



HAL
open science

Contributions to statistics and machine learning for physics-based attack detection in industrial systems

Guillaume Ansel

► **To cite this version:**

Guillaume Ansel. Contributions to statistics and machine learning for physics-based attack detection in industrial systems. Signal and Image Processing. Ecole nationale supérieure Mines-Télécom Atlantique, 2021. English. NNT : 2021IMTA0231 . tel-03363094

HAL Id: tel-03363094

<https://theses.hal.science/tel-03363094>

Submitted on 3 Oct 2021

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

THÈSE DE DOCTORAT DE

L'ÉCOLE NATIONALE SUPÉRIEURE
MINES-TÉLÉCOM ATLANTIQUE BRETAGNE
PAYS DE LA LOIRE – IMT ATLANTIQUE

ÉCOLE DOCTORALE N° 601
*Mathématiques et Sciences et Technologies
de l'Information et de la Communication*
Spécialité : *Signal, Image, Vision*

Par

Guillaume ANSEL

Contributions to Statistics and Machine Learning for Physics-based Attack Detection in Industrial Systems

Thèse présentée et soutenue à IMT Atlantique, le 17 février 2021
Unité de recherche : UMR CNRS 6285 Lab-STICC, Équipe MATRIX
Thèse N° : 2021IMTA0231

Rapporteurs avant soutenance :

Michel BARBEAU Professeur, School of Computer Science, Carleton University
Philippe FORSTER Professeur, Université Paris Nanterre

Composition du Jury :

Présidente :	Nora CUPPENS	Professeur, Polytechnique Montréal
Examineurs :	Michel BARBEAU	Professeur, School of Computer Science, Carleton University
	Philippe FORSTER	Professeur, Université Paris Nanterre
	Abdourrahmane ATTO	Maître de conférences, Polytech Annecy-Chambéry
	Frédéric CUPPENS	Professeur, Polytechnique Montréal
Dir. de thèse :	Dominique PASTOR	Professeur, IMT Atlantique

Invité :

Jérôme IDIER Directeur de recherche, Centrale Nantes

Remerciements

Je voudrais tout d'abord remercier mes directeurs et encadrants de thèse, Dominique Pastor, Frédéric Cuppens et Nora Cuppens, pour m'avoir donné cette opportunité et pour leur support, leur aide et leurs conseils tout au long de cette thèse.

Je remercie aussi le jury, Michel Barbeau, Philippe Forster, Abdourrahmane Atto et Jérôme Idier pour avoir accepté de prendre le temps de lire mes travaux et pour leurs retours.

Je voudrais également remercier tout le monde au département SC d'IMT Atlantique pour m'avoir accueilli et pour leur support pendant toutes ces années, avant et pendant ma thèse. En particulier Christophe, pour m'avoir permis de travailler avec lui et sans qui je n'aurais probablement pas pu suivre ce parcours. Thierry, pour ses conseils de méthodologie et d'organisation. Jean-Marc, Thierry, Christophe, Martine, Monique, Frédéric, Djalil, (et j'en oublie probablement...) pour tous les bons moments passés régulièrement en pause ensemble. Yassine et The Phuong pour ces années partagées ensemble dans le même bureau à discuter et s'entraider. Nicolas, pour ses conseils et l'opportunité de travailler avec lui après cette thèse. Erwan, mon frère de thèse, pour tous les moments passés au quotidien durant ces trois années de thèse. Et tous les autres doctorants, pour leur support mutuel durant ces années. Courage, vous y arriverez aussi !

Au-delà du département SC, je remercie aussi tous les professeurs et enseignants que j'ai eu l'occasion de rencontrer, que ce soit à IMT Atlantique, en prépa, au lycée, au collège, ou même à l'école primaire, pour m'avoir appris tant de choses au cours de toutes ces années.

Je voudrais aussi citer mes amis, Thomas, Pierre-Alexandre, Théo, Timothée, Phil, pour leur soutien et leur influence dans ma vie.

Enfin, je voudrais remercier toute ma famille, et tous particulièrement mes parents, pour tout ce qu'ils ont fait pour moi tous les jours, pour tout ce qu'ils m'ont appris, et pour m'avoir permis d'arriver jusqu'ici.

Merci pout tout !

Résumé long

Table des matières

1	Introduction	5
2	Random Distortion Testing	6
3	Invariance et RDT Généralisé	8
4	RDT Asymptotique	9
5	Apprentissage de comportements et détection d'anomalies sur signaux réels	10
6	Conclusion	11

1 Introduction

Avec la croissance rapide des réseaux informatiques et des infrastructures connectées, la cybersécurité est devenue un enjeu majeur des systèmes informatiques modernes, afin de protéger aussi bien les données personnelles de chacun que les infrastructures critiques dont nous dépendons. Le développement récent des villes intelligentes et de l'industrie 4.0 a encore fortement accru la nécessité de sécuriser ces infrastructures, car ces développements s'appuient sur une forte augmentation du nombre de systèmes connectés. Ces développements sont accompagnés d'une augmentation du nombre de vulnérabilités pouvant être exploitées pour attaquer ces infrastructures. Il y a donc un besoin croissant de nouvelles méthodes de détection d'attaques, capables de détecter ces nouvelles attaques rapidement.

La plupart des méthodes actuelles sont basées sur le principe de *misuse detection*, qui consiste à caractériser les attaques, et qui utilisent une base de données de signatures d'attaques. Ces approches nécessitent d'abord de découvrir et d'analyser manuellement les attaques. Il y a donc un délai entre l'apparition d'une nouvelle attaque et la possibilité de la détecter, délai durant lequel ces systèmes sont vulnérables. Nous allons donc nous intéresser dans cette thèse à des méthodes de type *détection d'anomalies*, permettant de détecter des anomalies comme étant des déviations du fonctionnement normal, donc potentiellement de détecter ces nouvelles attaques plus rapidement, puisqu'elles ne nécessitent pas de caractériser les attaques à détecter.

Pour cette thèse, nous avons décidé de travailler sur des systèmes industriels, type SCADA (Supervisory Control And Data Acquisition). Ce sont des systèmes critiques, parmi lesquels on peut trouver des systèmes comme des chaînes de production, des systèmes de production et de distribution d'énergie, d'eau, de gaz, etc. Il s'agit de systèmes dits cyber-physiques, qui interagissent avec un processus physique via des capteurs et des actionneurs, pilotés par des automates programmables. Tous ces composants sont en réseau et peuvent être surveillés par des opérateurs depuis un système de contrôle. Comme ces composants sont en réseau, on y retrouve des composants typiques pour la détection d'attaques tels que des NIDS (Network-based Intrusion Detection System) qui surveillent le trafic réseau. Par contre, ces types de systèmes de détection prennent rarement en compte la partie physique de ces systèmes, et ne sont pas forcément capables de surveiller les capteurs et les actionneurs du processus physique. Par exemple, ces systèmes

seraient incapables de détecter une attaque purement physique, où une personne agirait directement avec le processus, sans aucune interaction via le réseau.

Pour essayer de combler ce manque, nous allons nous intéresser dans cette thèse à des méthodes de type détection d'anomalies en surveillant les signaux physiques de ces systèmes industriels. Pour développer de telles méthodes, nous allons procéder en deux temps.

1. Tout d'abord, nous allons commencer par le développement d'un cadre théorique approprié pour l'apprentissage et la détection d'anomalies. Ce cadre théorique permettra ensuite de justifier les choix effectués pour le développement de méthodes de détection et expliquer les performances obtenues. Un aspect important de ce cadre théorique est qu'il doit être robuste, c'est-à-dire nécessiter peu d'hypothèses, et aussi maintenir des performances proches des performances théoriques, même si ces hypothèses ne sont pas tout à fait respectées. Un autre point essentiel est qu'il doit permettre de contrôler le taux de fausses alarmes, un point important en pratique qui influe fortement sur l'utilisabilité des méthodes développées.
2. Puis dans un deuxième temps, sur la base de ce cadre théorique, nous allons développer des méthodes de détection d'anomalies, c'est-à-dire des méthodes permettant d'apprendre le comportement normal du système, et ensuite de détecter des anomalies comme des déviations par rapport à ce comportement appris.

Pour cette thèse, nous avons choisi de nous appuyer sur l'approche RDT (Random Distortion Testing) [17], basée sur un problème de tests d'hypothèses, et qui nous semble approprié pour le développement de méthodes de détections robustes. Cette approche est présentée dans la section 2.

Les sections 3 et 4 présentent deux extensions de l'approche RDT développées durant cette thèse, respectivement le RDT Généralisé et le RDT Asymptotique. Le but de ces deux extensions est de lever certaines limitations de l'approche RDT par rapport aux propriétés du bruit considéré dans le problème RDT.

Enfin, la section 5 présente des applications de cette théorie sur des signaux réels issus du jeu de données SWaT (Secure Water Treatment) [2], un jeu de données enregistré sur un système de traitement des eaux, comportant des enregistrements durant deux semaines de l'ensemble des capteurs et des actionneurs du système. Ce jeu de données comporte également des attaques, nous permettant ainsi de tester les performances des méthodes développées. Nous présentons dans cette section deux méthodes de détection basées sur l'approche RDT : une première méthode de détection de discontinuités, puis une deuxième méthode de détection de changements de moyenne permettant de segmenter des signaux physiques, qui pourra potentiellement servir de base pour une méthode complète d'apprentissage et de détection d'anomalies.

2 Random Distortion Testing

L'approche RDT définit le problème de test d'hypothèses suivant :

$$\left\{ \begin{array}{l} \text{Observation : } \left\{ \begin{array}{l} Y(\omega) = \Theta(\omega) + X(\omega) \\ X \sim \mathcal{N}(0, C) \\ \Theta \text{ et } X \text{ indépendants} \end{array} \right. \\ \text{Hypothèse nulle } \mathcal{H}_0 : \nu_C(\Theta(\omega) - \theta_0) \leq \tau \\ \text{Hypothèse alternative } \mathcal{H}_1 : \nu_C(\Theta(\omega) - \theta_0) > \tau \end{array} \right. \quad (1)$$

On dispose d'une observation Y , qui est la somme d'une variable aléatoire Θ , représentant un phénomène d'intérêt, et de bruit Gaussien X . Cette modélisation est un modèle typique d'observation d'un phénomène physique à l'aide d'un capteur qui vient bruite l'observation. À partir de cette observation, on veut déterminer si la réalisation $\Theta(\omega)$ est suffisamment proche ou non d'un certain modèle déterministe

θ_0 défini par l'utilisateur. La proximité entre le signal Θ et le modèle θ_0 est définie par une tolérance, notée $\tau > 0$, elle-aussi définie par l'utilisateur.

Ce problème de test d'hypothèses diffère de ceux que l'on peut trouver en théorie statistique de la décision. En effet, dans des problèmes de décision classique, on cherche à obtenir une information sur la distribution qui a généré l'observation dont on dispose. Par exemple, on peut chercher à savoir si notre observation est issue d'une Gaussienne de moyenne 0 ou de moyenne 1. Plus généralement, ces problèmes considèrent une famille de distributions dépendant d'un paramètre θ , et on cherche à obtenir une information sur θ . Dans le problème RDT, l'information que l'on recherche est différente, car on cherche à savoir si la réalisation du signal que l'on a observée est proche ou non d'un certain modèle. On ne cherche pas à obtenir une information sur la distribution de Θ , qui reste complètement inconnue. À vrai dire, une même distribution pour Θ peut potentiellement rendre vraie une hypothèse ou l'autre, alors que pour les autres problèmes, chaque distribution est associée à une seule hypothèse.

Le problème RDT est intéressant pour plusieurs raisons. Tout d'abord, il n'y a besoin d'aucune information sur le signal Θ , sa distribution de probabilité est complètement inconnue. Ceci fait que l'on peut appliquer cette approche à un grand nombre de signaux, tant que le modèle d'observation est respecté.

Ensuite, il y a l'utilisation de la tolérance τ , qui permet d'introduire une robustesse par rapport au modèle θ_0 . En effet, il est rare que le signal soit exactement égal au modèle que l'on considère, et la plupart du temps, de légères déviations par rapport au modèle souhaité n'auront pas d'impact sur le signal. Cette tolérance permet donc d'autoriser un léger décalage avec le modèle de façon contrôlable.

Enfin, bien que l'on ait aucune information sur Θ , on peut trouver un test optimal pour ce problème. Le test en question est le test à seuil $\mathcal{J}_{\lambda_\gamma(\tau)}$, défini de la façon suivante :

$$\begin{aligned} \mathcal{J}_{\lambda_\gamma(\tau)}: \mathbb{R}^d &\rightarrow \{0, 1\} \\ y &\mapsto \begin{cases} 1 & \text{si } \nu_C(y - \theta_0) > \lambda_\gamma(\tau) \\ 0 & \text{sinon} \end{cases} \end{aligned} \quad (2)$$

où le seuil $\lambda_\gamma(\tau)$ vérifie l'équation $\mathbb{F}_{\chi_d^2(\tau^2)}(\lambda_\gamma(\tau)^2) = \gamma$, avec $\mathbb{F}_{\chi_d^2(\tau^2)}$ représentant la fonction de répartition d'une variable aléatoire suivant une loi du χ^2 non centrée à d degrés de liberté et de paramètre de non-centralité τ^2 .

Ce test $\mathcal{J}_{\lambda_\gamma(\tau)}$ est optimal au sens du critère γ -MCCP (Maximum Constant Conditional Power), un critère d'optimalité adapté au problème RDT, qui n'est pas sans rappeler le critère UMPI (Uniformly Most Powerful Invariant), utilisé pour des problèmes de décision faisant intervenir des propriétés d'invariance. Ce critère peut être décomposé en trois points :

1. Le test $\mathcal{J}_{\lambda_\gamma(\tau)}$ a niveau γ , autrement dit sa probabilité de fausse alarme est toujours inférieure ou égale à γ :

$$\sup_{\Theta: \mathbb{P}[\nu_C(\Theta - \theta_0) \leq \tau] \neq 0} \mathbb{P}[\mathcal{J}_{\lambda_\gamma(\tau)}(\Theta + X) = 1 \mid \nu_C(\Theta - \theta_0) \leq \tau] \leq \gamma \quad (3)$$

2. Pour toute variable aléatoire Θ et pour $\mathbb{P}_{\nu_C(\Theta - \theta_0)}$ -presque tout $\rho > \tau$, le test $\mathcal{J}_{\lambda_\gamma(\tau)}$ a puissance conditionnelle constante étant donné $\Theta \in \Upsilon_\rho$:

$$\forall \theta \in \Upsilon_\rho, \mathbb{P}[\mathcal{J}_{\lambda_\gamma(\tau)}(\Theta + X) = 1 \mid \Theta \in \Upsilon_\rho] = \beta_{\mathcal{J}_{\lambda_\gamma(\tau)}}(\theta) \quad (4)$$

Cette propriété est l'analogue de la propriété d'invariance du critère UMPI, mais elle porte ici sur la fonction puissance du test au lieu de s'appliquer sur le test.

3. Enfin, parmi tous les tests \mathcal{J} qui vérifient ces deux premières propriétés, le test $\mathcal{J}_{\lambda_\gamma(\tau)}$ maximise la probabilité de détection. Pour toute variable aléatoire $\Theta \in \Upsilon_\rho$ et pour $\mathbb{P}_{\nu_C(\Theta - \theta_0)}$ -presque tout $\rho > \tau$, on a :

$$\mathbb{P}[\mathcal{J}_{\lambda_\gamma(\tau)}(\Theta + X) = 1 \mid \Theta \in \Upsilon_\rho] \geq \mathbb{P}[\mathcal{J}(\Theta + X) = 1 \mid \Theta \in \Upsilon_\rho] \quad (5)$$

On dispose donc d'un test optimal pour ce problème, mais cette optimalité n'est valide que si les conditions du problème sont bien respectées. Il est notamment nécessaire de connaître parfaitement les propriétés du bruit X , et celui-ci doit être Gaussien, ce qui n'est pas nécessairement le cas en pratique. Il est notamment courant d'avoir besoin d'estimer la matrice de covariance du bruit, ce qui peut dégrader les performances du test. Pour cette raison, nous avons développé deux extensions de la théorie RDT prenant ces aspects en compte. La première, le RDT généralisé, étend l'approche pour des distributions de bruit non Gaussiennes, mais présentant des propriétés d'invariance. La seconde, le RDT asymptotique, permet de prendre en compte l'estimation de la variance du bruit σ_0^2 et du modèle θ_0 .

3 Invariance et RDT Généralisé

L'approche RDT s'applique uniquement au cas d'une observation $Y = \Theta + X$ où le bruit X a une distribution Gaussienne. Le but du RDT généralisé est d'étendre cette approche à d'autres distributions de bruit. Nous n'allons pas travailler ici avec une distribution quelconque pour le bruit, mais avec des distributions présentant des propriétés d'invariance. En effet, il s'avère que dans les développements théoriques du RDT, cette notion d'invariance est centrale afin d'obtenir la structure du test et en déduire l'optimalité [17]. Par exemple, dans le cas du problème RDT, la distribution Gaussienne du bruit est invariante par rotation, c'est-à-dire que si $X \sim \mathcal{N}(0, \sigma^2 I_d)$ est une variable aléatoire Gaussienne, et g est une rotation de \mathbb{R}^d , alors X et $g(X)$ ont la même distribution de probabilité.

Pour poser le problème GRDT, on considère maintenant une distribution de probabilité P , invariante sous l'action d'un groupe G de transformations linéaires de \mathbb{R}^d , et on considère également un maximal invariant $M: \mathbb{R}^d \rightarrow \mathbb{R}$ de G . Le problème GRDT est défini de la façon suivante :

$$\left\{ \begin{array}{l} \textbf{Observation : } \begin{cases} Y(\omega) = \Theta(\omega) + X(\omega) \\ X \sim P \\ \Theta \text{ et } X \text{ indépendants} \end{cases} \\ \textbf{Hypothèse nulle } \mathcal{H}_0 : M(\Theta(\omega)) \leq \tau \\ \textbf{Hypothèse alternative } \mathcal{H}_1 : M(\Theta(\omega)) > \tau \end{array} \right. \quad (6)$$

Pour ce problème, la question est maintenant de savoir s'il est possible de trouver un test optimal, au sens du critère γ -MCCP présenté précédemment. Nous avons pu identifier deux résultats théoriques qui peuvent aider à trouver un tel test. Le premier résultat, énoncé dans le Théorème 2.2.12, donne une condition suffisante pour qu'un test soit γ -MCCP dans le cas général.

Le deuxième résultat, énoncé dans le Théorème 2.2.16, donne un test optimal dans le cas spécifique où le bruit X est sphériquement invariant, ce qui signifie que le maximal invariant est la norme euclidienne (ceci inclut par exemple le cas du bruit Gaussien), et si la famille de distributions $\{\|\theta + X\|_2, \theta \in \mathbb{R}^d\}$ a un rapport de vraisemblance monotone. Dans ce cas, on peut trouver un seuil $\lambda_\gamma(\tau)$ tel que le test $\mathcal{T}_{\lambda_\gamma(\tau)}$ est γ -MCCP.

On peut donc trouver un test optimal sous ces conditions, mais il reste maintenant le problème de trouver des distributions de probabilités pour le bruit vérifiant ces contraintes. En effet, la propriété de rapport de vraisemblance monotone est très contraignante, car elle concerne une famille de distributions très particulière qui est intégralement définie par le choix de la distribution de X . On ne peut donc pas se contenter de partir d'une famille de distributions ayant rapport de vraisemblance monotone : il faut choisir une distribution de probabilité pour X , et vérifier si la famille résultante a rapport de vraisemblance monotone ou non. Pour le moment, à l'exception de la distribution Gaussienne, nous n'avons pas encore pu déterminer d'autres distributions vérifiant cette propriété. Il est également important de noter que nous avons principalement étudié le cas d'une distribution présentant une invariance sphérique, il faudrait également étudier d'autres maximaux invariants afin de voir s'il est possible d'établir des résultats similaires au Théorème 2.2.16 dans ces cas.

4 RDT Asymptotique

Après avoir étudié les propriétés d'invariance du bruit pour étendre l'approche RDT, nous allons maintenant considérer le problème de l'estimation des paramètres du modèle, et voir comment l'estimation de ces paramètres influe sur les performances. Dans beaucoup d'applications, il est courant de devoir estimer les propriétés du bruit car elles ne sont pas parfaitement connues. On peut alors se demander comment ceci affecte les performances du test $\mathcal{J}_{\lambda_\gamma(\tau)}$ et si on conserve l'optimalité dans ces conditions. Le but ici est de vérifier ces points, en considérant une estimation de la variance du bruit σ_0 , ainsi que du modèle θ_0 .

Cette étude est effectuée en deux temps. Nous commencerons par étudier le comportement asymptotique du test, lorsque les estimées de ces valeurs convergent vers les vraies valeurs de θ_0 et σ_0 . Puis nous étudierons le comportement non-asymptotique de ce test à travers des simulations.

Le problème ARDT est défini de la façon suivante :

$$\left\{ \begin{array}{l} \textbf{Observation :} \begin{cases} Y(\omega) = \Theta(\omega) + X(\omega) \\ X \sim \mathcal{N}(0, \sigma_0^2 I_d) \\ \Theta \text{ et } X \text{ indépendants} \end{cases} \\ \textbf{Hypothèse nulle } \mathcal{H}_0: \|\Theta(\omega) - \theta_0\|_2 \leq \sigma_0 \tau \\ \textbf{Hypothèse alternative } \mathcal{H}_1: \|\Theta(\omega) - \theta_0\|_2 > \sigma_0 \tau \end{array} \right. \quad (7)$$

avec θ_0 et σ_0 inconnus, et $\hat{\theta}_n$ et $\hat{\sigma}_n$ des estimateurs consistants de ces valeurs.

Afin de simplifier l'étude, nous considérons ici uniquement le cas de bruit blanc Gaussien de variance σ_0^2 , et non pas le cas général d'une matrice de covariance C quelconque. On a donc pour tout vecteur $y \in \mathbb{R}^d$, $\nu_C(\sigma_0^2 I_d)y = \|y\|_2 / \sigma_0^2$.

On suppose donc θ_0 et σ_0 inconnus, mais on dispose d'estimateurs consistants $\hat{\theta}_n$ et $\hat{\sigma}_n$ de ces deux valeurs, c'est-à-dire qu'ils convergent en probabilité vers θ_0 et σ_0 respectivement.

Étant donné que le test RDT $\mathcal{J}_{\lambda_\gamma(\tau)}$ est optimal lorsque θ_0 et σ_0 sont connus, il est naturel d'étudier le test obtenu lorsque l'on remplace ces valeurs par leurs estimées. Nous allons donc étudier le test $\tilde{\mathcal{J}}_{\lambda_\gamma(\tau)}$ défini par :

$$\begin{aligned} \tilde{\mathcal{J}}_{\lambda_\gamma(\tau)}: \mathbb{R}^d \times \mathbb{R}^d \times (0, \infty) &\rightarrow \{0, 1\} \\ (y, \theta, \sigma) &\mapsto \begin{cases} 1 & \text{si } \|y - \theta\|_2 > \sigma \lambda_\gamma(\tau) \\ 0 & \text{sinon} \end{cases} \end{aligned} \quad (8)$$

Pour ce test, on dispose de deux résultats théoriques concernant ses propriétés asymptotiques lorsque les estimateurs $\hat{\theta}_n$ et $\hat{\sigma}_n$ convergent. Ces deux résultats concernent le niveau asymptotique du test et un critère d'optimalité asymptotique.

Le premier résultat, énoncé dans le Théorème 3.2.2, indique que le test $\tilde{\mathcal{J}}_{\lambda_\gamma(\tau)}$ a niveau asymptotique γ lorsque $\hat{\theta}_n$ et $\hat{\sigma}_n$ convergent, si le signal Θ et les estimateurs sont indépendants :

$$\limsup_n \sup_{\Theta \text{ ind. de } (\hat{\theta}_n, \hat{\sigma}_n)} \mathbb{P}[\tilde{\mathcal{J}}_{\lambda_\gamma(\tau)}(\Theta + X, \hat{\theta}_n, \hat{\sigma}_n) = 1 \mid \|\Theta - \theta_0\|_2 \leq \sigma_0 \tau] \leq \gamma \quad (9)$$

Le deuxième résultat, énoncé dans le Théorème 3.2.4, donne l'optimalité asymptotique de ce test. Pour tout test $\tilde{\mathcal{F}}$ avec niveau γ si θ_0 et σ_0 sont supposés connus, avec puissance conditionnelle constante pour tout $\rho > \tau$, et dont la région critique est un $\mathbb{P}_{(Y, \theta_0, \sigma_0)}$ -continuity set, le test $\tilde{\mathcal{J}}_{\lambda_\gamma(\tau)}$ est asymptotiquement plus puissant que $\tilde{\mathcal{F}}$:

$$\begin{aligned} \limsup_n \left(\mathbb{P}[\tilde{\mathcal{J}}_{\lambda_\gamma(\tau)}(\Theta + X, \hat{\theta}_n, \hat{\sigma}_n) = 1 \mid \|\Theta - \theta_0\|_2 = \sigma_0 \rho] \right. \\ \left. - \mathbb{P}[\tilde{\mathcal{F}}(\Theta + X, \hat{\theta}_n, \hat{\sigma}_n) = 1 \mid \|\Theta - \theta_0\|_2 = \sigma_0 \rho] \right) \geq 0 \end{aligned} \quad (10)$$

Après ces résultats théoriques, nous avons étudié le comportement non-asymptotique du test $\tilde{\mathcal{F}}_{\lambda_\gamma(\tau)}$ par simulations. Ces simulations, présentées dans la Section 3.3, ont permis d'étudier plus en détail le comportement de ce test. Les premières simulations ont consisté à mesurer le taux de fausse alarme effectif du test, afin de vérifier le Théorème 3.2.2. Les résultats présentés sur la Figure 3.2 montre que d'une part, l'estimation de σ augmente le taux de fausse alarme effectif au dessus de γ , de façon très significative dans certains cas, et d'autre part ce taux de fausse alarme décroît lorsque l'estimation de σ_0 s'améliore, et semble bien tendre vers γ , ce qui confirme ce résultat théorique.

Après avoir vérifié ce résultat, nous avons développé une méthode permettant de retrouver le niveau γ lorsque σ_0 est estimé, en ajustant la tolérance τ . Cette tolérance vient alors compenser à la fois une déviation potentielle du signal par rapport au modèle et l'effet de l'estimation de σ_0 . Ceci ne dégrade pas les performances globales du test $\mathcal{F}_{\lambda_\gamma(\tau)}$, et on retrouve les mêmes courbes ROC (Receiver Operating Characteristic) lorsque l'on applique ce test à un problème de détection.

5 Apprentissage de comportements et détection d'anomalies sur signaux réels

Après ces aspects théoriques, nous pouvons maintenant voir comment s'appuyer sur ces résultats théoriques pour l'apprentissage et la détection sur signaux réels.

La première méthode développée à partir de ces outils permet de détecter des discontinuités sur des signaux qui sont supposés être continus. Nous avons effectivement constaté qu'un certain nombre d'attaques menées engendrent des discontinuités sur certains signaux, alors qu'en temps normal, ces signaux présentent des variations lentes, et sont donc continus.

Afin de détecter ces discontinuités, nous avons développé une méthode de détection s'appuyant sur une transformation en ondelettes, et effectuant la détection sur les coefficients de détail utilisant le test RDT $\tilde{\mathcal{F}}_{\lambda_\gamma(\tau)}$. La variance du bruit est estimée sur ces mêmes coefficients en utilisant l'estimateur MAD (Median Absolute Deviation), afin d'avoir une estimation robuste face aux pics pouvant être causés par ces discontinuités. Les paramètres γ et τ du test RDT peuvent être ajustés en s'appuyant sur un enregistrement normal du signal de façon à obtenir le taux de fausse alarme voulu. Avec cette méthode, nous avons pu détecter les neuf discontinuités présentes sur le signal testé, avec une seule fausse alarme sur un total de 224 965 échantillons testés.

La seconde méthode développée est une méthode de détection de changements. L'idée est que les composants d'un système industriel ont tendance à avoir un comportement régulier, avec quelques phases de fonctionnements qui sont répétées au cours d'un cycle. On voudrait donc détecter ces transitions d'une phase à la suivante, afin de caractériser le comportement normal du système. Pour cela, nous avons développé une méthode de détection de changements basée sur le test RDT, permettant de détecter des changements dans la moyenne d'un signal. Cette méthode consiste à estimer sur le signal la moyenne de la phase courante, puis tester par blocs d'échantillons si la moyenne du signal a changé significativement ou non par rapport à la valeur mesurée. Si oui, on considère que l'on a détecté une nouvelle phase du signal, et on estime alors la moyenne de la nouvelle phase. La détection est effectuée en utilisant le test RDT, dont la tolérance permet d'indiquer l'amplitude minimale des changements que l'on souhaite détecter. Le RDT asymptotique est utilisé ici afin de prendre en compte l'estimation de la moyenne actuelle du signal et de la variance du bruit. On obtient alors une segmentation du signal selon ses différentes phases de fonctionnement.

Sur la base de cette méthode de segmentation, nous envisageons de développer une méthode complète d'apprentissage et de détection d'anomalies. L'idée de cette méthode est d'identifier dans un premier temps les phases des différents signaux d'un système, et ensuite de caractériser ces phases (durée, valeur moyenne, variance, phases précédentes et suivantes...) sur la base de données enregistrées en fonctionnement normal. Ceci nous donne alors un modèle du comportement normal des différents composants du système. Sur la base de ce modèle, on peut alors détecter des anomalies en mesurant les caractéristiques du système, et en les comparant au modèle appris.

6 Conclusion

Cette thèse a permis d'explorer l'utilisation de l'approche RDT pour la détection d'anomalies. Cette approche nous a donné une base pour développer des méthodes de détection robustes et applicables sur des signaux dont la connaissance à priori est potentiellement très limitée. Nous avons développé deux extensions de cette approche, le RDT Généralisé et le RDT Asymptotique afin d'étendre davantage son cadre d'application, et nous avons pu développer des méthodes de détection sur cette base, avec des résultats encourageants lors d'essais sur signaux réels. Sur la base de ces résultats, nous envisageons de développer une méthode complète d'apprentissage et de détection, permettant de caractériser le comportement normal de systèmes industriels et de détecter des anomalies et des attaques menées contre ces systèmes de façon robuste.

Contents

Résumé long	5
1 Introduction	5
2 Random Distortion Testing	6
3 Invariance et RDT Généralisé	8
4 RDT Asymptotique	9
5 Apprentissage de comportements et détection d'anomalies sur signaux réels	10
6 Conclusion	11
Introduction	19
Context	19
A short primer on industrial systems and cybersecurity	20
Industrial Control Systems	20
Detection schemes: misuse detection vs. anomaly detection	21
Data	21
The SWaT testbed	22
Objectives and contributions	25
Thesis structure	26
Preliminary: mathematical notions and notations	29
1 Hypothesis Testing and Random Distortion Testing	31
1.1 Traditional hypothesis testing approaches	31
1.1.1 Non-Bayesian binary classification	31
1.1.2 Optimality	33
1.2 The RDT approach	36
1.2.1 Problem statement	36
1.2.2 Optimality: γ -MCCP tests	37
1.2.3 Thresholding tests and optimality	39
1.3 Robustness and hypothesis testing	40
1.3.1 Huber's approach to robust hypothesis testing	41
1.3.2 RDT and robust hypothesis testing	43
2 Invariance and Generalized Random Distortion Testing	45
2.1 Invariance in group theory	46
2.1.1 Group theory	46
2.1.2 Invariance	47
2.1.3 Orbits and maximal invariant	49
2.1.4 Invariance applied to probability distributions	51

2.2	Generalization of the RDT approach	53
2.2.1	Problem statement	54
2.2.2	Redefining notions for the GRDT problem	55
2.2.3	Preliminary results	56
2.2.4	Generalization when the maximal invariant is the euclidean norm	61
2.3	Conclusion and perspectives	64
3	Asymptotic Random Distortion Testing	65
3.1	Preliminary results	65
3.2	Asymptotic RDT	68
3.2.1	Problem statement	68
3.2.2	Asymptotic size and power	69
3.3	Simulation experiments	75
3.3.1	Level	75
3.3.2	Application to a detection problem	78
3.3.3	Recovering a test with level γ when estimating σ_0	86
3.4	Conclusion	90
4	Learning behaviors and detecting anomalies on real signals	91
4.1	Detecting discontinuities on continuous signals	92
4.2	Segmenting and learning phases, and detecting anomalies	96
4.2.1	Change-detection in time series	97
4.2.2	An RDT-based change-in-mean detection method	98
4.2.3	Application to real signals	101
4.2.4	Perspectives: towards learning a model and detecting anomalies	105
4.3	Conclusion	108
	Conclusion and perspectives	111
	Appendices	115
A	Deterministic Distortion Testing	117
B	Finding suitable families of TP-2 distributions for the GRDT problem	121
B.1	Notations and preliminary results	121
B.2	Calculations	122
C	Uniform continuity of the Generalized Marcum function $Q_{d/2}$	127
	Acronyms	135
	List of Publications	137
	Bibliography	139

List of Figures

Introduction	
1	Example of SCADA architecture 20
2	SWaT Process 22
3	Example of signals recorded under normal operating conditions in the SWaT testbed 23
4	Example of signals recorded in presence of attacks 24
1 Hypothesis Testing and Random Distortion Testing	
1.1	Threshold $\lambda_\gamma(\tau)$ against τ (left) and γ (right) for $d = 2$ 40
3 Asymptotic Random Distortion Testing	
3.1	Measured false alarm rate against γ with known σ_0 for different values of DNR 77
3.2	Measured false alarm rate against γ with estimated σ_0 for different values of DNR and numbers of samples N used to estimate σ_0 78
3.3	Measured false alarm rate against number of samples N used to estimate σ_0 for different values of γ with DNR = 5 dB 79
3.4	ROC curves obtained for the NP and RDT tests with no distortion and σ_0 known. 81
3.5	Measured false alarm rate against γ for the NP and RDT tests with no distortion and σ_0 known 81
3.6	ROC curves obtained for the NP and RDT tests with distortion and σ_0 known for SNR = 15 dB 82
3.7	Measured false alarm rate against γ for the NP and RDT tests with distortion and σ_0 known 83
3.8	ROC curves for the NP and RDT tests with estimation of σ_0 , SNR = 15 dB, DNR = 5 dB 84
3.9	Measured false alarm rate against γ for the Neyman-Pearson test with estimation of σ_0 85
3.10	Measured false alarm rate against γ for the RDT test with estimation of σ_0 85
3.11	Measured values of τ^* against γ for different values of N and DNR. 87
3.12	Measured false-alarm rate against γ using the adjusted tolerance τ^* 87
3.13	Measured values of τ^* against γ for different values of N and DNR with $\sigma_0 = 2$ 88
3.14	Measured values of τ^* against γ for different values of N and DNR with $\sigma_0 = 100$ 89
3.15	Comparison of ROC curves obtained using the tolerances τ and τ^* 89
4 Learning behaviors and detecting anomalies on real signals	
4.1	Examples of discontinuities observed on the water level signal LIT101 92
4.2	Wavelet transform on the signal LIT101 93
4.3	Wavelet transform on normal LIT101 signal 94
4.4	Threshold applied on the detail coefficients 95
4.5	Late attack detection 96
4.6	Sensor LIT101 and actuators MV101 and P101 during normal operation 97
4.7	Proposed RDT-based change-in-mean detection algorithm 100

4.8	Signal LIT101 and its derivative under normal operation.....	101
4.9	Tracked signal mean on the derivative of LIT101 ($N = 40, \gamma = 10^{-3}, \tau' = 0.05$).....	102
4.10	Tracked signal mean on the derivative of LIT101 ($N = 10, \gamma = 10^{-3}, \tau' = 0.05$).....	103
4.11	Tracked signal mean on the derivative of LIT101 ($N = 10, \gamma = 10^{-6}, \tau' = 0.05$).....	103
4.12	Tracked signal mean on the derivative of LIT101 with attacks ($N = 40, \gamma = 10^{-3}, \tau' = 0.05$).....	104
4.13	Tracked signal mean on the derivative of LIT101 with attacks ($N = 40, \gamma = 10^{-2}, \tau' = 0.01$).....	105
4.14	K-means output on normal LIT101 signal.....	106
4.15	Example of segmentation artifacts caused by short segments.....	106
4.16	K-means output on normal LIT101 signal (zoom around 0).....	107
4.17	K-means output on normal LIT101 signal with manual seeding.....	107
4.18	Comparison of slopes obtained on signal LIT101 both in normal and attack conditions... ..	107

List of Theorems, Lemmas and Definitions

1 Hypothesis Testing and Random Distortion Testing	
Definition 1.1.1	Test and critical region 32
Definition 1.1.2	Level, size and power function of a test 33
Definition 1.1.3	Uniformly Most Powerful (UMP) Test 33
Theorem 1.1.4	Neyman-Pearson Lemma 34
Definition 1.1.5	Monotone Likelihood Ratio 34
Theorem 1.1.6	Karlin-Rubin Theorem 35
Definition 1.2.1	Mahalanobis norm 36
Definition 1.2.2	RDT problem statement 36
Definition 1.2.3	Orbits Υ_ρ and set of orbits \mathfrak{F} 37
Definition 1.2.4	Conditional size and power 37
Lemma 1.2.5	Size and conditional size 38
Definition 1.2.6	Constant conditional power function given $\Theta \in \Upsilon_\rho$ 38
Definition 1.2.7	γ -MCCP (Maximum Constant Conditional Power) test 38
Definition 1.2.8	Thresholding tests 39
Lemma 1.2.9	Threshold $\lambda_\gamma(\tau)$ 40
Theorem 1.2.10	Optimality of the test $\mathcal{J}_{\lambda_\gamma(\tau)}$ 40
2 Invariance and Generalized Random Distortion Testing	
Definition 2.1.1	Group 46
Definition 2.1.2	Group homomorphism 46
Lemma 2.1.3	Properties of group homomorphisms 46
Definition 2.1.4	Group action 46
Lemma 2.1.5	Group action characterization 47
Definition 2.1.6	Invariant subset 48
Lemma 2.1.7	Invariant subset characterization 48
Lemma 2.1.8	Group action on a subset 48
Lemma 2.1.9	Transfer of a group action 49
Definition 2.1.10	Orbit 50
Definition 2.1.11	Invariant function 50
Definition 2.1.12	Maximal invariant 50
Lemma 2.1.13	Characterization of invariant functions 50
Lemma 2.1.14	Equivalent definition of orbits using a maximal invariant 51
Lemma 2.1.15	Group action on measures 51
Definition 2.1.16	Invariant measure 53
Lemma 2.1.17	Invariant measure and random variables 53
Definition 2.2.1	Generalized RDT problem statement 54
Definition 2.2.2	Conditional power and size 55
Definition 2.2.3	Thresholding tests 55

Definition 2.2.4	Constant conditional power function given $\Theta \in \Upsilon\rho$	56
Definition 2.2.5	γ -MCCP (Maximum Constant Conditional Power) test	56
Lemma 2.2.6	Distribution of $M(\theta + X)$	56
Lemma 2.2.7	Constant power function on a subset	57
Lemma 2.2.8	Power function of invariant tests	57
Lemma 2.2.9	Size and conditional size	58
Lemma 2.2.10	Constant condition power function on an orbit	59
Lemma 2.2.11	Constant conditional power function and constant power function	59
Theorem 2.2.12	Sufficient condition for a γ -MCCP test	60
Lemma 2.2.13	Probability density function of spherically invariant random variables	61
Lemma 2.2.14	Size of $\mathcal{T}_{\lambda_\gamma(\tau)}$	62
Lemma 2.2.15	62
Theorem 2.2.16	Optimality of $\mathcal{T}_{\lambda_\gamma(\tau)}$	63
3	Asymptotic Random Distortion Testing	
Lemma 3.1.1	Convergence in distribution to a deterministic constant	66
Lemma 3.1.2	Convergence in probability of random vectors	66
Lemma 3.1.3	67
Lemma 3.1.4	67
Definition 3.2.1	ARDT problem statement	68
Theorem 3.2.2	Asymptotic level of $\tilde{\mathcal{T}}_{\lambda_\gamma(\tau)}$	69
Lemma 3.2.3	71
Theorem 3.2.4	Asymptotic optimality of $\tilde{\mathcal{T}}_{\lambda_\gamma(\tau)}$	72
Lemma 3.2.5	74
Lemma 3.2.6	74
A	Deterministic Distortion Testing	
Definition A.0.1	DDT problem statement	117
Definition A.0.2	γ -MCP (Maximum Constant Power) test	118
Lemma A.0.3	γ -MCP tests and γ -MCCP tests	118
B	Finding suitable families of TP-2 distributions for the GRDT problem	
Lemma B.1.1	Spherically invariant random vectors	122
Lemma B.1.2	122
C	Uniform continuity of the Generalized Marcum function $Q_{d/2}$	
Lemma C.0.1	127
Lemma C.0.2	129
Lemma C.0.3	130
Theorem C.0.4	132

Introduction

Contents

Context	19
A short primer on industrial systems and cybersecurity	20
Industrial Control Systems	20
Detection schemes: misuse detection vs. anomaly detection	21
Data	21
The SWaT testbed	22
Objectives and contributions	25
Thesis structure	26

Context

Since the advent of the first computer networks back in the sixties, cybersecurity went from being nearly inexistent to becoming a major component of every computer system. As these systems are now integral parts of our lives and are present in major critical infrastructures, securing them is more crucial than ever. From someone's personal computer to banks to the internal systems power plants, every computer system can potentially be targeted by hackers with different goals: some may simply want to annoy people, some others try to extort money of out them — e.g. ransomwares, which target individuals and companies alike — but some actively attack specific infrastructures for possibly political or strategic reasons, for instance accessing sensitive information or damaging critical infrastructures. One example of such an attack is Stuxnet [1], discovered in 2010 and believed to have been developed by the NSA to stop the Iranian nuclear program.

Nowadays, with the advent of 5G networks and IoT (Internet of Things), these threats are omnipresent as the number of connected devices is rapidly increasing. Despite that, manufacturers do not necessarily take security into account when designing their products. While these technologies provide many useful services that can notably improve quality of life — for example Smart Cities aim to improve transportation networks, power management or healthcare among other things — they can also cause many security issues that a malicious actor could use to cause significant damage to these infrastructures. It is not hard to imagine how much damage one could cause by messing with traffic signals or water supplies for example, let alone a power plant.

In the industrial field we can see a similar trend developing with Industry 4.0 and Smart Factories. These factories rely on a large number of sensors to monitor every part of the production chain and make decisions autonomously to optimize the process, monitor anomalies, and help making decisions. These industrial systems are often referred to as ICSs (Industrial Control Systems). They build upon similar technologies to IoT, and as such inherit some of its vulnerabilities, as increased connectivity leads to more potentially exploitable vulnerabilities (zero-days).

Overall, despite many efforts to secure these systems, they are increasingly vulnerable to cyberattacks, especially new kinds of attacks that make use of these connected devices as an attack vector. Therefore it is more important than ever to develop new attack detection methods that can reliably detect these novel attacks in time before any damage occurs. For these reasons, we have chosen to work on developing new detection methods focusing mainly on industrial systems.

A short primer on industrial systems and cybersecurity

Industrial Control Systems

ICSs are cyber-physical systems, i.e. systems composed of physical components — such as sensors and actuators — that interact with a physical process and are networked together and controlled by computer systems. In these systems, the sensors and actuators that are part of the industrial process are connected to PLCs (Programmable Logic Controllers), which are industrial computers programmed to control the process. Multiple PLCs are used to control different parts of the system and ensure redundancy. These components are networked together and connected to control systems which are used by operators to manage and operate the process. Several network architectures exist for these systems, such as SCADAs (Supervisory Control And Data Acquisitions) or DCSs (Distributed Control Systems). An example of SCADA architecture in a water treatment process is shown in Fig. 1. ICSs can be found in many industrial sectors, including for example power grids, water distribution, telecommunication networks and manufacturing.

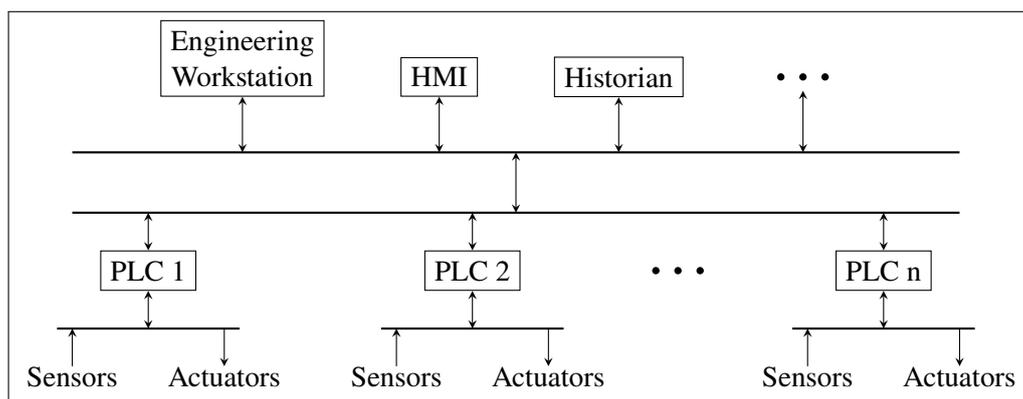


Figure 1: Example of SCADA architecture

Adapted by permission from Springer Nature: Jonathan Goh et al. "A Dataset to Support Research in the Design of Secure Water Treatment Systems". In: *Critical Information Infrastructures Security*. Ed. by Grigore Havarneanu et al. Lecture Notes in Computer Science. Cham: Springer International Publishing, 2017, pp. 88–99. ISBN: 978-3-319-71368-7. DOI: 10.1007/978-3-319-71368-7_8

Previously, most industrial systems used to be air-gapped, i.e. their components were on a separate network which was completely disconnected from any other network. Nowadays with the development of communication technologies, it is common to have these systems connected in some way to other networks. For example, the factory network can be connected to the corporate network of the company, so that managers have a view of the current state of the systems (productivity, required maintenance, etc.) to help their decision-making processes. However, that also means that they could be accessed from the Internet. An attacker could then potentially access the control systems and affect for example the control signals and disrupt the system. Wireless technologies are also increasingly used to help connect remote components. If someone can get in range of these connections, and if they are not properly secured, this can be used as an entry point to access the system. Another aspect of these systems is that some industries still use old components that were not designed to be networked in this way: older systems may not offer ways to validate integrity of the transmitted data or may not require authentication to perform critical operations. This would make them be easy targets in a modern network.

Detection schemes: misuse detection vs. anomaly detection

Most attack detection methods in cybersecurity fall into two main categories: *Misuse detection* and *Anomaly detection*.

Misuse detection, also referred to as *Signature-based detection*, refers to methods that detect attacks by defining what the attacks are. Everything else is then considered to be normal. The attacks to be detected are characterized using signatures, which are created by experts analyzing existing threats and integrated into signature databases. The detection system then attempts to recognize these signatures and raises an alarm whenever it finds a match. For example, for malware detection, these signatures contain part of the instructions present in the malware. For network intrusion detection, the signatures can contain sequences of packets or specific sequences of bytes in packets that have been determined to be malicious. However it is important to note that these detection mechanisms are reactive: they are unable to detect new forms of attacks for which no signature exists. This means that to detect a new threat, we have to wait until a new signature for this threat is made available, which means that experts first have to discover and analyze it. This analysis can take some time, and all systems remain vulnerable to this threat until new signatures are available.

Anomaly detection [3] is the opposite principle: instead of defining the attacks to be detected, we instead define what is normal, and everything else is considered an anomaly. This approach would allow detecting new attacks as long as they do not match any known normal behavior. However the main issue with this approach is learning the normal behavior of the system. With multiple users on a network, each one will have different habits and behaviors which all need to be learned accurately, while also considering that they can change over time. In addition to these human behaviors we also have to take into account all the automated and autonomous processes that might exist. As some of these are automated, they tend to be more regular, which makes discovering and learning them easier. However they still evolve over time: components may get replaced, or software may get upgraded which can change their behaviors. Because of this complexity, anomaly detection methods tend to present a high false alarm rate, making them difficult to use as many false alarms require a lot of time to manually confirm or infirm.

While learning the normal behavior of a system remains a daunting task, a shift to usable anomaly detection methods is required in order to detect new threats as soon as possible and minimize the damage they could cause. This is the focus of this work: developing methods to learn the normal behavior of a system while controlling false alarms rates. That is not to say that current signature-based methods should not be employed as well: detecting existing and known threats remains important so that they do not come back in the future. One way of achieving this is to combine both approaches in a hybrid detection method. Such work is beyond the scope of this study, but is a potential perspective to improve performances of the developed methods.

Data

As mentioned in the previous section, we will be working on anomaly detection methods which require learning the normal behavior of the studied system. As with every learning-based approach, we need to work on data captured on a system to learn its behavior. This section will discuss what kind of data is available in a cybersecurity context and the specific constraints that exist in this field.

Most attacks targeting a system are conducted remotely through a network. Therefore monitoring the network traffic is a very important part of cyberattack detection. This analysis is usually performed by IDSs (Intrusion Detection Systems), which scan all network packets that transit through them and raise an alarm whenever they detect suspicious activity. For industrial systems specifically, we also have access to the state of the physical process through the measurements of the sensors and the state of the actuators.

A major issue for cybersecurity research is data availability. As this field deals with sensitive information, releasing realistic datasets is difficult without compromising the security of actual systems. For general IDSs, a few datasets exist, such as DARPA [4] and NSL-KDD [5], which are network traffic captures in simulated environments and are commonly used to train and evaluate machine-learning-based

approaches. For industrial systems specifically, very few datasets are publicly available, which significantly hampers the research efforts to develop new adequate detection methods. Existing datasets for industrial systems include several datasets developed by the Mississippi State university for power, water and gas systems [6], and the SWaT dataset [2], which is described in more details in the next section.

An additional problem is that existing datasets can quickly become outdated because of technology development. For example, the previously mentioned DARPA and NSL-KDD were made from data captured in 1999. Since then, both technologies used in networks and existing threats have significantly evolved and thus the benefits from using these datasets today are very minor. It is also very important to have a properly labelled dataset: one can easily imagine what sort of problems could emerge if some attacks were mis-characterized as normal data for example.

The SWaT testbed

In the following, we will be using the SWaT (Secure Water Treatment) dataset [2] to conduct experiments and test our methods. The SWaT system is a testbed developed by Singapore University of Technology and Design. The goal of this testbed is to offer a platform that replicates a modern industrial system that is accessible to researchers, so that they can test new security systems under realistic conditions and conduct attacks on a real system. The physical state of the system is recorded as well as all network traffic between its components. This testbed replicates the different processes that compose a water treatment system, including water collection, quality assessment, and filtration. Figure 2 details the different steps of this process, which can be decomposed in six different sub-processes named from P1 to P6.

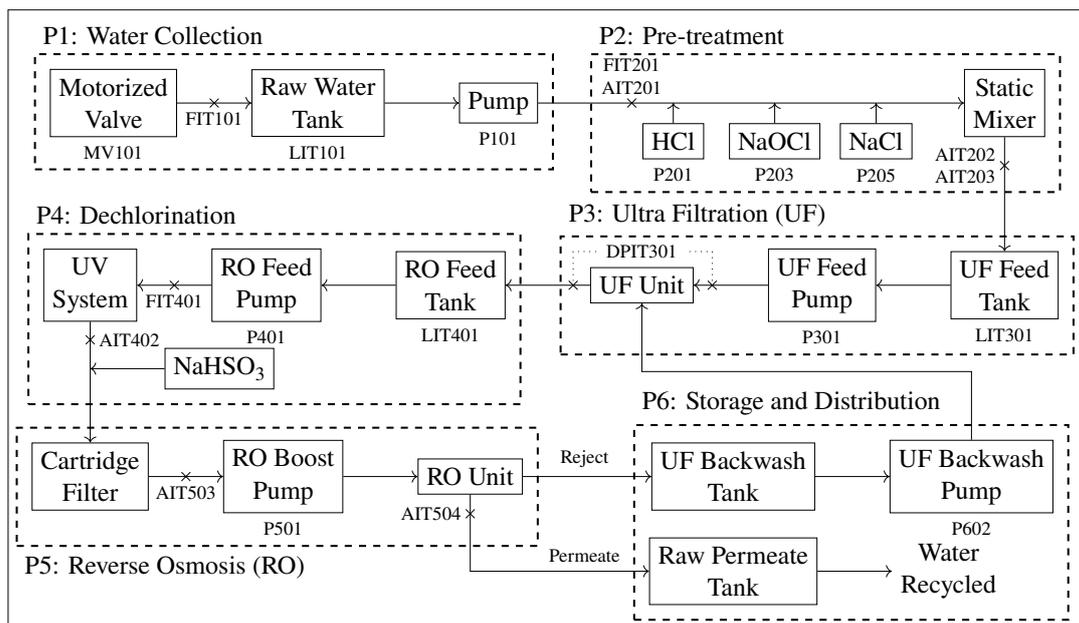


Figure 2: SWaT Process

Adapted by permission from Springer Nature: Jonathan Goh et al. "A Dataset to Support Research in the Design of Secure Water Treatment Systems". In: *Critical Information Infrastructures Security*. Ed. by Grigore Havarneanu et al. Lecture Notes in Computer Science. Cham: Springer International Publishing, 2017, pp. 88–99. ISBN: 978-3-319-71368-7. DOI: 10.1007/978-3-319-71368-7_8

This system is composed of 26 actuators and 25 sensors, among which we can find pumps, valves, flow meters, water level sensors, etc. The state of each actuator and the measurement of each sensor is recorded once per second. In the available dataset, the system has been recorded continuously over a period of 11 days. The system was first left running under normal operating conditions for 6 days. Then during the next 5 days, a total of 36 attacks were conducted on the system in an attempt to disrupt the process.

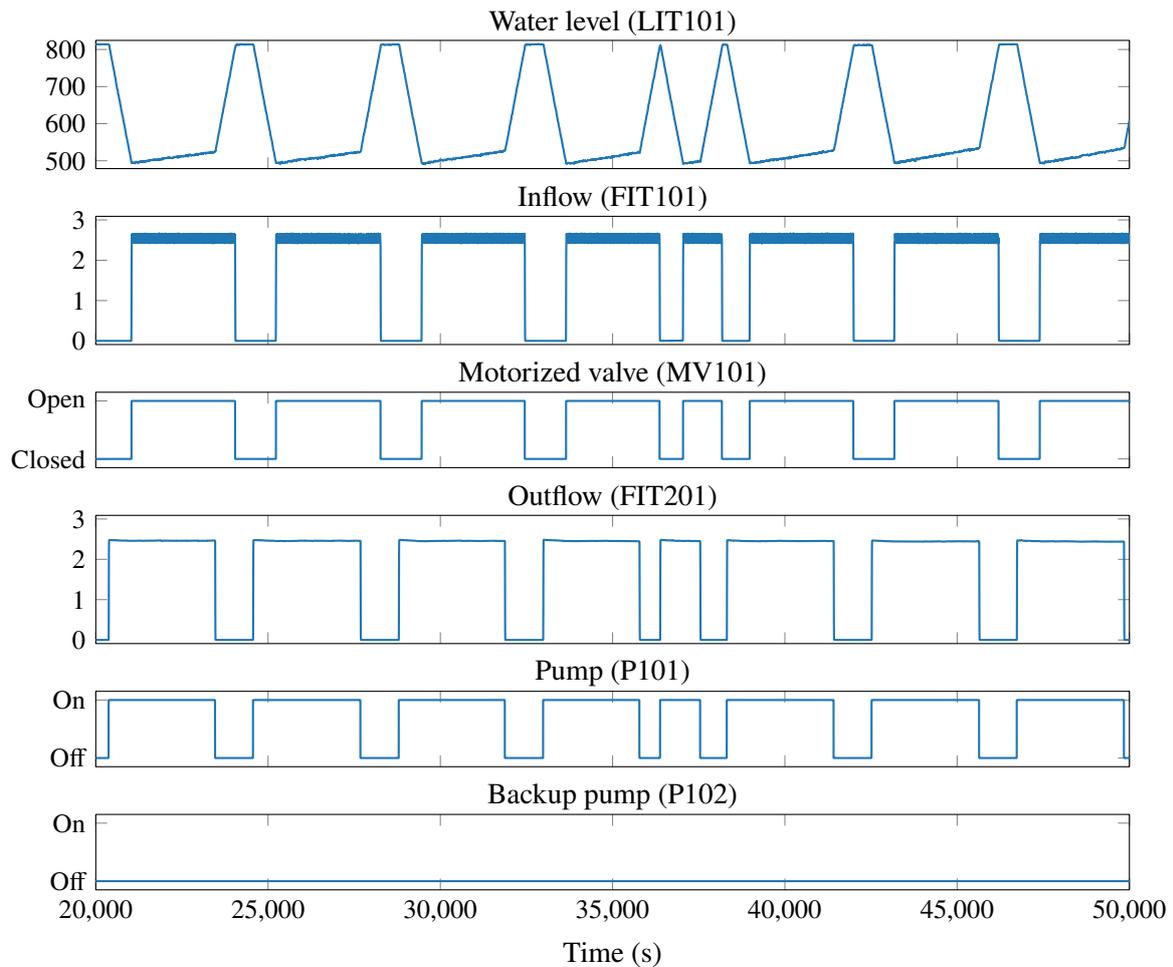


Figure 3: Example of signals recorded under normal operating conditions in the SWaT testbed

As an example, three sensors and three actuators of this system are displayed in Figs. 3 and 4. The three sensors respectively measure the water level (LIT101) in a tank of the first subprocess P1, the inflow rate (FIT101) to and the outflow rate (FIT201) from this tank. The inflow is controlled by a motorized valve (MV101) and the outflow is controlled by two pumps (P101 and P102). The second pump P102 is a backup pump which is only supposed to run if the main one fails.

Figure 3 displays the measurements of these sensors and the state of these actuators under normal operating conditions. We can observe here a regular cycle in each of these signals which is characteristic of the normal behavior of the system, with occasional variations: here for example we have a regular cycle that takes approximately one hour, with one notable exception between 35,000 s and 40,000 s where we can observe an instance of a much shorter cycle.

Figure 4 shows the same sensors and actuators, this time with three attacks conducted during the displayed time period. The areas marked with a red background indicate the instants during which attacks were conducted, as labelled in the dataset. As an aside, it is important to remember that while this red background displays the timestamps during which an attack was conducted, it does not mean that this attack affected all of or only the displayed sensors and actuators. It may affect a single sensor, or even none of the displayed ones, and it could also affect directly or indirectly other components of the system. It is difficult to describe the exact effects of an attack on the system as they can remain localized or propagate to other components depending on the nature of the attack. The description of the attacks given by the authors of the dataset state the intent of each attack and how it was conducted, but only give a short description of its consequences on the system, mainly whether the attack was successful or not, and sometimes a brief

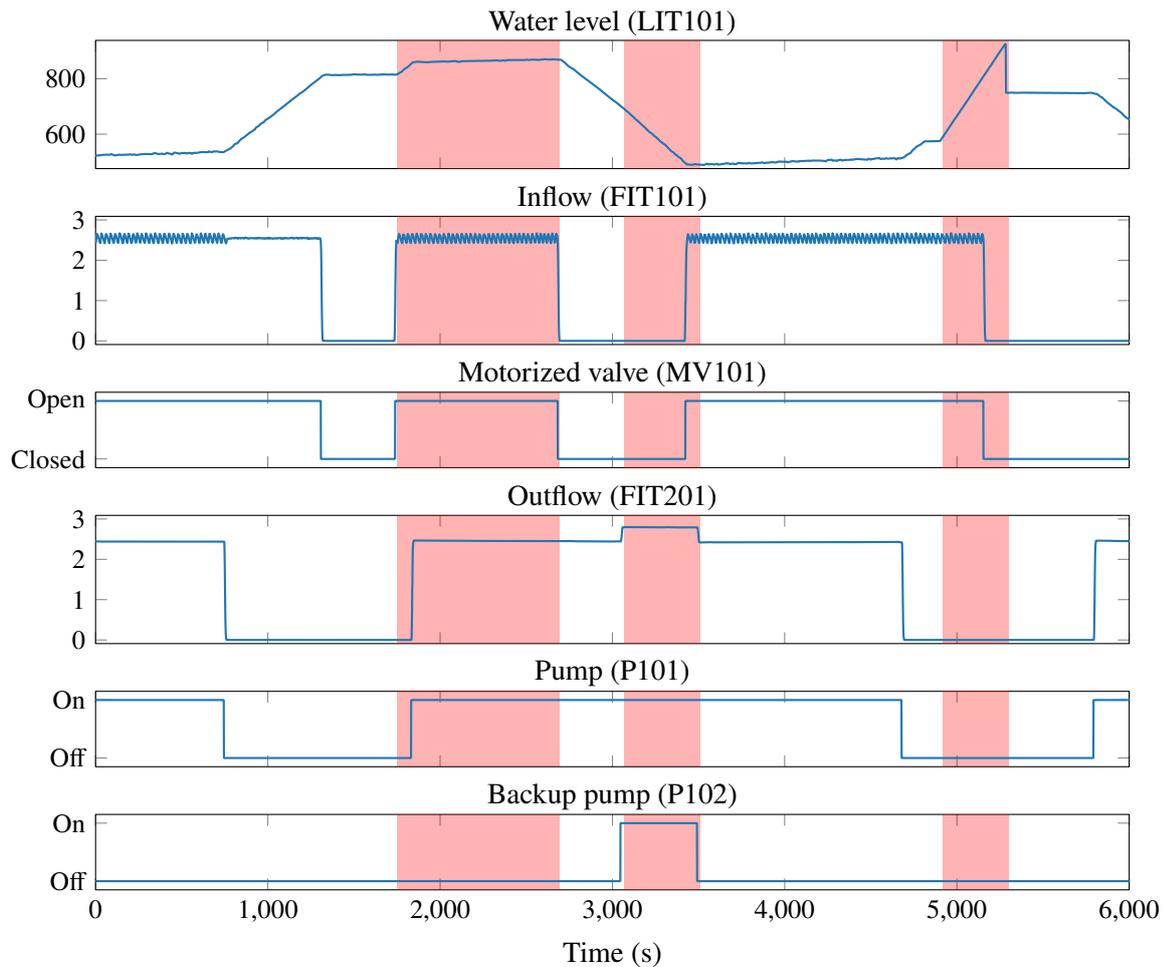


Figure 4: Example of signals recorded in presence of attacks

description of some of its side-effects.

To give an example of attacks that were conducted on this system, we give a brief description of the three attacks visible in Fig. 4:

1. The first attack (from 1,800 s to 2,700 s) consisted in forcing the inflow valve MV101 open, in an attempt to overfill the water tank.
2. The second attack (from 3,000 s to 3,500 s) consisted in forcing both pumps P101 and P102 to be on at the same time. This could cause excessive flow in the pipes and may damage both the pumps and the pipes.
3. This third attack (around 5,000 s) consisted in spoofing the sensor LIT101: instead of observing the real value of this sensor, the PLC receives a value of the attacker's choice, who in this case simulates a rapid increase in the water level.

This example displays two kinds of attack that can be conducted on such a system:

- Manipulating an actuator directly, causing it to potentially put the system in a dangerous state.
- Spoofing the value of a sensor, causing the system to have an erroneous view of the system and react wrongly.

In more complex scenarios, an attacker could combine these approaches in order to maximize damage or remain stealthy over a long period of time.

Objectives and contributions

Our initial ambition when we started this thesis was to develop a complete anomaly-based detection method, which would be able to learn the normal behavior of a system from a database of signals, and then use this knowledge to detect attacks by comparing the current state of the system to what we would expect from the learned model. This is the basic principle of *physics-based attack detection*, a survey of which can be found in [7].

One of the main issues when developing such a method is how to obtain an accurate model of the monitored system. Indeed, the detection performance of these methods will depend highly on the quality of the model. An inaccurate model may lead to a high number of false-alarms, making it unusable in practice, or it could also leave some attacks undetected. Since we are looking at physical systems, one option to build a model would be to use the underlying equations of physics that regulate the studied phenomena. However, this can quickly become intractable with many components, as their interactions become increasingly complex for larger systems. This renders this approach unusable beyond very small systems with few components.

An alternative approach consists in considering some mathematical model that uses a certain number of parameters, and fitting this model to data collected on the system to estimate these parameters. This is how *system identification* methods operate [8, 9], which can offer Auto-Regressive or State-Space representations of a system, among others. These methods usually rely on the ability to conduct dedicated experiments on the studied system, as they do not simply rely on data recorded during the operation of the system, but rather study how the system reacts to specific inputs to model its behavior. This might be conceivable for a system that is being designed and not yet operational, but it may not be possible to interrupt a currently operational process to conduct such experiments.

Yet another approach that can be considered is learning-based, i.e. using only data to build a model, without any prior knowledge of the system. Neural networks are a commonly encountered example of this principle, using a sufficiently large database to learn a model of some phenomenon through a training phase. A few examples of neural networks used to learn the normal behavior of industrial systems can be found in [10, 11, 12, 13]. While neural networks have proven to be a powerful tool in some fields (e.g. image classification, pattern recognition, among many others), they also present several downsides that can make their use problematic in some situations. A first important issue is that they require a large enough, properly labelled dataset to train. Indeed, the model obtained after training depends highly on the training dataset, and if this dataset is not exhaustive, the resulting model would be limited to the situations encountered during that training, without any guarantee of what may happen in other cases. A second important issue is explainability, which is not specific to neural networks, but applies to the entire field of artificial intelligence [14, 15]. Explainability refers to the possibility of understanding why a certain system gives a specific output for some input. For example, taking the examples of attacks shown on Fig. 4, a system trained to detect attacks might simply indicate that there is an attack occurring at a given timestamp without any other context. Without any explanation, it is up to the operator to figure out first whether the system is actually under attack, and then figure out what the attack is in order to take appropriate action. An explainable system could give more insight about the nature of the attack: it may pinpoint some sensors that are affected by the attack (for example sensor LIT101 in the case of the first attack), or even offer an interpretation of the attack (for example describing that the value of the sensor LIT101 exceeds its nominal value). Such explanations can therefore lead to faster and more appropriate responses to attacks, which can help minimize damage and downtime. Having these explanations may also help understand the cause of potential errors, which then indicate how to improve the system, whether through modifying its design or improving the dataset used for training. Several approaches have been proposed to attempt to extract explanations from different machine learning methods, which are summarized in [15]. However, deep neural networks remain challenging in that regard because of the large number of parameters and layers that make finding relevant features difficult.

Overall, obtaining an accurate and understandable model of a given system remains a complex problem

that is still largely unsolved. As such, we have to keep in mind that any model that we can obtain of a given system is a simplified representation of a complex physical system. In addition, “even if the model is reasonably good, our knowledge of the parameters in it, e.g., covariance functions, time constants, etc., may not be enough to justify a direct numerical evaluation of formulas derived from the model” [16]. To mitigate this mismatch between model and reality, we will be relying on the RDT (Random Distortion Testing) approach [17], which was designed with this issue in mind. The approach stems from statistical decision theory, and incorporates model mismatch in the problem statement itself, while also offering an optimality and the ability to control the false-alarm rate. More precisely, this approach consists in deciding whether some noisy signal drifts by too much from a certain deterministic model, without requiring any prior knowledge of the probability distribution of the signal. The only assumption made is that the observation results from some unknown signal in presence of additive and independent Gaussian noise, with perfectly known covariance matrix. However, having perfect knowledge of the noise properties is questionable in practice. Indeed, there are many cases where the noise covariance matrix is not known, but estimated instead using real data. In addition, while assuming that the noise is Gaussian is common practice, this is not always the case and may lead to losing the optimality and robustness properties of this approach. In our attempts to use the RDT approach on industrial signals, we have encountered both of these situations, which are not covered by the RDT approach [17] or its currently existing extensions [18, 19, 20]. We have therefore developed two extensions of the RDT approach, one designed to account for other noise distributions, leveraging the invariance properties on which the RDT framework is based on, and the other to take the noise variance estimation into account for Gaussian noise, when the noise components are independent and identically distributed.

We will then make use of the RDT framework and of these newly developed extensions to develop the basis for an anomaly detection method, applicable to the physical signals of an industrial system. This method consists in first learning the characteristics of the different phases of any given signal, and then detecting whenever the signal properties no longer match what was previously learned. Our main contribution in this regard is the development of a change-in-mean detection method, allowing one to segment a signal to find the different regime changes that are present. We will present this method and demonstrate it on signals from the SWaT dataset we presented earlier, then we will discuss some potential ideas for a complete anomaly detection method based on this change-detection algorithm.

The performance evaluation that will be presented on real signals will be limited to qualitative observations on the results yielded by our methods. Indeed, when performing our tests on real signals, we found that we were very limited by the dataset. This dataset only contains 36 attacks, which is not sufficient to get a good idea of the detection performance of our methods. In addition, we found that only having labels indicating whether an attack was in progress or not was quite limiting, notably to evaluate the quality of our segmentation obtained with our change detection method. Changes on the signals occur regardless of the presence of an attack, and the dataset does not provide any ground truth regarding these changes, making it difficult to evaluate the quality of the resulting segmentation. As such, we could only conduct limited experiments on real data, and we instead decided to focus more on the theoretical aspects of the RDT approach.

Thesis structure

Prior to the first chapter, a brief chapter will introduce some of the main notations that will be used throughout this thesis, and also recall a few important mathematical notions.

The first chapter will introduce the RDT approach, and compare it to usual hypothesis testing approaches to highlight the differences that exist between them, in terms of the problem structure and of the optimality criterion used in either case. We will also briefly discuss the robustness of tests and explain what makes the RDT approach a good candidate to find tests that are robust to model mismatches.

The second chapter will present our first extension of the RDT approach, which we named Generalized RDT, in which we attempt to make the RDT approach applicable to a wider class of noise distributions

beyond the usual Gaussian distribution. We will start this chapter with a discussion on invariance, to reintroduce this notion in the context of hypothesis testing, and then leverage it to redevelop the RDT approach assuming that the noise distribution is not necessary Gaussian, but invariant with respect to some group of transformations of the space \mathbb{R}^d .

The third chapter will then introduce Asymptotic RDT, our second extension of the RDT approach, which consists in accounting for the estimation of the parameters of the RDT test, namely the noise variance σ and the deterministic model θ_0 . We will present this extension in two steps. First we present theoretical results in the asymptotic case, showing that we can asymptotically reach the theoretical performance of the test as the estimators converge. Then we will present simulation results to study the behavior of the RDT test in the non-asymptotic case, in order to help quantify the performance loss incurred by the use of estimates instead of perfectly known values.

Finally, the fourth and final chapter will present applications of this work to cybersecurity on real signals from the SWaT dataset. We will present two detection methods based on the RDT framework that we developed and test them on the SWaT dataset. The first method consists in detecting discontinuities on physical signals that are supposed to be continuous, and offers a way to learn the appropriate parameters using a clean dataset. The second method is a change-in-mean detection method, allowing segmentation of signals to identify the different phases that are present. We will then discuss some ideas based on this method to build a complete learning-based anomaly detection method.

The conclusion will recapitulate the different results that we presented, and will also evoke some perspectives for our work. We will present some ideas to further extend the RDT approach, as well as possibilities for a full anomaly detection method based on the detection methods presented in the fourth chapter.

Several appendices are present after the conclusion. They include additional details on some aspects presented throughout this thesis.

Preliminary: mathematical notions and notations

This preliminary chapter gathers the notations that will be used throughout this thesis, and also contains a few reminders regarding some important mathematical notions.

Throughout this thesis, we consider a probability space $(\Omega, \Sigma, \mathbb{P})$ with respect to which all random variables and vectors are defined. For some set E , $\mathcal{M}(\Omega, E)$ designates the set of all measurable applications (i.e. random variables) defined on $(\Omega, \Sigma, \mathbb{P})$ that take their values in E .

For any integer $d \geq 1$, we denote by \mathcal{F}^d the Borel σ -algebra of \mathbb{R}^d .

Continuity set. For a random variable $X \in \mathcal{M}(\Omega, \mathbb{R}^d)$, a Borel set $B \in \mathcal{F}^d$ is said to be a $\mathbb{P}X^{-1}$ -continuity set if:

$$\mathbb{P}[X \in \partial B] = 0$$

where ∂B is the boundary of B : $\partial B = \bar{B} \setminus \mathring{B}$, where \bar{B} is the closure of B , and \mathring{B} is the interior of B .

Convergence of sequences of random variables. We recall below the three common modes of convergences for sequences of random variables. Let $X \in \mathcal{M}(\Omega, \mathbb{R}^d)$ be a random vector $(X_n)_{n \in \mathbb{N}} \in \mathcal{M}(\Omega, \mathbb{R}^d)^{\mathbb{N}}$ be a sequence of random vectors.

- The sequence (X_n) converges in distribution to X if for every $\mathbb{P}X^{-1}$ -continuity set A of \mathcal{F}^d we have:

$$\mathbb{P}[X_n \in A] \rightarrow \mathbb{P}[X \in A]$$

Convergence in distribution will be denoted as $X_n \xrightarrow{\mathcal{D}} X$.

- The sequence (X_n) converges in probability to X with respect to the probability measure \mathbb{P} if:

$$\forall \epsilon > 0, \mathbb{P}[\|X_n - X\|_2 \geq \epsilon] \rightarrow 0$$

Convergence in probability with respect to \mathbb{P} will be denoted as $X_n \xrightarrow{\mathbb{P}} X$.

- The sequence (X_n) converges almost surely to X if:

$$\mathbb{P}[\lim_{n \rightarrow \infty} X_n = X] = 1$$

or equivalently:

$$\mathbb{P}(\{\omega \in \Omega : X_n(\omega) \rightarrow X(\omega)\}) = 1$$

Almost sure convergence will be denoted as $X_n \xrightarrow{\text{a.s.}} X$.

We also remind that almost-sure convergence implies both convergence in probability and in distribution, and that convergence in probability implies convergence in distribution.

Conditional probability. Given an event $A \in \Sigma$ and a random variable $X \in \mathcal{M}(\Omega, \mathbb{R}^d)$, the *conditional probability* $\mathbb{P}(A \mid X = x)$ is the only element of $L^1(\mathbb{R}^d, \mathcal{F}^d, \mathbb{P}X^{-1})$ that satisfies:

$$\forall B \in \mathcal{F}^d, \mathbb{P}(A \cap [X \in B]) = \int_B \mathbb{P}(A \mid X = x) \mathbb{P}X^{-1}(dx) \quad (11)$$

Pushforward measure (or image measure). Given two measurable sets (X, \mathcal{X}) and (Y, \mathcal{Y}) , a measure μ defined on (X, \mathcal{X}) and a measurable function $f: (X, \mathcal{X}) \rightarrow (Y, \mathcal{Y})$, the pushforward measure (or image measure) μf^{-1} is the measure on (Y, \mathcal{Y}) defined by:

$$\begin{aligned} \mu f^{-1}: \mathcal{Y} &\rightarrow [0, +\infty) \\ B &\mapsto \mu f^{-1}(B) = \mu(f^{-1}(B)) \end{aligned} \quad (12)$$

We also recall that f is a measurable function from (X, \mathcal{X}) to (Y, \mathcal{Y}) if:

$$\forall B \in \mathcal{Y}, f^{-1}(B) \in \mathcal{X} \quad (13)$$

1 Hypothesis Testing and Random Distortion Testing

Contents

1.1	Traditional hypothesis testing approaches	31
1.1.1	Non-Bayesian binary classification	31
1.1.2	Optimality	33
1.2	The RDT approach	36
1.2.1	Problem statement	36
1.2.2	Optimality: γ -MCCP tests	37
1.2.3	Thresholding tests and optimality	39
1.3	Robustness and hypothesis testing	40
1.3.1	Huber's approach to robust hypothesis testing	41
1.3.2	RDT and robust hypothesis testing	43

In this first chapter, we introduce the RDT (Random Distortion Testing) theory on which the rest of this thesis is based. In short, the RDT approach considers the problem of testing whether a signal observed in presence of noise lies within a given ball or not. As such, this problem can be seen more as a topological testing problem, rather than a usual hypothesis testing problem. To explain this difference in more details, we will start by going through some basics of non-Bayesian hypothesis testing. We have chosen to restrict our study to non-Bayesian approaches, as these approaches do not require as much prior knowledge on the problem as Bayesian approaches do. For example, if we consider a problem that would consist in testing whether there is an attack or not on a system, a Bayesian approach would need to know the a priori probability that the system is under attack, which is pretty much impossible to reasonably quantify. Following this, we will then present the main elements of the RDT theory, as they were first introduced in [17]. We will explain why we think this is an interesting approach that can be used to build robust methods, as an approach that requires few parameters to set and few assumptions on the signal of interest. Since robustness is an important aspect to consider in the cybersecurity field, we will also take a look at Huber's approach to robust hypothesis testing [21] and compare that to the RDT approach.

1.1 Traditional hypothesis testing approaches

1.1.1 Non-Bayesian binary classification

An hypothesis testing problem can be formulated as follows: given an observation y (scalar or vector) which is a realization of some random variable Y , we want to get some information on the probability distribution P of Y . The information we get is whether this probability distribution belongs to some set of probability distributions. Here we only consider binary classification. We can formalize a non-bayesian decision problem as follows (\mathcal{P} : set of all probability distributions):

Non-bayesian decision problem

$$\left. \begin{array}{l}
 \textbf{Data model:} \\
 \exists \mathcal{F} \subset \mathcal{P}, \exists P \in \mathcal{F}, \exists Y \in \mathcal{M}(\Omega, \mathbb{R}^d), \\
 \left\{ \begin{array}{l} (Y \sim P) \\ \wedge (\forall y \in \mathbb{R}^d, \exists \omega \in \Omega, y = Y(\omega)) \end{array} \right. \\
 \textbf{Testing problem:} \\
 \text{Given one realization } y = Y(\omega), \text{ determine whether:} \\
 \left\{ \begin{array}{l} \mathcal{H}_0: P \in \mathcal{F}_0 \\ \text{or } \mathcal{H}_1: P \in \mathcal{F}_1 \end{array} \right. \\
 \text{with } \mathcal{F}_0 \cup \mathcal{F}_1 = \mathcal{F} \text{ and } \mathcal{F}_0 \cap \mathcal{F}_1 = \emptyset
 \end{array} \right\} \quad (1.1)$$

\mathcal{F} , \mathcal{F}_0 and \mathcal{F}_1 are known sets of probability distributions, while $P \in \mathcal{F}$ is an unknown probability distribution. \mathcal{H}_0 is called the *null hypothesis* and \mathcal{H}_1 the *alternative hypothesis*. If either of the sets \mathcal{F}_0 or \mathcal{F}_1 contains only a single probability distribution, the associated hypothesis is said to be a *simple hypothesis*. Otherwise it is said to be a *composite hypothesis*.

Instead of working directly with distributions, it is common to work with a parameter θ (either scalar or vector) associated to each possible distribution and define the problem using that parameter. For example, we could consider a problem where we consider all 1-dimensional Gaussian distributions with variance 1 and mean $\theta \in \mathbb{R}$, and we want to test whether $\theta \leq 0$ or $\theta > 0$.

Given such a problem, the issue is to figure out how to decide whether $P \in \mathcal{F}_0$ or $P \in \mathcal{F}_1$. This is done using a test function, which associates to any possible observation $y \in \mathbb{R}^d$ a decision.

Definition 1.1.1: Test and critical region

A *test* is any function defined on \mathbb{R}^d that can take the values 0 or 1:

$$\begin{array}{rcl}
 \mathcal{T}: \mathbb{R}^d & \rightarrow & \{0, 1\} \\
 y & \mapsto & \begin{cases} 0 & \text{if } \mathcal{H}_0 \text{ is accepted} \\ 1 & \text{if } \mathcal{H}_0 \text{ is rejected} \end{cases}
 \end{array} \quad (1.2)$$

For any test \mathcal{T} , the set $K_{\mathcal{T}} = \mathcal{T}^{-1}(\{1\})$ is called the *critical region* of \mathcal{T} .

Here we are actually considering a slightly simplified definition of tests, which does not apply to *randomized tests*. While randomized tests are important from a theoretical viewpoint — notably for the Neyman-Pearson lemma that will be presented in the following —, we will not introduce them here to simplify this brief presentation of non-bayesian hypothesis testing.

It is common to consider the null hypothesis \mathcal{H}_0 as a nominal situation, whereas the alternative hypothesis \mathcal{H}_1 tends to be a rare case that we want to detect. With this interpretation, we should keep in mind that the problem is asymmetrical: the order of the hypotheses is important and swapping them around can change the meaning of the problem, as well as the resulting tests depending on the criterion used to choose them.

For any given problem, there are four possible outcomes:

1. Accept \mathcal{H}_0 and \mathcal{H}_0 is true (true negative)
2. Accept \mathcal{H}_0 and \mathcal{H}_1 is true (Type II error, i.e. false negative or missed detection)
3. Reject \mathcal{H}_0 and \mathcal{H}_0 is true (Type I error, i.e. false positive or false alarm)
4. Reject \mathcal{H}_0 and \mathcal{H}_1 is true (true positive)

Being able to quantify the occurrences of each of these outcomes is important to evaluate the performance of a given test. This is usually done through the notions of *size*, *level*, and *power function*, which help measure the probability of each of these four outcomes.

Definition 1.1.2: Level, size and power function of a test

Let $\mathcal{T}: \mathbb{R}^d \rightarrow \{0, 1\}$ be a test.

The *size* $\alpha_{\mathcal{T}}$ and the *power function* $\beta_{\mathcal{T}}$ of \mathcal{T} are defined as:

$$\text{Power function: } \forall P \in \mathcal{F}, \beta_{\mathcal{T}}(P) = \mathbb{P}[\mathcal{T}(Y) = 1] \quad \text{with } Y \sim P$$

$$\text{Size: } \alpha_{\mathcal{T}} = \sup_{P \in \mathcal{F}_0} \beta_{\mathcal{T}}(P)$$

A test \mathcal{T} is said to have *level* γ if $\alpha_{\mathcal{T}} \leq \gamma$.

These metrics help characterize the behavior of any given test. The power function simply defines for a given test the probability to reject the null hypothesis \mathcal{H}_0 , for each probability distribution of interest. The size of a test can be interpreted as its maximum false-alarm rate, i.e. the maximal probability to wrongly reject the null hypothesis \mathcal{H}_0 . Ideally, we would want a test that can perfectly reject \mathcal{H}_0 with no false alarm whatsoever. Of course, finding such a test is not possible in the vast majority of cases. Instead, we have to balance the false-alarm and detection rates in order to find an appropriate equilibrium. This is usually done by defining an optimality criterion, and then finding an optimal test with regards to this criterion.

1.1.2 Optimality

A widely used criterion to find an optimal test for a given testing problem is the UMP (Uniformly Most Powerful) criterion.

Definition 1.1.3: Uniformly Most Powerful (UMP) Test

A test \mathcal{T}^* is said to be *Uniformly Most Powerful (UMP)* with level γ if:

1. \mathcal{T}^* has level γ .
2. For any other test \mathcal{T} with level γ , we have:

$$\forall P \in \mathcal{F}_1, \beta_{\mathcal{T}^*}(P) \geq \beta_{\mathcal{T}}(P)$$

UMP tests are desirable for two reasons:

1. First, it has a given level γ , which means it ensures that its false alarm-rate is controlled and will not exceed this value. This is important in contexts where false-alarms can be costly to deal with. To get back to the cybersecurity context we presented in the introduction, false alarms will usually require manual investigation to diagnose the source of the alarm and verify whether it is real or not. This manual intervention can be costly, in terms of both money and time, and having the ability to control the maximum false-alarm rate also means keeping these costs under control.
2. Second, it maximizes the detection rate for every possible distribution under \mathcal{H}_1 , meaning that no other test with level γ can offer a higher detection rate for any distribution under \mathcal{H}_1 .

In this sense, UMP tests are the best tests when one is willing to accept a false-alarm rate up to γ . However, while UMP tests are desirable, they are rare in practice. In fact, for most problems, no such test exists, and we therefore have to rely on suboptimal tests. Note that a suboptimal test does not mean that it is bad, but depending on the circumstances, it may not perform as well as another one.

Nonetheless, there are two important cases where a UMP test exists. These cases are covered by the Neyman-Pearson lemma [22] and the Karlin-Rubin theorem [23]. We will start by presenting the Neyman-Pearson lemma, which gives a most powerful test for any binary hypothesis testing problem composed of two simple hypotheses. As a reminder, we do not consider randomized tests in this presentation, hence the expression of the test given here will be slightly simplified.

Theorem 1.1.4: Neyman-Pearson Lemma

Let P_0 and P_1 be two probability distributions with respective probability density functions f_0 and f_1 . Let $Y \in \mathcal{M}(\Omega, \mathbb{R}^d)$ be a random variable with unknown probability distribution $P \in \{P_0, P_1\}$. Let $\gamma \in (0, 1)$.

Consider the following hypothesis problem:

$$\begin{cases} \mathcal{H}_0: P = P_0 \\ \mathcal{H}_1: P = P_1 \end{cases} \quad (1.3)$$

Let \mathcal{J}_{NP} be the test defined by:

$$\begin{aligned} \mathcal{J}_{NP}: \mathbb{R}^d &\rightarrow \{0, 1\} \\ y &\mapsto \begin{cases} 1 & \text{if } f_1(y) > \lambda_{NP} f_0(y) \\ 0 & \text{otherwise} \end{cases} \end{aligned} \quad (1.4)$$

where $\lambda_{NP} \geq 0$ is a threshold chosen such that $\mathbb{P}[\mathcal{J}_{NP}(Y) = 1] = \gamma$ when $Y \sim P_0$.

The test \mathcal{J}_{NP} is most powerful with level γ for testing $\mathcal{H}_0: P = P_0$ against $\mathcal{H}_1: P = P_1$.

As a short note, in the statement of the Neyman-Pearson lemma, we mentioned that the test \mathcal{J}_{NP} is “most powerful” instead of UMP. This is completely equivalent, as the term “uniformly” signify that the test is most powerful for every possible distribution under \mathcal{H}_1 . Since both hypotheses are simple, the term “uniformly” is commonly dropped as it does not add anything in this particular case.

For composite hypotheses, the Karlin-Rubin theorem presents one particular case where it is possible to find a UMP test. Before stating this theorem, we first need to define the notion of *Monotone likelihood ratio*, which plays a key role.

Definition 1.1.5: Monotone Likelihood Ratio

Let $\mathcal{F} = \{P_\theta, \theta \in \Theta\}$ be a family of probability distributions defined on \mathbb{R} , indexed by $i \in \Theta \subset \mathbb{R}$, each distribution P_θ having probability density function f_θ . Let $T: \mathbb{R} \rightarrow \mathbb{R}$ be a function. The family \mathcal{F} is said to have a *monotone likelihood ratio in T* if for any pair $(\theta, \theta') \in \Theta^2$ such that $\theta < \theta'$, the function $f_{\theta'}/f_\theta$ is a strictly increasing function of T :

$$\forall (\theta, \theta') \in \Theta^2, \theta < \theta' \Rightarrow (\exists g: \mathbb{R} \rightarrow \mathbb{R}, (g \text{ is increasing}) \wedge (f_{\theta'}/f_\theta = g \circ T)) \quad (1.5)$$

which we can rewrite as:

$$\forall (\theta, \theta') \in \Theta^2, \forall (x, x') \in \mathbb{R}^2, (\theta < \theta' \wedge T(x) < T(x')) \Rightarrow \frac{f_{\theta'}(x)}{f_\theta(x)} \leq \frac{f_{\theta'}(x')}{f_\theta(x')} \quad (1.6)$$

An interpretation of this property is that as $T(x)$ increases, higher values of the parameter $\theta \in \Theta$ become increasingly likely in terms of the likelihood ratio.

The monotone likelihood ratio property can also be found under the name *total positivity of order 2* (TP-2) when applied to functions of two parameters, which are not necessarily probability density functions [24]. For a strictly positive function $k: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ of two parameters, k is TP-2 if it verifies:

$$\forall (x_1, x_2, y_1, y_2) \in \mathbb{R}^4, (x_1 < x_2 \wedge y_1 < y_2) \Rightarrow k(x_1, y_1)k(x_2, y_2) \geq k(x_1, y_2)k(x_2, y_1) \quad (1.7)$$

We can see that this is a particular case of Eq. (1.6) with the two parameters being θ and x , and using the identity function for T .

Here are a few examples of families of distribution with monotone likelihood ratio, parametrized by θ :

- The family of 1-dimensional Gaussian distributions $\mathcal{N}(\theta, \sigma^2)$ with known variance σ^2 , for $\theta \in \mathbb{R}$.
- The family of 1-dimensional Gaussian distributions $\mathcal{N}(\mu_0, \theta^2)$ with known mean μ_0 , for $\theta > 0$.
- The family of non-central Chi-squared distributions $\chi_{d,\theta}^2$ with d degrees of freedom and non-centrality parameter $\theta \geq 0$
- The one-parameter exponential families of distributions [25, Corollary 3.4.1], which are families of distribution having a density p_θ of the form:

$$p_\theta(x) = C(\theta)e^{Q(\theta)T(x)}h(x) \tag{1.8}$$

where Q is a non-decreasing function. These families of distributions have a monotone likelihood ratio in T .

With this definition, we can now state the Karlin-Rubin theorem, which presents another case where we can find a UMP test, this time for composite hypotheses.

Theorem 1.1.6: Karlin-Rubin Theorem

Let $\mathcal{F} = \{P_\theta, \theta \in \Theta\}$ be a family of probability distributions defined on \mathbb{R} , indexed by $\theta \in \Theta \subset \mathbb{R}$. Assume that the family \mathcal{F} has monotone likelihood ratio in a function $T: \mathbb{R} \rightarrow \mathbb{R}$. Let $Y \in \mathcal{M}(\Omega, \mathbb{R})$ be a random variable with unknown probability distribution $P_\theta \in \mathcal{F}$. Let $\theta_0 \in \Theta$ and $\gamma \in (0, 1)$. Consider the following hypothesis testing problem:

$$\begin{cases} \mathcal{H}_0: \theta \leq \theta_0 \\ \mathcal{H}_1: \theta > \theta_0 \end{cases} \tag{1.9}$$

Let \mathcal{J}_{KR} be the test defined by:

$$\begin{aligned} \mathcal{J}_{KR}: \mathbb{R} &\rightarrow \{0, 1\} \\ y &\mapsto \begin{cases} 1 & \text{if } T(y) > \lambda_{KR} \\ 0 & \text{otherwise} \end{cases} \end{aligned} \tag{1.10}$$

where $\lambda_{KR} \in \mathbb{R}$ is a threshold chosen such that $\mathbb{P}[\mathcal{J}_{KR}(Y) = 1] = \gamma$ when $Y \sim P_{\theta_0}$.

1. The test \mathcal{J}_{KR} is UMP with level γ for testing $\mathcal{H}_0: \theta \leq \theta_0$ against $\mathcal{H}_1: \theta > \theta_0$.
2. Its power function $\beta_{\mathcal{J}_{KR}}$ is strictly increasing for any $\theta \in \Theta$ such that $\beta_{\mathcal{J}_{KR}}(\theta) \notin \{0, 1\}$.

In most other cases however, it is impossible to find a UMP test. One way to find tests is to restrict the search to a subset of all sets. Taking invariance properties of the problem into account is one way to do so, resulting in a UMPI (Uniformly Most Powerful Invariant) test. Invariance will be discussed in more details in the next chapter.

If invariance is not sufficient to find an optimal test, then we may have to rely on suboptimal tests. In that case, the choice of test is highly dependent on the problem at hand and the desired properties of the test. It is important to remember that a suboptimal test is not necessarily a poorly performing test. The two notions are actually disconnected from one another. Indeed, a test being suboptimal only means that it is not optimal for a given criterion. For the UMP optimality criterion, a suboptimal test can still have level γ and offer adequate performance. We just have to keep in mind that in some circumstances, other tests may perform better with the same level.

A common problem that does not necessarily have a UMP test consists in testing $\mathcal{H}_0: \theta = \theta_0$ against $\mathcal{H}_1: \theta \neq \theta_0$ for a given θ_0 . Whether a UMP test exists or not depends on the probability distributions associated to each value of θ . In a large number of cases, no such test exists. However, there exists a family of three commonly used test for this particular problem, commonly referred to as the “holy trinity”. These three tests are the Wald test [26], the Rao test [27] and the GLRT (Generalized Likelihood Ratio Test). In short, these three tests rely in some way on the maximum-likelihood estimation of the parameter θ . The Wald and Rao tests rely on asymptotic properties of these estimators, related to the Fisher information matrix. They slightly differ from one another in that the Wald test requires computing the maximum-likelihood estimate $\hat{\theta}$ of θ under \mathcal{H}_1 , whereas this is not necessary for the Rao test. The GLRT reuses the principle of the likelihood ratio test used in the Neyman-Pearson lemma, and simply replaces the densities used to compute the likelihood ratio by the densities that maximize the likelihood under each hypothesis. An interesting property of these three tests is that they are asymptotically equivalent, meaning that as the number of available observations increases and approaches infinity, the three tests will take the same decision with probability 1 [28].

1.2 The RDT approach

1.2.1 Problem statement

After presenting some general notions regarding typical binary hypothesis testing problems, we will now introduce the RDT approach. This approach differs in its objective. As we have seen previously, most hypothesis testing problems consist in testing whether the distribution that generated an observation belongs to one set or the other. In the RDT approach, we do not want to determine whether the distribution that generated the observation belong to one set or the other, but whether the phenomenon of interest lies close enough to some model.

In order to present the RDT problem, we first need to introduce the *Mahalanobis norm*, which is an integral part of this problem.

Definition 1.2.1: Mahalanobis norm

Let C be a $d \times d$ definite positive symmetric matrix.
The *Mahalanobis norm* with matrix C is defined by:

$$\forall y \in \mathbb{R}^d, \nu_C(y) = \sqrt{y^T C^{-1} y} \quad (1.11)$$

We now introduce the RDT problem [17]:

Definition 1.2.2: RDT problem statement

Data model:

$$\left. \begin{aligned} &\exists Y \in \mathcal{M}(\Omega, \mathbb{R}^d), \exists \Theta \in \mathcal{M}(\Omega, \mathbb{R}^d), \exists X \in \mathcal{M}(\Omega, \mathbb{R}^d), \\ &\left\{ \begin{array}{l} (X \sim \mathcal{N}(0, C) \text{ with } C \text{ a } d \times d \text{ definite positive matrix}) \\ \wedge (\Theta \text{ and } X \text{ are independent}) \\ \wedge (Y = \Theta + X) \\ \wedge (\forall y \in \mathbb{R}^d, \exists \omega \in \Omega, y = Y(\omega)) \end{array} \right. \end{aligned} \right\} \quad (1.12)$$

Testing problem:

Given one realization $y = Y(\omega) = \Theta(\omega) + X(\omega)$, determine whether:

$$\left\{ \begin{array}{l} \mathcal{H}_0: \nu_C(\Theta(\omega) - \theta_0) \leq \tau \\ \text{or} \\ \mathcal{H}_1: \nu_C(\Theta(\omega) - \theta_0) > \tau \end{array} \right.$$

with $\tau > 0$ and $\theta_0 \in \mathbb{R}^d$

Let us compare this problem to the general decision problem stated in Eq. (1.1). The data model is an application of what we presented before, with the set of distributions of interest being the set of distributed obtained when convolving any probability distribution with a Gaussian distribution. If we define $\mathcal{F} = \{\mu * f_X, \mu \in \mathcal{P}\}$, we can rephrase the data model like this:

$$\exists P \in \mathcal{F}, \exists Y \in \mathcal{M}(\Omega, \mathbb{R}^d), \begin{cases} (Y \sim P) \\ \wedge (\forall y \in \mathbb{R}^d, \exists \omega \in \Omega, y = Y(\omega)) \end{cases} \quad (1.13)$$

The difference between these problems appears in the testing problem: we do not want to know whether the distribution that generated the observation belongs to one set or another. In fact, we do not even define these subsets of \mathcal{F} . Instead, what we want to determine is whether the realization of the signal lies close enough to some model θ_0 or not. An important consequence of this is that for the same distribution, either hypothesis may be true depending on ω , whereas with traditional hypothesis testing problems, each distribution is associated to exactly one hypothesis.

The noise structure plays an important role in the definition of the RDT problem: it is taken into account through the Mahalanobis norm using the noise covariance matrix C . We give more details about this in the next chapter regarding invariance.

Many approaches to conventional testing problems rely on *likelihood*, whether it is through the likelihood ratio for the Neyman-Pearson lemma, the Karlin-Rubin theorem or the Generalized Likelihood Ratio test for example, or through the properties of the likelihood function as in the Wald and Rao tests for example. Using the likelihood function requires having a family of probability distributions that can be parameterized by some scalar or vector parameter θ , which has always been the case so far in the shown examples prior to this section. For the RDT problem however, the family of distributions that we consider is very large and cannot be parametrized appropriately. Consequently, the likelihood is not an relevant approach to solve the RDT problem.

1.2.2 Optimality: γ -MCCP tests

After introducing the RDT problem and how it differs from usual hypothesis testing problems, we now present the main theoretical results in order to introduce an optimal test. Finding an optimal test requires of first defining an optimality criterion, which in this case is the γ -MCCP (Maximum Constant Conditional Power) criterion, defined in Definition 1.2.7. This optimality criterion relies on notions of size and power which are similar to the ones defined previously, but adapted to the problem at hand.

In the following, we will rely on the notion of *orbits*, which are closely linked to the invariance properties of the RDT problem. We will give a simple definition of these orbits here, a more thorough study will be conducted in Chapter 2.

Definition 1.2.3: Orbits Υ_ρ and set of orbits \mathfrak{F}

For every $\rho \geq 0$, we define the *orbit* Υ_ρ as the set $\Upsilon_\rho = \{y \in \mathbb{R}^d : \nu_C(y - \theta_0) = \rho\}$. We denote by \mathfrak{F} the set of all orbits Υ_ρ : $\mathfrak{F} = \{\Upsilon_\rho, \rho \geq 0\}$.

We can now define the relevant notions of size and power for the RDT problem, which we refer to as *conditional size* and *conditional power*.

Definition 1.2.4: Conditional size and power

Let $\mathcal{T}: \mathbb{R}^d \rightarrow \{0, 1\}$ be a test, $\Theta \in \mathcal{M}(\Omega, \mathbb{R}^d)$ a random variable and $\rho > 0$. We define the *conditional power* and *conditional size* of \mathcal{T} as:

$$\begin{aligned} \text{Conditional power:} & \quad \mathbb{P}[\mathcal{T}(\Theta + X) = 1 \mid \Theta \in \Upsilon_\rho] \\ \text{Conditional size:} & \quad \sup_{\Theta \in \mathcal{M}(\Omega, \mathbb{R}^d): \mathbb{P}[\nu_C(\Theta - \theta_0) \leq \tau] \neq 0} \mathbb{P}[\mathcal{T}(\Theta + X) = 1 \mid \nu_C(\Theta - \theta_0) \leq \tau] \end{aligned} \quad (1.14)$$

Similarly to the definitions of size and power given in Definition 1.1.2, the conditional size and power are probabilities of rejecting the null hypothesis based on the received observation. Notably, the conditional size still represents the maximum false-alarm rate of a given test for the RDT problem. The main difference with the usual definition of size and power function is the use of conditional probabilities. Remember that, in the RDT problem, choosing a certain probability distribution for Θ does not preclude either hypothesis from being true. Conditional probabilities are therefore used to distinguish between the two hypotheses. Since we are using conditional probabilities to study the behavior of tests, it means we are looking at local properties of these tests on each orbit Υ_ρ , whereas the power function used in regular testing problems only considers the global behavior of the test for each distribution of interest. Nonetheless, we still define a power function $\beta_{\mathcal{T}}$ and a size $\alpha_{\mathcal{T}}$ for each test \mathcal{T} in the sense of Definition 1.1.2, as it will prove to be useful and convenient later:

$$\begin{aligned} \text{Power function: } \quad & \forall \theta \in \mathbb{R}^d, \beta_{\mathcal{T}}(\theta) = \mathbb{P}[\mathcal{T}(\theta + X) = 1] \\ \text{Size: } \quad & \alpha_{\mathcal{T}} = \sup_{\theta \in \mathbb{R}^d: \nu_C(\theta - \theta_0) \leq \tau} \beta_{\mathcal{T}}(\theta) \end{aligned} \quad (1.15)$$

These notions of size and power are natural if we consider the problem of testing $\nu_C(\theta - \theta_0) \leq \tau$ against $\nu_C(\theta - \theta_0) > \tau$, where $\theta \in \mathbb{R}^d$ is a deterministic value. The next result shows that this notion of size is relevant to control the conditional size in the RDT problem, and therefore the false-alarm rate.

Lemma 1.2.5: Size and conditional size

For any test $\mathcal{T}: \mathbb{R}^d \rightarrow \{0, 1\}$, we have:

$$\sup_{\Theta \in \mathcal{M}(\Omega, \mathbb{R}^d): \mathbb{P}[\nu_C(\Theta - \theta_0) \leq \tau] \neq 0} \mathbb{P}[\mathcal{T}(\Theta + X) = 1 \mid \nu_C(\Theta - \theta_0) \leq \tau] = \alpha_{\mathcal{T}}$$

Just as previously, we will then say that a test has *level* γ for any $\gamma \in (0, 1)$ if $\alpha_{\mathcal{T}} \leq \gamma$.

Similarly, we can also link the notions of power and conditional power. This is done through the notion of *constant conditional power function*, which will then be used to define our optimality criterion.

Definition 1.2.6: Constant conditional power function given $\Theta \in \Upsilon_\rho$

Let $\Theta \in \mathcal{M}(\Omega, \mathbb{R}^d)$ be a random variable independent of $X \sim \mathcal{N}(0, C)$ and $\rho \geq 0$.

A test \mathcal{T} is said to have *constant conditional power function given $\Theta \in \Upsilon_\rho$* if:

$$\forall \theta \in \Upsilon_\rho, \mathbb{P}[\mathcal{T}(\Theta + X) = 1 \mid \Theta \in \Upsilon_\rho] = \beta_{\mathcal{T}}(\theta)$$

With this definition, we can now present the γ -MCCP optimality criterion:

Definition 1.2.7: γ -MCCP (Maximum Constant Conditional Power) test

Let $\mathcal{T}^*: \mathbb{R}^d \rightarrow \{0, 1\}$ be a test, and let $\gamma \in (0, 1)$.

The test \mathcal{T}^* is said to have level γ and Maximum Constant Conditional Power over \mathfrak{F} — and we simply say that \mathcal{T}^* is γ -MCCP — if:

- (i) \mathcal{T}^* has level γ .
- (ii) For any random variable $\Theta \in \mathcal{M}(\Omega, \mathbb{R}^d)$ and for $\mathbb{P}(\nu_C(\Theta - \theta_0))^{-1}$ -almost every $\rho > \tau$, \mathcal{T}^* has constant conditional power function given $\Theta \in \Upsilon_\rho$.
- (iii) For any random variable $\Theta \in \mathcal{M}(\Omega, \mathbb{R}^d)$, for $\mathbb{P}(\nu_C(\Theta - \theta_0))^{-1}$ -almost every $\rho > \tau$ and for any test \mathcal{T} with level γ and constant conditional power function given $\Theta \in \Upsilon_\rho$, we have:

$$\mathbb{P}[\mathcal{T}^*(\Theta + X) = 1 \mid \Theta \in \Upsilon_\rho] \geq \mathbb{P}[\mathcal{T}(\Theta + X) = 1 \mid \Theta \in \Upsilon_\rho] \quad (1.16)$$

At first glance, this criterion might seem very daunting, but in fact, it is actually similar to the UMPI criterion. Indeed, a UMPI test has to fulfill three criterions: it needs to have level γ , be invariant, and be more powerful than any other invariant test with level γ . These three criterions match the three points described in Definition 1.2.7:

1. The first point regarding the level is straightforward.
2. The second point can be interpreted as a notion of invariance of the power function of the test, instead of the test itself.
3. The third point describes the “most powerful” aspect of the test. Here, we consider the class of tests that have constant conditional power function given $\Theta \in \Upsilon_\rho$, for at least one random vector $\Theta \in \mathcal{M}(\Omega, \mathbb{R}^d)$ and at least one $\rho > \tau$. For each test \mathcal{T} that satisfies these conditions, a γ -MCCP test has greater conditional power on the orbits Υ_ρ on which \mathcal{T} has constant conditional power function.

1.2.3 Thresholding tests and optimality

Now that we have introduced the relevant optimality criterion, we will introduce an optimal test in that sense. The tests we consider here are *thresholding tests* on the Mahalanobis norm ν_C .

Definition 1.2.8: Thresholding tests

For any $t \geq 0$, we define the *thresholding test* \mathcal{T}_t by:

$$\begin{aligned} \mathcal{T}_t: \mathbb{R}^d &\rightarrow \{0, 1\} \\ y &\mapsto \begin{cases} 1 & \text{if } \nu_C(y - \theta_0) > t \\ 0 & \text{otherwise} \end{cases} \end{aligned} \quad (1.17)$$

The reason we are choosing to look at these tests in particular is because of the structure of the RDT problem. Indeed, this problem is invariant through the use of the Mahalanobis norm. We will go more in depth regarding the subject of invariance in the next chapter, but as a brief example, consider the case where we have white Gaussian noise, i.e. $C = \sigma I_d$ with $\sigma > 0$. In this case, the distribution of the noise is spherically invariant, and the Mahalanobis norm — which is then simply the euclidean norm, up to a multiplicative factor — is also invariant by rotation. Since the RDT problem consists in detecting whether the Mahalanobis distance between $\Theta(\omega)$ and the model θ_0 exceeds a certain threshold, it is not unreasonable to consider tests where we put a threshold on the Mahalanobis distance between the observation y and θ_0 .

Of course, we need to choose the threshold t appropriately, and then verify whether the resulting test \mathcal{T}_t is indeed γ -MCCP. The choice of that threshold is directed by the desired level γ and the distribution of the noise X , and more accurately by the distribution of $\nu_C(X)$. To find the optimal threshold that can offer a level γ , we will use the *generalized Marcum function* [29], which is defined by:

$$\begin{aligned} Q_{d/2}: [0, \infty) \times [0, \infty) &\rightarrow \mathbb{R} \\ (\rho, \eta) &\mapsto 1 - \mathbb{F}_{\chi_d^2(\rho^2)}(\eta^2) \end{aligned} \quad (1.18)$$

where $\mathbb{F}_{\chi_d^2(\rho^2)}$ is the cumulative distribution function of the non-central χ^2 distribution with d degrees of freedom and non-centrality parameter ρ^2 . An important consequence of the definition of this function is that for any $\theta \in \mathbb{R}^d$, and for any random variable $X \in \mathcal{M}(\Omega, \mathbb{R}^d)$ that follows the Gaussian distribution $\mathcal{N}(\theta, I_d)$, we have:

$$\forall \lambda \geq 0, \mathbb{P}[\|X\|_2 \geq \lambda] = Q_{d/2}(\|\theta\|_2, \lambda) \quad (1.19)$$

The following lemma defines the threshold $\lambda_\gamma(\tau)$ using this function $Q_{d/2}$. This is the threshold of interest that will allow us to find a γ -MCCP test for the RDT problem. This lemma also describes a few properties of this threshold with regards to the parameters γ and τ .

Lemma 1.2.9: Threshold $\lambda_\gamma(\tau)$

- For any $\gamma \in (0, 1)$, there exists a single positive real $\lambda_\gamma(\tau) \geq 0$ that verifies the equation $Q(\tau, \lambda_\gamma(\tau)) = 1 - \gamma$.
- For any $\gamma \in (0, 1)$, the function $\lambda_\gamma: \tau \in [0, \infty) \mapsto \lambda_\gamma(\tau) \in [0, \infty)$ is continuous everywhere and strictly increasing, and we have $\lim_{\tau \rightarrow +\infty} \lambda_\gamma(\tau) = +\infty$
- For any $\tau \in [0, \infty)$, the function $\gamma \in (0, 1) \mapsto \lambda_\gamma(\tau) \in [0, \infty)$ is continuous everywhere and strictly decreasing, and we have $\lim_{\gamma \rightarrow 1} \lambda_\gamma(\tau) = 0$ and $\lim_{\gamma \rightarrow 0} \lambda_\gamma(\tau) = +\infty$

To give an idea of how the parameters γ and τ affect the threshold $\lambda_\gamma(\tau)$, Figure 1.1 shows $\lambda_\gamma(\tau)$ as a function of γ and τ , using $d = 2$ as an example.

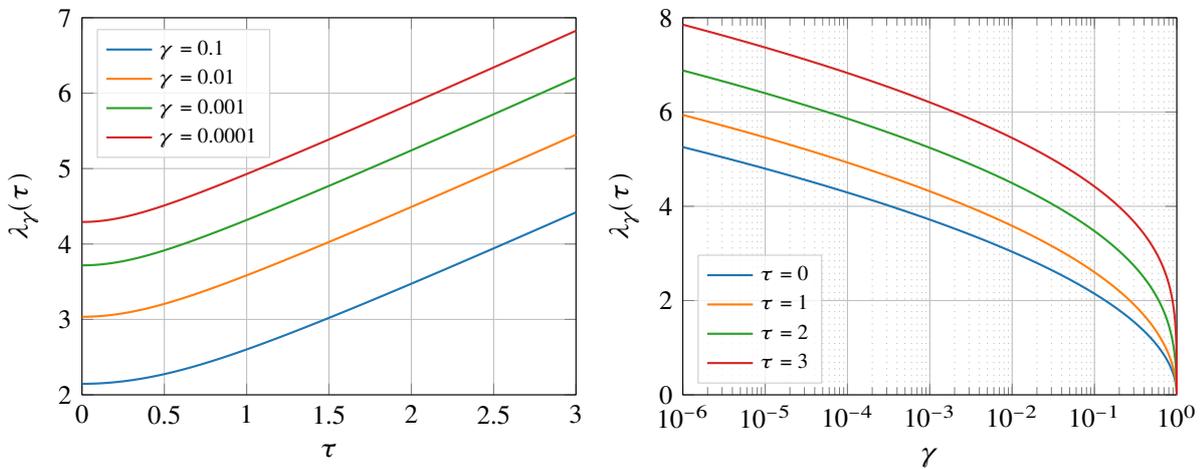


Figure 1.1: Threshold $\lambda_\gamma(\tau)$ against τ (left) and γ (right) for $d = 2$

Theorem 1.2.10: Optimality of the test $\mathcal{J}_{\lambda_\gamma(\tau)}$

The thresholding test $\mathcal{J}_{\lambda_\gamma(\tau)}$ has size γ and is γ -MCCP.

1.3 Robustness and hypothesis testing

When applying hypothesis testing methods to a real problem, we need an accurate model of reality. In the case of traditional testing methods that we presented in Section 1.1, this means having perfect knowledge of all the potential probability distributions that our observation may follow. In practice however, assuming that we can have such knowledge of all involved distributions is questionable. Indeed, “mathematical models are often significant simplifications/idealizations of complex physical problems” [16], and mismatches between these models and reality may have adverse effects on the performance of the test used if they are not accounted for. Therefore, it is important to consider approaches able to take these potential mismatches into account, so that we have robust methods that can guarantee a certain level of performance in real scenarios.

Huber [30, 21] presents such an approach, which consists in determining whether the probability distribution that generated a given observation is close enough or not to given models. We will first explain

this approach, then we will consider the RDT problem and see how this approach is well-suited for robust hypothesis testing.

1.3.1 Huber’s approach to robust hypothesis testing

Following Levy [21], to explain Huber’s approach to robust hypothesis testing, we will consider a testing problem as an example: consider that, given an observation y of a random variable Y with unknown probability distribution P , we want to test $\mathcal{H}_0: P = P_0$ against $\mathcal{H}_1: P = P_1$, where both P_0 and P_1 are known. This is a typical problem where both hypotheses are simple. In such a problem, we would derive the Neyman-Pearson test \mathcal{J}_{NP} based on the likelihood ratio of P_0 and P_1 , and this test is known to be optimal with a given level γ in this case.

Now imagine that both probability distributions P_0 and P_1 are approximations of the real probability distributions that may be encountered in practice, obtained, for example, after observation of a sufficiently large number of samples from each real distribution. This means that P_0 and P_1 are not the actual possible probability distributions of Y , but only approximations of the real distributions. Therefore, the Neyman-Pearson test of level γ computed using P_0 and P_1 is not actually optimal in this scenario, and it also may not actually have level γ in practice.

However, if the distributions P_0 and P_1 are decent approximations of the real distributions, we can hope that the actual distributions we will encounter are close enough to these. If this is indeed the case, and we want to find a test with level γ , one solution consists in testing $\mathcal{H}_0: P \approx P_0$ against $\mathcal{H}_1: P \approx P_1$. Of course, we have to define the meaning of $P \approx P_0$ and $P \approx P_1$: what does it mean for a probability distribution to be “close enough” to another one? Defining this notion of closeness means defining a neighborhood of each distribution P_0 and P_1 , i.e. two sets \mathcal{F}_0 and \mathcal{F}_1 , that contain P_0 and P_1 respectively. These two sets represent all of the possible real probability distributions that we may encounter in practice, that are close enough to P_0 and P_1 .

How do we choose these neighborhoods? Or in other words, how do we define that two distributions are close enough? There are several answers to this question, of which we give a couple examples. In the following, we consider that we want to define a neighborhood \mathcal{F} for some probability distribution P_0 that admits a probability density function f_0 . To define \mathcal{F} , we can consider for example the following approaches:

Contamination model. Let $\varepsilon \in (0, 1)$. A contamination model consists in assuming that the observation is drawn either with probability $1 - \varepsilon$ from P_0 , or with probability ε from some other unknown probability distribution. The set \mathcal{F} is here defined by:

$$\mathcal{F} = \{(1 - \varepsilon)P_0 + \varepsilon P, P \in \mathcal{P}\} \quad (1.20)$$

Distance-based approaches. Several distances can be defined between probability distributions. For two probability distributions P_1 and P_2 , we can for example use the following distances:

- The *Kolmogorov distance* is defined as the maximum distance between the cumulative distribution functions F_1 and F_2 of P_1 and P_2 respectively:

$$d_K(P_1, P_2) = \sup_{y \in \mathbb{R}^d} |F_1(y) - F_2(y)| \quad (1.21)$$

- The *Total variation distance*, defined by

$$d_{TV}(P_1, P_2) = \int |f_1(y) - f_2(y)| dy \quad (1.22)$$

where f_1 and f_2 are the respective probability density functions of P_1 and P_2 .

- The *Kullback-Leibler divergence*, which is not actually a distance, but is nonetheless used as a metric to measure mismatches between distributions. It is defined by:

$$D(P_2 | P_1) = \int \ln\left(\frac{f_2(y)}{f_1(y)}\right) f_2(y) dy \quad (1.23)$$

Using any of these distances d (or pseudo-distances), we can then define a neighborhood \mathcal{F} of P_0 by:

$$\mathcal{F} = \{P \in \mathcal{P} \mid d(P, P_0) \leq \varepsilon\} \quad (1.24)$$

for some chosen real ε .

Now let us return to the problem of testing $\mathcal{H}_0: P \approx P_0$ against $\mathcal{H}_1: P \approx P_1$. Using any of the previously presented methods we can define two neighborhoods \mathcal{F}_0 and \mathcal{F}_1 of P_0 and P_1 respectively, which we assume to be disjoint. Using these neighborhoods, our problem becomes testing $\mathcal{H}_0: P \in \mathcal{F}_0$ against $\mathcal{H}_1: P \in \mathcal{F}_1$. We now need a criterion to choose a test for this problem.

Before doing so, we introduce the following notations:

- Probability of false alarm: for any test \mathcal{T} and any probability distribution $P \in \mathcal{F}_0$ under \mathcal{H}_0 , we define \mathbb{P}_{FA} by:

$$\mathbb{P}_{\text{FA}}(\mathcal{T}, P) = \int \mathcal{T}(y) P(dy) \quad (1.25)$$

- Probability of missed detection: for any test \mathcal{T} and any probability distribution $P \in \mathcal{F}_1$ under \mathcal{H}_1 , we define \mathbb{P}_{MD} by:

$$\mathbb{P}_{\text{MD}}(\mathcal{T}, P) = \int (1 - \mathcal{T}(y)) P(dy) \quad (1.26)$$

Huber defines the *robust Neyman-Pearson testing problem*, which uses the following criterion to choose a test:

$$\mathcal{T}_{RNP} = \arg \min_{\mathcal{T} \in K_\gamma} \max_{P \in \mathcal{F}_1} \mathbb{P}_{\text{MD}}(\mathcal{T}, P) \quad (1.27)$$

where K_γ designates the set of all tests \mathcal{T} that have level γ :

$$K_\gamma = \{\mathcal{T} \mid \sup_{P \in \mathcal{F}_0} \mathbb{P}_{\text{FA}}(\mathcal{T}, P) \leq \gamma\} \quad (1.28)$$

This test \mathcal{T}_{RNP} is a test with level γ , that minimizes the maximum missed-detection rate for all distributions of \mathcal{F}_1 . This test is not necessarily optimal in the UMP sense, as it does not seek to maximize the detection rate for every distribution of \mathcal{F}_1 .

From this criterion, how can we find this test \mathcal{T}_{RNP} ? One known result [21, Sec 6.3.1, p.240] is that the test \mathcal{T}_{RNP} is a likelihood ratio test using two distributions $g_0 \in \mathcal{F}_0$ and $g_1 \in \mathcal{F}_1$:

$$\forall y \in \mathbb{R}^d, \mathcal{T}_{RNP}(y) = \begin{cases} 1 & \text{if } \frac{g_1(y)}{g_0(y)} > \lambda \\ 1 & \text{otherwise} \end{cases} \quad (1.29)$$

where λ is chosen such that $\mathbb{P}_{\text{FA}}(\mathcal{T}_{RNP}, g_0) = \gamma$. In addition, both of these distributions can be considered as the least-favorable distributions, in the sense that we have the following inequalities:

$$\forall P \in \mathcal{F}_0, \mathbb{P}_{\text{FA}}(\mathcal{T}_{RNP}, P) \leq \mathbb{P}_{\text{FA}}(\mathcal{T}_{RNP}, g_0) \quad (1.30)$$

$$\forall P \in \mathcal{F}_1, \mathbb{P}_{\text{MD}}(\mathcal{T}_{RNP}, P) \leq \mathbb{P}_{\text{MD}}(\mathcal{T}_{RNP}, g_1) \quad (1.31)$$

Beyond these inequalities however, there is no real method available to identify the test \mathcal{T}_{RNP} : “the most effective approach is often to guess the solution and then verify that it works” [21, p.241]. This means having to rely on a trial-and-error approach, where we pick two distributions $g_0 \in \mathcal{F}_0$ and $g_1 \in \mathcal{F}_1$,

compute the appropriate likelihood ratio test, and verify if Eqs. (1.30) and (1.31) are verified for a certain number of distributions of \mathcal{F}_0 and \mathcal{F}_1 .

This approach is therefore quite impractical to use, as there is no general way to find the test that satisfies the robust Neyman-Pearson testing problem (Eq. (1.27)), and can also be highly affected by the choices made to define the neighborhoods \mathcal{F}_0 and \mathcal{F}_1 . In addition, the definition of these neighborhoods is not very intuitive. Indeed, it is fairly difficult to represent these notions of closeness between probability distributions, which means that we cannot easily represent what kinds of differences we are taking into account between our nominal models (P_0 and P_1) and reality when designing the robust test \mathcal{J}_{RNP} .

1.3.2 RDT and robust hypothesis testing

We will now explain how the RDT approach is an appropriate choice for robust hypothesis testing. The robustness of the RDT approach comes from the problem statement, which incorporates the notion of model mismatch. As a reminder, the RDT problem consists in deciding whether a realization of some signal Θ is close enough to a model θ_0 , using a noisy observation Y of Θ . In a sense, we can interpret this problem as wanting to know whether θ_0 is a reasonable model of the signal $\Theta(\omega)$ or not. This is done by determining whether we have $\nu_C(\Theta(\omega) - \theta_0) \leq \tau$ or $\nu_C(\Theta(\omega) - \theta_0) > \tau$, which is more robust than wanting to test $\Theta(\omega) = \theta_0$ against $\Theta(\omega) \neq \theta_0$. Indeed, in many cases, it does not make much sense to test whether $\Theta(\omega) = \theta_0$, as this is often impossible in practice, and not necessarily relevant. It may not make much sense to attempt to detect every deviation of Θ from θ_0 , including very minor ones. The use of the tolerance τ therefore gives this approach a certain robustness to these minor changes, and also allows the user to define the amplitude of the deviations that should be detected.

In addition, the signal model used in the RDT problem makes very few assumptions on the observation Y . The model used assumes that the observation Y consists of the signal Θ , whose distribution is completely unknown, in presence of some additive independent Gaussian noise X . Because of this, there is no possible model mismatch regarding the signal Θ , since it does not require any knowledge. As such, the only possible source of mismatches has to do with the noise, whose distribution needs to be known perfectly. However, we will see in Chapter 3 that the parameters of this Gaussian noise can be estimated, and that we still retain an asymptotic optimality by doing so.

Overall, these different aspects make the RDT approach a good candidate to build robust detection methods: the nature of the problem itself deals with robustness, and the few assumptions made on the signal model allow it to be used regardless of the distribution of the signal of interest.

2 Invariance and Generalized Random Distortion Testing

Contents

2.1	Invariance in group theory	46
2.1.1	Group theory	46
2.1.2	Invariance	47
2.1.3	Orbits and maximal invariant	49
2.1.4	Invariance applied to probability distributions	51
2.2	Generalization of the RDT approach	53
2.2.1	Problem statement	54
2.2.2	Redefining notions for the GRDT problem	55
2.2.3	Preliminary results	56
2.2.4	Generalization when the maximal invariant is the euclidean norm	61
2.3	Conclusion and perspectives	64

We mentioned several times so far that the RDT approach is closely related to the notion of invariance, regarding the properties of the Gaussian noise which is part of the problem statement. The article which introduced RDT [17] already made this connection back in 2013, but did not make use of it. Indeed, it was not needed for any theoretical development, as we could instead work directly with the well-known properties of the Gaussian noise. In this chapter, we will revisit the work that was done back in 2013 using the invariance point of view. The goal here is to explicit the role of invariance in this approach, and attempt to generalize it to other forms of noise that also present invariance properties.

The notion of invariance has been studied for hypothesis testing, as a method to find an optimal test among a restricted class of tests. For example, we may be able to find a UMPI test when a UMP test does not exist. Invariance is notably presented by Lehmann in [25]. However, we were unsatisfied with the presentations of invariance that we found in the literature for hypothesis testing. Indeed, we found these presentations do not clarify how invariance in hypothesis testing is related to the original notion of invariance found in group theory [31]. As such, we would like to start by formalizing the notion of invariance and its application to statistical decision theory. After that, we will discuss our attempts to generalize the RDT approach using invariance. These attempts resulted in interesting theoretical results, but unfortunately did not allow us to discover any other applicable noise distributions so far.

In this chapter, we will provide proofs for most of the results presented here, even ones that are commonly found in the literature, especially in the first section regarding invariance. We have made this choice in order to have a presentation of invariance and of the Generalized RDT as self-contained as possible, and also to palliate some gaps that we found in the literature regarding the use of invariance in hypothesis testing. We notably found that Lehmann's presentation of invariance in [25] was sometimes lacking, with notably very unclear proofs. As such, we wanted to palliate some of these shortcomings, and present a hopefully clearer picture of some of these concepts.

2.1 Invariance in group theory

In this first section, we offer a formalized path from group theory to invariance for hypothesis testing. We will start by recalling some required notions of group theory, most notably the concept of group action, then show how invariance is defined. We will also properly introduce the notions of orbits and maximal invariant, which we will need in our generalization of the RDT approach. After this, we will see how we can connect these notions to what is commonly presented in the literature about hypothesis testing.

2.1.1 Group theory

We start this section by recalling a few basic notions of group theory which will be useful later.

Definition 2.1.1: Group

Let G be a set and let $*$ be a binary operation defined for every pair of elements of G . $(G, *)$ is a *group* if it verifies the following properties:

1. Closure: $\forall (g_1, g_2) \in G^2, g_1 * g_2 \in G$
2. Associativity: $\forall (g_1, g_2, g_3) \in G^3, (g_1 * g_2) * g_3 = g_1 * (g_2 * g_3)$
3. Neutral element: $\exists ! e \in G, \forall g \in G, e * g = g * e = g$
4. Inverse element: $\forall g \in G, \exists ! g^{-1} \in G, g * g^{-1} = g^{-1} * g = e$

Definition 2.1.2: Group homomorphism

Let $(G, *)$ and (H, \cdot) be two groups, and let $f: G \rightarrow H$ be a function. The function f is a *group homomorphism* from $(G, *)$ to (H, \cdot) if it verifies:

$$\forall (g_1, g_2) \in G^2, f(g_1 * g_2) = f(g_1) \cdot f(g_2)$$

From these two definitions follow two immediate properties of a group homomorphism.

Lemma 2.1.3: Properties of group homomorphisms

Let f be a group homomorphism from $(G, *)$ to (H, \cdot) . Then:

1. *Neutral element compatibility:* if e_G and e_H are the neutral elements of $(G, *)$ and (H, \cdot) respectively, we have: $f(e_G) = e_H$.
2. *Inverse element compatibility:* for any element $g \in G$, we have: $f(g^{-1}) = f(g)^{-1}$.

Given a set \mathcal{X} , we will denote by $\text{Perm}(\mathcal{X})$ the set of all permutations of \mathcal{X} , i.e. the set of all \mathcal{X} -valued bijections defined on \mathcal{X} (also known as transformations). The set $\text{Perm}(\mathcal{X})$ with the composition operation \circ is a group.

Definition 2.1.4: Group action [31]

Let $(G, *)$ be a group and \mathcal{X} be a set.

We say that a function $\Pi : (G, *) \rightarrow (\text{Perm}(\mathcal{X}), \circ)$ is a *group action* of G on \mathcal{X} if Π is a group homomorphism between $(G, *)$ and $(\text{Perm}(\mathcal{X}), \circ)$.

The idea of a group action is to associate to each element of the group G a transformation of the set \mathcal{X} , effectively transferring the group structure to a certain subset of transformations of \mathcal{X} .

We now introduce an equivalent definition of a group action, used notably by Eaton in [24].

Lemma 2.1.5: Group action characterization

Let $\Pi: G \rightarrow \mathcal{X}^{\mathcal{X}}$.

Π is a group action of G on \mathcal{X} if and only if both following properties are true:

$$(i) \quad \forall (g, g') \in G \times G, \forall x \in \mathcal{X}, \Pi(g)(\Pi(g')(x)) = \Pi(g * g')(x)$$

$$(ii) \quad \forall x \in \mathcal{X}, \Pi(e)(x) = x$$

where e is the identity element of G . The second property can also be written $\Pi(e) = \text{Id}_{\mathcal{X}}$.

Proof. If Π is a group action, then Π verifies both (i) and (ii) since (i) is the definition of a group homomorphism and (ii) is the first property of Lemma 2.1.3.

Conversely, assume that Π verifies (i) and (ii).

Note that this does not assume that Π is a permutation of \mathcal{X} , this is in fact what we need to prove. The fact that Π is a group homomorphism will then directly come as a consequence of (i).

Let $g \in G$.

1. Injectivity:

Let $(x, x') \in \mathcal{X}^2$ such that $\Pi(g)(x) = \Pi(g)(x')$.

By applying $\Pi(g^{-1})$ we get:

$$\Pi(g^{-1})(\Pi(g)(x)) = \Pi(g^{-1})(\Pi(g)(x'))$$

We have:

$$\begin{aligned} \Pi(g^{-1})(\Pi(g)(x)) &= \Pi(g^{-1} * g)(x) && \text{from (i)} \\ &= \Pi(e)(x) \\ &= x && \text{from (ii)} \end{aligned}$$

With the exact same reasoning, we can also prove that $\Pi(g^{-1})(\Pi(g)(x')) = x'$.

Therefore $x = x'$, meaning that $\Pi(g)$ is injective.

2. Surjectivity:

Let $y \in \mathcal{X}$. We want to show that there exists $x \in \mathcal{X}$ such that $\Pi(g)(x) = y$.

Let $x = \Pi(g^{-1})(y)$. We have :

$$\begin{aligned} \Pi(g)(x) &= \Pi(g)(\Pi(g^{-1})(y)) \\ &= \Pi(g * g^{-1})(y) && \text{from (i)} \\ &= \Pi(e)(y) \\ &= y && \text{from (ii)} \end{aligned}$$

Therefore $\Pi(g)$ is surjective.

$\Pi(g)$ is thus bijective and $\Pi(g^{-1})$ is its inverse function, hence: $\Pi(g) \in \text{Perm}(\mathcal{X})$.

Since Π also verifies (i), Π is a group homomorphism. \square

2.1.2 Invariance

We have now introduced the basic elements that are required to introduce the notion of invariance. This is done through the definition of invariant subsets of \mathcal{X} , given a certain group action Π .

Definition 2.1.6: Invariant subset

Let Π be a group action of G on \mathcal{X} , and let $E \subset \mathcal{X}$ be a subset of \mathcal{X} . E is an *invariant subset* of \mathcal{X} under the action Π of G on \mathcal{X} if:

$$\forall x \in E, \forall g \in G, \Pi(g)(x) \in E$$

We can also equivalently write:

$$\forall g \in G, \Pi(g)(E) \subset E$$

where $\Pi(g)(E) = \{\Pi(g)(x), x \in E\}$

Lemma 2.1.7: Invariant subset characterization

A subset $E \subset \mathcal{X}$ is invariant under the action of Π of G on \mathcal{X} if and only if:

$$\forall g \in G, \Pi(g)(E) = E$$

Proof.

$$\Pi(g)(E) = E \Leftrightarrow (\Pi(g)(E) \subset E \wedge \Pi(g)(E) \supset E)$$

- If we have $\Pi(g)(E) = E$ for all $g \in G$, then E is trivially invariant under the action of Π .
- Conversely, assume that E is invariant under the action of Π .
Let $g \in G$. We know by definition of an invariant subset that $\Pi(g)(E) \subset E$. All we need to prove is that $\Pi(g)(E) \supset E$.
Let $y \in E$. Since G is a group, $g^{-1} \in G$ exists. Let $x = \Pi(g^{-1})(y)$. We have:

$$\begin{aligned} \Pi(g)(x) &= \Pi(g)(\Pi(g^{-1})(y)) \\ &= \Pi(g * g^{-1})(y) \\ &= \Pi(e)(y) \\ &= y \end{aligned}$$

In addition to that, by definition of x , we have $x \in \Pi(g^{-1})(E)$, and since we assumed that E is invariant under the action of Π , we have $\Pi(g^{-1})(E) \subset E$, meaning that $x \in E$.
Therefore we have found an element $x \in E$ such that $y = \Pi(g)(x)$, meaning that $y \in \Pi(g)(E)$.
Hence $\Pi(g)(E) = E$. □

Lemma 2.1.8: Group action on a subset

If E is an invariant subset of \mathcal{X} under the action Π of G on \mathcal{X} , then Π_E defined by

$$\begin{aligned} \Pi_E: (G, *) &\rightarrow (\text{Perm}(E), \circ) \\ g &\mapsto \Pi(g)|_E \end{aligned}$$

is a group action of G on E , where $\Pi(g)|_E$ is the restriction of $\Pi(g)$ to the set E :

$$\begin{aligned} \Pi(g)|_E: E &\rightarrow E \\ x &\mapsto \Pi(g)(x) \end{aligned}$$

Proof. Let $g \in G$.

First, $\Pi(g)|_E$ is properly defined, since E is invariant under the action of Π , and therefore for any $x \in E$ we have $\Pi(g)|_E(x) = \Pi(g)(x) \in E$.

We will show that $\Pi(g)|_E \in \text{Perm}(E)$.

- (*Injectivity*) $\Pi(g)|_E$ is the restriction of $\Pi(g)$ which is an injective function. Therefore $\Pi(g)|_E$ is injective.
- (*Surjectivity*) Let $y \in E$. We want to find $x \in E$ such that $\Pi(g)|_E(x) = y$. Let $x = \Pi(g^{-1})(y)$. Since $y \in E$ and $g^{-1} \in G$, $x \in E$ because E is invariant.

$$\begin{aligned} \Pi(g)|_E(x) &= \Pi(g)(x) \\ &= \Pi(g)(\Pi(g^{-1})(y)) \\ &= \Pi(g * g^{-1})(y) \\ &= \Pi(e)(y) \\ &= y \end{aligned}$$

Therefore $\Pi(g)|_E$ is surjective.

Hence $\Pi(g)|_E \in \text{Perm}(E)$. From its definition, Π_E trivially verifies the properties (i) and (ii), therefore it is a group action of G on E . \square

One last important notion before moving on to applications to probability distributions is the notion of transfer of a group action from one set to another.

Lemma 2.1.9: Transfer of a group action

Let \mathcal{X} and \mathcal{Y} be two sets, $f: \mathcal{X} \rightarrow \mathcal{Y}$ a bijection between \mathcal{X} and \mathcal{Y} and Π a group action of G on \mathcal{X} . Let γ be the application defined by:

$$\begin{aligned} \gamma: (G, *) &\rightarrow (\text{Perm}(\mathcal{Y}), \circ) \\ g &\mapsto f \circ \Pi(g) \circ f^{-1} \end{aligned}$$

γ is a group action of G on \mathcal{Y} .

Proof. Let ζ be the application defined for any $h \in \text{Perm}(\mathcal{X})$ by $\zeta(h) = f \circ h \circ f^{-1}$. Note that with this definition of ζ , we have $\gamma = \zeta \circ \Pi$.

We will now prove that ζ is a group homomorphism between $(\text{Perm}(\mathcal{X}), \circ)$ and $(\text{Perm}(\mathcal{Y}), \circ)$:

- For all $h \in \text{Perm}(\mathcal{X})$, $\zeta(h)$ is defined on \mathcal{Y} and takes its values in \mathcal{Y} , therefore $\zeta(h) \in \text{Perm}(\mathcal{Y})$
- For all $h \in \text{Perm}(\mathcal{X})$, $\zeta(h)$ is bijective since it is a composition of bijective functions, therefore $\zeta(h) \in \text{Perm}(\mathcal{Y})$
- $\forall h, h' \in \text{Perm}(\mathcal{X}), \zeta(h) \circ \zeta(h') = f \circ h \circ f^{-1} \circ f \circ h' \circ f^{-1} = f \circ h \circ h' \circ f^{-1} = \zeta(h \circ h')$

Therefore ζ is a group homomorphism between $(\text{Perm}(\mathcal{X}), \circ)$ and $(\text{Perm}(\mathcal{Y}), \circ)$. Hence we can deduce that $\gamma = \zeta \circ \Pi$ is a group homomorphism between $(G, *)$ and $(\text{Perm}(\mathcal{Y}), \circ)$ as the composition of two group homomorphisms, which means that γ is a group action of G on \mathcal{Y} . \square

2.1.3 Orbits and maximal invariant

We now introduce the notions of maximal invariant and orbits which play an important part to introduce invariance in hypothesis testing problems. In the following, we consider a group action Π of a group G on a set \mathcal{X} .

From this group action, we can define an equivalence relation between elements of \mathcal{X} by:

$$\forall (x, x') \in \mathcal{X}^2, x \equiv x' \text{ if } \exists g \in G, x' = \Pi(g)(x) \tag{2.1}$$

The fact that this is indeed an equivalence relation results directly from the properties of a group action. For every $x \in \mathcal{X}$, the equivalence class of x is the set $\{\Pi(g)(x), g \in G\}$ and is called the *orbit* of x .

Definition 2.1.10: Orbit

For an element $x \in \mathcal{X}$, the *orbit* of x is the set $\text{orb}(x) \subset \mathcal{X}$ defined by:

$$\text{orb}(x) = \{\Pi(g)(x), g \in G\} \tag{2.2}$$

Definition 2.1.11: Invariant function

Let f be a function defined on \mathcal{X} .

The function f is said to be *invariant under* Π if:

$$\forall x \in \mathcal{X}, \forall g \in G, f(\Pi(g)(x)) = f(x)$$

Equivalently, a function f defined on \mathcal{X} is invariant if and only if it is constant on each orbit of \mathcal{X} .

Definition 2.1.12: Maximal invariant

Let M be a function defined on \mathcal{X} .

The function M is a *maximal invariant* of the group action Π if it verifies both following properties:

- (i) M is invariant under Π
- (ii) $\forall (x, x') \in \mathcal{X}^2, M(x) = M(x') \Rightarrow \exists g \in G, x' = \Pi(g)(x)$

To put it simply, a maximal invariant is an invariant function defined on \mathcal{X} that also takes different values on different orbits. We can note that given a group action, it is always possible to find at least a maximal invariant. For example, we can always consider the function $M: x \in \mathcal{X} \mapsto \text{orb}(x) \in 2^{\mathcal{X}}$. This function M is indeed a maximal invariant since the orbits are the equivalence classes of the relation defined in Eq. (2.1).

For a more concrete example, consider $\mathcal{X} = \mathbb{R}^d$, $G = O(d)$ the orthogonal group of \mathbb{R}^d with the composition operation, and Π the trivial group action defined for every $g \in G$ by $\Pi(g) = g$. A maximal invariant of this group action is the euclidean norm $\|\cdot\|_2$, and the orbits generated by this group action are the $d - 1$ -spheres of \mathbb{R}^d . In \mathbb{R}^2 and \mathbb{R}^3 , the orthogonal group $O(d)$ contains notably the rotations around 0. It is fairly easy to see in this case that an element on a sphere remains on that same sphere after applying a rotation, and therefore that its norm is preserved.

Lemma 2.1.13: Characterization of invariant functions [24, Proposition 7.7]

Let f be a function defined on \mathcal{X} and let M be a maximal invariant of Π .

f is invariant under Π if and only if it is a function of the maximal invariant M , i.e. there exists a function \bar{f} defined on $M(\mathcal{X})$ such that:

$$\forall x \in \mathcal{X}, f(x) = \bar{f}(M(x)) \tag{2.3}$$

In the following, we will be using an alternate definition of the orbits, which uses a maximal invariant. This is done to simplify some notations throughout this thesis.

Lemma 2.1.14: Equivalent definition of orbits using a maximal invariant

Given a maximal invariant M of Π , let Υ_ρ be the sets defined by:

$$\forall \rho \in M(\mathcal{X}), \Upsilon_\rho = \{x \in \mathcal{X} : M(x) = \rho\} = M^{-1}(\{\rho\}) \quad (2.4)$$

We have:

$$\forall x \in \mathcal{X}, \text{orb}(x) = \Upsilon_{M(x)} \quad (2.5)$$

which we can also write:

$$\forall x \in \mathcal{X}, \forall \rho \in M(\mathcal{X}), \text{orb}(x) = \Upsilon_\rho \quad (2.6)$$

This lemma is simply a matter of notation, and its proof comes directly from the definition of a maximal invariant and of the orbits. We introduce here the notation Υ_ρ , which references each orbit by the value taken by a given maximal invariant M on this orbit. The chosen maximal invariant is implicit, but we will only work with a single maximal invariant at a time in the following, so there should be no confusion.

2.1.4 Invariance applied to probability distributions

In the following, we consider $\mathcal{X} = \text{Meas}_+(\mathbb{R}^d)$ the set of all positive measures defined on \mathcal{F}^d . This set includes of course the set of all probability distributions.

For every $g \in G$, we define the function Π_g by:

$$\begin{aligned} \Pi_g: \text{Meas}_+(\mathbb{R}^d) &\rightarrow \text{Meas}_+(\mathbb{R}^d) \\ \mu &\mapsto \mu g^{-1} \end{aligned}$$

where μg^{-1} is the pushforward measure of μ by g . We also define the function Π by:

$$\forall g \in G, \Pi(g) = \Pi_g$$

Lemma 2.1.15: Group action on measures

Π is a group action of G on $\text{Meas}_+(\mathbb{R}^d)$.

Proof. We first need to prove that, for every $g \in G$, Π_g is bijective, i.e. that it is both injective and surjective. Let $g \in G$.

- **Injectivity:** Let $\mu, \nu \in \text{Meas}_+(\mathbb{R}^d)^2$ be two measures. We have:

$$\begin{aligned} \Pi_g(\mu) = \Pi_g(\nu) &\Leftrightarrow \forall B \in \mathcal{B}(\mathbb{R}^d), \mu g^{-1}(B) = \nu g^{-1}(B) \\ &\Leftrightarrow \forall B \in \mathcal{B}(\mathbb{R}^d), \mu(g^{-1}(B)) = \nu(g^{-1}(B)) \\ &\Leftrightarrow \forall B \in g^{-1}(\mathcal{B}(\mathbb{R}^d)), \mu(B) = \nu(B) \end{aligned}$$

with $g^{-1}(\mathcal{B}(\mathbb{R}^d)) = \{g^{-1}(B), B \in \mathcal{B}(\mathbb{R}^d)\}$. We will show that $g^{-1}(\mathcal{B}(\mathbb{R}^d)) = \mathcal{B}(\mathbb{R}^d)$.

- Let $B \in g^{-1}(\mathcal{B}(\mathbb{R}^d))$. There is $B' \in \mathcal{B}(\mathbb{R}^d)$ such that $B = g^{-1}(B')$. Since $g^{-1} \in G$, g^{-1} is a measurable transformation, therefore $B \in \mathcal{B}(\mathbb{R}^d)$: $g^{-1}(\mathcal{B}(\mathbb{R}^d)) \subset \mathcal{B}(\mathbb{R}^d)$.
- Let $B \in \mathcal{B}(\mathbb{R}^d)$ and $B' = g(B)$. We have $g^{-1}(B') = g^{-1}(g(B)) = B$. Therefore $B \in g^{-1}(\mathcal{B}(\mathbb{R}^d))$: $g^{-1}(\mathcal{B}(\mathbb{R}^d)) \supset \mathcal{B}(\mathbb{R}^d)$.

Hence $g^{-1}(\mathcal{B}(\mathbb{R}^d)) = \mathcal{B}(\mathbb{R}^d)$, which means:

$$\begin{aligned} \Pi_g(\mu) = \Pi_g(\nu) &\Leftrightarrow \forall B \in g^{-1}(\mathcal{B}(\mathbb{R}^d)), \mu(B) = \nu(B) \\ &\Leftrightarrow \forall B \in \mathcal{B}(\mathbb{R}^d), \mu(B) = \nu(B) \\ &\Leftrightarrow \mu = \nu \end{aligned}$$

Therefore Π_g is injective.

- **Surjectivity:** Let $\mu \in \text{Meas}_+(\mathbb{R}^d)$ and $\nu = \mu g \in \text{Meas}_+(\mathbb{R}^d)$ the image measure of μ by g^{-1} . We have:

$$\begin{aligned} \forall B \in \mathcal{B}(\mathbb{R}^d), \Pi_g(\nu)(B) &= \nu g^{-1}(B) \\ &= \nu(g^{-1}(B)) \\ &= \mu g(g^{-1}(B)) \\ &= \mu(g(g^{-1}(B))) \\ &= \mu(B) \end{aligned}$$

Therefore $\mu = \Pi_g(\nu)$: Π_g is surjective.

Hence Π_g is bijective: $\Pi_g \in \text{Perm}(\mathcal{X})$.

We can now verify that Π is a group action of G on $\text{Meas}_+(\mathbb{R}^d)$. Let $g, h \in G^2$. We have:

$$\begin{aligned} \forall \mu \in \text{Meas}_+(\mathbb{R}^d), \Pi_g \circ \Pi_h(\mu) &= \Pi_g(\Pi_h(\mu)) \\ &= \Pi_g(\mu h^{-1}) \\ &= \mu h^{-1} g^{-1} \\ &= \mu(g \circ h)^{-1} \\ &= \Pi_{g \circ h}(\mu) \end{aligned}$$

Therefore $\Pi_g \circ \Pi_h = \Pi_{g \circ h}$: Π is a group action of G on \mathcal{X} . □

Using the group action Π , we can now define what an invariant hypothesis testing problem is (for usual testing problems, not in the RDT case). Consider a family of probability distributions $\mathcal{F} = \{P_\theta, \theta \in \Theta\}$ indexed by a parameter $\theta \in \Theta$ in some set Θ , such that no two values of $\theta \in \mathcal{D}$ refer to the same probability distribution:

$$\forall (\theta, \theta') \in \mathcal{D}^2, \theta \neq \theta' \Rightarrow P_\theta \neq P_{\theta'}$$

and let Θ_0 and Θ_1 be two disjoint subsets of Θ such that $\Theta_0 \cup \Theta_1 = \Theta$, and let $\mathcal{F}_0 = \{P_\theta, \theta \in \Theta_0\}$ and $\mathcal{F}_1 = \{P_\theta, \theta \in \Theta_1\}$. Given some observation Y that with probability distribution P , the problem of testing $P \in \mathcal{F}_0$ against $P \in \mathcal{F}_1$ is said to be invariant under Π if both families \mathcal{F}_0 and \mathcal{F}_1 are invariant under the action Π . In this case, the set \mathcal{F} is also invariant under the action Π .

As mentioned in Chapter 1, it may be easier to work with the parameter θ instead of working with the probability distribution P . Here, this means testing $\theta \in \Theta_0$ against $\theta \in \Theta_1$. For this, we introduce the function f defined by:

$$\begin{aligned} f: \mathcal{F} &\rightarrow \Theta \\ P_\theta &\mapsto \theta \end{aligned}$$

Since each distribution of \mathcal{F} is uniquely indexed by a value of Θ , f is bijective. Therefore, from Lemma 2.1.9, the application γ defined by:

$$\begin{aligned} \gamma: G &\rightarrow \text{Perm}(\mathcal{D}) \\ g &\mapsto f \circ \Pi(g) \circ f^{-1} \end{aligned}$$

is a group action of G on Θ . For $i \in \{0, 1\}$, we have:

$$\begin{aligned} \forall g \in G, \forall \theta \in \Theta_i, \gamma(g)(\theta) &= f(\Pi(g)(f^{-1}(\theta))) \\ &= f(\Pi(g)(P_\theta)) \\ &= f(P_\theta g^{-1}) \\ &= f(P_{\theta_g}) && \text{for some } \theta_g \in \Theta_i, \text{ because } \mathcal{F}_i \text{ is invariant under } \Pi \\ &= \theta_g \end{aligned}$$

We can therefore deduce that the sets Θ_0 and Θ_1 are also invariant under the action γ , and therefore the problem of testing $\theta \in \Theta_0$ against $\theta \in \Theta_1$ is also invariant.

In the literature regarding invariance of hypothesis testing problems (see for example [25]), the function $\gamma(g)$ is often referred to as \bar{g} , and the set of functions $\{\gamma(g), g \in G\}$ is referred to as \bar{G} . However, these functions \bar{g} are typically not defined through a group action transfer, but instead as the function that maps to any $\theta \in \Theta$ the parameter $\theta_g \in \Theta$, defined here in our calculation of $\gamma(g)(\theta)$, that result from the invariance of \mathcal{F} .

Before introducing the Generalized RDT problem, we will present the notion of *invariant measure*, and a few associated properties. One difference regarding invariance between usual testing problems and RDT is that instead of considering invariant sets of probability distributions, we will work with a single invariant measure.

Definition 2.1.16: Invariant measure

Let $\mu \in \mathcal{X}$ be a measure.

We say that μ is an *invariant measure* if the set $\{\mu\}$ is invariant under the action of Π , i.e.:

$$\forall g \in G, \mu g^{-1} = \mu$$

The following lemma is a direct consequence of this definition.

Lemma 2.1.17: Invariant measure and random variables

Let P be an invariant probability distribution. If $X \in \mathcal{M}(\Omega, \mathbb{R}^d)$ is a random variable with distribution P , then for any $g \in G$, X and $g(X)$ have the same distribution:

$$\forall g \in G, \mathbb{P}X^{-1} = \mathbb{P}g(X)^{-1} = P \quad (2.7)$$

Proof. We have:

$$\begin{aligned} \forall B \in \mathcal{B}^d, \mathbb{P}g(X)^{-1}(B) &= \mathbb{P}[g(X) \in B] \\ &= \mathbb{P}[X \in g^{-1}(B)] \\ &= \mathbb{P}X^{-1}(g^{-1}(B)) \\ &= P(g^{-1}(B)) \\ &= P(B) && \text{because } P \text{ is invariant} \\ &= \mathbb{P}X^{-1}(B) \end{aligned}$$

Therefore $\mathbb{P}g(X)^{-1} = \mathbb{P}X^{-1}$, which means that X and $g(X)$ share the same probability distribution. \square

If the probability distribution of a random variable $X \in \mathcal{M}(\Omega, \mathbb{R}^d)$ is invariant, we will also simply say that X is invariant.

2.2 Generalization of the RDT approach

After explaining the notion of invariance, we will now see its application to the RDT problem. The main goal of this section is to offer a generalization of the RDT problem where the distribution of the noise X is not necessarily Gaussian, but instead has an invariant distribution with respect to some group G , allowing us to offer an optimal test in more situations.

The work presented in this section was conducted in great part with the assistance of Sabrina Bourmani, who initiated the work on this generalization in her thesis [32].

We will consider here a group (G, \circ) of *linear* transformations defined on \mathbb{R}^d . The linearity of the functions $g \in G$ is important to keep in mind, and is adapted to the problem we consider, which is linear

in the sense that we consider an observation that is the sum of some independent signal Θ and noise X . We also consider an \mathbb{R} -valued maximal invariant $M: \mathbb{R}^d \rightarrow \mathbb{R}$ of G . For every $\rho \in \mathbb{R}$, Υ_ρ is an orbit:

$$\Upsilon_\rho = \{y \in \mathbb{R}^d, M(y) = \rho\} \quad (2.8)$$

We call \mathfrak{F} the set of all these orbits: $\mathfrak{F} = \{\Upsilon_\rho, \rho \in M(\mathbb{R}^d)\}$.

We start by stating the Generalized RDT problem, adapted from the RDT problem presented in Chapter 1, and explain the differences that exist between these two problems. We will then redefine the essential elements of the RDT problem (size, power, optimality, etc.) for the GRDT case. After that, we will present several results which are completely independent from the distribution of the noise X . Finally we will focus on the case where the maximal invariant is the euclidean norm $\|\cdot\|_2$, in which case we can derive a test $\mathcal{J}_{\lambda, \tau}$ similarly to the RDT case, and present results regarding its optimality.

2.2.1 Problem statement

We start by stating the GRDT (Generalized Random Distortion Testing) testing problem, and compare it to the RDT testing problem presented in Definition 1.2.2. The Generalized RDT problem is the following hypothesis testing problem:

Definition 2.2.1: Generalized RDT problem statement

$$\left\{ \begin{array}{l} \textbf{Data model:} \\ \exists Y \in \mathcal{M}(\Omega, \mathbb{R}^d), \exists \Theta \in \mathcal{M}(\Omega, \mathbb{R}^d), \exists X \in \mathcal{M}(\Omega, \mathbb{R}^d), \exists P \in \mathcal{P}, \\ \left\{ \begin{array}{l} (X \sim P) \\ \wedge (P \text{ is invariant under } G) \\ \wedge (\Theta \text{ and } X \text{ are independent}) \\ \wedge (Y = \Theta + X) \\ \wedge (\forall y \in \mathbb{R}^d, \exists \omega \in \Omega, y = Y(\omega)) \end{array} \right. \\ \textbf{Testing problem:} \\ \text{Given one realization } y = Y(\omega) = \Theta(\omega) + X(\omega), \text{ determine whether:} \\ \left\{ \begin{array}{l} \mathcal{H}_0: M(\Theta(\omega)) \leq \tau \\ \text{or} \\ \mathcal{H}_1: M(\Theta(\omega)) > \tau \end{array} \right. \\ \text{with } \tau \in \mathbb{R} \end{array} \right. \quad (2.9)$$

For comparison, here is the RDT problem which was stated in Definition 1.2.2:

Definition 1.2.2: RDT problem statement

$$\left\{ \begin{array}{l} \textbf{Data model:} \\ \exists Y \in \mathcal{M}(\Omega, \mathbb{R}^d), \exists \Theta \in \mathcal{M}(\Omega, \mathbb{R}^d), \exists X \in \mathcal{M}(\Omega, \mathbb{R}^d), \\ \left\{ \begin{array}{l} (X \sim \mathcal{N}(0, C) \text{ with } C \text{ a } d \times d \text{ definite positive matrix}) \\ \wedge (\Theta \text{ and } X \text{ are independent}) \\ \wedge (Y = \Theta + X) \\ \wedge (\forall y \in \mathbb{R}^d, \exists \omega \in \Omega, y = Y(\omega)) \end{array} \right. \\ \textbf{Testing problem:} \\ \text{Given one realization } y = Y(\omega) = \Theta(\omega) + X(\omega), \text{ determine whether:} \\ \left\{ \begin{array}{l} \mathcal{H}_0: \nu_C(\Theta(\omega) - \theta_0) \leq \tau \\ \text{or} \\ \mathcal{H}_1: \nu_C(\Theta(\omega) - \theta_0) > \tau \end{array} \right. \\ \text{with } \tau > 0 \text{ and } \theta_0 \in \mathbb{R}^d \end{array} \right. \quad (1.12)$$

The main difference between these two problems is the distribution of the noise X , which is no longer necessarily Gaussian, but still presents invariance properties with respect to a given group G .

The GRDT problem actually encompasses the RDT problem, where X follows a Gaussian distribution with positive definite covariance matrix C . The maximal invariant M is then the Mahalanobis norm using the matrix C , and the associated orbits are ellipsoids. To explicit the associated group G , we first need to do a little bit of work with the matrix C : consider an eigendecomposition of $C = U\Delta U^T$, where $\Delta = \text{diag}(\xi_1, \dots, \xi_d)$ is a diagonal matrix containing the eigenvalues of C , and U is a $d \times d$ orthogonal matrix. Let $\Phi = \Delta^{-1/2}U^T$ where $\Delta^{-1/2} = \text{diag}(\xi_1^{-1/2}, \dots, \xi_d^{-1/2})$. The group G of interest for the RDT problem is the group of linear transforms g defined by:

$$\forall y \in \mathbb{R}^d, g(y) = \Phi^{-1}R\Phi y \quad (2.10)$$

where R is any $d \times d$ orthogonal matrix. If we consider the case where $C = I_d$, then we find that the maximal invariant is the euclidean norm, and the group G is the orthogonal group $O(d)$ of \mathbb{R}^d . In the case of \mathbb{R}^2 and \mathbb{R}^3 , this group notably contains the rotations of the space centered on the origin. It is fairly easy in that case to see that a Gaussian distribution with mean 0 remains unaffected by such a rotation.

We can also note that in the GRDT case, we decided to only consider the case where $\theta_0 = 0$. This is only to simplify calculations, and can be done without loss of generality. One can also notice that we now consider a tolerance $\tau \in \mathbb{R}$ which is not necessarily positive. Indeed, unlike the Mahalanobis norm, the maximal invariant is not necessarily positive.

2.2.2 Redefining notions for the GRDT problem

Before describing the main theoretical results, we first need to adapt several definitions given in Chapter 1 for the GRDT problem. These adaptations are all straightforward and should not require any additional explanations. The only differences are the use of the maximal invariant M instead of the Mahalanobis norm ν_C , the fact that we consider $\theta_0 = 0$, and the fact that the tolerance τ is not necessarily positive.

Definition 2.2.2: Conditional power and size

Let $\mathcal{J}: \mathbb{R}^d \rightarrow \{0, 1\}$ be a test, $\Theta \in \mathcal{M}(\Omega, \mathbb{R}^d)$ a random variable and $\rho \in \mathbb{R}$. We define the *conditional power* and *conditional size* of \mathcal{J} for the GRDT problem as:

$$\text{Conditional power:} \quad \mathbb{P}[\mathcal{J}(\Theta + X) = 1 \mid \Theta \in \Upsilon_\rho] \quad (2.11)$$

$$\text{Conditional size:} \quad \sup_{\Theta \in \mathcal{M}(\Omega, \mathbb{R}^d): \mathbb{P}[M(\Theta) \leq \tau] \neq 0} \mathbb{P}[\mathcal{J}(\Theta + X) = 1 \mid M(\Theta) \leq \tau] \quad (2.12)$$

We also define the size and power of \mathcal{J} as:

$$\text{Power function:} \quad \forall \theta \in \mathbb{R}^d, \beta_{\mathcal{J}}(\theta) = \mathbb{P}[\mathcal{J}(\theta + X) = 1] \quad (2.13)$$

$$\text{Size:} \quad \alpha_{\mathcal{J}} = \sup_{\theta \in \mathbb{R}^d: M(\theta) \leq \tau} \beta_{\mathcal{J}}(\theta) \quad (2.14)$$

Like previously, a test \mathcal{J} is said to have *level* $\gamma \in (0, 1)$ for the GRDT problem if $\alpha_{\mathcal{J}} \leq \gamma$. We note \mathcal{K}_γ the set of all tests of level γ .

Definition 2.2.3: Thresholding tests

For any $t \in \mathbb{R}$, the *thresholding test* \mathcal{J}_t is defined by:

$$\begin{aligned} \mathcal{J}_t: \mathbb{R}^d &\rightarrow \{0, 1\} \\ y &\mapsto \begin{cases} 1 & \text{if } M(y) > t \\ 0 & \text{otherwise} \end{cases} \end{aligned} \quad (2.15)$$

Note that these thresholding tests are invariant functions, as they depend on the observation y only through $M(y)$.

Definition 2.2.4: Constant conditional power function given $\Theta \in \Upsilon_\rho$

Let $\rho \in \mathbb{R}$. Let $\Theta \in \mathcal{M}(\Omega, \mathbb{R}^d)$ be a random variable independent of $X \sim \mathcal{N}(0, C)$. A test \mathcal{J} is said to have constant conditional power function given $\Theta \in \Upsilon_\rho$ if:

$$\forall \theta \in \Upsilon_\rho, \mathbb{P}[\mathcal{J}(\Theta + X) = 1 \mid \Theta \in \Upsilon_\rho] = \beta_{\mathcal{J}}(\theta)$$

Definition 2.2.5: γ -MCCP (Maximum Constant Conditional Power) test

Let $\mathcal{F}^*: \mathbb{R}^d \rightarrow \{0, 1\}$ be a test, and let $\gamma \in (0, 1)$.

The test \mathcal{F}^* is said to have level γ and Maximum Constant Conditional Power over \mathfrak{F} — and we simply say that \mathcal{F}^* is γ -MCCP — if:

- (i) \mathcal{F}^* has level γ .
- (ii) For any random variable $\Theta \in \mathcal{M}(\Omega, \mathbb{R}^d)$ and for $\mathbb{P}(M(\Theta))^{-1}$ -almost every $\rho > \tau$, \mathcal{F}^* has constant conditional power function given $\Theta \in \Upsilon_\rho$.
- (iii) For any random variable $\Theta \in \mathcal{M}(\Omega, \mathbb{R}^d)$, for $\mathbb{P}(M(\Theta))^{-1}$ -almost every $\rho > \tau$ and for any test \mathcal{J} with level γ and constant conditional power function given $\Theta \in \Upsilon_\rho$, we have:

$$\mathbb{P}[\mathcal{F}^*(\Theta + X) = 1 \mid \Theta \in \Upsilon_\rho] \geq \mathbb{P}[\mathcal{J}(\Theta + X) = 1 \mid \Theta \in \Upsilon_\rho] \quad (2.16)$$

Now that we have defined the γ -MCCP optimality criterion for the GRDT problem, the question is thus: do such tests exist?

2.2.3 Preliminary results

We start the study of the GRDT problem with several preliminary results that will be helpful in the following. The results presented in this section do not depend on the noise distribution X , and therefore do not depend on the group G or the maximal invariant M .

Lemma 2.2.6: Distribution of $M(\theta + X)$

- (i) The distribution of $M(\theta + X)$ depends on θ only through $M(\theta)$, i.e. for any pair $(\theta_1, \theta_2) \in \mathbb{R}^d \times \mathbb{R}^d$ such that $M(\theta_1) = M(\theta_2)$, the random variables $M(\theta_1 + X)$ and $M(\theta_2 + X)$ have the same probability distribution.
- (ii) For any $\rho \in \mathbb{R}$ and any random variable $\Xi \in \mathcal{M}(\Omega, \mathbb{R}^d)$ such that $\Xi \in \Upsilon_\rho$ almost surely, $M(\Xi + X)$ has the same probability distribution as $M(\theta + X)$ for any $\theta \in \Upsilon_\rho$.

Proof of (i). Let $(\theta_1, \theta_2) \in \mathbb{R}^d \times \mathbb{R}^d$ such that $M(\theta_1) = M(\theta_2)$. Since M is a maximal invariant, there exists a linear transformation $g \in G$ such that $\theta_2 = g(\theta_1)$. We have:

$$\begin{aligned} M(\theta_2 + X) &= M(g(\theta_1) + X) \\ &= M(g(\theta_1 + g^{-1}(X))) && \text{by linearity of } g \\ &= M(\theta_1 + g^{-1}(X)) && \text{because } M \text{ is a maximal invariant} \end{aligned}$$

Therefore, we have:

$$\begin{aligned}
 \forall B \in \mathcal{B}, \mathbb{P}M(\theta_2 + X)^{-1}(B) &= \mathbb{P}[M(\theta_2 + X) \in B] \\
 &= \mathbb{P}[M(\theta_1 + g^{-1}(X)) \in B] \\
 &= \mathbb{P}[M(\theta_1 + X) \in B] && \text{because } X \text{ is invariant under } G \\
 &= \mathbb{P}M(\theta_1 + X)^{-1}(B)
 \end{aligned}$$

Therefore $M(\theta_1 + X)$ and $M(\theta_2 + X)$ have the same distribution.

Proof of (ii). Let $\rho \in \mathbb{R}$ and $\Xi \in \mathcal{M}(\Omega, \mathbb{R}^d)$ such that $\Xi \in \Upsilon_\rho$ almost surely. We have:

$$\begin{aligned}
 \forall B \in \mathcal{B}, \mathbb{P}M(\Xi + X)^{-1}(B) &= \int \mathbb{P}[M(\Xi + X) \in B \mid \Xi = \xi] \mathbb{P}\Xi^{-1}(d\xi) \\
 &= \int_{\Upsilon_\rho} \mathbb{P}[M(\xi + X) \in B] \mathbb{P}\Xi^{-1}(d\xi) && \text{because } \Xi \in \Upsilon_\rho \text{ a.s.}
 \end{aligned}$$

From (i), $\xi \mapsto \mathbb{P}[M(\xi + X) \in B]$ is constant on Υ_ρ , hence for any $\theta \in \Upsilon_\rho$:

$$\begin{aligned}
 \forall B \in \mathcal{B}, \mathbb{P}M(\Xi + X)^{-1}(B) &= \mathbb{P}[M(\theta + X) \in B] \int_{\Upsilon_\rho} \mathbb{P}\Xi^{-1}(d\xi) \\
 &= \mathbb{P}[M(\theta + X) \in B] \quad \square
 \end{aligned}$$

Lemma 2.2.7: Constant power function on a subset

Let $\mathcal{J}: \mathbb{R}^d \rightarrow \{0, 1\}$ be a test. Let E be a subset of \mathbb{R}^d . The test \mathcal{J} has constant power function on E — which simply means that the function $\beta_{\mathcal{J}}$ is constant on E — if and only if for any random variable $\Xi \in \mathcal{M}(\Omega, \mathbb{R}^d)$, independent of X , such that $\Xi \in E$ almost surely, we have:

$$\forall \theta \in E, \mathbb{P}[\mathcal{J}(\Xi + X) = 1] = \beta_{\mathcal{J}}(\theta) \quad (2.17)$$

Proof. If we have Eq. (2.17), then we can immediately deduce that \mathcal{J} has constant power function on E . Conversely, note that for any random variable $\Xi \in \mathcal{M}(\Omega, \mathbb{R}^d)$, by definition of a conditional probability, we have:

$$\begin{aligned}
 \mathbb{P}[\mathcal{J}(\Xi + X) = 1] &= \int \mathbb{P}[\mathcal{J}(\Xi + X) = 1 \mid \Xi = \xi] \mathbb{P}\Xi^{-1}(d\xi) \\
 &= \int \mathbb{P}[\mathcal{J}(\xi + X) = 1] \mathbb{P}\Xi^{-1}(d\xi) \\
 &= \int \beta_{\mathcal{J}}(\xi) \mathbb{P}\Xi^{-1}(d\xi)
 \end{aligned}$$

Thus, if $\Xi \in E$ almost surely, we have:

$$\mathbb{P}[\mathcal{J}(\Xi + X) = 1] = \int_E \beta_{\mathcal{J}}(\xi) \mathbb{P}\Xi^{-1}(d\xi)$$

If \mathcal{J} has constant power function on E , then this equation becomes:

$$\forall \theta \in E, \mathbb{P}[\mathcal{J}(\Xi + X) = 1] = \beta_{\mathcal{J}}(\theta) \int_E \mathbb{P}\Xi^{-1}(d\xi) = \beta_{\mathcal{J}}(\theta) \quad \square$$

Lemma 2.2.8: Power function of invariant tests

Any invariant test \mathcal{J} has constant power function on every orbit $\Upsilon_\rho \in \mathfrak{F}$.

Proof. Let \mathcal{T} be an invariant test, $\rho \in \mathbb{R}$ and $(\theta_1, \theta_2) \in \Upsilon_\rho$. Since θ_1 and θ_2 belong to the same orbit Υ_ρ , there exists $g \in G$ such that $\theta_2 = g(\theta_1)$. We have:

$$\begin{aligned}
 \beta_{\mathcal{T}}(\theta_2) &= \mathbb{P}[\mathcal{T}(\theta_2 + X) = 1] \\
 &= \mathbb{P}[\mathcal{T}(g(\theta_1) + X) = 1] \\
 &= \mathbb{P}[\mathcal{T}(g(\theta_1 + g^{-1}(X))) = 1] && \text{by linearity of } g \\
 &= \mathbb{P}[\mathcal{T}(\theta_1 + g^{-1}(X)) = 1] && \text{because } \mathcal{T} \text{ is invariant} \\
 &= \mathbb{P}[\mathcal{T}(\theta_1 + X) = 1] && \text{because } X \text{ is invariant} \\
 &= \beta_{\mathcal{T}}(\theta_1)
 \end{aligned}$$

Therefore \mathcal{T} has constant power function on every orbit $\Upsilon_\rho \in \mathfrak{F}$. \square

A direct consequence of this is that any thresholding test \mathcal{T}_η has constant power function on every orbit, since thresholding tests are invariant.

The following lemma is the equivalent of Lemma 1.2.5, expressing that the size and the conditional size of any test are identical for the GRDT problem.

Lemma 2.2.9: Size and conditional size

For any test $\mathcal{T}: \mathbb{R}^d \rightarrow \{0, 1\}$, we have:

$$\sup_{\Theta \in \mathcal{M}(\Omega, \mathbb{R}^d): \mathbb{P}[M(\Theta) \leq \tau] \neq 0} \mathbb{P}[\mathcal{T}(\Theta + X) = 1 \mid M(\Theta) \leq \tau] = \alpha_{\mathcal{T}} \quad (2.18)$$

Proof. Let \mathcal{T} be some test.

First, we recall the definition of $\alpha_{\mathcal{T}}$, given previously in Eq. (2.14):

$$\alpha_{\mathcal{T}} = \sup_{\theta \in \mathbb{R}^d: M(\theta) \leq \tau} \beta_{\mathcal{T}}(\theta) \quad (2.19)$$

Let $\theta \in \mathbb{R}^d$ such that $M(\theta) \leq \tau$. For any $\Theta \in \mathcal{M}(\Omega, \mathbb{R}^d)$ such that $\Theta = \theta$ almost surely, we have $\mathbb{P}[\mathcal{T}(\Theta + X) = 1 \mid M(\Theta) \leq \tau] = \beta_{\mathcal{T}}(\theta)$. Since this is true for any θ such that $M(\theta) \leq \tau$, we have:

$$\alpha_{\mathcal{T}} \leq \sup_{\Theta \in \mathcal{M}(\Omega, \mathbb{R}^d): \mathbb{P}[M(\Theta) \leq \tau] \neq 0} \mathbb{P}[\mathcal{T}(\Theta + X) = 1 \mid M(\Theta) \leq \tau]$$

Conversely, for any $\Theta \in \mathcal{M}(\Omega, \mathbb{R}^d)$ independent of $X \sim \mathcal{N}(0, C)$, any $\theta \in \mathbb{R}^d$, and any borel set B of \mathbb{R} , we have the following:

$$\mathbb{P}\left([\mathcal{T}(\Theta + X) = 1] \cap [M(\Theta) \in B] \mid \Theta = \theta\right) = \mathbb{P}\left([\mathcal{T}(\theta + X) = 1] \cap [M(\theta) \in B]\right)$$

In the following, $\mathbb{1}_B$ is the indicator function of B , defined by:

$$\begin{aligned}
 \mathbb{1}_B: \mathbb{R} &\rightarrow \{0, 1\} \\
 x &\mapsto \begin{cases} 1 & \text{if } x \in B \\ 0 & \text{otherwise} \end{cases}
 \end{aligned}$$

There are two cases to consider:

- If $M(\theta) \in B$, then $\mathbb{1}_B(M(\theta))\beta_{\mathcal{T}}(\theta) = \beta_{\mathcal{T}}(\theta)$ and:

$$\mathbb{P}\left([\mathcal{T}(\Theta + X) = 1] \cap [M(\Theta) \in B] \mid \Theta = \theta\right) = \beta_{\mathcal{T}}(\theta)$$

- Otherwise, if $M(\theta) \notin B$, then $\mathbb{1}_B(M(\theta))\beta_{\mathcal{T}}(\theta) = 0$ and:

$$\mathbb{P}\left([\mathcal{T}(\Theta + X) = 1] \cap [M(\Theta) \in B] \mid \Theta = \theta\right) = 0$$

Combining these two cases yields:

$$\mathbb{P}\left([\mathcal{J}(\Theta + X) = 1] \cap [M(\Theta) \in B] \mid \Theta = \theta\right) = \mathbb{1}_B(M(\theta))\beta_{\mathcal{J}}(\theta) \quad (2.20)$$

In particular, if Θ is such that $\mathbb{P}[M(\Theta) \leq \tau] \neq 0$ and is independent of X , then it follows from the definition of a conditional probability and Eq. (2.20) with $B = (-\infty, \tau]$ that:

$$\mathbb{P}\left([\mathcal{J}(\Theta + X) = 1] \cap [M(\Theta) \leq \tau]\right) = \int_{\{\theta \in \mathbb{R}^d, M(\theta) \leq \tau\}} \beta_{\mathcal{J}}(\theta) \mathbb{P}_{\Theta}(\mathrm{d}\theta)$$

The right hand side of this equality is less than or equal to $\alpha_{\mathcal{J}}\mathbb{P}[M(\Theta) \leq \tau]$. From Bayes' rule, we get that:

$$\sup_{\Theta \in \mathcal{M}(\Omega, \mathbb{R}^d): \mathbb{P}[M(\Theta) \leq \tau] \neq 0} \mathbb{P}[\mathcal{J}(\Theta + X) = 1 \mid M(\Theta) \leq \tau] \leq \alpha_{\mathcal{J}}$$

which completes the proof. \square

Lemma 2.2.10: Constant condition power function on an orbit

Let $\rho \in \mathbb{R}$ and $\Theta \in \mathcal{M}(\Omega, \mathbb{R}^d)$. If a test \mathcal{J} has constant conditional power function given $\Theta \in \Upsilon_{\rho}$, then \mathcal{J} has constant power function on Υ_{ρ} .

Proof. This is a direct consequence of the definition of a constant conditional power function (see Definition 2.2.4). \square

Lemma 2.2.11: Constant conditional power function and constant power function

(i) A test \mathcal{J} has constant power function on every orbit $\Upsilon_{\rho} \in \mathfrak{F}$ if and only if there exists a function $Q_{\mathcal{J}}^*: \mathbb{R} \rightarrow [0, 1]$ such that for any $\Theta \in \mathcal{M}(\Omega, \mathbb{R}^d)$ and for $\mathbb{P}M(\Theta)^{-1}$ -almost every $\rho \in \mathbb{R}$:

$$Q_{\mathcal{J}}^*(\rho) = \mathbb{P}[\mathcal{J}(\Theta + X) = 1 \mid \Theta \in \Upsilon_{\rho}] \quad (2.21)$$

(ii) If a test \mathcal{J} has constant power function on every orbit $\Upsilon_{\rho} \in \mathfrak{F}$ and $Q_{\mathcal{J}}^*: \mathbb{R} \rightarrow [0, 1]$ satisfies Eq. (2.21), then:

$$\forall \rho \in \mathbb{R}, \forall \theta \in \Upsilon_{\rho}, \beta_{\mathcal{J}}(\theta) = Q_{\mathcal{J}}^*(\rho)$$

(iii) A test \mathcal{J} has constant power function on every orbit $\Upsilon_{\rho} \in \mathfrak{F}$ if and only if, for any $\Theta \in \mathcal{M}(\Omega, \mathbb{R}^d)$, the test \mathcal{J} has constant conditional power given $\Theta \in \Upsilon_{\rho}$ for $\mathbb{P}M(\Theta)^{-1}$ -almost every $\rho \in \mathbb{R}$.

Proof. Let $\mathcal{J}: \mathbb{R}^d \rightarrow \{0, 1\}$ be a test. We will start by proving statement (i).

- First, assume that \mathcal{J} has constant power function on every orbit $\Upsilon_{\rho} \in \mathfrak{F}$. We can define the function $Q_{\mathcal{J}}^*: \mathbb{R} \rightarrow [0, 1]$ such that:

$$\forall \rho \in \mathbb{R}, \forall \theta \in \Upsilon_{\rho}, Q_{\mathcal{J}}^*(\rho) = \beta_{\mathcal{J}}(\theta)$$

Let $\Theta \in \mathcal{M}(\Omega, \mathbb{R}^d)$ be a random variable independent of $X \sim \mathcal{N}(0, C)$ and B a Borel set of \mathbb{R} . From Eq. (2.20), we get that:

$$\forall \theta \in \mathbb{R}^d, \mathbb{P}\left([\mathcal{J}(\Theta + X) = 1] \cap [M(\Theta) \in B] \mid \Theta = \theta\right) = \mathbb{1}_B(M(\theta))Q_{\mathcal{J}}^*(M(\theta))$$

Using the standard change of variable formula [33, Theorem 16.13], we have:

$$\int \mathbb{1}_B(M(\theta))Q_{\mathcal{J}}^*(M(\theta))\mathbb{P}^{\Theta^{-1}}(\mathrm{d}\theta) = \int_B Q_{\mathcal{J}}^*(\rho)\mathbb{P}M(\Theta)^{-1}(\mathrm{d}\rho)$$

Therefore:

$$\mathbb{P}([\mathcal{J}(\Theta + X) = 1] \cap [M(\Theta) \in B]) = \int_B Q_{\mathcal{J}}^*(\rho) \mathbb{P}M(\Theta)^{-1}(d\rho)$$

On the other hand, we have:

$$\mathbb{P}([\mathcal{J}(\Theta + X) = 1] \cap [M(\Theta) \in B]) = \int_B \mathbb{P}[\mathcal{J}(\Theta + X) = 1 \mid M(\Theta) = \rho] \mathbb{P}M(\Theta)^{-1}(d\rho)$$

Since this is true for any Borel set B of \mathbb{R} , from the definition of a conditional probability, we can deduce that for $\mathbb{P}M(\Theta)^{-1}$ -almost every $\rho \in \mathbb{R}$ we have:

$$Q_{\mathcal{J}}^*(\rho) = \mathbb{P}[\mathcal{J}(\Theta + X) = 1 \mid \Theta \in \Upsilon_{\rho}]$$

- Conversely, assume that there is a function $Q_{\mathcal{J}}^*: \mathbb{R} \rightarrow [0, 1]$ which satisfies Eq. (2.21). Given any $\rho \in \mathbb{R}$, any $\theta \in \Upsilon_{\rho}$ and any random variable $\Theta \in \mathcal{M}(\Omega, \mathbb{R}^d)$ such that $\Theta = \theta$ almost surely, we have:

$$\beta_{\mathcal{J}}(\theta) = \mathbb{P}[\mathcal{J}(\Theta + X) = 1 \mid \Theta \in \Upsilon_{\rho}]$$

Therefore we have:

$$\forall \rho \in \mathbb{R}, \forall \theta \in \Upsilon_{\rho}, \beta_{\mathcal{J}}(\theta) = Q_{\mathcal{J}}^*(\rho)$$

Hence \mathcal{J} has constant power function on every orbit $\Upsilon_{\rho} \in \mathfrak{F}$. This also proves statement (ii).

Statement (iii) is a direct consequence of statements (i) and (ii). \square

The following theorem expresses a sufficient condition to have a γ -MCCP test. Assuming that we can find a test that fulfills the first two criteria of Definition 2.2.5 — meaning that we have a test with level γ and constant conditional power function —, this theorem gives a sufficient condition for this test to be most powerful, and therefore γ -MCCP.

Theorem 2.2.12: Sufficient condition for a γ -MCCP test

Let \mathcal{T}^ be a test with level γ for the GRDT problem such that, for any $\Theta \in \mathcal{M}(\Omega, \mathbb{R}^d)$, \mathcal{T}^* has constant conditional power function given $\Theta \in \Upsilon_{\rho}$ for $\mathbb{P}M(\Theta)^{-1}$ -almost every $\rho > \tau$. If, given any $\rho > \tau$, there exists two random variables Ξ_0 and Ξ_1 , with $\Xi_0 \in \Upsilon_{\tau}$ and $\Xi_1 \in \Upsilon_{\rho}$ almost surely, such that the test \mathcal{T}^* is most powerful with level γ for testing $\mathcal{H}_0: Y = \Xi_0 + X$ against $\mathcal{H}_1: Y = \Xi_1 + X$, then \mathcal{T}^* is γ -MCCP for the GRDT problem.*

Proof. Assume that for any $\rho > \tau$ there exists two random variables Ξ_0 and Ξ_1 with $\Xi_0 \in \Upsilon_{\tau}$ and $\Xi_1 \in \Upsilon_{\rho}$ almost surely, such that the test \mathcal{T}^* is most powerful with level γ for testing \mathcal{H}_0 against \mathcal{H}_1 . Let $\Theta \in \mathcal{M}(\Omega, \mathbb{R}^d)$. Let $\rho > \tau$ and let \mathcal{J} be a test with level γ for the GRDT problem and constant conditional power function given $\Theta \in \Upsilon_{\rho}$. By hypothesis, there exists Ξ_0 and Ξ_1 such that \mathcal{T}^* is most powerful with level γ for testing \mathcal{H}_0 against \mathcal{H}_1 . We have:

$$\begin{aligned} \mathbb{P}[\mathcal{J}(\Xi_0 + X) = 1] &= \mathbb{P}[\mathcal{J}(\Xi_0 + X) = 1 \mid \Xi_0 \in \Upsilon_{\tau}] && \text{because } \Xi_0 \in \Upsilon_{\tau} \text{ almost surely} \\ &\leq \sup_{\Theta \in \mathcal{M}(\Omega, \mathbb{R}^d): \mathbb{P}[M(\Theta) \leq \tau] \neq 0} \mathbb{P}[\mathcal{J}(\Theta + X) = 1 \mid M(\Theta) \leq \tau] \\ &\leq \alpha_{\mathcal{J}} && \text{from Lemma 2.2.9} \\ &\leq \gamma && \text{because } \mathcal{J} \text{ has level } \gamma \text{ for GRDT} \end{aligned}$$

Therefore \mathcal{J} has level γ for \mathcal{H}_0 against \mathcal{H}_1 . Since \mathcal{T}^* is most powerful with level γ for testing \mathcal{H}_0 against \mathcal{H}_1 , we have:

$$\mathbb{P}[\mathcal{T}^*(\Xi_1 + X) = 1] \geq \mathbb{P}[\mathcal{J}(\Xi_1 + X) = 1]$$

Since \mathcal{T} and \mathcal{T}^* both have constant conditional power function given $\Theta \in \Upsilon_\rho$, we have:

$$\begin{aligned} \forall \theta \in \Upsilon_\rho, \mathbb{P}[\mathcal{T}(\Theta + X) = 1 \mid \Theta \in \Upsilon_\rho] &= \beta_{\mathcal{T}}(\theta) \\ \text{and } \forall \theta \in \Upsilon_\rho, \mathbb{P}[\mathcal{T}^*(\Theta + X) = 1 \mid \Theta \in \Upsilon_\rho] &= \beta_{\mathcal{T}^*}(\theta) \end{aligned}$$

From Lemma 2.2.10, since \mathcal{T} and \mathcal{T}^* both have constant conditional power function given $\Theta \in \Upsilon_\rho$, they both have constant power function on Υ_ρ . From Lemma 2.2.7, since $\Xi_1 \in \Upsilon_\rho$ almost surely, we also have:

$$\begin{aligned} \forall \theta \in \Upsilon_\rho, \mathbb{P}[\mathcal{T}(\Xi_1 + X) = 1] &= \beta_{\mathcal{T}}(\theta) \\ \text{and } \forall \theta \in \Upsilon_\rho, \mathbb{P}[\mathcal{T}^*(\Xi_1 + X) = 1] &= \beta_{\mathcal{T}^*}(\theta) \end{aligned}$$

Hence:

$$\begin{aligned} \mathbb{P}[\mathcal{T}(\Xi_1 + X) = 1] &= \mathbb{P}[\mathcal{T}(\Theta + X) = 1 \mid \Theta \in \Upsilon_\rho] \\ \text{and } \mathbb{P}[\mathcal{T}^*(\Xi_1 + X) = 1] &= \mathbb{P}[\mathcal{T}^*(\Theta + X) = 1 \mid \Theta \in \Upsilon_\rho] \end{aligned}$$

Therefore:

$$\mathbb{P}[\mathcal{T}^*(\Theta + X) = 1 \mid \Theta \in \Upsilon_\rho] \geq \mathbb{P}[\mathcal{T}(\Theta + X) = 1 \mid \Theta \in \Upsilon_\rho]$$

This proves that \mathcal{T}^* is γ -MCCP. □

2.2.4 Generalization when the maximal invariant is the euclidean norm

Until now, the only assumption we made on the noise X is that the group G is composed of linear transformations g . However, we were unable so far to find a γ -MCCP test using only these very limited assumptions. Therefore, as a starting point, we decided to focus on using the euclidean norm $\|\cdot\|_2$ as the maximal invariant, since we know that we can find a γ -MCCP test in the Gaussian case. This generalization is not as broad as we originally intended, but would still let us consider other spherically invariant probability distributions, such as generalized Gaussian distributions.

The following lemma provides an important relation between the probability density function of a spherically invariant variable and of its norm. This is an important lemma, as it is what allows us to continue our reasoning using the euclidean norm. For other maximal invariants, we do not necessarily have a similar relation, which is likely why we have trouble working in the general case.

Lemma 2.2.13: Probability density function of spherically invariant random variables

For any spherically invariant random variable $Z \in \mathcal{M}(\Omega, \mathbb{R}^d)$ that has a probability density function f_Z , $\|Z\|_2$ also has a probability density function $f_{\|Z\|_2}$ and we have:

$$\forall z \in \mathbb{R}^d, f_{\|Z\|_2}(\|z\|_2) = m(S^{d-1}) \|z\|_2^{d-1} f_Z(z)$$

where $m(S^{d-1})$ is the surface measure of the unit sphere S^{d-1} in \mathbb{R}^d :

$$m(S^{d-1}) = \frac{2\pi^{d/2}}{\Gamma(d/2)}$$

Proof. This is a consequence of [24, Proposition 7.15] □

From Lemma 2.2.6, for any $\theta \in \mathbb{R}^d$ the distribution of $\|\theta + X\|_2$ only depends on θ through $\|\theta\|_2$. We denote by P_ρ the probability distribution of $\|\theta + X\|_2$ with $\rho = \|\theta\|_2$, and f_ρ the associated probability density function.

In order to find an optimal test, we need to make an additional assumption, which is that the family of densities $\{f_\rho, \rho \geq 0\}$ has a monotone likelihood ratio (cf. Definition 1.1.5), i.e.:

$$\forall \rho_1 > \rho_0, \forall x_1 > x_0, \frac{f_{\rho_1}(x_1)}{f_{\rho_0}(x_1)} \geq \frac{f_{\rho_1}(x_0)}{f_{\rho_0}(x_0)}$$

For any $\rho > 0$, let $\lambda_\gamma(\rho)$ be such that $\mathbb{P}[\|\theta + X\|_2 > \lambda_\gamma(\rho)] = \gamma$ for any $\theta \in \Upsilon_\rho$.

Lemma 2.2.14: Size of $\mathcal{J}_{\lambda_\gamma(\tau)}$

Under the assumption that the family of densities $\{f_\rho, \rho \geq 0\}$ has a monotone likelihood ratio, the thresholding test $\mathcal{J}_{\lambda_\gamma(\tau)}$ has size γ for the GRDT problem.

Proof. Consider the following hypothesis problem:

$$(*) \begin{cases} \text{Observation: } \exists \rho \geq 0, Z \sim P_\rho \\ \mathcal{H}_0: \rho \leq \tau \\ \mathcal{H}_1: \rho > \tau \end{cases}$$

From the Karlin-Rubin theorem (cf. Theorem 1.1.6), since the family of distributions $\{f_\rho, \rho \geq 0\}$ has monotone likelihood ratio, the test $\bar{\mathcal{J}}_{\lambda_\gamma(\tau)}$ defined by

$$\bar{\mathcal{J}}_{\lambda_\gamma(\tau)}: \mathbb{R} \rightarrow \{0, 1\}$$

$$z \mapsto \begin{cases} 0 & \text{if } z \leq \lambda_\gamma(\tau) \\ 1 & \text{if } z > \lambda_\gamma(\tau) \end{cases}$$

has size γ for (*):

$$\sup_{\rho \in \mathbb{R}: \rho \leq \tau} \beta_{\bar{\mathcal{J}}_{\lambda_\gamma(\tau)}}(\rho) = \gamma$$

By definition of $\mathcal{J}_{\lambda_\gamma(\tau)}$, we have $\mathcal{J}_{\lambda_\gamma(\tau)} = \bar{\mathcal{J}}_{\lambda_\gamma(\tau)} \circ \|\cdot\|_2$, and:

$$\begin{aligned} \forall \theta \in \mathbb{R}^d, \beta_{\mathcal{J}_{\lambda_\gamma(\tau)}}(\theta) &= \mathbb{P}[\mathcal{J}_{\lambda_\gamma(\tau)}(\theta + X) = 1] \\ &= \mathbb{P}[\bar{\mathcal{J}}_{\lambda_\gamma(\tau)}(\|\theta + X\|_2) = 1] \\ &= \mathbb{P}[\bar{\mathcal{J}}_{\lambda_\gamma(\tau)}(Z) = 1] && \text{with } Z \sim P_{\|\theta\|_2} \\ &= \beta_{\bar{\mathcal{J}}_{\lambda_\gamma(\tau)}}(\|\theta\|_2) \end{aligned}$$

Therefore:

$$\begin{aligned} \alpha_{\mathcal{J}_{\lambda_\gamma(\tau)}} &= \sup_{\theta \in \mathbb{R}^d: \|\theta\|_2 \leq \tau} \beta_{\mathcal{J}_{\lambda_\gamma(\tau)}}(\theta) \\ &= \sup_{\theta \in \mathbb{R}^d: \|\theta\|_2 \leq \tau} \beta_{\bar{\mathcal{J}}_{\lambda_\gamma(\tau)}}(\|\theta\|_2) \\ &= \sup_{\rho \in [0, \tau]} \beta_{\bar{\mathcal{J}}_{\lambda_\gamma(\tau)}}(\rho) \\ &= \gamma \end{aligned}$$

Hence the test $\mathcal{J}_{\lambda_\gamma(\tau)}$ has size γ for the RDT problem. \square

Lemma 2.2.15

Let $\rho_1 > \rho_0$. Let Ξ_0 and Ξ_1 be two random variables uniformly distributed on $\rho_0 S^{d-1}$ and $\rho_1 S^{d-1}$ respectively. Given any $\gamma \in (0, 1)$, the test $\mathcal{J}_{\lambda_\gamma(\rho_0)}$ is most powerful with size γ for testing $\mathcal{H}_0: Y = \Xi_0 + X$ against $\mathcal{H}_1: Y = \Xi_1 + X$.

Proof. The Neyman-Pearson test \mathcal{J}_{NP} of level γ for testing $\mathcal{H}_0: Y = \Xi_0 + X$ against $\mathcal{H}_1: Y = \Xi_1 + X$

is given by:

$$\forall x \in \mathbb{R}^d, \mathcal{J}_{NP}(x) = \begin{cases} 0 & \text{if } \frac{f_{\Xi_1+X}(x)}{f_{\Xi_0+X}(x)} \leq \lambda \\ 1 & \text{if } \frac{f_{\Xi_1+X}(x)}{f_{\Xi_0+X}(x)} > \lambda \end{cases}$$

with λ determined by $\mathbb{P}\left[\frac{f_{\Xi_1+X}(\Xi_0+X)}{f_{\Xi_0+X}(\Xi_0+X)} > \lambda\right] = \gamma$. Since X , Ξ_0 and Ξ_1 are spherically invariant, $\Xi_0 + X$ and $\Xi_1 + X$ are also spherically invariant, and from Lemma 2.2.13 we have:

$$\forall x \in \mathbb{R}^d, \frac{f_{\Xi_1+X}(x)}{f_{\Xi_0+X}(x)} = \frac{f_{\|\Xi_1+X\|_2}(\|x\|_2)}{f_{\|\Xi_0+X\|_2}(\|x\|_2)}$$

From Lemma 2.2.6, the probability density functions of $\|\Xi_0 + X\|_2$ and $\|\Xi_1 + X\|_2$ are f_{ρ_0} and f_{ρ_1} respectively, hence:

$$\forall x \in \mathbb{R}^d, \frac{f_{\Xi_1+X}(x)}{f_{\Xi_0+X}(x)} = \frac{f_{\rho_1}(\|x\|_2)}{f_{\rho_0}(\|x\|_2)}$$

Since the family of densities $\{f_\rho, \rho > 0\}$ has a monotone likelihood ratio, we can rewrite the test \mathcal{J}_{NP} as:

$$\forall x \in \mathbb{R}^d, \mathcal{J}_{NP}(x) = \begin{cases} 0 & \text{if } \|x\|_2 \leq \mu \\ 1 & \text{if } \|x\|_2 > \mu \end{cases}$$

with μ determined by $\mathbb{P}[\|\Xi_0 + X\|_2 > \mu] = \gamma$. Since $\Xi_0 \in \Upsilon_{\rho_0}$, from Lemma 2.2.6, we have $\mathbb{P}[\|\Xi_0 + X\|_2 > \mu] = \mathbb{P}[\|\theta + X\|_2 > \mu]$ for any $\theta \in \Upsilon_{\rho_0}$. By definition of $\lambda_\gamma(\rho_0)$, we then have $\mathbb{P}[\|\Xi_0 + X\|_2 > \lambda_\gamma(\rho_0)] = \gamma$. Therefore $\mathcal{J}_{NP} = \mathcal{J}_{\lambda_\gamma(\rho_0)}$. Hence $\mathcal{J}_{\lambda_\gamma(\rho_0)}$ is most powerful for testing $\mathcal{H}_0: Y = \Xi_0 + X$ against $\mathcal{H}_1: Y = \Xi_1 + X$. \square

Theorem 2.2.16: Optimality of $\mathcal{J}_{\lambda_\gamma(\tau)}$

Under the assumption that the family of densities $\{f_\rho, \rho \geq 0\}$ has a monotone likelihood ratio, the test $\mathcal{J}_{\lambda_\gamma(\tau)}$ is γ -MCCP.

Proof. From Lemma 2.2.14, $\mathcal{J}_{\lambda_\gamma(\tau)}$ has size γ for the RDT problem. From Lemma 2.2.8, $\mathcal{J}_{\lambda_\gamma(\tau)}$ has constant power function on every orbit $\Upsilon_\rho \in \mathfrak{F}$, therefore from Lemma 2.2.11, it has constant conditional power function given any $\Theta \in \mathcal{M}(\Omega, \mathbb{R}^d)$ for $\mathbb{P}\|\Theta\|_2^{-1}$ -almost every $\rho > 0$. Applying Lemma 2.2.15 with $\rho_0 = \tau$ and every $\rho_1 > \tau$ and Theorem 2.2.12, we get that $\mathcal{J}_{\lambda_\gamma(\tau)}$ is γ -MCCP for the RDT problem. \square

To summarize our findings, we are able to find a γ -MCCP test under the following conditions:

1. The noise distribution is spherically invariant (i.e. the maximal invariant is the euclidean norm).
2. The family of densities $\{f_\rho, \rho > 0\}$ of $\|\theta + X\|_2$ where $\theta \in \Upsilon_\rho$ has a monotone likelihood ratio.

These two constraints are actually quite restrictive, and make it very difficult to find applicable distributions other than the Gaussian one. The second requirement is particularly constraining, because we cannot simply pick a family of densities that has a monotone likelihood ratio, and work from there. Indeed, the densities of this family have to correspond to the densities of $\|\theta + X\|_2$, which is completely determined by the distribution chosen for X . Therefore we have to first choose a distribution for X , and then verify whether the resulting family of distributions $\{f_\rho, \rho > 0\}$ has a monotone likelihood ratio. We describe in Appendix B how the calculations unfold in the Gaussian case. Outside of the Gaussian case, we have not been able to find other distributions that fulfill both of these requirements.

It is not impossible that the Gaussian distribution is the only one that respects these conditions, which could possibly hint that the RDT approach is intrinsically linked to the Gaussian noise. However, we have to keep in mind that the results presented in this chapter only offer sufficient conditions to find an optimal test, and we have not established many necessary conditions for such a test. Therefore, there may exist an alternate reasoning that could lead to an optimal test for other distributions.

2.3 Conclusion and perspectives

While our generalization of the RDT approach was not as successful as we had hoped for, we still managed to present some interesting aspects in this chapter. First, we offer a more complete picture of invariance regarding its use in hypothesis testing, compared to what is commonly found in the literature, which we hope can help clarify its use and the notions involved. Then we attempted to generalize the RDT approach to other noise distributions. While we were not able to find other applicable distributions using the results presented here, we managed nonetheless to revisit the original RDT approach under the invariance viewpoint, minimizing the reliance on the properties of the Gaussian distribution. Doing so also allowed us to see which properties of this distribution are essential to find an optimal test, notably the importance of the monotone likelihood ratio.

While the works presented here did not allow us to generalize the RDT to other distributions, there are a few other approaches we can consider to find other applicable distributions:

- We restricted the latter part of our study to spherically invariant distributions, meaning that the maximal invariant is the euclidean norm. We did so because in this case, we can express easily the density of the norm using Lemma 2.2.13. For other maximal invariants, we would need to figure out if we can find a similar relation, and see if this could lead to an optimal test. Lemma 2.2.13 comes from [24, Proposition 7.15], which gives the density of the maximal invariant $M(X)$ depending on the density of X . However, this density is not given with respect to the usual Lebesgue measure of \mathbb{R} , but to the image measure of the Lebesgue measure of \mathbb{R}^d using M .
- Instead of finding other applicable distributions, an idea could be to use Gaussian densities, expressed with respect to a measure that is not the Lebesgue measure, but a different one.

In the following chapter, we will explore a different idea to extend the RDT approach, which consists in taking into account the effects of estimating the parameters of the RDT model.

3 | Asymptotic Random Distortion Testing

Contents

3.1	Preliminary results	65
3.2	Asymptotic RDT	68
3.2.1	Problem statement	68
3.2.2	Asymptotic size and power	69
3.3	Simulation experiments	75
3.3.1	Level	75
3.3.2	Application to a detection problem	78
3.3.3	Recovering a test with level γ when estimating σ_0	86
3.4	Conclusion	90

In the previous chapter, we attempted to generalize the RDT approach to more noise distributions by studying the role of invariance in the RDT problem. In this chapter, we study a different aspect which is the effect of parameter estimation on the RDT test $\mathcal{J}_{\lambda, \gamma(\tau)}$. So far, we always assumed that the parameter of the problem are perfectly known. But in some scenarios, this hypothesis may be questionable. For example, in many concrete problems, we do not have perfect knowledge of the properties of the noise, and we may need to estimate its covariance matrix. Depending on the quality of this estimation, the effective performance of the test may be negatively affected, resulting in a potential loss of optimality: the test may no longer have the desired level γ , or it may also no longer be most powerful.

With this in mind, we decided to study how estimating the model θ_0 and the noise variance σ_0 (assuming that the noise is Gaussian and white) affects the RDT test. The reason why we decided to focus on the estimation of these two parameters will become apparent in the next chapter, when we present one of the methods we developed to help detect anomalies on signals.

We will first study the asymptotic properties of the RDT test $\mathcal{J}_{\lambda, \gamma(\tau)}$, assuming that we have consistent estimators of θ_0 and σ_0 . The goal here is to find out to what extent the optimality of the RDT test is preserved when we estimate these parameters. In doing so, we will exhibit an appropriate asymptotic optimality criterion. The results of this section make up what we refer to as ARDT (Asymptotic Random Distortion Testing).

After presenting these theoretical results, we will then conduct simulations to not only confirm these results, but also try to quantify the performance loss that occurs when using coarse estimates. We will also present simulations applying the RDT test to a detection problem, compare it to the Neyman-Pearson test, and show how both tests are affected when an unknown distortion is added to the model.

3.1 Preliminary results

We start this chapter by stating a few preliminary results, mainly related to properties of convergence of sequences of random variables.

Lemma 3.1.1: Convergence in distribution to a deterministic constant

Let $(X_n)_{n \in \mathbb{N}} \in \mathcal{M}(\Omega, \mathbb{R}^d)^{\mathbb{N}}$ be a sequence of random vectors and $c \in \mathbb{R}^d$ a deterministic vector. If (X_n) converges in distribution to c , then (X_n) converges in probability to c .

While this result is stated in [34, Corollary 1.5.4 B, p.19], there is no associated proof, and we could not find any in the literature. Therefore we have chosen to present a proof of this result.

Proof. Let $\epsilon > 0$. We want to prove that $\mathbb{P}[\|X_n - c\|_2 \geq \epsilon] \xrightarrow{n \rightarrow \infty} 0$, or equivalently that $\mathbb{P}[\|X_n - c\|_2 < \epsilon] \xrightarrow{n \rightarrow \infty} 1$. Let $B(c, \epsilon)$ be the open ball with center c and radius ϵ . Since $X_n \xrightarrow{D} c$, and since $B(c, \epsilon)$ is an open set, from the Portmanteau theorem [35, Theorem 2.1, p.16], we have:

$$\liminf \mathbb{P}[X_n \in B(c, \epsilon)] \geq \mathbb{P}[c \in B(c, \epsilon)] = 1$$

Since we always have $0 \leq \mathbb{P}[X_n \in B(c, \epsilon)] \leq 1$, we therefore have $\liminf \mathbb{P}[X_n \in B(c, \epsilon)] = 1$. In addition to that, we also have:

$$1 \geq \limsup \mathbb{P}[X_n \in B(c, \epsilon)] \geq \liminf \mathbb{P}[X_n \in B(c, \epsilon)] = 1$$

Hence $\limsup \mathbb{P}[X_n \in B(c, \epsilon)] = \liminf \mathbb{P}[X_n \in B(c, \epsilon)] = 1$. Therefore $\lim \mathbb{P}[X_n \in B(c, \epsilon)]$ exists and $\lim \mathbb{P}[X_n \in B(c, \epsilon)] = 1$. \square

Lemma 3.1.2: Convergence in probability of random vectors

Let $(X_n)_{n \in \mathbb{N}} = (X_n^1, \dots, X_n^d)_{n \in \mathbb{N}} \in \mathcal{M}(\Omega, \mathbb{R}^d)^{\mathbb{N}}$ be a sequence of random vectors and $X = (X^1, \dots, X^d) \in \mathcal{M}(\Omega, \mathbb{R}^d)$ be a random vector. We have:

$$X_n \xrightarrow{\mathbb{P}} X \Leftrightarrow \forall i \in \{1, \dots, d\}, X_n^i \xrightarrow{\mathbb{P}} X^i$$

In other words, convergence in probability of a sequence of random vectors is equivalent to convergence in probability of each component.

While this result may appear fairly trivial, we were unable to find a proof of this lemma anywhere in the literature. Therefore we decided to present a proof here.

Proof.

- First, assume that $X_n \xrightarrow{\mathbb{P}} X$. Let $i \in \{1, \dots, d\}$ and $\epsilon > 0$. We have $|X_n^i - X^i| \leq \|X_n - X\|_2$. Therefore $\mathbb{P}[|X_n^i - X^i| \geq \epsilon] \leq \mathbb{P}[\|X_n - X\|_2 \geq \epsilon]$. Hence $\mathbb{P}[|X_n^i - X^i| \geq \epsilon] \rightarrow 0$.
- Conversely, assume that $\forall i \in \{1, \dots, d\}, X_n^i \xrightarrow{\mathbb{P}} X^i$. Let $\epsilon > 0$. We have:

$$\mathbb{P}[\|X_n - X\|_2 \geq \epsilon] \leq \mathbb{P}\left(\bigcup_{i=1}^d \left[|X_n^i - X^i| \geq \frac{\epsilon}{d}\right]\right) \quad (3.1)$$

Indeed, if for every $i \in \{1, \dots, d\}$ we have $|X_n^i - X^i| < \frac{\epsilon}{d}$, then $\|X_n - X\|_2 < \epsilon$. Therefore by contraposition we get Eq. (3.1). From Eq. (3.1) we get:

$$\mathbb{P}[\|X_n - X\|_2 \geq \epsilon] \leq \sum_{i=1}^d \mathbb{P}\left[|X_n^i - X^i| \geq \frac{\epsilon}{d}\right] \rightarrow 0$$

Therefore $X_n \xrightarrow{\mathbb{P}} X$. \square

Lemma 3.1.3

Let $\Omega_0 \in \Sigma$ such that $\mathbb{P}(\Omega_0) \neq 0$, and let \mathbb{P}_{Ω_0} be the probability measure defined by:

$$\begin{aligned} \mathbb{P}_{\Omega_0}: \Sigma &\rightarrow [0, 1] \\ A &\mapsto \frac{\mathbb{P}(A \cap \Omega_0)}{\mathbb{P}(\Omega_0)} \end{aligned}$$

We have:

$$X_n \xrightarrow{\mathbb{P}_{\Omega_0}} X$$

Proof. Let $\epsilon > 0$. We have:

$$\begin{aligned} \forall n \in \mathbb{N}, \mathbb{P}_{\Omega_0}(\|X_n - X\|_2 \geq \epsilon) &= \frac{\mathbb{P}(\|X_n - X\|_2 \geq \epsilon \cap \Omega_0)}{\mathbb{P}(\Omega_0)} \\ &\leq \frac{\mathbb{P}(\|X_n - X\|_2 \geq \epsilon)}{\mathbb{P}(\Omega_0)} \end{aligned}$$

Since $X_n \xrightarrow{\mathbb{P}} X$, we have $\mathbb{P}(\|X_n - X\|_2 \geq \epsilon) \rightarrow 0$. Therefore $\mathbb{P}_{\Omega_0}(\|X_n - X\|_2 \geq \epsilon) \rightarrow 0$. Hence $X_n \xrightarrow{\mathbb{P}_{\Omega_0}} X$. \square

Lemma 3.1.4

Let $(f_n)_{n \in \mathbb{N}}$ be a sequence of \mathbb{R} -valued measurable functions defined on a measure space $(\Omega, \mathcal{B}, \mu)$, with μ a finite measure. If there exists $C > 0$ such that for all $n \in \mathbb{N}$, we have $|f_n| \leq C$, and if we have:

$$\forall A \in \mathcal{B}, \lim_n \int_A f_n d\mu \geq 0$$

then:

$$\limsup_n f_n \geq 0 \quad \mu - a.e. \quad (3.2)$$

Proof. For every $n \in \mathbb{N}$, let $g_n: x \mapsto \sup\{f_k(x) \mid k \geq n\}$ (pointwise supremum). Let $g = \limsup_n f_n = \lim_n g_n$. Since the functions f_n are measurable and uniformly bounded by C , the functions g_n and g are also measurable and uniformly bounded by C . The constant function $x \in \Omega \mapsto C$ is μ -integrable since μ is a finite measure, therefore from the Dominated Convergence Theorem, we have:

$$\forall A \in \mathcal{B}, \int_A g_n d\mu \xrightarrow{n} \int_A g d\mu \quad (3.3)$$

By definition of g_n , we have:

$$\forall n \in \mathbb{N}, \forall k \geq n, \forall x \in \Omega, g_n(x) \geq f_k(x)$$

Therefore:

$$\forall n \in \mathbb{N}, \forall k \geq n, \forall A \in \mathcal{B}, \int_A g_n d\mu \geq \int_A f_k d\mu$$

From Eq. (3.2), we can deduce that:

$$\forall A \in \mathcal{B}, \forall n \in \mathbb{N}, \int_A g_n d\mu \geq 0$$

Therefore with Eq. (3.3), we get:

$$\forall A \in \mathcal{B}, \int_A g d\mu \geq 0$$

From this we can deduce that $g \geq 0$ μ -almost everywhere, i.e.:

$$\limsup_n f_n \geq 0 \quad \mu - \text{a.e.} \quad \square$$

3.2 Asymptotic RDT

3.2.1 Problem statement

Let us start by stating the problem we consider in this chapter.

Definition 3.2.1: ARDT problem statement

$$\left. \begin{array}{l}
 \textbf{Data model:} \\
 \exists Y \in \mathcal{M}(\Omega, \mathbb{R}^d), \exists \Theta \in \mathcal{M}(\Omega, \mathbb{R}^d), \exists X \in \mathcal{M}(\Omega, \mathbb{R}^d), \\
 \left\{ \begin{array}{l}
 (X \sim \mathcal{N}(0, \sigma I_d) \text{ with } \sigma > 0 \text{ unknown}) \\
 \wedge (\Theta \text{ and } X \text{ are independent}) \\
 \wedge (Y = \Theta + X) \\
 \wedge (\forall y \in \mathbb{R}^d, \exists \omega \in \Omega, y = Y(\omega))
 \end{array} \right. \\
 \textbf{Testing problem:} \\
 \text{Given one realization } y = Y(\omega) = \Theta(\omega) + X(\omega), \text{ determine whether:} \\
 \left\{ \begin{array}{l}
 \mathcal{H}_0: \|\Theta(\omega) - \theta_0\|_2 \leq \sigma_0 \tau \\
 \text{or} \\
 \mathcal{H}_1: \|\Theta(\omega) - \theta_0\|_2 > \sigma_0 \tau
 \end{array} \right. \\
 \text{with } \tau > 0 \text{ known, } \theta_0 \in \mathbb{R}^d \text{ unknown,} \\
 \text{given consistent estimators } \hat{\theta}_n \text{ and } \hat{\sigma}_n \text{ of } \theta_0 \text{ and } \sigma_0 \text{ respectively}
 \end{array} \right\} \quad (3.4)$$

The main difference with the original RDT problem presented in Definition 1.2.2 is of course the fact that θ_0 and σ_0 are unknown. We instead assume that we have access to consistent estimators $\hat{\theta}_n$ and $\hat{\sigma}_n$ of these values. As a reminder, the fact that these estimators are consistent simply means that $\hat{\theta}_n$ and $\hat{\sigma}_n$ converge in probability to θ_0 and σ_0 respectively as n increases.

Another important difference to notice is the fact that the noise X is white, i.e. its covariance matrix is of the form $\sigma_0^2 I_d$. The reason why we are not considering any covariance matrix C is related to the invariance properties of the noise. Indeed, first consider the family of Gaussian distributions with covariance matrix $\sigma^2 I_d$ for all $\sigma > 0$. All of these distributions are spherically invariant. Therefore replacing σ_0 with any estimate $\hat{\sigma}$ will not affect the invariance properties of the noise.

Now consider the case of a family of Gaussian distributions with any covariance matrix C . In this case, while all of these probability distributions are invariant, it is not the same group that leaves every distribution invariant. The group depends on the covariance matrix C (see Eq. (2.10) for the definition of this group), and therefore replacing C with an estimate means we do not necessarily preserve the invariance structure of the noise.

With this in mind, we decided to only consider the case $C = \sigma_0^2 I_d$ for now, which will slightly simplify the study. While this is not exhaustive, it remains an interesting and useful case, notably in the 1-dimensional case. We can also easily extend this study to the case when the covariance matrix has the form $\sigma_0^2 C$ with C known and σ_0 estimated. Note that, as a consequence of this, the Mahalanobis norm $\nu_C(y)$ of any vector $y \in \mathbb{R}^d$ becomes $\|y\|_2 / \sigma_0$.

To study the behavior of a test taking into account the estimates $\hat{\theta}$ and $\hat{\sigma}$, we will consider a slightly modified definition of a test compared to the definition we gave in Definition 1.1.1. We will denote such test $\tilde{\mathcal{T}}$, and they will be functions of three parameters (the observation, the model estimate and the noise standard deviation estimate) defined on $\mathbb{R}^d \times \mathbb{R}^d \times (0, +\infty)$ and taking their values in $\{0, 1\}$. As an example,

we can consider for any $t > 0$ the thresholding test $\tilde{\mathcal{T}}_t$ defined by:

$$\begin{aligned} \tilde{\mathcal{T}}_t: \mathbb{R}^d \times \mathbb{R}^d \times (0, \infty) &\rightarrow \{0, 1\} \\ (y, \theta, \sigma) &\mapsto \begin{cases} 1 & \text{if } \frac{\|y - \theta\|_2}{\sigma} > t \\ 0 & \text{otherwise} \end{cases} \end{aligned} \quad (3.5)$$

These test $\tilde{\mathcal{T}}_t$ are a direct adaptation of the thresholding tests defined in Definition 1.2.8, in which the model θ_0 and the noise standard deviation σ_0 were both assumed to be known. We simply substituted θ_0 and σ_0 with the parameters θ and σ respectively. In addition to this, for each test $\tilde{\mathcal{T}}$, we also define its critical region $K_{\tilde{\mathcal{T}}}$ by:

$$K_{\tilde{\mathcal{T}}} = \{(y, \theta, \sigma) \in \mathbb{R}^d \times \mathbb{R}^d \times (0, \infty) : \tilde{\mathcal{T}}(y, \theta, \sigma) = 1\} \quad (3.6)$$

Note that in the following, we may shorten $\mathbb{R}^d \times \mathbb{R}^d \times (0, \infty)$ to \mathbb{R}^{2d+1} to simplify notation. With this notation, when we consider for example $(y, \theta, \sigma) \in \mathbb{R}^{2d+1}$, it is implied that y , θ and σ belong to \mathbb{R}^d , \mathbb{R}^d , and $(0, \infty)$ respectively.

In the following, we consider two consistent estimators $\hat{\theta}_n$ and $\hat{\sigma}_n$ of θ_0 and σ_0 respectively, i.e. $\hat{\theta}_n \xrightarrow{\mathbb{P}} \theta_0$ and $\hat{\sigma}_n \xrightarrow{\mathbb{P}} \sigma_0$, with $\hat{\sigma}_n > 0$ for every $n \in \mathbb{N}$.

3.2.2 Asymptotic size and power

In this section, we will now derive asymptotic results regarding the size and the power of the test $\tilde{\mathcal{T}}_{\lambda_\gamma(\tau)}$, which is adapted from the test $\mathcal{T}_{\lambda_\gamma(\tau)}$, known to be γ -MCCP for the RDT problem. The thresholding test $\tilde{\mathcal{T}}_{\lambda_\gamma(\tau)}$ is defined by:

$$\begin{aligned} \tilde{\mathcal{T}}_{\lambda_\gamma(\tau)}: \mathbb{R}^d \times \mathbb{R}^d \times (0, \infty) &\rightarrow \{0, 1\} \\ (y, \theta, \sigma) &\mapsto \begin{cases} 1 & \text{if } \frac{\|y - \theta\|_2}{\sigma} > \lambda_\gamma(\tau) \\ 0 & \text{otherwise} \end{cases} \end{aligned} \quad (3.7)$$

with the threshold $\lambda_\gamma(\tau)$ defined in Lemma 1.2.9.

From Chapter 1, we know that the test $\mathcal{T}_{\lambda_\gamma(\tau)}$ is optimal for the RDT problem. Since $\hat{\theta}_n$ and $\hat{\sigma}_n$ converge to θ_0 and σ_0 respectively, we may expect that $\tilde{\mathcal{T}}_{\lambda_\gamma(\tau)}$ is asymptotically optimal, i.e. it should asymptotically have level γ and be asymptotically most powerful.

We start with the asymptotic level of $\tilde{\mathcal{T}}_{\lambda_\gamma(\tau)}$.

Theorem 3.2.2: Asymptotic level of $\tilde{\mathcal{T}}_{\lambda_\gamma(\tau)}$

Let $\mathcal{S} = \{\Xi \in \mathcal{M}(\Omega, \mathbb{R}^d) : \forall n \in \mathbb{N}, \Xi \text{ and } (\hat{\theta}_n, \hat{\sigma}_n) \text{ are independent}\}$. We have:

$$\limsup_n \sup_{\Xi \in \mathcal{S} : \mathbb{P}[\|\Xi - \theta_0\|_2 \leq \sigma_0 \tau] \neq 0} \mathbb{P}[\tilde{\mathcal{T}}_{\lambda_\gamma(\tau)}(\Xi + X, \hat{\theta}_n, \hat{\sigma}_n) = 1 \mid \|\Xi - \theta_0\|_2 \leq \sigma_0 \tau] \leq \gamma$$

A noteworthy difference compared to the RDT case is that while $\mathcal{T}_{\lambda_\gamma(\tau)}$ has exactly size γ , as stated in Theorem 1.2.10, here in the ARDT case, we only know that the test $\tilde{\mathcal{T}}_{\lambda_\gamma(\tau)}$ has asymptotic level γ , i.e. its asymptotic size is lower than or equal to γ .

Proof. Let $\Xi_0 \in \mathcal{S}$ such that $\mathbb{P}[\|\Xi_0 - \theta_0\|_2 \leq \sigma_0 \tau] \neq 0$.

$$\begin{aligned} \mathbb{P}[\tilde{\mathcal{T}}_{\lambda_\gamma(\tau)}(\Xi_0 + X, \hat{\theta}_n, \hat{\sigma}_n) = 1 \mid \|\Xi_0 - \theta_0\|_2 \leq \sigma_0 \tau] \\ &= \int \mathbb{P}[\tilde{\mathcal{T}}_{\lambda_\gamma(\tau)}(\Xi_0 + X, \theta, \sigma) = 1 \mid \|\Xi_0 - \theta_0\|_2 \leq \sigma_0 \tau] \mathbb{P}(\hat{\theta}_n, \hat{\sigma}_n)^{-1}(d\theta, d\sigma) \\ &\leq \int \sup_{\Xi \in \mathcal{M}(\Omega, \mathbb{R}^d): \mathbb{P}[\|\Xi - \theta_0\|_2 \leq \sigma_0 \tau] \neq 0} \mathbb{P}[\tilde{\mathcal{T}}_{\lambda_\gamma(\tau)}(\Xi + X, \theta, \sigma) = 1 \mid \|\Xi - \theta_0\|_2 \leq \sigma_0 \tau] \mathbb{P}(\hat{\theta}_n, \hat{\sigma}_n)^{-1}(d\theta, d\sigma) \end{aligned}$$

From Lemma 1.2.5, since $\tilde{\mathcal{T}}_{\lambda_\gamma(\tau)}(\cdot, \theta, \sigma)$ is a test for any $\theta \in \mathbb{R}^d$ and any $\sigma > 0$, we have:

$$\sup_{\Xi \in \mathcal{M}(\Omega, \mathbb{R}^d): \mathbb{P}[\|\Xi - \theta_0\|_2 \leq \sigma_0 \tau] \neq 0} \mathbb{P}[\tilde{\mathcal{T}}_{\lambda_\gamma(\tau)}(\Xi + X, \theta, \sigma) = 1 \mid \|\Xi - \theta_0\|_2 \leq \sigma_0 \tau] = \sup_{\xi \in \mathbb{R}^d: \|\xi - \theta_0\|_2 \leq \sigma_0 \tau} \mathbb{P}[\tilde{\mathcal{T}}_{\lambda_\gamma(\tau)}(\xi + X, \theta, \sigma) = 1]$$

Therefore:

$$\begin{aligned} \mathbb{P}[\tilde{\mathcal{T}}_{\lambda_\gamma(\tau)}(\Xi_0 + X, \hat{\theta}_n, \hat{\sigma}_n) = 1 \mid \|\Xi_0 - \theta_0\|_2 \leq \sigma_0 \tau] \\ &\leq \int \sup_{\xi \in \mathbb{R}^d: \|\xi - \theta_0\|_2 \leq \sigma_0 \tau} \mathbb{P}[\tilde{\mathcal{T}}_{\lambda_\gamma(\tau)}(\theta + X, \theta, \sigma) = 1] \mathbb{P}(\hat{\theta}_n, \hat{\sigma}_n)^{-1}(d\theta, d\sigma) \\ &\leq \int \sup_{\xi \in \mathbb{R}^d: \|\xi - \theta_0\|_2 \leq \sigma_0 \tau} \mathbb{P}\left[\frac{\|\xi + X - \theta\|_2}{\sigma} > \lambda_\gamma(\tau)\right] \mathbb{P}(\hat{\theta}_n, \hat{\sigma}_n)^{-1}(d\theta, d\sigma) \end{aligned}$$

Let $W = \xi + X - \theta \sim \mathcal{N}(\xi - \theta, \sigma_0^2 I_d)$. Using Eq. (1.19), we get:

$$\begin{aligned} \mathbb{P}\left[\frac{\|W\|_2}{\sigma} > \lambda_\gamma(\tau)\right] &= \mathbb{P}\left[\frac{\|W\|_2}{\sigma_0} > \frac{\sigma}{\sigma_0} \lambda_\gamma(\tau)\right] \\ &= Q_{d/2}\left(\frac{\|\xi - \theta\|_2}{\sigma_0}, \frac{\sigma}{\sigma_0} \lambda_\gamma(\tau)\right) \end{aligned}$$

The function $Q_{d/2}$ is continuous and increases with its first argument, therefore:

$$\begin{aligned} \sup_{\xi \in \mathbb{R}^d: \|\xi - \theta_0\|_2 \leq \sigma_0 \tau} Q_{d/2}\left(\frac{\|\xi - \theta\|_2}{\sigma_0}, \frac{\sigma}{\sigma_0} \lambda_\gamma(\tau)\right) &= Q_{d/2}\left(\sup_{\xi \in \mathbb{R}^d: \|\xi - \theta_0\|_2 \leq \sigma_0 \tau} \frac{\|\xi - \theta\|_2}{\sigma_0}, \frac{\sigma}{\sigma_0} \lambda_\gamma(\tau)\right) \\ &= Q_{d/2}\left(\tau + \frac{\|\theta - \theta_0\|_2}{\sigma_0}, \frac{\sigma}{\sigma_0} \lambda_\gamma(\tau)\right) \end{aligned}$$

Hence:

$$\begin{aligned} \mathbb{P}[\tilde{\mathcal{T}}_{\lambda_\gamma(\tau)}(\Xi_0 + X, \hat{\theta}_n, \hat{\sigma}_n) = 1 \mid \|\Xi_0 - \theta_0\|_2 \leq \sigma_0 \tau] \\ &\leq \int Q_{d/2}\left(\tau + \frac{\|\theta - \theta_0\|_2}{\sigma_0}, \frac{\sigma}{\sigma_0} \lambda_\gamma(\tau)\right) \mathbb{P}(\hat{\theta}_n, \hat{\sigma}_n)^{-1}(d\theta, d\sigma) \end{aligned}$$

This inequality is valid for any $\Xi_0 \in \mathcal{S}$ such that $\mathbb{P}[\|\Xi_0 - \theta_0\|_2 \leq \sigma_0 \tau] \neq 0$, and the right-hand side does not depend on Ξ_0 . Therefore:

$$\begin{aligned} \sup_{\Xi \in \mathcal{S}: \mathbb{P}[\|\Xi - \theta_0\|_2 \leq \sigma_0 \tau] \neq 0} \mathbb{P}[\tilde{\mathcal{T}}_{\lambda_\gamma(\tau)}(\Xi + X, \hat{\theta}_n, \hat{\sigma}_n) = 1 \mid \|\Xi - \theta_0\|_2 \leq \sigma_0 \tau] \\ &\leq \int Q_{d/2}\left(\tau + \frac{\|\theta - \theta_0\|_2}{\sigma_0}, \frac{\sigma}{\sigma_0} \lambda_\gamma(\tau)\right) \mathbb{P}(\hat{\theta}_n, \hat{\sigma}_n)^{-1}(d\theta, d\sigma) \end{aligned}$$

From Appendix C, the Generalized Marcum function $Q_{d/2}$ is uniformly continuous on $[0, +\infty) \times [0, +\infty)$. It is also bounded by definition. Therefore, from the Portmanteau theorem [35], since $(\hat{\theta}_n)_{n \in \mathbb{N}}$ and $(\hat{\sigma}_n)_{n \in \mathbb{N}}$ converges in distribution to θ_0^2 and σ_0^2 respectively, and by definition of $\lambda_\gamma(\tau)$, we have:

$$\begin{aligned} \lim_n \int Q_{d/2} \left(\tau + \frac{\|\theta - \theta_0\|_2}{\sigma_0}, \frac{\sigma}{\sigma_0} \lambda_\gamma(\tau) \right) \mathbb{P}(\hat{\theta}_n, \hat{\sigma}_n)^{-1}(d\theta, d\sigma) &= Q_{d/2}(\tau, \lambda_\gamma(\tau)) \\ &= \gamma \end{aligned}$$

Hence:

$$\limsup_n \sup_{\Xi \in \mathcal{S}: \mathbb{P}[\|\Xi - \theta_0\|_2 \leq \sigma_0 \tau] \neq 0} \mathbb{P}[\tilde{\mathcal{J}}_{\lambda_\gamma(\tau)}(\Xi + X, \hat{\theta}_n, \hat{\sigma}_n) = 1 \mid \|\Xi - \theta_0\|_2 \leq \sigma_0 \tau] \leq \gamma$$

which is the desired result. \square

Now that we have our result regarding the asymptotic level of the test $\tilde{\mathcal{J}}_{\lambda_\gamma(\tau)}$, we can now study its asymptotic power.

Let $Z = (Y, \theta_0, \sigma_0)$ and for every $n \in \mathbb{N}$, let $Z_n = (Y, \hat{\theta}_n, \hat{\sigma}_n)$. From Lemma 3.1.2, Z_n converge in probability to Z , since $\hat{\theta}_n$ and $\hat{\sigma}_n$ converge in probability to θ_0 and σ_0 respectively.

For every $\rho \geq 0$, every $A \in \mathcal{R}^{2d+1}$ and every $n \in \mathbb{N}$, let Π_n and Π be the functions defined by:

$$\Pi_n(\rho, A) = \mathbb{P}[Z_n \in A \mid \|\Theta - \theta_0\|_2 = \rho] \quad (3.8)$$

$$\Pi(\rho, A) = \mathbb{P}[Z \in A \mid \|\Theta - \theta_0\|_2 = \rho] \quad (3.9)$$

Lemma 3.2.3

For any $\mathbb{P}Z^{-1}$ -continuity set $A \in \mathcal{R}^{2d+1}$, we have:

$$\forall B \in \mathcal{R}, \int_B \Pi_n(\rho, A) \mathbb{P}\|\Theta - \theta_0\|_2^{-1}(d\rho) \xrightarrow{n} \int_B \Pi(\rho, A) \mathbb{P}\|\Theta - \theta_0\|_2^{-1}(d\rho)$$

Proof. Let $B \in \mathcal{R}$. By definition of a conditional probability, we have:

$$\mathbb{P}([Z_n \in A] \cap [\|\Theta - \theta_0\|_2 \in B]) = \int_B \Pi_n(\rho, A) \mathbb{P}\|\Theta - \theta_0\|_2^{-1}(d\rho) \quad (3.10)$$

as well as:

$$\mathbb{P}([Z \in A] \cap [\|\Theta - \theta_0\|_2 \in B]) = \int_B \Pi(\rho, A) \mathbb{P}\|\Theta - \theta_0\|_2^{-1}(d\rho) \quad (3.11)$$

There are two possibilities: either $\mathbb{P}[\|\Theta - \theta_0\|_2 \in B] = 0$ or $\mathbb{P}[\|\Theta - \theta_0\|_2 \in B] \neq 0$.

- If $\mathbb{P}[\|\Theta - \theta_0\|_2 \in B] = 0$, then $\mathbb{P}([Z_n \in A] \cap [\|\Theta - \theta_0\|_2 \in B]) = \mathbb{P}([Z \in A] \cap [\|\Theta - \theta_0\|_2 \in B]) = 0$. Therefore we trivially get:

$$\int_B \Pi_n(\rho, A) \mathbb{P}\|\Theta - \theta_0\|_2^{-1}(d\rho) \xrightarrow{n} \int_B \Pi(\rho, A) \mathbb{P}\|\Theta - \theta_0\|_2^{-1}(d\rho)$$

- Otherwise, if $\mathbb{P}[\|\Theta - \theta_0\|_2 \in B] \neq 0$, we have from Bayes' rule:

$$\mathbb{P}([Z_n \in A] \cap [\|\Theta - \theta_0\|_2 \in B]) = \mathbb{P}[Z_n \in A \mid \|\Theta - \theta_0\|_2 \in B] \mathbb{P}[\|\Theta - \theta_0\|_2 \in B]$$

Recall that we consider a probability space $(\Omega, \Sigma, \mathbb{P})$ on which all random variables are defined. Let $\Omega_0 = [\|\Theta - \theta_0\|_2 \in B] \in \Sigma$ and let \mathbb{P}_{Ω_0} be the probability measure defined by:

$$\begin{aligned} \mathbb{P}_{\Omega_0}: \quad \Sigma &\rightarrow [0, 1] \\ M &\mapsto \frac{\mathbb{P}(M \cap \Omega_0)}{\mathbb{P}(\Omega_0)} \end{aligned}$$

From Lemma 3.1.3, since $\mathbb{P}(\Omega_0) \neq 0$, we have $Z_n \xrightarrow{\mathbb{P}_{\Omega_0}} Z$. Since A is a $\mathbb{P}Z^{-1}$ -continuity set, it is also a $\mathbb{P}_{\Omega_0}Z^{-1}$ -continuity set. Indeed, we have:

$$\begin{aligned} \mathbb{P}_{\Omega_0}[Z \in \partial A] &= \frac{\mathbb{P}([Z \in \partial A] \cap \Omega_0)}{\mathbb{P}(\Omega_0)} \\ &\leq \frac{\mathbb{P}[Z \in \partial A]}{\mathbb{P}(\Omega_0)} \\ &= 0 \end{aligned} \quad \text{because } A \text{ is a } \mathbb{P}Z^{-1}\text{-continuity set}$$

Therefore from the Portmanteau theorem [35, Theorem 2.1], we have:

$$\mathbb{P}_{\Omega_0}[Z_n \in A] \xrightarrow{n} \mathbb{P}_{\Omega_0}[Z \in A]$$

Hence:

$$\mathbb{P}([Z_n \in A] \cap [\|\Theta - \theta_0\|_2 \in B]) \xrightarrow{n} \mathbb{P}([Z \in A] \cap [\|\Theta - \theta_0\|_2 \in B])$$

which yields the desired result using Eqs. (3.10) and (3.11). \square

We can now present the result regarding the asymptotic power of $\tilde{\mathcal{T}}_{\lambda_\gamma(\tau)}$, which is our optimality criterion in the ARDT problem.

Theorem 3.2.4: Asymptotic optimality of $\tilde{\mathcal{T}}_{\lambda_\gamma(\tau)}$

For any test $\tilde{\mathcal{T}}: \mathbb{R}^{2d+1} \rightarrow \{0, 1\}$ such that the test $\tilde{\mathcal{T}}(\cdot, \theta_0, \sigma_0)$ has level γ for the RDT problem, and has constant conditional power function given $\|\Theta - \theta_0\|_2 = \rho$ for $\Pr \|\Theta - \theta_0\|_2^{-1}$ -almost every ρ , if $K_{\tilde{\mathcal{T}}}$ is a $\mathbb{P}Z^{-1}$ -continuity set, then we have:

$$\limsup_n (\Pi_n(\rho, K_{\tilde{\mathcal{T}}_{\lambda_\gamma(\tau)}}) - \Pi_n(\rho, K_{\tilde{\mathcal{T}}})) \geq 0 \quad (3.12)$$

Proof. As a reminder, $K_{\mathcal{J}_{\lambda_\gamma(\tau)}}$ and $K_{\tilde{\mathcal{T}}_{\lambda_\gamma(\tau)}}$, the respective critical regions of $\mathcal{J}_{\lambda_\gamma(\tau)}$ and $\tilde{\mathcal{T}}_{\lambda_\gamma(\tau)}$, are defined by:

$$\begin{aligned} K_{\mathcal{J}_{\lambda_\gamma(\tau)}} &= \left\{ y \in \mathbb{R}^d : \frac{\|y - \theta_0\|_2}{\sigma_0} > \lambda_\gamma(\tau) \right\} \\ K_{\tilde{\mathcal{T}}_{\lambda_\gamma(\tau)}} &= \left\{ (y, \theta, \sigma) \in \mathbb{R}^{2d+1} : \frac{\|y - \theta\|_2}{\sigma} > \lambda_\gamma(\tau) \right\} \end{aligned}$$

Let $\tilde{\mathcal{T}}: \mathbb{R}^{2d+1} \rightarrow \{0, 1\}$ be a test such that the test \mathcal{T} defined by

$$\begin{aligned} \mathcal{T}: \mathbb{R}^d &\rightarrow \{0, 1\} \\ y &\mapsto \tilde{\mathcal{T}}(y, \theta_0, \sigma_0) \end{aligned}$$

has level γ for the RDT problem and has constant conditional power function given $\|\Theta - \theta_0\|_2 = \rho$ for $\mathbb{P}\|\Theta - \theta_0\|_2^{-1}$ -almost every ρ .

Since $\mathcal{J}_{\lambda_\gamma(\tau)}$ is γ -MCCP, we have by definition for $\mathbb{P}\|\Theta - \theta_0\|_2^{-1}$ -almost every $\rho > \tau$:

$$\mathbb{P}[\mathcal{J}_{\lambda_\gamma(\tau)}(Y) = 1 \mid \|\Theta - \theta_0\|_2 = \rho] \geq \mathbb{P}[\mathcal{T}(Y) = 1 \mid \|\Theta - \theta_0\|_2 = \rho] \quad (3.13)$$

We can equivalently write that for $\mathbb{P}\|\Theta - \theta_0\|_2^{-1}$ -almost every $\rho > \tau$, we have:

$$\mathbb{P}[Y \in K_{\mathcal{J}_{\lambda_\gamma(\tau)}} \mid \|\Theta - \theta_0\|_2 = \rho] \geq \mathbb{P}[Y \in K_{\mathcal{T}} \mid \|\Theta - \theta_0\|_2 = \rho] \quad (3.14)$$

By definition of Z (reminder: $Z = (Y, \theta_0, \sigma_0)$), we have:

$$Y \in K_{\mathcal{F}_{\lambda_\gamma(\tau)}} \Leftrightarrow \frac{\|y - \theta_0\|_2}{\sigma_0} > \lambda_\gamma(\tau) \Leftrightarrow Z \in K_{\tilde{\mathcal{F}}_{\lambda_\gamma(\tau)}}$$

and $Y \in K_{\mathcal{F}} \Leftrightarrow \mathcal{J}(y) = 1 \Leftrightarrow \tilde{\mathcal{F}}(y, \theta_0, \sigma_0) = 1 \Leftrightarrow Z \in K_{\tilde{\mathcal{F}}}$

Therefore, for $\mathbb{P}\|\Theta - \theta_0\|_2^{-1}$ -almost every $\rho > \tau$:

$$\mathbb{P}[Z \in K_{\tilde{\mathcal{F}}_{\lambda_\gamma(\tau)}} \mid \|\Theta - \theta_0\|_2 = \rho] \geq \mathbb{P}[Z \in K_{\tilde{\mathcal{F}}} \mid \|\Theta - \theta_0\|_2 = \rho] \quad (3.15)$$

Equation (3.15) also means that for $\mathbb{P}\|\Theta - \theta_0\|_2^{-1}$ -almost every $\rho > \tau$:

$$\Pi(\rho, K_{\tilde{\mathcal{F}}_{\lambda_\gamma(\tau)}}) \geq \Pi(\rho, K_{\tilde{\mathcal{F}}}) \quad (3.16)$$

Therefore:

$$\forall B \in \mathcal{B}((\tau, \infty)), \int_B \Pi(\rho, K_{\tilde{\mathcal{F}}_{\lambda_\gamma(\tau)}}) \mathbb{P}\|\Theta - \theta_0\|_2^{-1}(\mathrm{d}\rho) \geq \int_B \Pi(\rho, K_{\tilde{\mathcal{F}}}) \mathbb{P}\|\Theta - \theta_0\|_2^{-1}(\mathrm{d}\rho) \quad (3.17)$$

If $K_{\tilde{\mathcal{F}}_{\lambda_\gamma(\tau)}}$ and $K_{\tilde{\mathcal{F}}}$ are $\mathbb{P}Z^{-1}$ -continuity sets, then from Lemma 3.2.3 we have:

$$\forall B \in \mathcal{B}((\tau, \infty)), \lim_n \int_B \Pi_n(\rho, K_{\tilde{\mathcal{F}}_{\lambda_\gamma(\tau)}}) \mathbb{P}\|\Theta - \theta_0\|_2^{-1}(\mathrm{d}\rho) \geq \lim_n \int_B \Pi_n(\rho, K_{\tilde{\mathcal{F}}}) \mathbb{P}\|\Theta - \theta_0\|_2^{-1}(\mathrm{d}\rho) \quad (3.18)$$

We can rewrite Eq. (3.18) as:

$$\forall B \in \mathcal{B}((\tau, \infty)), \lim_n \int_B (\Pi_n(\rho, K_{\tilde{\mathcal{F}}_{\lambda_\gamma(\tau)}}) - \Pi_n(\rho, K_{\tilde{\mathcal{F}}})) \mathbb{P}\|\Theta - \theta_0\|_2^{-1}(\mathrm{d}\rho) \geq 0 \quad (3.19)$$

By applying Lemma 3.1.4, we get that for $\mathbb{P}\|\Theta - \theta_0\|_2^{-1}$ -almost every $\rho > \tau$:

$$\limsup_n (\Pi_n(\rho, K_{\tilde{\mathcal{F}}_{\lambda_\gamma(\tau)}}) - \Pi_n(\rho, K_{\tilde{\mathcal{F}}})) \geq 0 \quad (3.20)$$

which is the sought result. \square

In this theorem, we specify that we consider the class of tests $\tilde{\mathcal{F}}$ whose critical region $K_{\tilde{\mathcal{F}}}$ is a $\mathbb{P}Z^{-1}$ -continuity set. We can wonder what this constraint implies on the test $\tilde{\mathcal{F}}$. If we consider a Borel set $A \in \mathcal{R}^{2d+1}$, what does $A \in \mathcal{R}^{2d+1}$ being a $\mathbb{P}Z^{-1}$ -continuity set mean? Going back to the definition of a continuity set, this means that $\mathbb{P}[Z \in \partial A] = 0$. Since $Z = (Y, \theta_0, \sigma_0)$, this means that $\mathbb{P}[(Y, \theta_0, \sigma_0) \in (\partial A)] = 0$. In the following, for any set $B \in \mathcal{R}^{2d+1}$, we will denote by B_{θ_0, σ_0} its (θ_0, σ_0) -section, which is defined by:

$$(\partial B)_{\theta_0, \sigma_0} = \{y \in \mathbb{R}^d : (y, \theta_0, \sigma_0) \in \partial B\} \quad (3.21)$$

With this notation, if A is a $\mathbb{P}Z^{-1}$ -continuity set, then we have $\mathbb{P}[Y \in (\partial A)_{\theta_0, \sigma_0}] = 0$, where $(\partial A)_{\theta_0, \sigma_0}$ denotes the (θ_0, σ_0) -section of ∂A , the boundary of A . For a test $\tilde{\mathcal{F}}$, this means that we have $\mathbb{P}[Y \in (\partial K_{\tilde{\mathcal{F}}})_{\theta_0, \sigma_0}] = 0$.

One aspect that may be interesting to study is whether this constraint tells us anything on the set $\partial(K_{\tilde{\mathcal{F}}}_{\theta_0, \sigma_0})$. Indeed, let \mathcal{J} be the test $\tilde{\mathcal{F}}(\cdot, \theta_0, \sigma_0)$. The set $K_{\tilde{\mathcal{F}}}_{\theta_0, \sigma_0}$ is defined by:

$$\begin{aligned} K_{\tilde{\mathcal{F}}}_{\theta_0, \sigma_0} &= \{y \in \mathbb{R}^d : (y, \theta_0, \sigma_0) \in K_{\tilde{\mathcal{F}}}\} \\ &= \{y \in \mathbb{R}^d : \tilde{\mathcal{F}}(y, \theta_0, \sigma_0) = 1\} \\ &= \{y \in \mathbb{R}^d : \mathcal{J}(y) = 1\} \\ &= K_{\mathcal{J}} \end{aligned} \quad (3.22)$$

Therefore studying the boundary of $K_{\tilde{\mathcal{F}}_{\theta_0, \sigma_0}}$ means studying the boundary of $K_{\mathcal{F}}$, the critical region of a test defined on \mathbb{R}^d , which is potentially easier to work with.

In short, the question is whether, for some set $A \in \mathcal{F}^{2d+1}$, we can link $(\partial A)_{\theta_0, \sigma_0}$ and $\partial(A_{\theta_0, \sigma_0})$. We first present some useful relations between sets and sections in Lemma 3.2.5, and then Lemma 3.2.6 shows the result we obtained.

Lemma 3.2.5

For two sets A and B , we have:

$$(A \cap B)_{\theta_0, \sigma_0} = A_{\theta_0, \sigma_0} \cap B_{\theta_0, \sigma_0} \quad (3.23)$$

$$(A^c)_{\theta_0, \sigma_0} = (A_{\theta_0, \sigma_0})^c \quad (3.24)$$

$$\overline{A_{\theta_0, \sigma_0}} \subset \overline{A}_{\theta_0, \sigma_0} \quad (3.25)$$

Proof. Proof of Eq. (3.23):

$$\begin{aligned} x \in (A \cap B)_{\theta_0, \sigma_0} &\Leftrightarrow (x, \theta_0, \sigma_0) \in (A \cap B) \\ &\Leftrightarrow (x, \theta_0, \sigma_0) \in A \text{ and } (x, \theta_0, \sigma_0) \in B \\ &\Leftrightarrow x \in A_{\theta_0, \sigma_0} \text{ and } x \in B_{\theta_0, \sigma_0} \\ &\Leftrightarrow x \in A_{\theta_0, \sigma_0} \cap B_{\theta_0, \sigma_0} \end{aligned}$$

Proof of Eq. (3.24):

$$\begin{aligned} x \in (A^c)_{\theta_0, \sigma_0} &\Leftrightarrow (x, \theta_0, \sigma_0) \in A^c \\ &\Leftrightarrow (x, \theta_0, \sigma_0) \notin A \\ &\Leftrightarrow x \notin A_{\theta_0, \sigma_0} \\ &\Leftrightarrow x \in (A_{\theta_0, \sigma_0})^c \end{aligned}$$

Proof of Eq. (3.25):

$$\begin{aligned} x \in \overline{A_{\theta_0, \sigma_0}} &\Leftrightarrow \exists (x_n)_{n \in \mathbb{N}} \in (A_{\theta_0, \sigma_0})^{\mathbb{N}}, x_n \xrightarrow{n} x \\ &\Leftrightarrow \exists (x_n)_{n \in \mathbb{N}} \in (A_{\theta_0, \sigma_0})^{\mathbb{N}}, (x_n, \theta_0, \sigma_0) \xrightarrow{n} (x, \theta_0, \sigma_0) \\ &\Rightarrow \exists (u_n)_{n \in \mathbb{N}} \in A^{\mathbb{N}}, u_n \xrightarrow{n} (x, \theta_0, \sigma_0) \\ &\Rightarrow (x, \theta_0, \sigma_0) \in \overline{A} \\ &\Rightarrow x \in \overline{A}_{\theta_0, \sigma_0} \end{aligned} \quad \square$$

Lemma 3.2.6

For any set $A \in \mathcal{F}^{2d+1}$, we have $\partial(A_{\theta_0, \sigma_0}) \subset (\partial A)_{\theta_0, \sigma_0}$.

Proof. Using the definition of the boundary of a set and the previous lemma, we have:

$$\begin{aligned} (\partial A)_{\theta_0, \sigma_0} &= (\overline{A} \setminus \overset{\circ}{A})_{\theta_0, \sigma_0} = (\overline{A} \cap (\overset{\circ}{A})^c)_{\theta_0, \sigma_0} = \overline{A}_{\theta_0, \sigma_0} \cap ((\overset{\circ}{A})^c)_{\theta_0, \sigma_0} = \overline{A}_{\theta_0, \sigma_0} \cap ((\overset{\circ}{A})_{\theta_0, \sigma_0})^c \\ \partial(A_{\theta_0, \sigma_0}) &= \overline{A_{\theta_0, \sigma_0}} \setminus \overset{\circ}{A_{\theta_0, \sigma_0}} = \overline{A_{\theta_0, \sigma_0}} \cap \overset{\circ}{A_{\theta_0, \sigma_0}}{}^c \end{aligned}$$

For any set B , by definition of the interior of a set, we have $\overset{\circ}{B} = \overline{B^c}$. Thus:

$$\begin{aligned}\overline{A_{\theta_0, \sigma_0}}^c &= \overline{(A_{\theta_0, \sigma_0})^c} = \overline{(A^c)_{\theta_0, \sigma_0}} \\ ((\overset{\circ}{A})_{\theta_0, \sigma_0})^c &= \left(\overline{(A^c)_{\theta_0, \sigma_0}} \right)^c = \overline{A^c}_{\theta_0, \sigma_0}\end{aligned}$$

By applying Eq. (3.25) to A and A^c , we get:

$$\begin{aligned}\overline{A_{\theta_0, \sigma_0}} &\subset \overline{A}_{\theta_0, \sigma_0} \\ \overline{A_{\theta_0, \sigma_0}}^c &\subset ((\overset{\circ}{A})_{\theta_0, \sigma_0})^c\end{aligned}$$

Therefore $\partial(A_{\theta_0, \sigma_0}) \subset (\partial A)_{\theta_0, \sigma_0}$. □

What can we conclude from this result? Applying Lemma 3.2.6 to $K_{\tilde{\mathcal{F}}}$, the critical region of $\tilde{\mathcal{T}}$, we obtain $\partial K_{\tilde{\mathcal{F}}} \subset (\partial K_{\tilde{\mathcal{F}}})_{\theta_0, \sigma_0}$. Therefore, if $\mathbb{P}[Y \in (\partial K_{\tilde{\mathcal{F}}})_{\theta_0, \sigma_0}] = 0$, then we necessarily have $\mathbb{P}[Y \in \partial K_{\tilde{\mathcal{F}}}] = 0$. From this, we know that the class of tests we consider in Theorem 3.2.4 is composed of tests whose border is a $\mathbb{P}Y^{-1}$ -continuity set when θ_0 and σ_0 are known.

3.3 Simulation experiments

After presenting our theoretical results regarding Asymptotic RDT in the previous section, we will now confirm some of them through simulations. As previously stated, being able to ensure a certain level γ is an important aspect to consider when designing a test for a given problem. Ensuring that we verify this property will be an integral part of these simulations.

Since these simulations are the first ones that we present in this thesis, we will also take this opportunity to see how the RDT test performs in a detection problem in order to demonstrate a realistic use case. We will compare it to the typical Neyman-Pearson test to exhibit differences between them, both in terms of performance and assumptions made.

3.3.1 Level

We start off by studying the level of the RDT test $\mathcal{J}_{\lambda_\gamma(\tau)}$. The goal of this first section is twofold. On one hand, we will confirm the results described in Theorem 3.2.2 which, as a reminder, states that the test $\mathcal{J}_{\lambda_\gamma(\tau)}$ has asymptotic level γ as the estimates $\hat{\theta}_n$ and $\hat{\sigma}_n$ improve. On the other hand, we will also see how the RDT test $\mathcal{J}_{\lambda_\gamma(\tau)}$ behaves in the non-asymptotic case. This will allow us to measure the performance degradation caused by poor parameter estimation, and also to see at which point we approach the theoretical results.

We run Monte-Carlo simulations to estimate the effective false-alarm rate (which will be denoted as \mathbb{P}_{FA}) for given sets of parameters. The parameters that we have to choose for these simulations are:

- The dimension of the problem d ,
- The desired level γ ,
- The distribution of the signal Θ ,
- The tolerance τ ,
- The noise standard deviation σ_0 ,
- The model θ_0 ,
- The noise standard deviation estimator $\hat{\sigma}$,
- The model estimator $\hat{\theta}$.

After identifying the parameters to set, the next question is how to set them, as to ensure that the results we get are relevant. Among all these parameters, the one that requires the most attention is most likely

the distribution of Θ . Indeed, there are countless possible probability distributions, and so we have to be careful picking one that is representative of what we want to show.

We have already stated that we want to confirm results regarding the level of the test. This means that we want $\|\Theta\|_2 \leq \sigma_0 \tau$. For the sake of efficiency, it would be ideal to find a distribution for Θ that maximizes the false-alarm rate, which would then represent the worst-case scenario. From the properties of the power function of the RDT test, this is the case whenever we have $\|\Theta\|_2 = \sigma_0 \tau$ almost surely, meaning that Θ is almost surely on the sphere of radius $\sigma_0 \tau$. This already significantly narrows down our choice of distribution, but there are still many left to choose from. We chose a uniform spherical distribution with radius $\sigma_0 \tau$, since it does not favor any particular situation, and also because generating samples from such a distribution can be done easily (see for example [36]). Indeed, we can simply generate a spherical d -dimensional Gaussian vector and normalize it to get a uniformly distributed random vector on the unit sphere. With this choice, we expect the measured false-alarm rate to be equal to γ when the noise variance and the model are perfectly known, since the RDT test has exactly size γ in these circumstances.

Now that we have explained our choice of probability distribution for Θ , we will go through the choices that we made for the other parameters:

- We arbitrarily chose a 2-dimensional problem ($d = 2$).
- The level γ was set to values spreading from 10^{-6} to 1, which should cover a representative range of realistic cases.
- The noise standard deviation σ_0 was set to 1 for the sake of simplicity.
- For the same reason, we set the model θ_0 to 0.
- To choose the tolerance τ , we decided to express it relatively to the noise standard deviation σ_0 using the DNR (Maximum-Distortion-to-Noise ratio), which is simply the ratio between τ and σ_0 . In the following, we will express it in decibels: $\text{DNR}_{\text{dB}} = 20 \log(\tau/\sigma_0)$. We chose a few values of DNR: $-\infty$ dB (no distortion), -5 dB, 0 dB, and 5 dB.
- The noise standard deviation estimate is generated using the typical unbiased standard deviation estimator applied to N independent and identically distributed scalar Gaussian samples with mean 0 and variance σ_0^2 . We then need to choose the parameter N , which will dictate the quality of the noise standard deviation estimate. We have decided to pick values ranging from $N = 10$ samples to $N = 1000$. We will also have cases when we assume σ_0 to be known as reference.
- The model is assumed to be known: $\hat{\theta} = \theta_0$.

With the parameters defined, we can now present the simulation procedure that we used for each set of parameters. The following steps detail one simulation iteration, this process is then repeated enough times to get a good estimate of the detection rate:

1. Generate a signal sample $\Theta \sim \mathbb{P}_\Theta$.
2. Generate a noise sample $X \sim \mathcal{N}(0, \sigma_0^2 I_d)$.
3. From these samples, generate the observation $Y = \Theta + X$.
4. Generate an estimate $\hat{\sigma}$ of σ_0 (if needed, since some simulations assume σ_0 to be known):
 - (a) Generate N independent and identically distributed Gaussian scalar samples $(X_1, \dots, X_N) \in (\mathcal{M}(\Omega, \mathbb{R}))^N$ with mean 0 and variance σ_0^2 .
 - (b) Compute the unbiased standard deviation estimate of the samples (X_1, \dots, X_N) using the usual maximum-likelihood estimator:

$$\hat{\sigma}_N = \sqrt{\frac{1}{N-1} \sum_{i=1}^N \left(X_i - \frac{\sum_{j=1}^N X_j}{N} \right)^2} \quad (3.26)$$

5. Apply the test $\tilde{\mathcal{F}}_{\lambda_\gamma(\tau)}(y, \hat{\sigma})$ and count a detection if $\tilde{\mathcal{F}}_{\lambda_\gamma(\tau)}(y, \hat{\sigma}) = 1$.

The measured false alarm rate is then obtained by counting the number of detections and dividing it by the number of iterations.

Now that we have detailed our simulation procedure, we can finally start presenting some results. We claimed just before that if θ_0 and σ_0 are both known, then the measured false alarm rate should be equal to γ with these parameters, so let us start by verifying this claim. In Fig. 3.1, we plot the measured false alarm rate against the desired level γ for different values of DNR. As we can see on this figure, it seems that we indeed always have $\mathbb{P}_{\text{FA}} = \gamma$ regardless of the signal amplitude. Keep in mind that modifying the DNR means changing both the tolerance τ in the definition of the RDT problem and the amplitude of the generated signal Θ (since we chose Θ such that $\|\Theta\|_2 = \sigma_0 \tau$), so that they are always equal.

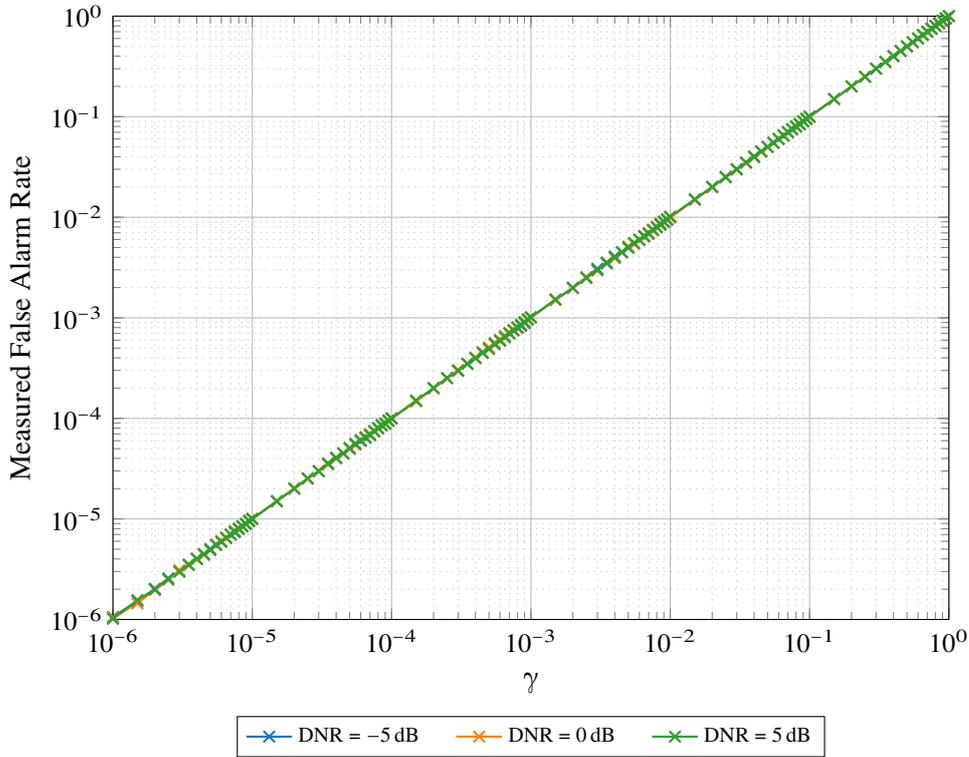


Figure 3.1: Measured false alarm rate against γ with known σ_0 for different values of DNR

After verifying this first claim, we can now introduce the estimation of σ_0 to see its effect on the false-alarm rate. These simulation results are shown in Fig. 3.2 for multiple values of DNR and N which, as a reminder, stands for the number of samples used to estimate σ_0 . On this figure, we can see that replacing σ_0 with an estimate $\hat{\sigma}$ results in an increase of the false-alarm rate, such that the RDT test no longer has level γ . The observed increase depends of course on the quality of the estimation. For example, this increase is barely noticeable for $N = 1000$ which we could have expected since this is a fairly good estimate of σ_0 , whereas with $N = 10$, we only get a mediocre estimate and we can see that it is clearly insufficient, especially for low values of γ . For instance, with $\gamma = 10^{-6}$ and $N = 10$, we can observe a measured false alarm rate greater than 10^{-3} , meaning that it would not be usable in this case since the effective false-alarm rate would be up to three orders of magnitude greater than expected in this case (up to, because remember that the signal distribution that we considered is the worst-case scenario, other distributions for Θ may result in a lower false-alarm rate).

Another aspect worth mentioning is the influence of the tolerance τ . We can see that increasing the tolerance results in a higher false-alarm rate in the worst-case scenario. This increase appears to be accentuated by low values of γ and N . At the time of writing, we do not have an explanation for this phenomenon.

The results displayed in this figure are coherent with Theorem 3.2.2. Indeed, if we consider the evolution of \mathbb{P}_{FA} as N increases with fixed values of γ and DNR, we can see that \mathbb{P}_{FA} decreases and approaches γ . Asymptotically, it seems likely we would indeed have $\mathbb{P}_{\text{FA}} \leq \gamma$, as stated in Theorem 3.2.2.

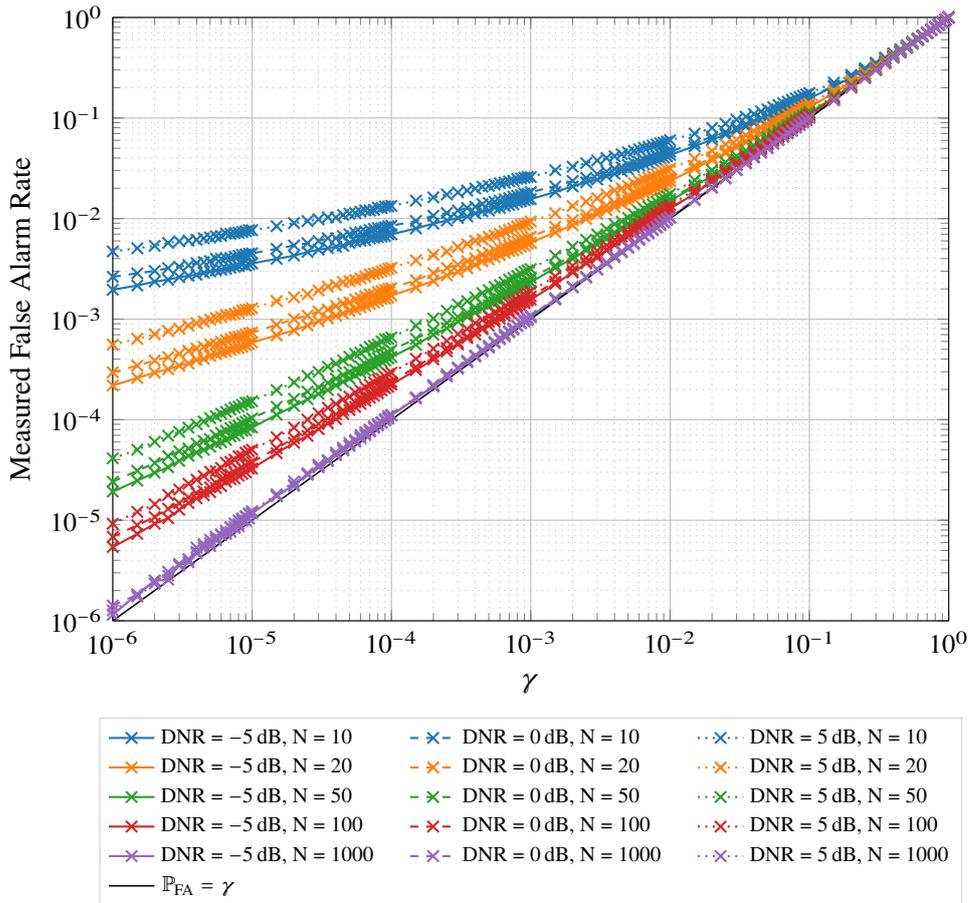


Figure 3.2: Measured false alarm rate against γ with estimated σ_0 for different values of DNR and numbers of samples N used to estimate σ_0

In addition, it appears that we actually have $\mathbb{P}_{\text{FA}} = \gamma$ asymptotically in this case. This would indicate that the RDT test $\tilde{\mathcal{T}}_{\lambda_\gamma(\tau)}$ does not only have asymptotic level γ as we have proven in Theorem 3.2.2, but also asymptotic size γ . We can observe the same behavior on Fig. 3.3, which presents the same simulation results, but this time plotting \mathbb{P}_{FA} against N for several values of γ . We can also clearly see here that lower values of γ require a more accurate estimate of σ_0 in order for the effective false alarm rate to approach γ .

3.3.2 Application to a detection problem

Now that we have somewhat confirmed the results regarding the asymptotic level of the RDT test, we will now look at an application to a detection problem. The idea here is to try to detect a signal that is affected by some unknown distortion. We will see how the estimation of σ_0 affects the performance of this test for such a problem, and we will also compare it to the Neyman-Pearson test. Of course, these two tests are very different, notably in terms of the knowledge of the problem required to use them. However we think this comparison is still useful to some extent in order to highlight these differences, see what assumptions are required to use them and find some of their limitations. The Neyman-Pearson test will be used sub-optimally here, since the test used here does not account for the distortion, which is assumed to be unknown. This is fairly common in practice, the Neyman-Pearson test is often used in circumstances where the models used under both hypotheses do not perfectly match reality. This of course affects the performance of test, but if the models are good enough, we can hope that the performance penalty remains small.

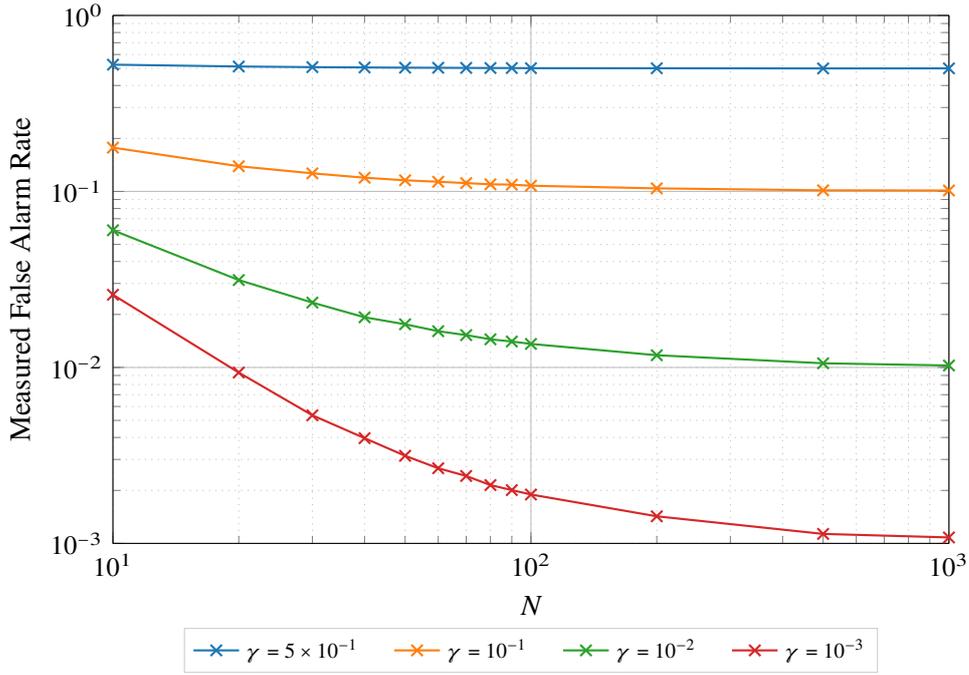


Figure 3.3: Measured false alarm rate against number of samples N used to estimate σ_0 for different values of γ with DNR = 5 dB

Let us start by defining the detection problem of interest. We consider a problem in which we want to test the presence of a fixed signal $\xi \in \mathbb{R}^d$ in presence of white Gaussian noise $X \sim \mathcal{N}(0, \sigma_0^2 I_d)$. This signal ξ is affected by some unknown additive random distortion $\Delta \in \mathcal{M}(\Omega, \mathbb{R}^d)$, independent from X . This distortion is assumed to have a bounded amplitude such that $\|\Delta\|_2 \leq \sigma_0 \tau$, with $\tau > 0$ being known. This detection problem can be summarized as follows:

Detection problem

Observation:

$$Y = \varepsilon \xi + \Delta + X$$

where :

$\varepsilon \in \{0, 1\}$ is unknown

$\xi \in \mathbb{R}^d$ is known

$\Delta \in \mathcal{M}(\Omega, \mathbb{R}^d)$ is unknown, with $\|\Delta\|_2 \leq \tau$ for a known $\tau \geq 0$

$X \sim \mathcal{N}(0, \sigma_0^2 I_d)$

Δ and X are independent

(3.27)

Testing problem:

Given one realization $y = Y(\omega)$, determine whether:

$\mathcal{H}_0: \varepsilon = 0$

$\mathcal{H}_1: \varepsilon = 1$

The presence of an unknown and bounded distortion Δ means that the RDT test $\tilde{\mathcal{J}}_{\lambda, \gamma(\tau)}$ is likely a good candidate for this problem. Indeed, under \mathcal{H}_0 , the generated signal $\Delta + X$ matches the kind of signal we consider in the RDT problem (see Definition 1.2.2) under \mathcal{H}_0 . Its performance will likely depend on the signal ξ and the maximal distortion amplitude τ considered.

Similarly to what we did previously, the tolerance τ and the signal ξ will be determined using the DNR (Maximum-Distortion-to-Noise ratio) and SNR (Signal-to-Noise ratio) respectively. These two values will be expressed in decibels and are defined by:

$$\begin{aligned} \text{DNR}_{\text{dB}} &= 20 \log(\tau / \sigma_0) \\ \text{SNR}_{\text{dB}} &= 20 \log(\|\xi\|_2 / \sigma_0) \end{aligned} \quad (3.28)$$

If we ignored the distortion ($\Delta = 0$), we could simply find and use the appropriate Neyman-Pearson test, since both hypotheses are simple. In this case, the Neyman-Pearson test of level γ is given by:

$$\begin{aligned} \mathcal{T}_{NP}: \mathbb{R}^d &\rightarrow \{0, 1\} \\ y &\mapsto \begin{cases} 1 & \text{if } \frac{y^T \xi}{\|\xi\|_2} > \sigma_0 \Phi^{-1}(\gamma) \\ 0 & \text{otherwise} \end{cases} \end{aligned} \quad (3.29)$$

where Φ stands for the cumulative distribution function of the standard Gaussian distribution. Since this problem also considers an additive distortion, this test \mathcal{T}_{NP} is likely not optimal, and may not even have level γ . In practice however, if the distortion is small enough, it is not uncommon to ignore it and use this test in hopes that the performance penalty is acceptable. We will take this simple approach for the upcoming simulations and see how this influences the results.

We will be assessing the performances of each test using ROC (Receiver Operating Characteristic) curves. These curves show the estimated detection probability against the estimated false alarm probability. Each ROC curve will be obtained by running simulations with the parameter γ spanning the range from 10^{-6} to 1 and every other parameter fixed. In addition to ROC curves we will also verify whether the desired level γ is respected, as this remains an important property to consider for a test to be usable in practice. We will also study how the estimation of σ_0 affects the performance of each test. The Neyman-Pearson test $\tilde{\mathcal{T}}_{NP}^*$ taking the noise standard deviation estimation into account is simply obtained by replacing σ_0 in its definition by the estimate $\hat{\sigma}$:

$$\begin{aligned} \tilde{\mathcal{T}}_{NP}^*: \mathbb{R}^d \times (0, \infty) &\rightarrow \{0, 1\} \\ (y, \hat{\sigma}) &\mapsto \begin{cases} 1 & \text{if } \frac{y^T \xi}{\|\xi\|_2} > \hat{\sigma} \Phi^{-1}(\gamma) \\ 0 & \text{otherwise} \end{cases} \end{aligned} \quad (3.30)$$

The simulations will be presented in the following order:

1. First we will assess both tests with no distortion ($\Delta = 0$) and σ_0 perfectly known.
2. We will then introduce the distortion Δ , while σ_0 remains known.
3. Finally, we will introduce the estimate $\hat{\sigma}_N$ in addition to the distortion Δ .

Simulations with no distortion and known noise standard deviation σ_0

First, we present simulation results in ideal conditions to set a baseline. There is no distortion ($\Delta = 0$) and the noise standard deviation is perfectly known ($\hat{\sigma} = \sigma_0$). We compare the performance of the tests \mathcal{T}_{NP} and $\mathcal{T}_{\lambda_\gamma(0)}$, since there is no distortion ($\tau = 0$). The ROC curves obtained for both tests through simulations are presented in Fig. 3.4. This figure shows that in these conditions, the NP test has better performance than the RDT test, meaning that it offers a better detection rate for the same false-alarm rate. This was expected since the NP test is known to be most powerful under these conditions (Theorem 1.1.4). Regarding the level of each tests, we can see on Fig. 3.5 that both tests have exactly size γ , which is also expected from their respective properties.

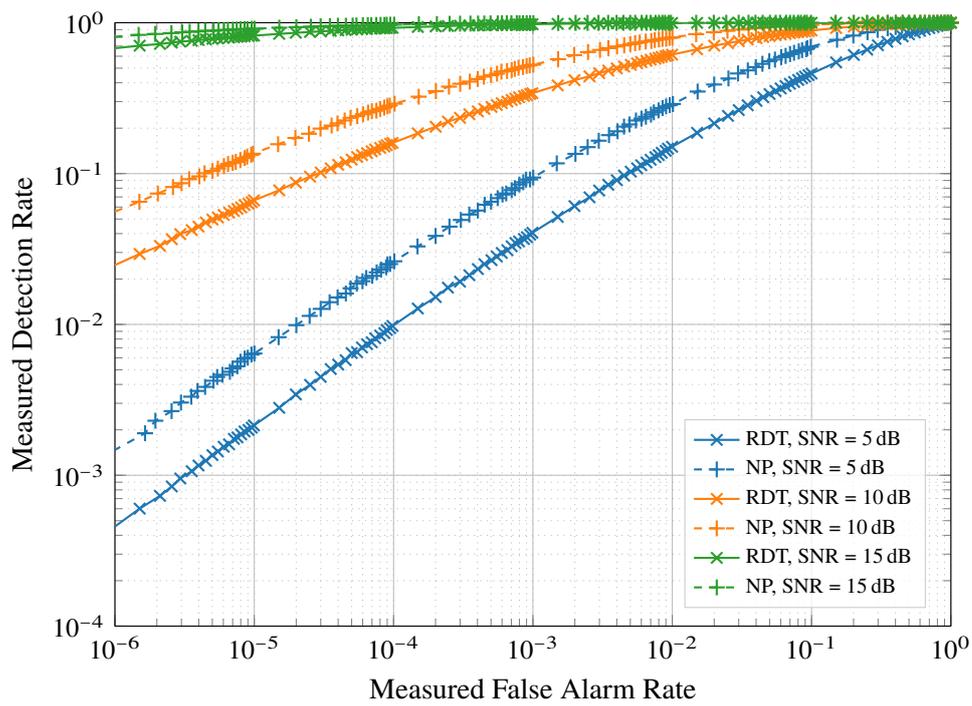


Figure 3.4: ROC curves obtained for the NP and RDT tests with no distortion and σ_0 known

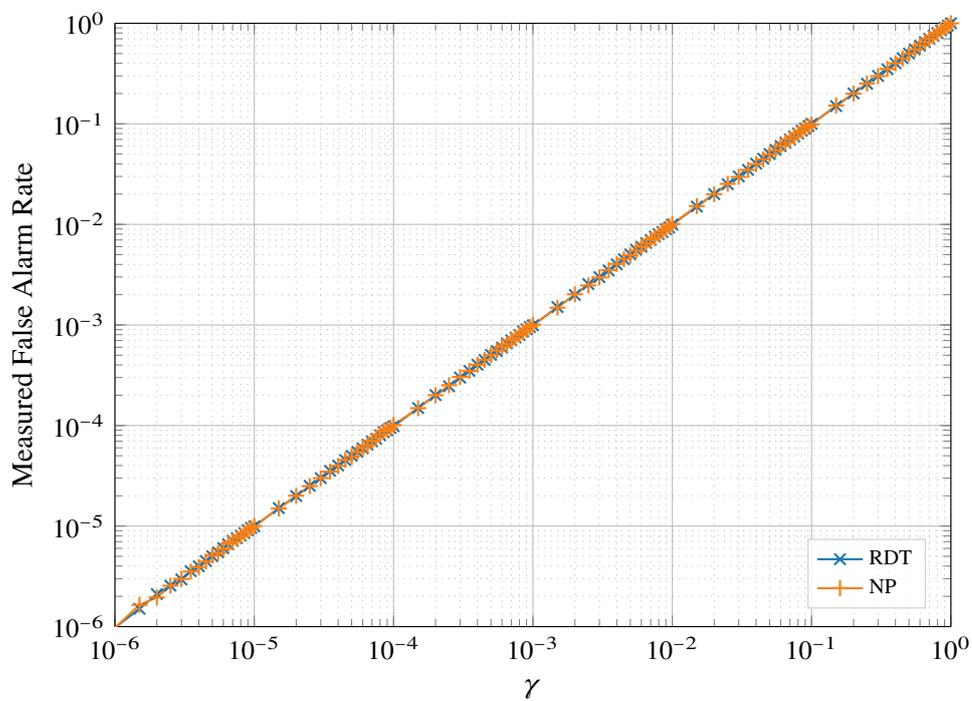


Figure 3.5: Measured false alarm rate against γ for the NP and RDT tests with no distortion and σ_0 known

Simulations with distortion Δ and known noise standard deviation σ_0

We now introduce the distortion Δ and see how it influences each test. As in the previous section, we use a uniformly distributed spherical distortion with radius $\sigma_0\tau$. Remember that this distortion Δ is considered to be unknown; only its maximum amplitude $\sigma_0\tau$ is assumed to be known. We will therefore now compare the tests \mathcal{T}_{NP} and $\mathcal{T}_{\lambda_\gamma(\tau)}$. We conserve the same Neyman-Pearson test, since it does not offer any way to take the distortion Δ into account. As such, we expect to observe some degradation in its performance, since it no longer matches the problem that we are considering.

Figure 3.6 shows the ROC curves obtained with a fixed SNR of 15 dB to see how increasing the DNR affects the performances of both tests. Judging from these curves, we can see that the Neyman-Pearson test seems to continue to outperform the RDT test as the distortion amplitude increases. Indeed, for the same measured false-alarm rate, the Neyman-Pearson test offers a higher detection rate.

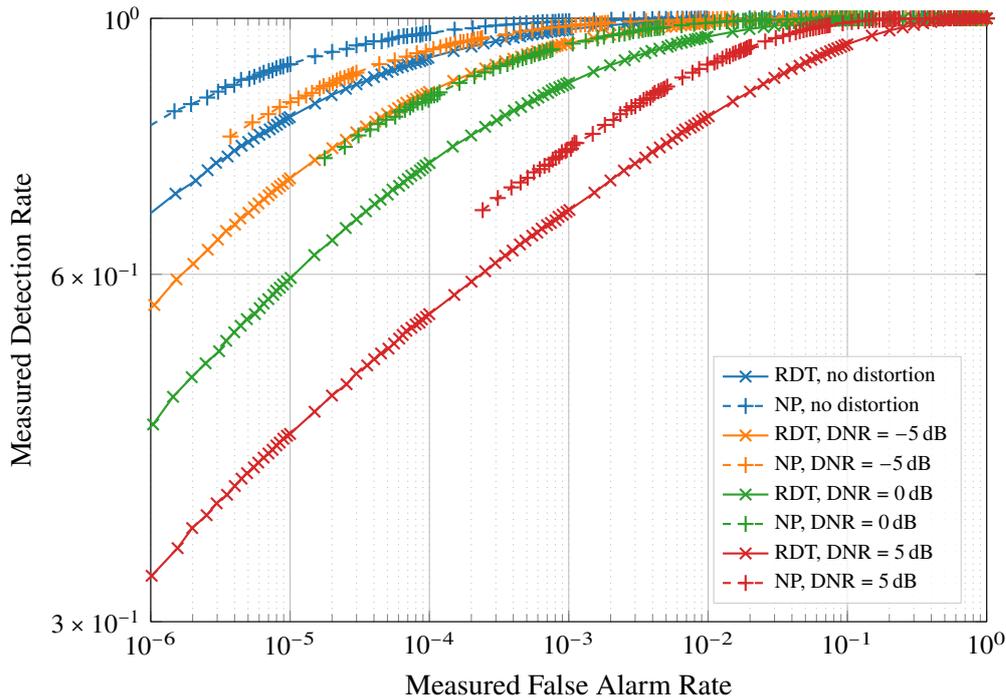


Figure 3.6: ROC curves obtained for the NP and RDT tests with distortion and σ_0 known for SNR = 15 dB

However, we can notice that as the DNR increases, the minimum false-alarm rate of the Neyman-Pearson test increases as well. For example, with a DNR of 5 dB, its false-alarm rate only goes as low as 2×10^{-4} , when using $\gamma = 10^{-6}$. What we notice here is that the Neyman-Pearson test no longer has level γ for this problem. We can confirm this by looking at Fig. 3.7, which shows the measured false-alarm rate against γ for both tests. We can see there that as soon as we introduce a distortion, the NP test no longer has level γ , and that the effective false-alarm rate appears to increase with the distortion amplitude. This makes it troublesome to use it in real scenarios where controlling the false-alarm rate is crucial to maintain proper operation. In contrast, the RDT test maintains its level γ regardless of the distortion amplitude, as we have already seen in Fig. 3.1. This was expected, since the RDT test is designed to take this kind of distortion into account.

Looking closely at this figure, one can notice that our claim about the NP test no longer having level γ is not quite correct: for values of γ between approximately 0.5 and 1 (shown in the magnified area on Fig. 3.7) we can see that the effective false-alarm rate appears to decrease slightly as the distortion amplitude increases, such that it is actually inferior to γ in that range. We currently do not have any explanation for this behavior. However, this range of values for γ tends to not be commonly used in practice, as they correspond to very high false-alarm rates. Therefore we have chosen to ignore this

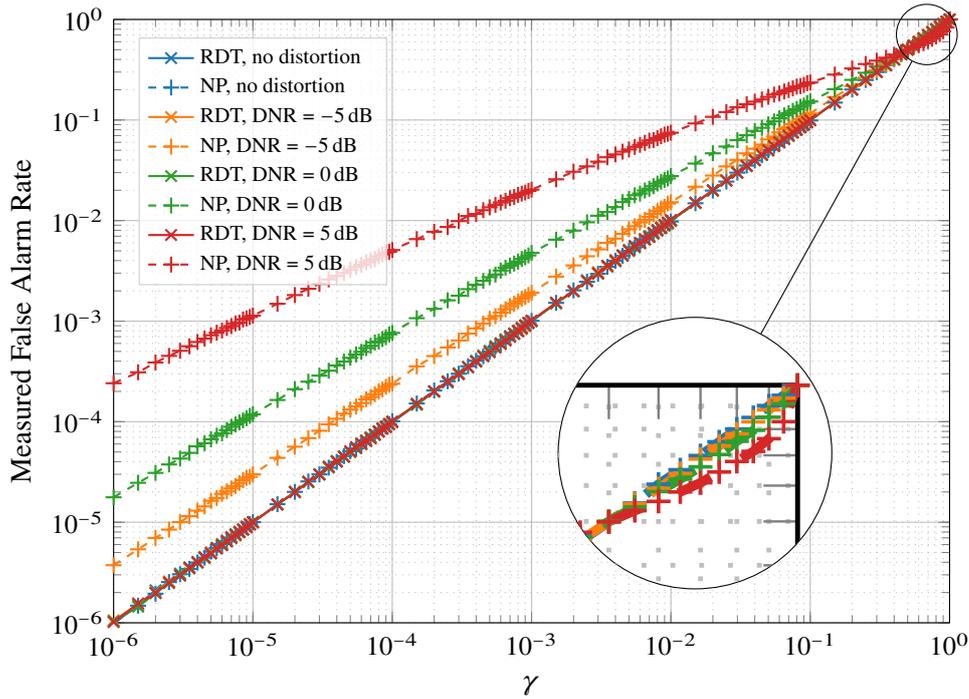


Figure 3.7: Measured false alarm rate against γ for the NP and RDT tests with distortion and σ_0 known

behavior for the time being.

We could consider modifying the Neyman-Pearson test to try to recover a test with level γ . One way to do so would be to adjust the parameter γ , and replace it with a value γ' such that the effective false-alarm rate is equal to γ . We can always find such a value γ' for any given γ , since the effective false-alarm rate decreases to 0 as γ tends to 0. For example, if the desired level is 10^{-3} and we have a DNR of 5 dB, then we can see from Fig. 3.7 that this level is effectively attained when choosing $\gamma \approx 10^{-5}$. However, we have to keep in mind that these simulations were performed with one specific distortion distribution. From the problem statement, we want to find a test that has level γ for any bounded distortion Δ . For the RDT test, we know that any distortion Δ such that we have $\|\Delta\|_2 = \sigma_0 \tau$ almost surely maximizes the false-alarm rate. However, we do not know whether this is the case for the NP test. Therefore, we would need further study to determine which distortion Δ would maximize the false-alarm rate, so that we can find the value of γ' that would yield a test with level γ for any distortion.

Simulations with distortion Δ and estimated noise standard deviation σ_0

After seeing the impact of the distortion on the performance of each test, we will now study the effect of the estimation of the noise standard deviation $\hat{\sigma}$, in addition to the distortion Δ . Figure 3.8 presents the ROC curves obtained for each test with SNR = 15 dB and DNR = 5 dB. We can see that, like previously, the NP test appears to outperform the RDT test, since it offers a higher detection rate at equal false-alarm rate. Comparing these curves to the ones where σ_0 is known shows that the main effect of the estimation of σ_0 is an increased false-alarm rate for both tests. This is particularly noticeable for low values of N . For example, with $N = 10$, the minimum measured false-alarm rate is approximately 5×10^{-4} for both NP and RDT tests, reached with $\gamma = 10^{-6}$.

Figures 3.9 and 3.10 show the measured false-alarm rate against γ for the NP and RDT tests respectively for different distortion amplitudes. We can clearly see on both of these figures the degradation on the false-rate rate caused by the estimation of σ_0 . In these circumstances, neither test has level γ . This was already the case for the NP test with only the distortion Δ , and the estimation the σ_0 accentuates it. A noteworthy difference in the behavior of these tests is that as N increases, the false-alarm rate of the RDT

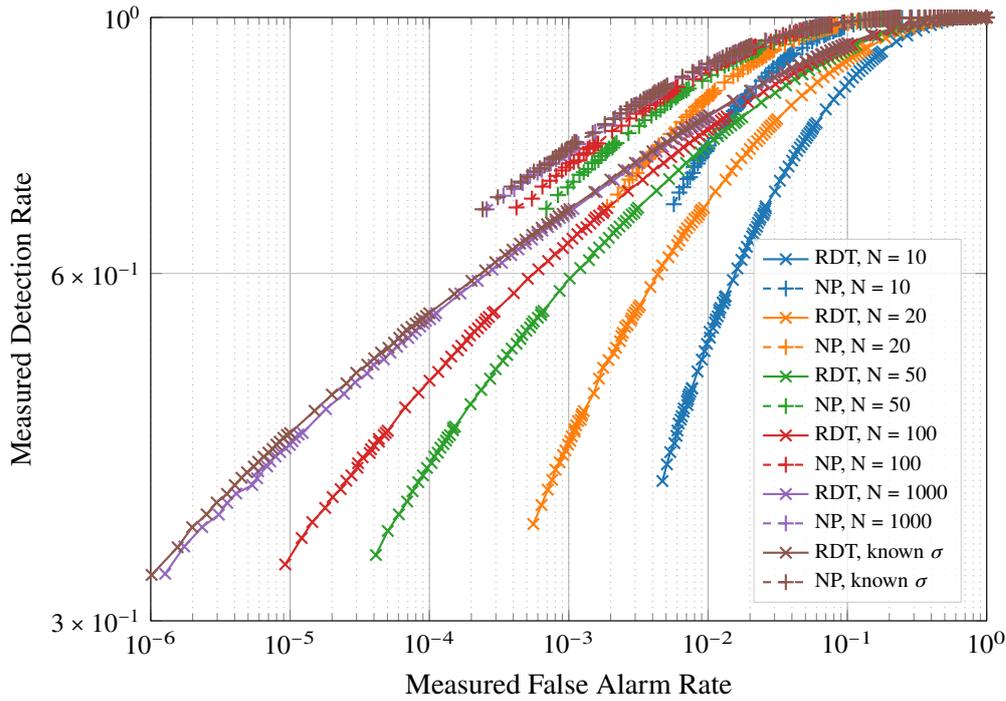


Figure 3.8: ROC curves for the NP and RDT tests with estimation of σ_0 , SNR = 15 dB, DNR = 5 dB

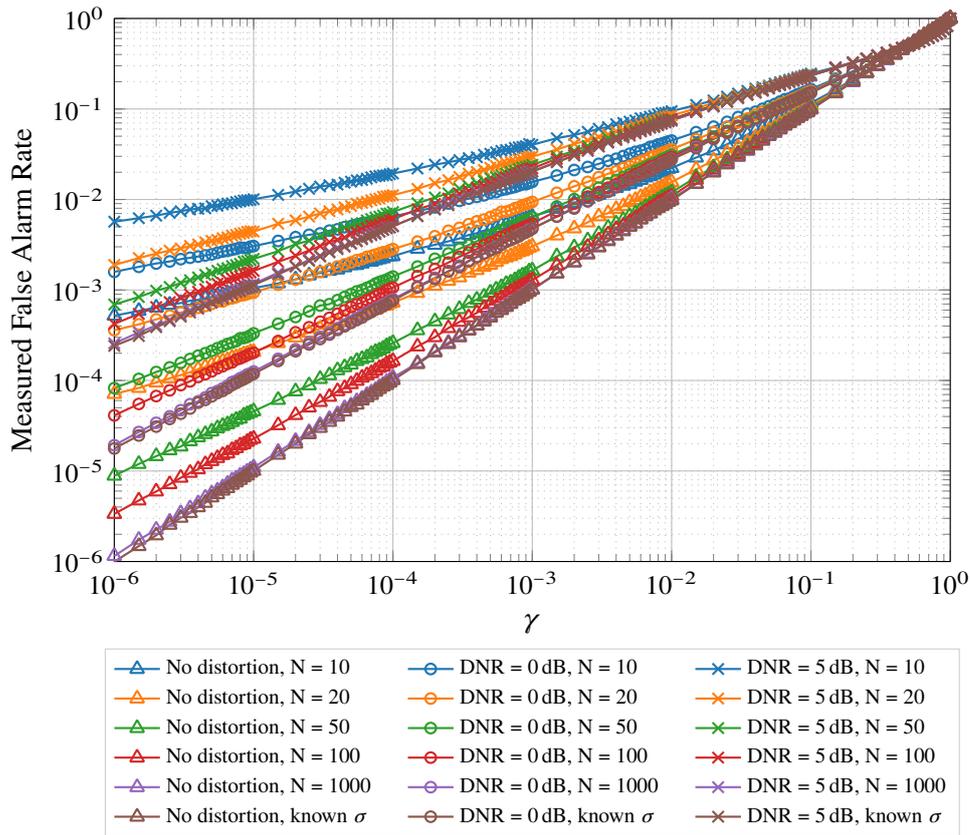
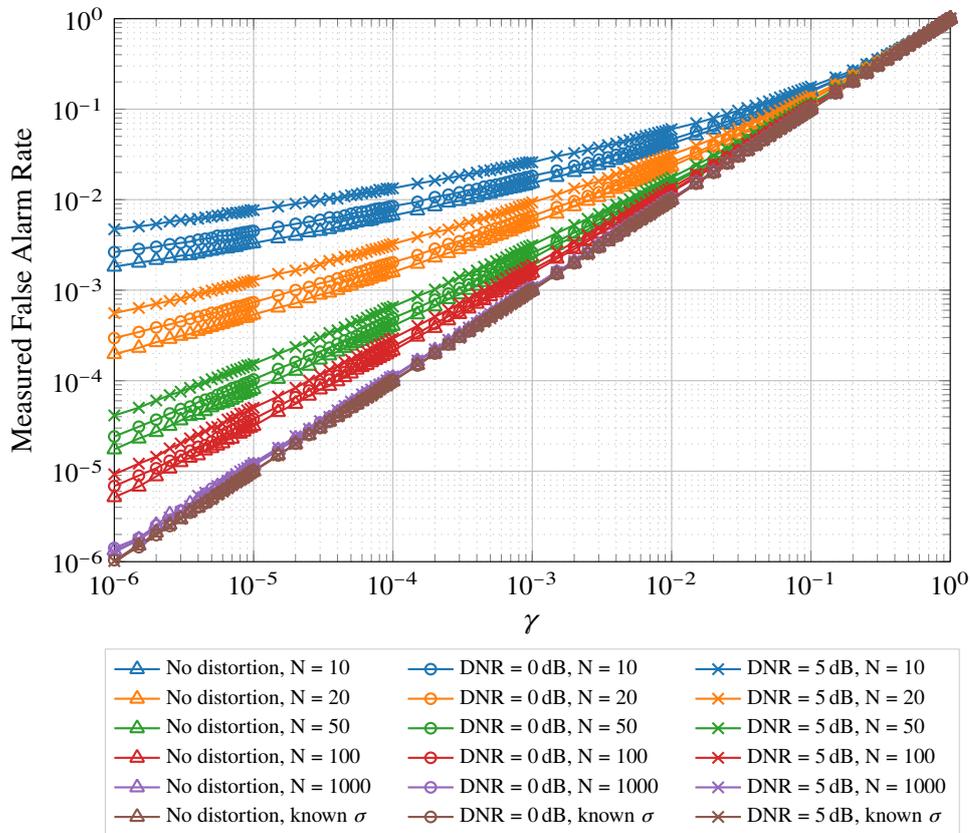
test approaches γ , but this is not the case for the NP test. For the latter, its false-alarm rate can only reach the false-alarm rate that we measured in presence of the distortion Δ , which is greater than γ . Therefore we do not have an asymptotic level guarantee for the NP test in presence of a distortion.

We can also use these figures to pinpoint how large N needs to be so that we can assume that we have reached the asymptotic regime, i.e. the performance loss caused by the estimation of σ_0 is negligible. Since the NP test cannot have level γ in presence of a distortion Δ , we will only consider the RDT test here. We can see on both Figs. 3.8 and 3.10 that the value of N required to approach γ seems to depend on γ . Indeed, it appears that we need a better estimate of σ_0 as γ decreases. For example, for $\gamma = 10^{-3}$, $N = 100$ can be considered sufficiently close to the case with known σ_0 . However, if we are targeting $\gamma = 10^{-6}$, using $N = 100$ can increase the false-alarm rate tenfold with a DNR of 5 dB, and therefore we would need to increase N further. Looking at the ROC curves of Fig. 3.8, it appears that N has little effect on the detection rate of the RDT test, and seems to mainly affect the false-alarm rate.

For information, we give in Table 3.1 the area under the ROC curve for the RDT test.

SNR	DNR	N (samples)					σ known
		10	20	50	100	1000	
5 dB	-5 dB	72.28 %	73.48 %	74.17 %	74.41 %	74.61 %	74.64 %
5 dB	0 dB	68.02 %	69.07 %	69.68 %	69.88 %	70.05 %	70.09 %
5 dB	5 dB	60.30 %	61.20 %	61.79 %	61.99 %	62.16 %	62.20 %
10 dB	-5 dB	91.19 %	92.80 %	93.67 %	93.94 %	94.18 %	94.20 %
10 dB	0 dB	87.21 %	88.91 %	89.87 %	90.18 %	90.44 %	90.48 %
10 dB	5 dB	76.62 %	78.27 %	79.20 %	79.52 %	79.78 %	79.81 %
15 dB	-5 dB	99.31 %	99.78 %	99.91 %	99.93 %	99.95 %	99.95 %
15 dB	0 dB	98.69 %	99.44 %	99.69 %	99.75 %	99.80 %	99.80 %
15 dB	5 dB	95.53 %	97.11 %	97.82 %	98.03 %	98.20 %	98.22 %

Table 3.1: Area under the ROC curve for the RDT test

Figure 3.9: Measured false alarm rate against γ for the Neyman-Pearson test with estimation of σ_0 Figure 3.10: Measured false alarm rate against γ for the RDT test with estimation of σ_0

Conclusion of these simulations

We have seen throughout these simulations that the NP and RDT tests behave fairly differently in presence of a distortion and when replacing σ_0 with an estimate $\hat{\sigma}$. The ROC curves show that the NP test always offers a higher detection rate for identical false-alarm rates, regardless of the distortion or the noise standard deviation estimation. However, as soon as a distortion is introduced, this test is no longer able to respect a given level γ , unlike the RDT test which is robust to a distortion, and also has level γ asymptotically.

3.3.3 Recovering a test with level γ when estimating σ_0

In the previous simulations we strongly insisted on the level of the tests, as we believe that this property is important to preserve. A nice feature of the RDT test is the use of a tolerance that can be used to compensate a lack of knowledge, in our case here the fact that σ_0 is not perfectly known. In this section, we describe a possible way to adjust the RDT test in order to recover a test with level γ using this tolerance τ . We decided to focus on the RDT test here since we only have to account for the effect of the noise variance estimation, whereas working with the NP test would require compensating both the distortion and the estimation of σ_0 .

The idea here is to adjust the tolerance τ so that it encompasses the effects of both the distortion and the estimation of the noise variance. This means that for a given tolerance τ and a given estimator $\hat{\sigma}$, we want to find an adjusted tolerance τ^* such that the test $\tilde{\mathcal{T}}_{\lambda_\gamma(\tau^*)}$ has size γ . Of course, this assumes that such a tolerance does exist. And then, if it exists, we have to figure out a way to find it.

First, let us start by addressing the existence of this tolerance τ^* . Judging from the previous simulations, the estimation of σ_0 increases the false-alarm rate beyond γ . From the definition of the test $\tilde{\mathcal{T}}_{\lambda_\gamma(\tau)}$ (see Eq. (3.7)), this means that the test statistics $\|y\|_2/\hat{\sigma}$ exceeds the threshold $\lambda_\gamma(\tau)$ too frequently. In order to lower the false-alarm rate back down to γ , we then need to increase this threshold $\lambda_\gamma(\tau)$. From Lemma 1.2.9, we know that the threshold $\lambda_\gamma(\tau)$ increases continuously with τ to $+\infty$. Therefore, by replacing τ in the test $\tilde{\mathcal{T}}_{\lambda_\gamma(\tau)}$ with a greater value, we should be able to find a value τ^* such that the test $\tilde{\mathcal{T}}_{\lambda_\gamma(\tau^*)}$ has size γ . Keep in mind that the tolerance τ used to define the RDT problem does not change: we still consider a distortion with maximum amplitude τ . We only modify the RDT test, keeping the same test structure and only replacing the tolerance used with a different value τ^* .

Now that we know that this tolerance τ^* exists, we need to be able to compute it. We present here a way to estimate it through simulations. We will proceed using the same simulations that we used in Section 3.3.1 to estimate the false-alarm rate of the RDT test. The idea is simple: we want to find τ^* such that $\mathbb{P}[\tilde{\mathcal{T}}_{\lambda_\gamma(\tau^*)}(\Theta + X, \hat{\sigma}) = 1] = \gamma$. Therefore, we estimate $\mathbb{P}[\tilde{\mathcal{T}}_{\lambda_\gamma(\tau^*)}(\Theta + X, \hat{\sigma}) = 1]$ for different values of τ^* , until the measured probability is close enough to γ .

Since these simulations can take some time (especially when γ is small), we have to choose the values τ^* somewhat carefully to complete this process in a reasonable time. We first determine a closed interval in which τ^* is contained, then we perform a binary search within that interval to get a good estimate of τ^* . At the beginning, we only know that τ^* is contained in the interval $[\tau, +\infty)$, which is not useable for a binary search. Therefore the first step consists in finding an upper bound for τ^* . To do so, we have chosen to estimate the false-alarm rate using the tolerance $k\tau$ for $k = 2, 3, 4, \dots$, until we find k such that the measured false-alarm rate is lower than γ . At this point, we know that τ^* is contained in the interval $[(k-1)\tau, k\tau]$, on which we then perform a binary search to refine our estimate of τ^* .

Figure 3.11 displays the values of τ^* we obtained through this process. From this figure, we can see that τ^* appears to be affected by the level γ , the estimator $\hat{\sigma}$ and the DNR. An expected takeaway from these results is that coarser estimates of σ_0 require increasing the tolerance more to compensate its effect. Another interesting aspect to note is that lower values of γ also require a larger increase to compensate the effect of the estimation of σ_0 . We also need to verify that the false-alarm rate we get using this tolerance τ^* is effectively γ . Figure 3.12 shows that this is indeed the case for every tested set of parameters.

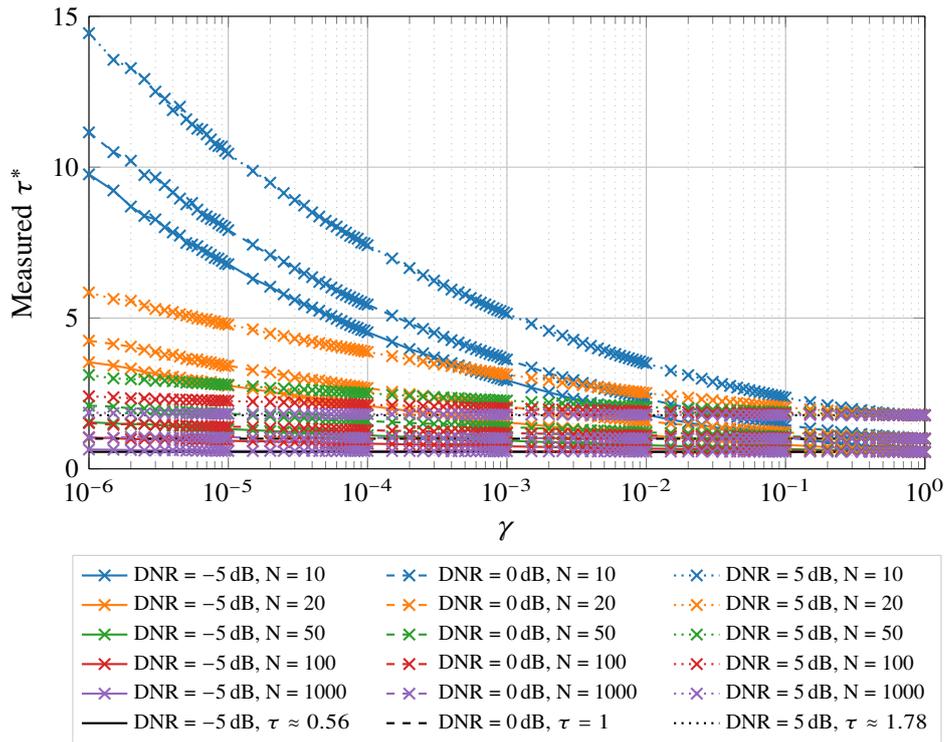


Figure 3.11: Measured values of τ^* against γ for different values of N and DNR

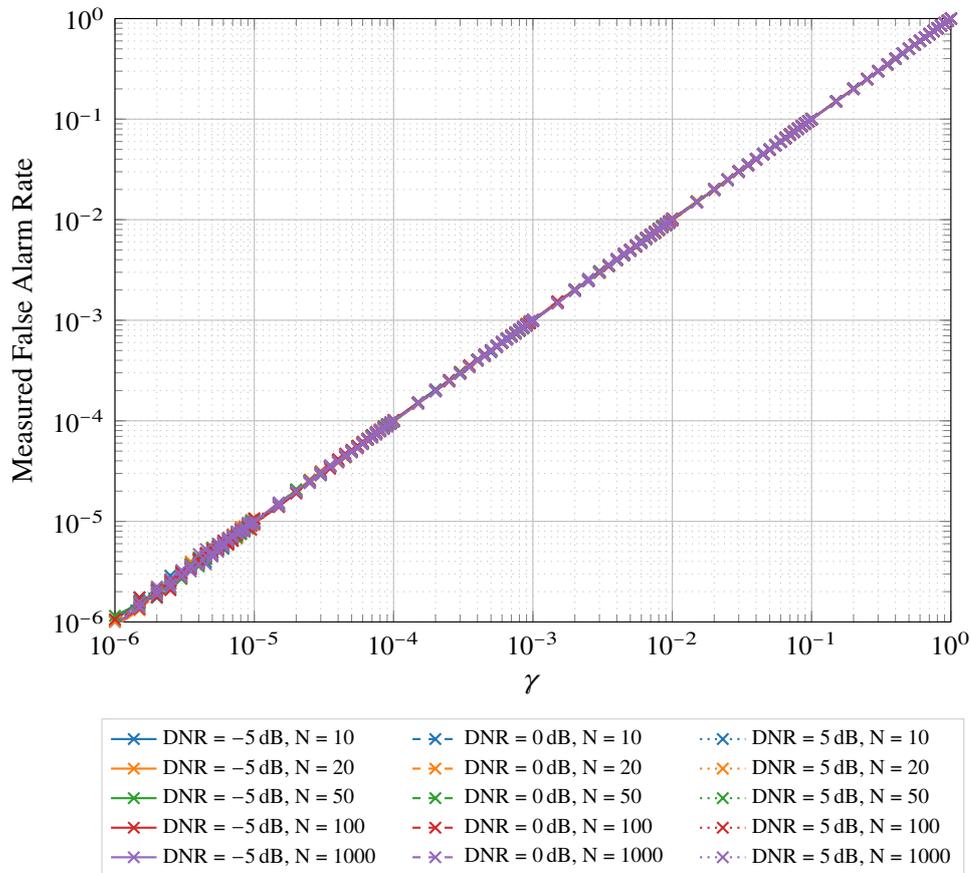


Figure 3.12: Measured false-alarm rate against γ using the adjusted tolerance τ^*

For our approach to be valid, we also need to verify whether τ^* depends on the noise standard deviation σ_0 . Indeed, recall that in practice we only have access to an estimate of σ_0 . If there is a dependency between σ_0 and τ^* , then simply replacing σ_0 with an estimate to find the corresponding value of τ^* may not necessarily result in a test with level γ . To verify whether such a dependency exists, we ran additional simulations using $\sigma_0 = 2$ and $\sigma_0 = 100$, whereas all previous simulations used $\sigma_0 = 1$. Figures 3.13 and 3.14 show the obtained values of τ^* obtained for these two values of σ_0 respectively. Comparing these results to the ones shown in Fig. 3.11 for $\sigma_0 = 1$, we can observe that we always get the same values τ^* regardless of the value of σ_0 .

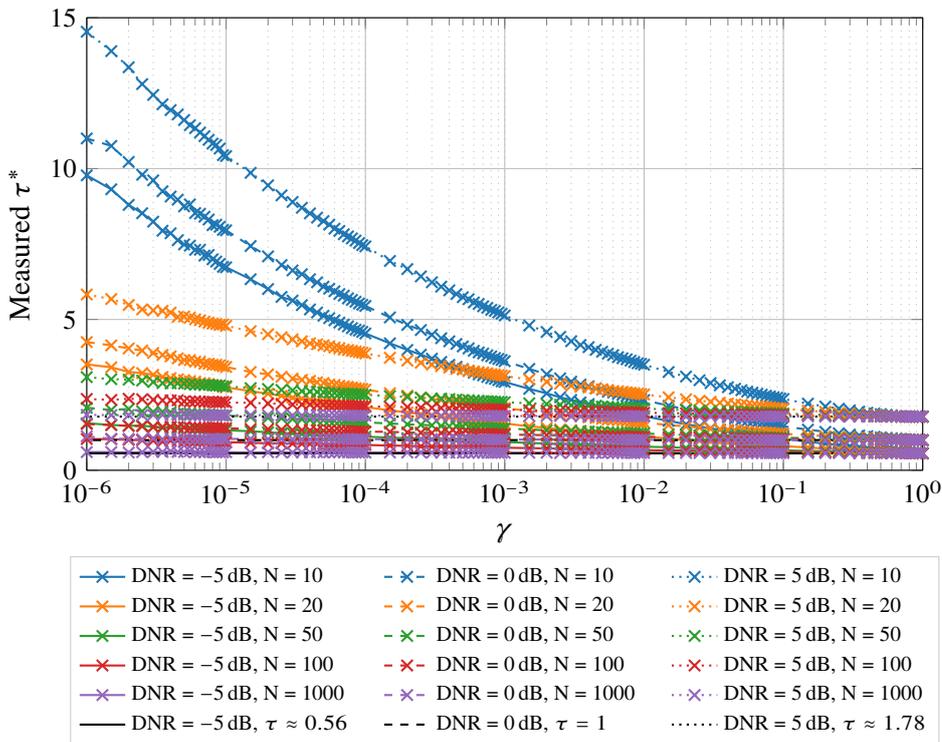


Figure 3.13: Measured values of τ^* against γ for different values of N and DNR with $\sigma_0 = 2$

Now that we know that the test $\tilde{\mathcal{T}}_{\lambda_\gamma(\tau^*)}$ has level γ , we will now compare it to the test $\tilde{\mathcal{T}}_{\lambda_\gamma(\tau)}$ in the detection problem that we studied in the previous section. Figure 3.15 shows the ROC curves obtained for each test, using a SNR of 15 dB and a DNR of 5 dB. We can see there that the ROC curves obtained using τ and τ^* are identical: both tests offer the same detection rate for the same false-alarm rate. The only difference is that each point of the ROC curve is obtained using a different value for γ depending on whether we are using τ or τ^* .

To understand this behavior, we first have to recall the definition of the tests we consider:

$$\begin{aligned} \tilde{\mathcal{T}}_{\lambda_\gamma(\tau)}: \mathbb{R}^d \times (0, +\infty) &\rightarrow \{0, 1\} \\ (y, \hat{\sigma}) &\mapsto \begin{cases} 1 & \text{if } \frac{\|y\|_2}{\hat{\sigma}} > \lambda_\gamma(\tau) \\ 0 & \text{otherwise} \end{cases} \end{aligned} \quad (3.31)$$

These are thresholding tests using the threshold $\lambda_\gamma(\tau)$, which is computed using the two parameters γ and τ . The ROC curve is obtained by measuring the false-alarm rate and detection rate for a fixed tolerance τ and for γ spanning the range $(0, 1)$. In other terms, a ROC curve characterizes a family of tests of the form $\{\tilde{\mathcal{T}}_{\lambda_\gamma(\tau)}, \gamma \in (0, 1)\}$, for a given tolerance τ .

What we need to study now is how the choice of τ affects this family of test, and therefore the ROC curve. First, let us consider a single test $\tilde{\mathcal{T}}_{\lambda_{\gamma_1}(\tau_1)}$, which is characterized by the pair of values (γ_1, τ_1) . Now

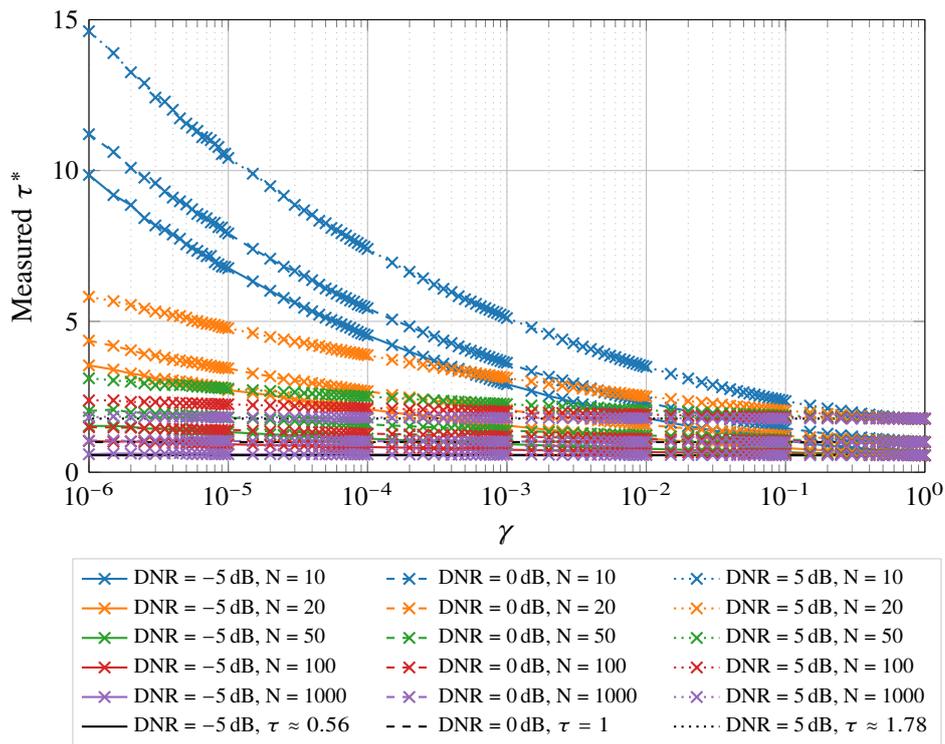


Figure 3.14: Measured values of τ^* against γ for different values of N and DNR with $\sigma_0 = 100$

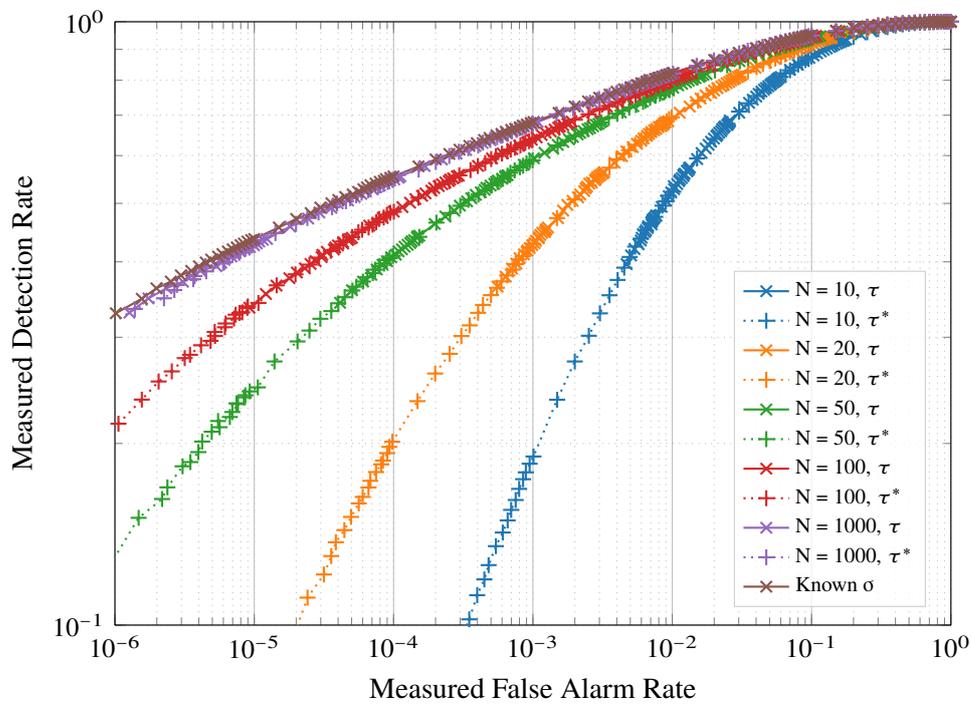


Figure 3.15: Comparison of ROC curves obtained using the tolerances τ and τ^*

consider a second tolerance τ_2 . From Lemma 1.2.9, the function $\gamma \mapsto \lambda_\gamma(\tau_2)$ is continuous and strictly decreasing from $+\infty$ to 0. Therefore there exists a unique value $\gamma_2 \in (0, 1)$ such that $\lambda_{\gamma_2}(\tau_2) = \lambda_{\gamma_1}(\tau_1)$. This means that the tests $\tilde{\mathcal{T}}_{\lambda_{\gamma_1}(\tau_1)}$ and $\tilde{\mathcal{T}}_{\lambda_{\gamma_2}(\tau_2)}$ are actually the same test, since they use the same threshold.

Now consider the families of tests $\mathcal{T}_1 = \{\tilde{\mathcal{T}}_{\lambda_\gamma(\tau_1)}, \gamma \in (0, 1)\}$ and $\mathcal{T}_2 = \{\tilde{\mathcal{T}}_{\lambda_\gamma(\tau_2)}, \gamma \in (0, 1)\}$. With the same reasoning, for each $\gamma \in (0, 1)$, we can find $\gamma' \in (0, 1)$ such that $\lambda_\gamma(\tau_1) = \lambda_{\gamma'}(\tau_2)$. Therefore the families \mathcal{T}_1 and \mathcal{T}_2 are actually composed of the same tests, only indexed differently. Since these families are identical, the resulting ROC curves for the tolerances τ_1 and τ_2 are also identical, which means that the choice of tolerance does not affect the overall performance of the test.

This is an interesting feature of the RDT test: the tolerance makes it possible to compensate both the distortion and the estimation of σ_0 without affecting the overall performance of the test, in order to maintain the guarantee on the level.

3.4 Conclusion

In this chapter we developed an extension of the RDT approach which improves upon the base RDT theory by taking parameter estimation into account, namely the model and the noise variance. We showed that the RDT test $\tilde{\mathcal{T}}_{\lambda_\gamma(\tau)}$ is asymptotically optimal, with an appropriate optimality criterion derived from the γ -MCCP criterion, and we then confirmed the results regarding its asymptotic level through simulations. We also performed some comparisons with the Neyman-Pearson test, allowing us to see notably how they compare and how this relates to their requirements to use them. We have also suggested a way to compensate the effect of the noise variance estimation on the level of the RDT test by adjusting the tolerance τ , which can be used to compensate a lack of knowledge of the problem without affecting the performance of the test.

As we mentioned throughout this chapter, there are further aspects of this asymptotic test that could use further study. First, we restricted our study to the case of a covariance matrix of the form $\sigma_0^2 I_d$ where σ_0 is estimated. For a more general approach, we would need to consider the estimation of any covariance matrix C , and verify if we still have an asymptotic optimality in this case. Regarding our simulations, we only focused on the estimation of σ_0 and the model θ_0 remained known all throughout. We would need to conduct additional simulations to see the effect of the estimation of θ_0 on the performance of the test. Finally, another aspect that would deserve further work is the adjustment of the tolerance τ to compensate the effect of the parameter estimation. We have shown that we can find a tolerance τ^* to recover a test with level γ without affecting its performance, and described a procedure to estimate this adjusted tolerance τ^* through simulations. However, this procedure is fairly slow, as it requires many iterations to get an accurate estimate of τ^* , especially for low values of γ . With more study of the behavior of τ^* , we may be able to devise a more efficient method to estimate this parameter.

In the next chapter, we will be using the aspects we studied in this chapter to build detection methods using the RDT test, which we will then test on real signals from the SWaT dataset.

4 Learning behaviors and detecting anomalies on real signals

Contents

4.1	Detecting discontinuities on continuous signals	92
4.2	Segmenting and learning phases, and detecting anomalies	96
4.2.1	Change-detection in time series	97
4.2.2	An RDT-based change-in-mean detection method	98
4.2.3	Application to real signals	101
4.2.4	Perspectives: towards learning a model and detecting anomalies	105
4.3	Conclusion	108

Until now, we have been working on theoretical aspects of hypothesis testing and the RDT approach, and also performed a few simulations in artificial scenarios. In this chapter, we will now make use of what we studied so far in order to build methods capable of working on real signals. As discussed in the introduction, the end goal of this work is to propose learning-based methods capable of detecting cyberattacks in industrial systems using the sensors and actuators present to first learn their normal behavior, then detect anomalies. Unlike the previous chapters which were mostly theoretical, this chapter is focused mainly on describing methods and experimenting on real signals. We will mainly use the SWaT dataset that we presented in the introduction as a support for our experiments.

Using these signals, we will present two methods developed using the tools described before:

- A first method that can detect discontinuities on signals that are supposed to be continuous. This method is fairly simple, and uses normal operation data to configure the parameters of interest. Since there are quite a few attacks that cause discontinuities in our dataset, we found this approach to be worthwhile to develop and present here.
- A second, more complete method that consists in segmenting signals in order to discover their different phases, then learn the characteristics of each phase. Based on this model of a signal, the goal is to detect whenever the behavior of the observed signal deviates from this model. This method is based on a change detection method in order to discover the different phases of the signal.

In the following we will describe these two methods and test them out on signals from the SWaT dataset in order to illustrate how they work. Note that the methods presented here are mainly intended as proofs of concepts. They should not be considered ready for use in real-world scenarios as there are still many aspects that require further development. However, the results presented below should give an idea of what is realistically achievable.

Before presenting the methods we developed, we should briefly mention the issue of performance evaluation. In the following, we will make qualitative comments on the performance of the described methods, but we will not present any quantitative performance metrics, such as precision, recall or the F-score which are commonly found metrics for performance evaluation, especially on real data. We found

that the SWaT dataset does not offer the required data to compute these metrics for our methods. The only labels present in this dataset indicate whether an attack is occurring at any given time, but our methods look at other aspects of the studied signals — such as transitions from one phase to another — which are not labeled here. Therefore we were not able to assess the performance of our methods quantitatively, and we had to resort to a limited qualitative evaluation.

4.1 Detecting discontinuities on continuous signals

In this first section we look at the issue of detecting discontinuities in signals that are supposed to be continuous. We observed that on some signals, numerous attacks that spoofed a sensor caused significant discontinuities at the beginning and at the end of the attack, i.e. a very large change between two consecutive samples. Two examples of such attacks are shown on Fig. 4.1, taken from the water level sensor LIT101 in the SWaT dataset. The two attacks around 361,000 s and 390,000 s create two discontinuities each, one at their start when the attacker starts spoofing the sensor, and one at the end when we can once again observe the true value of the sensor.

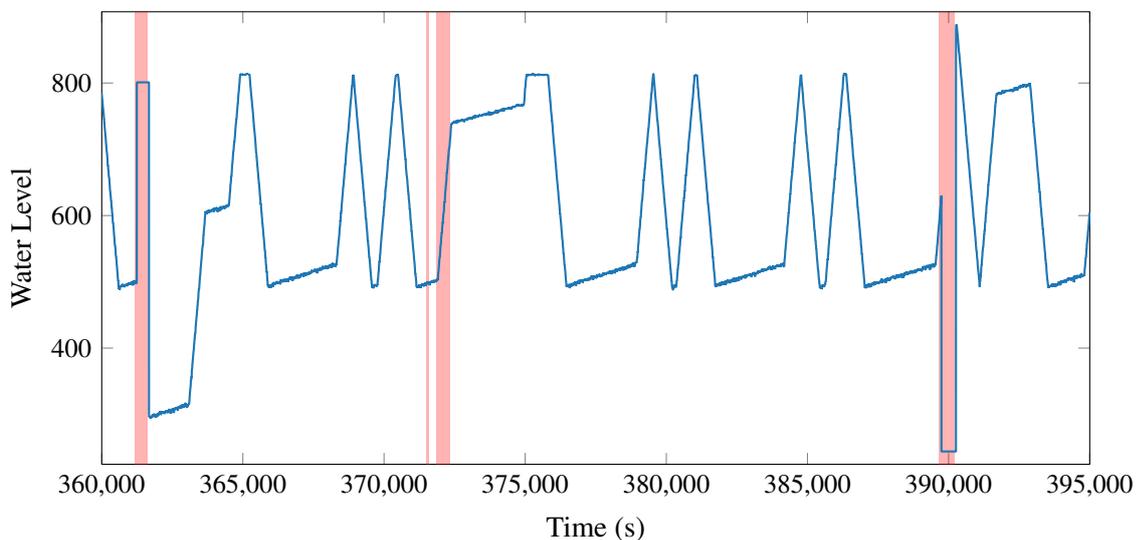


Figure 4.1: Examples of discontinuities observed on the water level signal LIT101

In this case, the observed signal is supposed to be continuous. Indeed, the measured water level varies slowly over time, and the sampling rate of 1 Hz is high enough compared to these variations. Discontinuities on such a signal are therefore anomalies that we can try to detect. Of course, this is not applicable to every type of signal. For example, some electrical signals may require a much higher sampling rate to be considered continuous, since electrical phenomena can be much faster than the mechanical ones that we measure here. As such, the detection method we propose here can only target specific signals that have slow variations compared to their sample rate. This method is an adaptation of a change-point detection method presented in [37], where it was used to segment respiratory signals.

This method uses a wavelet transform as a base to transform the signal into a representation where the discontinuities are easily detectable. In short, the wavelet transform is a tool that allows studying different scales of a signal by separating its low-frequency and high-frequency components. It operates by filtering the signal using two different filters, and downsampling each filter output by only keeping every other sample. Since discontinuities are very fast transitions, we should be able to detect them by analyzing the high frequencies of the signal. We can do so using a wavelet transform and looking for outliers in the detail coefficients. Figure 4.2 shows an example of this on the same signal LIT101 with a wavelet transform using the Daubechies 6 wavelet [38], with a single level of decomposition. We can see on this figure how each discontinuity creates a large spike in the detail coefficients.

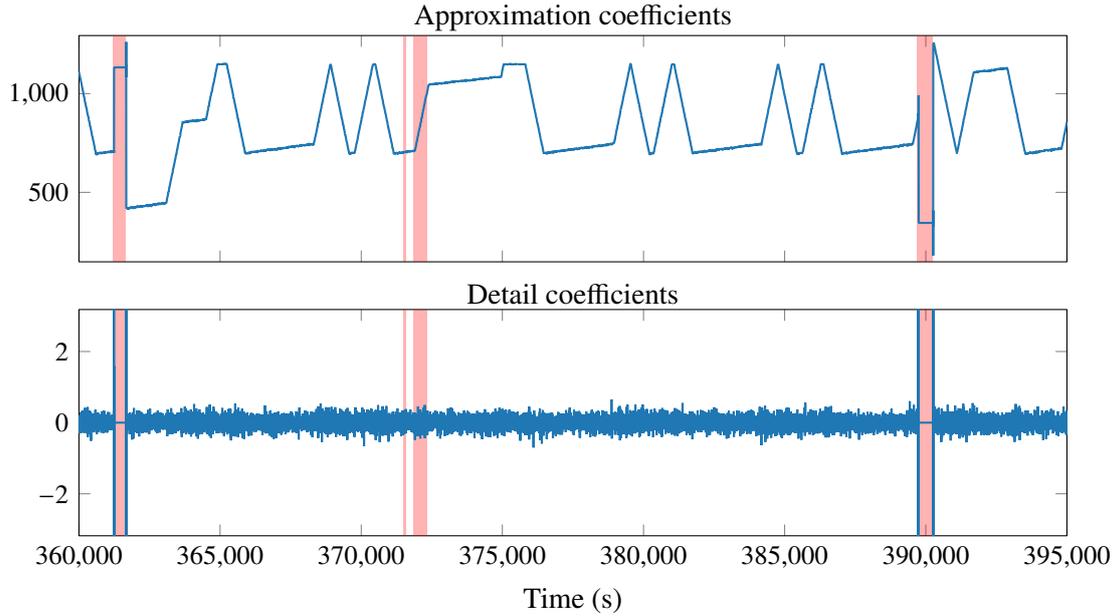


Figure 4.2: Wavelet transform on the signal LIT101

To detect these spikes, we can simply apply a threshold on the amplitude of the detail coefficients, and trigger a detection whenever this threshold is exceeded. Of course the question we have to address now is how to choose this threshold. For the experiments conducted here, we went for a fairly straightforward approach: we have been working with the RDT threshold $\lambda_\gamma(\tau)$ this entire time, so we also decided to use it to define our detection threshold, which leads us to the issue of choosing the parameters γ and τ .

To choose these parameters properly, let us start by setting the problem. We model the detail coefficients (y_1, \dots, y_N) as independent realizations of a random variable Y , which we decompose as the sum of some random variable $\Theta \in \mathcal{M}(\Omega, \mathbb{R})$ and a Gaussian noise $X \sim \mathcal{N}(0, \sigma^2)$, independent from each other. Now, recall that the RDT problem addresses the issue of testing $\mathcal{H}_0: \|Y\|_2/\sigma \leq \tau$ against $\mathcal{H}_1: \|Y\|_2/\sigma > \tau$, and the test $\mathcal{T}_{\lambda_\gamma(\tau)}$ is the answer we gave to that problem. The idea is then to set the tolerance τ to define the maximum acceptable amplitude of the detail coefficients. By definition of our problem, this maximum amplitude is therefore $\sigma\tau$.

However, describing the problem this way can cause an issue: if the noise variance varies over time, then our definition of a discontinuity will also change, which may not be a desirable behavior. This is the case in Fig. 4.2, where the variance of the detail coefficients changes slightly over time, depending on the current phase. Some attacks may also affect it. For instance, during the two attacks that spoofed the sensor (around 361,000 s and 390,000 s), there is no noise whatsoever between the two spikes: the detail coefficients are all equal to 0, since the base signal LIT101 is constant during that time.

It might be more practical to set this maximum amplitude directly, meaning that our problem becomes testing $\mathcal{H}_0: \|Y\|_2 \leq \tau'$ against $\mathcal{H}_1: \|Y\|_2 > \tau'$ where τ' is the desired maximum amplitude. Adjusting the test for this problem is easy: we simply have to consider the threshold $\lambda_\gamma(\tau'/\sigma)$, resulting in the thresholding test $\mathcal{T}_{\lambda_\gamma(\tau'/\sigma)}$ defined as:

$$\begin{aligned} \mathcal{T}_{\lambda_\gamma(\tau'/\sigma)}: \mathbb{R} &\rightarrow \{0, 1\} \\ y &\mapsto \begin{cases} 1 & \text{if } y \geq \sigma \lambda_\gamma(\tau'/\sigma) \\ 0 & \text{otherwise} \end{cases} \end{aligned} \quad (4.1)$$

By defining the tolerance τ' , we now have a constant definition of the discontinuities that we want to detect, regardless of the noise variance. This test takes the variations of the noise variance into account by adjusting the threshold appropriately. This noise variance is estimated on the detail coefficients using a

robust estimator. We cannot use the usual ML variance estimator because of the spikes that are present when discontinuities occur, since this estimator is not robust to such outliers, and would yield absurd results when used in these conditions. Instead, we have chosen to use the MAD (Median Absolute Deviation) estimator [39] which is defined for N samples $(y_1, \dots, y_N) \in \mathbb{R}^N$ by:

$$\text{MAD}(y_1, \dots, y_N) = 1.4826 \times \text{median}(|y_1 - \bar{y}|, \dots, |y_N - \bar{y}|) \quad (4.2)$$

with $\bar{y} = \text{median}(y_1, \dots, y_N)$, assuming that the samples (y_1, \dots, y_N) follow a Gaussian distribution. For simplicity, we decided to ignore the fact that the noise variance changes slightly over time and simply used a single estimate of σ using all samples, giving us a single threshold used for the entire signal. While the noise variance varies over time depending on the current state of the system, these variations are small, and using a single noise variance estimate computed using the entire signal is sufficient for demonstration purposes.

In order to use this test, we have to choose the parameters γ and τ' . As a reminder, the level γ represents the maximal desired false-alarm rate, and τ' the maximal acceptable amplitude of the detail coefficients. We can set γ with some considerations regarding an acceptable time between false alarms. For example, choosing $\gamma = 10^{-4}$ means that we can tolerate a false alarm every 10,000 tested samples on average. With a sample rate of 0.5 Hz ¹, this would translate to one false-alarm every 5.5 h on average.

Choosing the tolerance τ' is trickier, since we need some knowledge of the system in order to determine a normal range for the detail coefficients. An alternative that we will use here is using normal operation data to define this normal range. In addition to the data that contains attacks, we also have access to a one-week-long continuous recording of the same signal LIT101 under normal conditions. Using this recording, and performing the same wavelet transform, we can determine what the normal range of these detail coefficients should be. We are then effectively learning the parameter τ' , using normal observations of the signal only. The result of this wavelet transform on normal data is shown on Fig. 4.3.

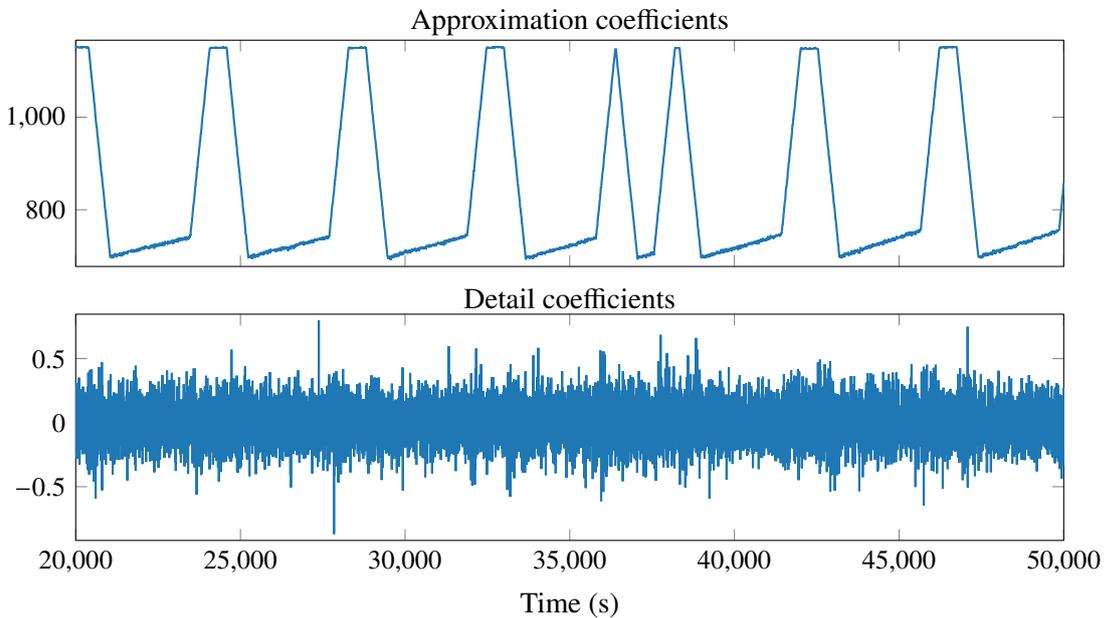


Figure 4.3: Wavelet transform on normal LIT101 signal

From the detail coefficients shown on this figure, we now want to determine a tolerance τ' . Since every data point here is normal, we want to choose a tolerance large enough so that all these samples are considered normal when using the test $\mathcal{J}_{\lambda_\gamma(\tau'/\sigma)}$. On the other hand, this tolerance should remain as small

¹Recall that we are testing the detail coefficients from the first level of wavelet decomposition, hence the sampling rate of the detail coefficients is half of that of the original signal.

as possible, because the samples observed here are the basis of our definition of normal behavior: samples with higher amplitude than the ones observed here should be treated as anomalies. Therefore we should choose the smallest possible tolerance τ' such that testing all samples of the normal signal yields no false alarm. We also have to consider the fact that we are willing to accept a false-alarm rate up to γ . Hence it may actually be more appropriate to choose a tolerance such that we get an effective false-alarm rate of γ on our training dataset.

Using a value of $\gamma = 10^{-4}$, and estimating the noise variance with the MAD estimator, we obtained on this training dataset a tolerance $\tau' = 0.358$. Now we can use these parameters to detect discontinuities on the signal with attacks, after re-estimating the noise variance σ on this signal. Detection is simply performed by applying the test $\mathcal{J}_{\lambda_{\gamma}(\tau'/\sigma)}$ to each detail coefficient, which simply consists in comparing the amplitude of each sample to the threshold $\sigma \lambda_{\gamma}(\tau'/\sigma)$. Figure 4.4 shows the threshold applied on the detail coefficients of the signal LIT101 with attacks.

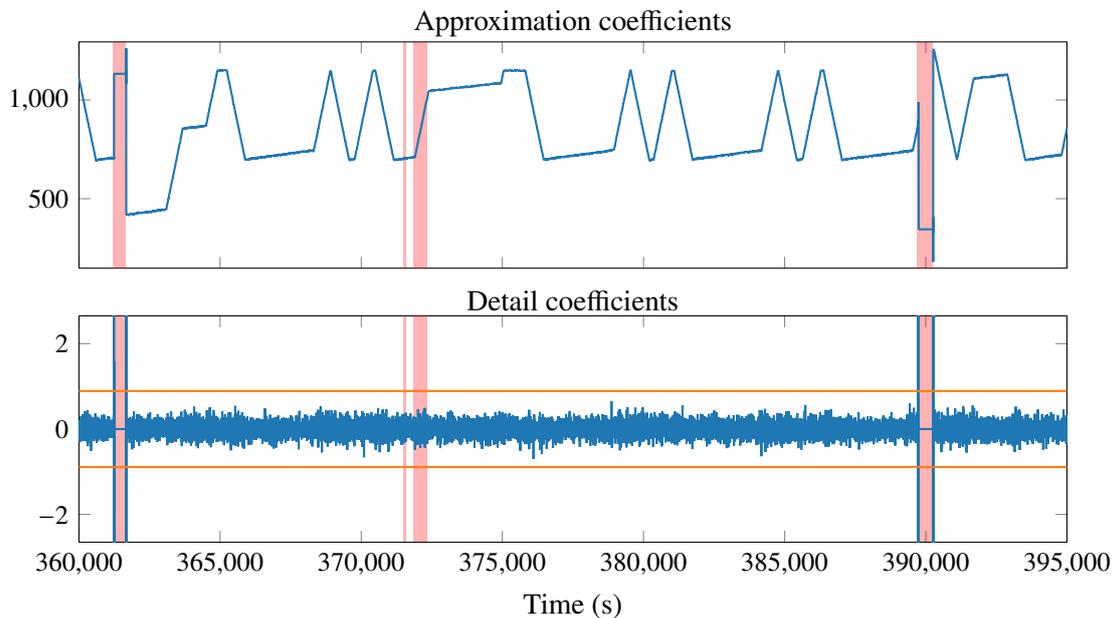


Figure 4.4: Threshold applied on the detail coefficients

Using this approach, we were able to detect all nine discontinuities present in this signal. These nine discontinuities are caused by five different attacks. Of these five attacks, four of them create a discontinuity at the start and at the end of the attack. We can therefore consider that each of these four attacks can be detected in real time at their start, which should be early enough to react. However, the fifth attack creates a discontinuity only when it ends, once the sensor is no longer spoofed, as can be seen in Fig. 4.5. In this case, the discontinuity is indeed detected, but that means we can only detect this attack once it is over, which is too late for a proper real-time response.

We also need to look at the effective false-alarm rate, which we expect to be lower than the level $\gamma = 10^{-4}$. In this case, a single sample triggered a false alarm among the 224,965 tested samples, yielding an effective false-alarm rate of approximately 5×10^{-6} . This false-alarm rate is therefore acceptable with regard to the set level.

We demonstrated the use of this method on a recorded signal, but adapting it to be usable in real time should be straightforward. Indeed, the wavelet transform should not cause any problem, since it simply consists in filtering the signal with a FIR (Finite Impulse Response) filter and downsampling its output. These two operations can be carried out in real time without any problem. Estimating the noise variance σ in real time should also not pose any problem, since we can adapt the MAD estimator to use a finite rolling window on the signal.

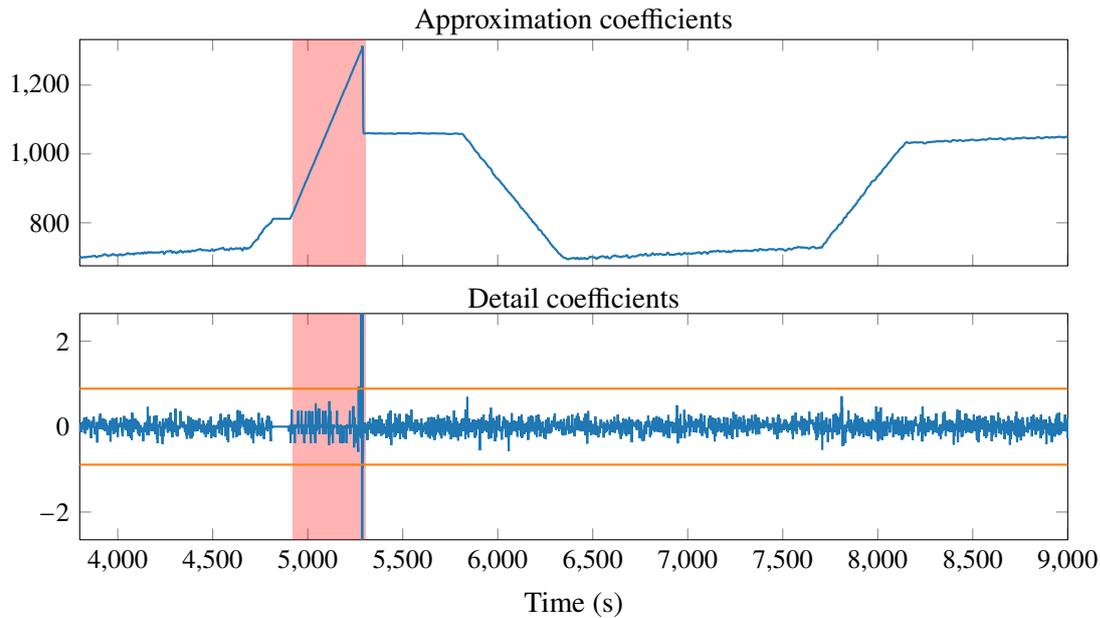


Figure 4.5: Late attack detection

Overall, this is a simple method to detect discontinuities in signals which can be used with little preliminary setup. The two parameters γ and τ' can be easily configured with a reference dataset containing only normal data, and testing this method on real data shows that we get good detection performance while properly controlling the false-alarm rate. The limited testing on real data could be completed with tests on artificial signal, which would allow a finer evaluation of the performance of this method. For example we could measure the influence of the amplitude of a discontinuity on the detection rate. It would also be interesting from a theoretical viewpoint to link the amplitude of a discontinuity on the signal to the amplitude of the spike created in the detail coefficient for a given wavelet. This would allow defining a tolerance based directly on a specification, without necessarily requiring normal operation data.

4.2 Segmenting and learning phases, and detecting anomalies

We now present in this section our second approach, which consists in learning the different phases present in a physical signal. For example, going back to the LIT101 signal that we have been using as an example so far, we can identify a succession of four phases, depending on the state of the actuators that are directly related, in this case the pumps and valves that control the inflow and outflow connected to the water tank in which the water level sensor LIT101 is located. Based on this observation, we decided to develop a method to discover these phases using normal operation data, then characterize these discovered phases in order to get a model of the system under normal conditions. Using this model, we then want to verify whether the current state of the system matches one of these discovered phases or not. If not, we can then consider that we detected an anomaly and raise an alarm. We can summarize this process in the following four steps:

1. Segmentation: using clean data without any anomalies or attacks, segment the signal according to the state changes of the system.
2. Feature extraction: for each segment, extract relevant characteristics that can be used to identify each phase.

3. **Training:** group the different segments using the extracted features to identify the existing phases, and define normal ranges of the features of interest for each phase. The output of this training phase constitutes our model of the system.
4. **Anomaly detection:** Using the model defined in the previous step, analyze the current state of the system and verify whether it matches any of the previously learned states.

We will focus on presenting the first part of this process, which is supported by the theoretical work presented in the previous chapters. As the rest of the learning and detection process will rely on this first step, it seemed important to us to offer a proper presentation and justification for the choices made to develop this segmentation method. Afterwards, we will briefly touch on ways to build a model based on this segmentation and how to perform detection.

4.2.1 Change-detection in time series

As we mentioned back in the introduction, components in an industrial systems are often characterized by a succession of phases. For example, let us consider the LIT101 sensor signal along with the directly related actuators MV101 and P101, which are represented on Fig. 4.6 in normal operating conditions. As we can see there, the LIT101 signal is piecewise linear with four different slopes, matching the four possible combinations of states of the actuators regulating the inflow and outflow. These relations are summarized in the following table:

Valve (MV101)	Pump (P101)	Water level (LIT101)
Closed	Off	Constant
Open	Off	Fast increase
Closed	On	Fast decrease
Open	On	Slow increase

Table 4.1: Relations between the states of the sensor LIT101 and the actuators MV101 and P101

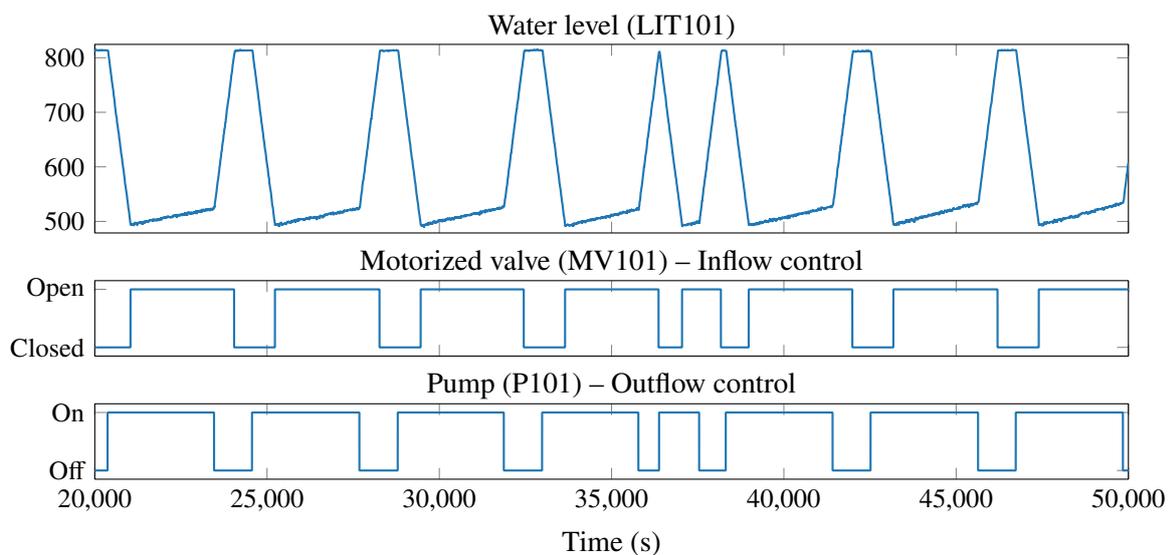


Figure 4.6: Sensor LIT101 and actuators MV101 and P101 during normal operation

Using the information from these actuators, we can easily segment the LIT101 signal into its different phases: a new segment simply starts whenever the state of one of the actuators changes. However, it would

be useful to be able to do this using the LIT101 signal on its own. For example, an attack could try to stealthily manipulate one of the actuators, forcing it into one state and also spoofing it so that this change is not visible on the transmitted actuator signals. Working on the LIT101 signal would in this case allow detecting discrepancies between this signal and the observed state of the actuators. As such, we want to use a method that would allow us to detect phase transitions, based solely on the signal LIT101.

Finding such transitions in signals is nothing new: the problem of detecting changes in time series can be found in many fields, such as quality control in manufacturing, fault detection, telecommunications, etc. The literature on this topic is vast (see for example [40, 41, 42]). We will present it only briefly.

In short, we can formalize the change detection problem as follows: consider a time series $(Y_n)_{n \in \mathbb{N}}$, where each sample Y_n has probability distribution P_{Y_n} . This time series is assumed to be piecewise stationary, meaning that the distribution of these samples remains constant for a certain duration and only changes at certain rupture points $n_i \in \mathbb{N}$, such that $P_{Y_{n_i}} \neq P_{Y_{n_i+1}}$. The goal of change detection methods is to find these rupture points and estimate them as well as possible from a realization of the time series.

However, detecting any change in the probability distribution of the series is a complex task, and may also lead to detecting a very large number of irrelevant changes. We instead usually consider detecting changes in some of its interesting properties. The most widely addressed problem is detecting changes in the mean, for which many solutions have been proposed. The Shewhart chart [43] and the CUSUM (Cumulative Sum) [44] are two of the most well-known approaches proposed to address it. The Shewhart chart was first introduced for quality control in manufacturing in order to verify that produced components were within specifications. In short, the Shewhart chart operates on blocks of N samples of the time series and assumes that all samples in each block are independent and identically distributed, sharing a common expectation μ , and also that we know the distribution of the samples both before and after a change occurs. For each block, it tests whether this mean μ is equal to a known nominal value μ_0 , meaning that there was no change, or whether the mean μ is now equal to another known value μ_1 , signifying that the mean of the process has changed. This test is simply performed using the appropriate Neyman-Pearson test since both hypotheses are simple by definition of this problem. The CUSUM is another change-detection method which can be seen as an adaptation of the SPRT (Sequential Probability Ratio Test) [45]. Unlike the Shewhart chart, this method does not operate on blocks of samples but processes the signal one sample at a time. Despite this difference, it relies on very similar hypotheses and also requires perfect knowledge of the distribution of the samples before and after change.

This distribution knowledge requirement can be found in many change detection approaches, and can severely limit their applicability to real-world scenarios where good data models are often not available. This observation in hypothesis testing is what led to the development of the RDT approach in the first place, and we therefore decided to develop our own change-detection method based on it in an attempt to inherit its properties. The goal here is then to develop a change-in-mean detection method that is applicable to real signals and makes as few assumptions as possible on the nature of the observed signals. As a reminder, the RDT approach only assumes that the observed signal Y is the result of the observation of some unknown phenomenon Θ in presence of additive and independent Gaussian noise X .

One may wonder if the wavelet-based method we developed in the previous section could be applied here to detect changes in the signal. Indeed, we are looking to detect changes in the slope of the signal, which are irregularities in its derivative, and wavelets are well-known to be an appropriate tool to detect such irregularities. However, after trying it out, we found that it was not actually applicable in this specific context. We will explain why this is the case later, when we start applying the method we are about to present to real signals.

4.2.2 An RDT-based change-in-mean detection method

The method that we developed is inspired from the Block-RDT control chart presented in [46], which is a block-based method that detects changes-in-mean in time series using the RDT test. However this method assumes that both the current mean of the signal and the noise variance are perfectly known. This is not applicable in our case since the goal of this segmentation is to then be able to learn these

values depending on the state of the system. We address this in our method by estimating both the current mean and the noise variance, and then improving these estimates as the current phase continues. This is where our study of the Asymptotic RDT extension in Chapter 3 comes into play, ensuring that we will asymptotically reach the performance we would get without estimation. Our method requires choosing the following parameters:

- The processing block size N
- The tolerance τ
- The desired false-alarm rate γ

In addition to these parameters, one can choose which estimators to use for the mean and the noise variance. For the mean, we will simply use the classic empirical mean of the samples, and for the noise variance we will use the MAD estimator (see Eq. (4.2)). We will explain why we chose the MAD estimator over the usual maximum-likelihood estimator in Section 4.2.3, where we apply this change-detection method to real signals. Note that both estimators are consistent (see [47] for the MAD), which is a requirement for the Asymptotic RDT to be applicable.

As previously, we assume that each sample Y_n of the time series can be written $Y_n = \Theta_n + X_n$ where $\Theta_n \in \mathcal{M}(\Omega, \mathbb{R})$ is a signal with unknown distribution and $X \sim \mathcal{N}(0, \sigma_n^2)$ is the Gaussian noise independent from Θ_n . Our method processes the signal in fixed-size blocks of N samples, and consists in testing block after block whether the empirical mean of the signal Θ within a block of N samples has changed significantly or not compared to that of the current phase.

In the following, for a time series $(A)_{n \in \mathbb{N}}$, we will denote by $\langle A \rangle_m^p$ its empirical mean between samples $m \in \mathbb{N}$ and $p \in \mathbb{N}$:

$$\langle A \rangle_m^p = \frac{1}{p - m + 1} \sum_{i=m}^p A_i \quad (4.3)$$

If the start and end samples m and p are not mentioned, we are implicitly talking about the samples contained in the currently processed block: if we are working on the k -th block (with k starting from 1), then $m = (k - 1)N + 1$ and $p = kN$.

Before we present our change-detection method, we have to mention several points regarding the RDT test we will be using:

- Just like we did in the previous section with the wavelet-based method, we will define the tolerance relatively to the signal norm, without normalizing it by the noise standard deviation. As a reminder, this means that we use a tolerance τ' to test whether $\|\Theta(\omega) - \theta_0\|_2 \leq \tau'$ or $\|\Theta(\omega) - \theta_0\|_2 > \tau'$, and instead of the regular test $\mathcal{J}_{\lambda_\gamma(\tau)}$, we use the test $\mathcal{J}_{\lambda_\gamma(\tau'/\sigma_0)}$.
- Since we are working on the mean of a block of N samples, the threshold to be used for a given tolerance τ and level γ is not $\lambda_\gamma(\tau)$, but $\lambda_\gamma(\tau\sqrt{N})/\sqrt{N}$ instead [46], with the test being applied to the empirical mean of the observations. The factor \sqrt{N} present here translates the decrease in noise variance that occurs when averaging N independent samples.
- Of course, we will also be considering estimates of the model θ_0 and the noise variance σ_0 , and as such we will be using a test of the form $\tilde{\mathcal{J}}_{\lambda_\gamma(\tau)}$ that takes these estimates into consideration, as defined in Chapter 3 (see Eq. (3.5) for the definition of these tests).

Putting all this together, the test that we will use is a test of the form $\tilde{\mathcal{J}}_{\lambda_\gamma(\tau'\sqrt{N}/\hat{\sigma})/\sqrt{N}}(\langle Y \rangle, \hat{\theta}, \hat{\sigma})$.

Now we have everything we need to present our change-detection method. It can be summarized in the following steps:

1. First we get an initial estimate $(\hat{\theta}, \hat{\sigma})$ of the initial signal parameters (θ_0, σ_0) using the first N samples.
2. Then, using the N next samples, we want to test whether the empirical mean of the signal Θ has changed significantly or not, meaning that we want to test whether $\|\langle \Theta(\omega) \rangle - \theta_0\|_2 \leq \tau'$ or

$\|\langle \Theta(\omega) \rangle - \theta_0\|_2 > \tau'$. We do this by using the estimates $\hat{\theta}$ and $\hat{\sigma}$ as substitutes for θ_0 and σ_0 respectively and applying the test $\tilde{\mathcal{F}}_{\lambda_\gamma(\tau' \sqrt{N}/\hat{\sigma})/\sqrt{N}}(\langle Y(\omega) \rangle, \hat{\theta}, \hat{\sigma})$.

3. What we do next depends on the result of the test:

- If $\tilde{\mathcal{F}}_{\lambda_\gamma(\tau' \sqrt{N}/\hat{\sigma})/\sqrt{N}}(\langle Y(\omega) \rangle, \hat{\theta}, \hat{\sigma}) = 0$, then there has not been any significant change and we consider that the current block is part of the same phase of the signal. In that case, we update our estimates of the mean and of the noise variance using the tested samples, and then go back to step 2 with the next block of N samples.
- Otherwise, if $\tilde{\mathcal{F}}_{\lambda_\gamma(\tau' \sqrt{N}/\hat{\sigma})/\sqrt{N}}(\langle Y(\omega) \rangle, \hat{\theta}, \hat{\sigma}) = 1$, then we decide that a significant enough change has occurred and that the signal has entered a new phase. In that case, we can report that a change has been detected within this block, and restart from step 1 using the next block of N samples to initialize the next phase. Note that we do not use the current block to initialize the next phase since we consider that there is a change within this block, which would compromise the initial parameter estimation.

This process is also described in algorithmic form in Fig. 4.7.

Input: Time series Y_n , block size N , tolerance τ' , desired false alarm probability γ

Output: Estimated instants of change

$\hat{\theta} \leftarrow \text{mean}(Y_0, \dots, Y_{N-1})$ (Initial mean estimate)

$\hat{\sigma} \leftarrow \text{std}(Y_0, \dots, Y_{N-1})$ (Initial variance estimate)

$n \leftarrow N$ (Number of samples in the current segment)

$i_s \leftarrow 0$ (Index of the first sample of the current segment)

repeat

$s \leftarrow \text{mean}(Y_{i_s+n}, \dots, Y_{i_s+n+N-1})$ (Current block mean)

$z \leftarrow |s - \hat{\theta}|/\hat{\sigma}$ (Test statistic)

$T \leftarrow \lambda_\gamma(\tau' \sqrt{N}/\hat{\sigma})/\sqrt{N}$ (Threshold)

if $z \leq T$ **then**

$n \leftarrow n + N$ (Extend the current segment)

$\hat{\theta} \leftarrow \text{mean}(Y_{i_s}, \dots, Y_{i_s+n+N-1})$ (Update the mean estimate)

$\hat{\sigma} \leftarrow \text{std}(Y_{i_s}, \dots, Y_{i_s+n+N-1})$ (Update the variance estimate)

else

Notify that a change has been detected between indices $i_s + n$ and $i_s + n + N - 1$

$i_s \leftarrow i_s + n + N$ (Start the next segment after the end of the current block)

$\hat{\theta} \leftarrow \text{mean}(Y_{i_s}, \dots, Y_{i_s+N-1})$ (Initial mean estimate)

$\hat{\sigma} \leftarrow \text{std}(Y_{i_s}, \dots, Y_{i_s+N-1})$ (Initial variance estimate)

$n \leftarrow N$

end if

until end of Y_n reached

Figure 4.7: Proposed RDT-based change-in-mean detection algorithm

Before moving on to application examples of this method, we should address the problem of setting the three required parameters that we mentioned, which as a reminder are the block size N , the tolerance τ and the desired false alarm rate γ :

- The tolerance τ represents the amplitude of changes that should be detected. Therefore setting it appropriately requires some prior knowledge of the system in order to choose a sensible value, taking into account the amplitude of the transitions that may occur.

- The block size N has to be carefully chosen by considering several aspects, including notably the duration of the phases in the signal, the sample rate of the signal, the desired resolution for estimating the change instants, and the overall performance of the method (detection and false-alarm rates).

Working with large blocks means that we quickly get good estimates of the model and of the noise variance. Since we use the mean of the signal over a block, it also yields a better SNR and therefore better performance. However, this also means that we assume that the signal changes slowly enough, otherwise we would easily miss short segments. In addition, we also get less precision in estimating the change instants, since we only know that they occur within a block of N samples.

We could work with faster signals by lowering the block size N , but this also means we would not get as much of an SNR increase from averaging samples, and therefore the overall performance of the method would be negatively affected. This effect is compounded by the coarser initial estimate of the model and noise variance, which is also performed on one block of N samples at the beginning of each detected phase.

- The level γ will usually be set last, and used to adjust the threshold based on initial results.

4.2.3 Application to real signals

After describing our proposed method, we will now show the results yielded by applying this method on real signals, using once again the LIT101 water level signal that we used previously, under both normal and attack conditions. We have seen that this signal is piecewise linear and as we discussed earlier, its slope is an information of interest. Therefore, instead of applying this algorithm directly to our signal, we will apply it to its derivative, which we expect to be piecewise constant with added noise.

We can see this on Fig. 4.8 where we show both the signal in normal conditions and its derivative, which is obtained simply by computing the difference between consecutive samples. The derivative may appear surprisingly noisy compared to the original signal. This is due to the fact that the variations of the original signal are very slow and we are observing it here over a very long period of time (a little over 8 hours). Therefore the steps in the derivative caused by the slope changes are actually fairly small compared to the noise level. This is why our wavelet-based detection method that we presented earlier does not work here to detect these discontinuities: unlike the ones we tried to detect earlier, these discontinuities are very small and end up being completely hidden in the surrounding noise, making them very difficult to detect in this manner.

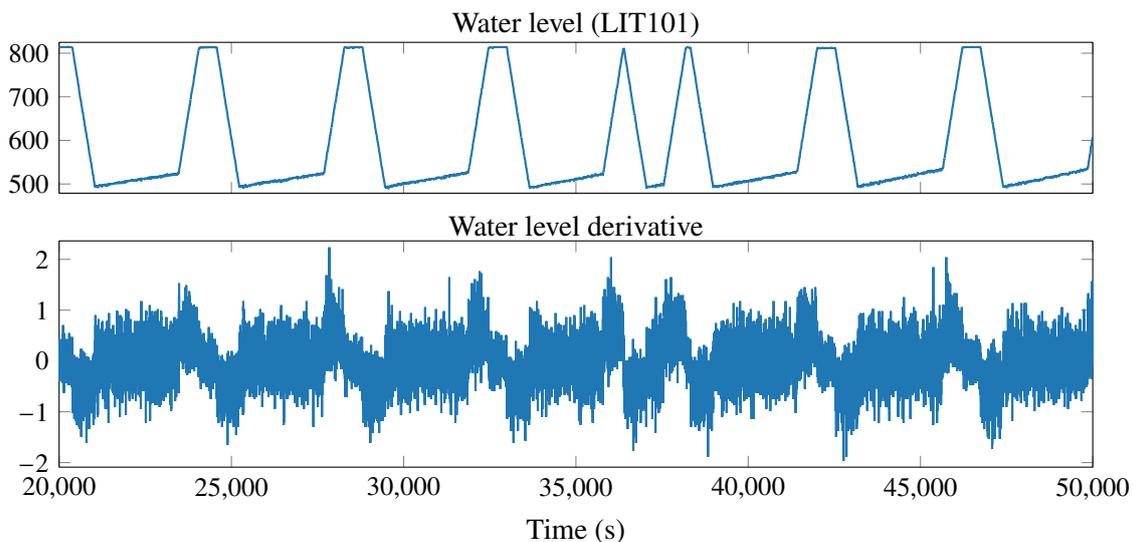


Figure 4.8: Signal LIT101 and its derivative under normal operation

Applying our change-detection method on this derivative yields the result visible on Fig. 4.9, where the tracked mean is shown on the bottom figure in orange. The following parameters were used:

- $N = 40$
- $\gamma = 10^{-3}$
- $\tau' = 0.05$

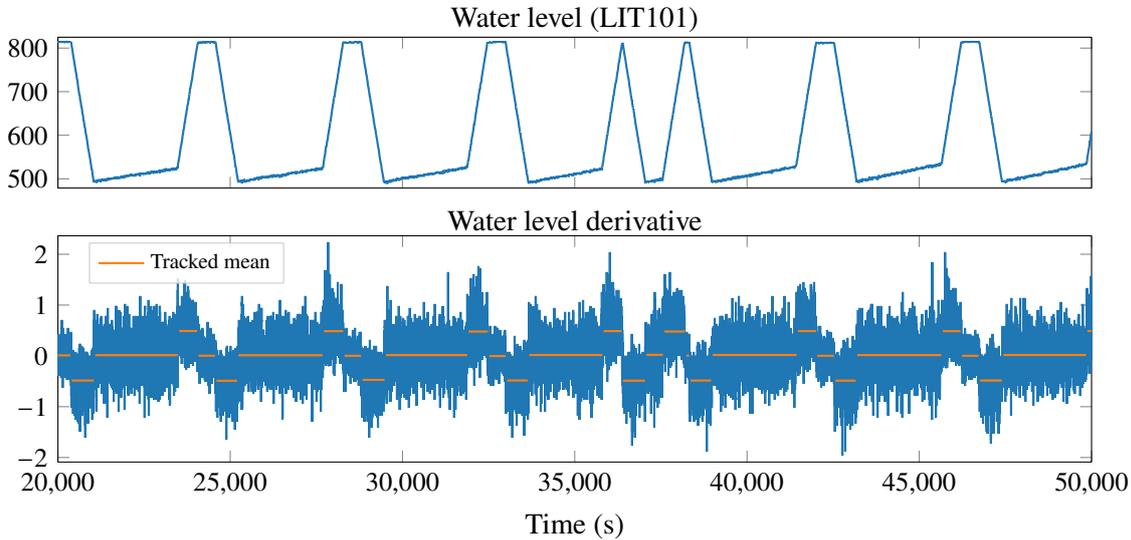


Figure 4.9: Tracked signal mean on the derivative of LIT101 ($N = 40$, $\gamma = 10^{-3}$, $\tau' = 0.05$)

31 segments were found by our method using these parameters, and they seem to visually match up with the different phases present on the original signal. In this example, we can consider that we were able to detect all changes without any false-alarm. Of course, this is only a qualitative evaluation on a portion of this signal, and does not necessarily reflect the actual performance of this method. Quantitative performance evaluation is tricky, since we do not have access to the exact change instants in this dataset. We would need to manually label the changes in this dataset, which would require manually identifying every change point in the entire signal with good accuracy, which is not easy considering the very slow variations of the signal. An alternative could be relying on actuator data, but the issue with using this data here is that we have to account for potential delays between a change happening on the actuators, and their effects being visible on the sensor, which would skew the performance evaluation.

We can also see gaps between each discovered segment, which correspond to the blocks of N samples in which changes were detected. As a reminder, this is by design since this method cannot find the change instant more accurately within each of these gaps, and we do not want to risk having a compromised initial estimate by using samples from two different phases. Since this signal has a sample rate of 1 Hz, this means that there is an uncertainty of up to 40 s in the location of each change point, and we also have to account for potential detection delays, for example a change occurring late in the previous block which was not detected at that time, but only in the next block.

We could improve this resolution by reducing the block size N , but this may cause other issues. An example of this is shown in Fig. 4.10, where we set the block size N to 10 samples. Comparing this result to the previous one, the changes of interest are still detected, but we also have in addition to these a considerable number of very short segments: from 31 segments found in the previous figure, we now have 206 segments, many of which do not match any change of interest.

We can attempt to circumvent this issue by lowering the parameter γ , which should make the method less sensitive to noise and therefore lower the number of false-alarms. We can see the result of this on Fig. 4.11, where we reduced γ from 10^{-3} to 10^{-6} , while keeping the other parameters unchanged. While this does significantly reduce the false-alarm rate, there are still a few left: we indeed have 43 segments instead of the expected 31. In addition to these remaining false-alarms, we can also see that some changes

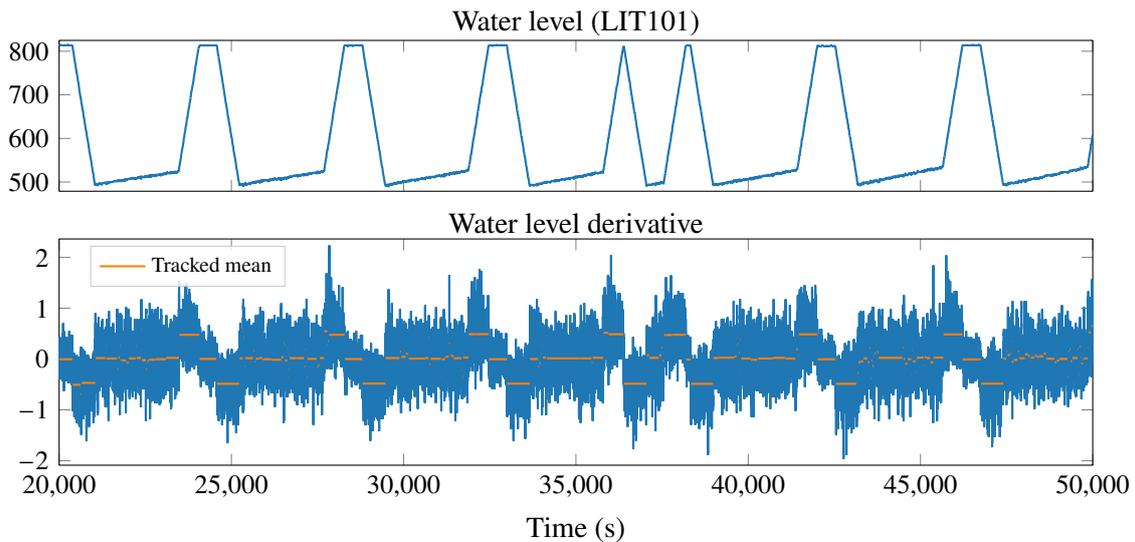


Figure 4.10: Tracked signal mean on the derivative of LIT101 ($N = 10$, $\gamma = 10^{-3}$, $\tau' = 0.05$)

were also missed. For example, changes around 24,000 s and 28,000 s were missed, and in each of these two instances, two expected segments were merged into a single one.

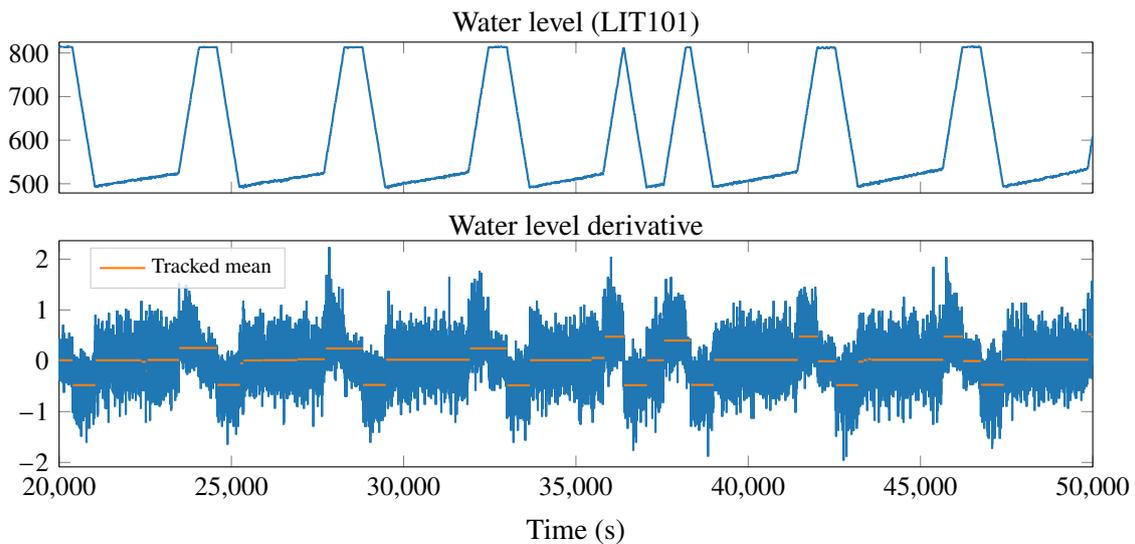


Figure 4.11: Tracked signal mean on the derivative of LIT101 ($N = 10$, $\gamma = 10^{-6}$, $\tau' = 0.05$)

One potential way to improve the resolution without compromising performance would consist in using a sequential test. This is fairly analogous to what is done in the CUSUM method, based on the SPRT sequential test. The basic idea of a sequential test is to process samples one by one, until a decision can be taken. For each tested sample, there are three possible outcomes:

1. We consider that a change has been detected.
2. We consider that there is no change.
3. There is not yet enough information to decide whether there was a change or not, and we continue with the next sample.

The advantage of this approach is a finer detection of the change-points by working one sample at a time instead of processing entire blocks of samples. An RDT-based sequential test exists ([19, 20]) and could be considered here as a replacement for the block-based test used. However there has not yet been any study of this sequential test when used with estimates of the model and of the noise variance, and therefore we cannot yet offer the same performance guarantees that we have here with the block approach.

We can also see how this change-detection method behaves in presence of attacks on Fig. 4.12. On the

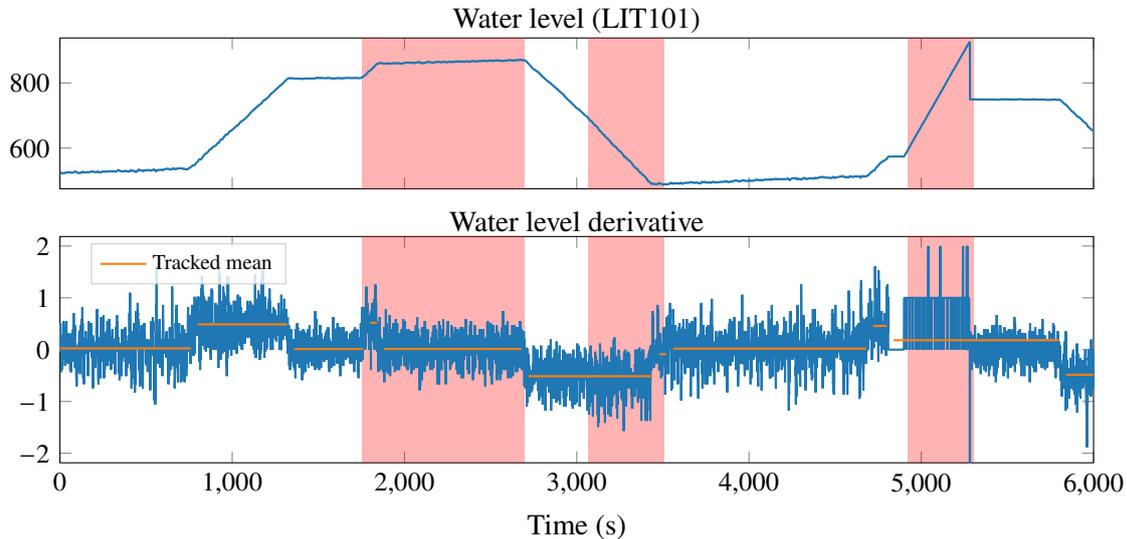


Figure 4.12: Tracked signal mean on the derivative of LIT101 with attacks ($N = 40$, $\gamma = 10^{-3}$, $\tau' = 0.05$)

three attacks that are displayed on that figure, we can note the following:

- In the first attack (around 2,000 s), we can see that we properly detect the two changes that occur at the beginning of the attack. These steps have comparable amplitudes to the ones we observed previously under normal conditions, therefore we expected to be able to detect these.
- In the second attack, there is a very slight change in the slope of the signal, which goes undetected with the parameters used here. We can detect it by lowering the threshold used (by either increase γ or decreasing τ'), but this will also result in a higher false-alarm rate. An example of this is shown on Fig. 4.13, where we lowered the detection threshold by increasing γ and lowering τ' .
- The third attack, which consists in spoofing the sensor, causes very unnatural behavior on the sensor derivative. Indeed, during this attack, there are very few different values observed, as if the signal was very coarsely quantized, and there is no noise whatsoever. This notably means that during this time, our model of a signal observed in independent Gaussian noise does not hold. We can detect a change at the beginning of the attack, but not at its end, as we have one segment extending from the start of the attack to the end of the following normal phase.

The behavior encountered in that third attack is the main reason why we are using the MAD as the noise variance estimator instead of the more commonly found maximum-likelihood estimator (see Eq. (3.26)). Indeed, looking at the end of this attack, where the spoofed signal suddenly goes back to normal, we can see that the derivative has a sudden negative spike. This spike is composed of a single sample with value -176 , whereas all other samples shown here are contained in the range $[-2, 2]$. If this sample ends up being used to estimate the noise variance, then using ML-estimator would result in a extremely large estimate for $\hat{\sigma}$, since that estimator is not robust to outliers. Using this estimation of the noise standard deviation would then lead to being completely incapable of detecting any more changes afterwards, because the threshold computed using this estimate also becomes very large. Replacing the ML-estimator with the

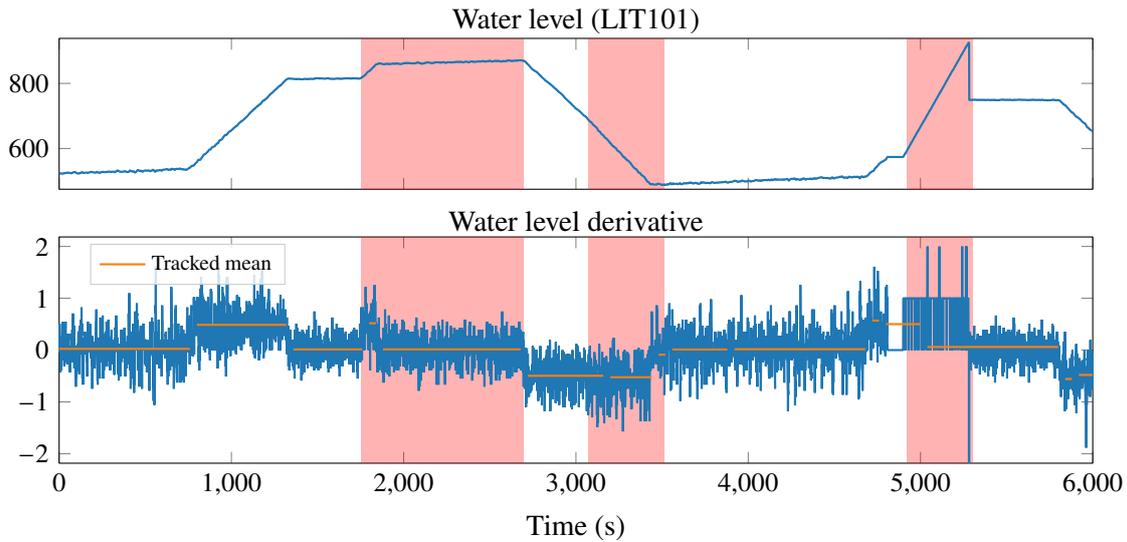


Figure 4.13: Tracked signal mean on the derivative of LIT101 with attacks ($N = 40$, $\gamma = 10^{-2}$, $\tau' = 0.01$)

MAD prevents such samples from significantly affecting the estimation of σ , and therefore allows the change-detection method to continue working as expected.

An issue that currently remains open is evaluating the performance of this segmentation. We have assessed the resulting segmentation qualitatively on a few examples, but we did not provide any actual metric that could be used to evaluate this segmentation quantitatively. The main issue here is that we do not have a ground truth to which we could compare the change points that we found. Indeed, the only labels available in the SWaT dataset indicate whether there is an attack happening at a given time, but no information is given regarding the normal operation of the system. We could try to estimate the real change points ourselves with the help of the actuator signals, but we found that this task would be very error-prone, and such labels may be considered dubious, since we are labelling the dataset ourselves after seeing the segmentation yielded by our own algorithm. We also noticed that there is a variable delay between a change in the state of an actuator, and its effect on the monitored sensor, making it very difficult to automate this task without manual adjustments. For all these reasons, we could not provide any performance measurements that would allow us to compare our method to others.

4.2.4 Perspectives: towards learning a model and detecting anomalies

After a fairly detailed presentation of our change-in-mean detection method, we will now present some in-progress ideas to use the previously performed segmentation as a base to learn characteristics of the signal of interest, and then detect anomalies from this basis. This section will be fairly short, as we were only able to conduct very limited testing using the signals at our disposal given the time constraints. We will start by considering a clustering approach to identify the phases of the signal and show how a basic K-means approach can perform as an example, reusing the segmentation shown in Fig. 4.9. We will then discuss some methods we considered to detect anomalies from this basis.

Using the previously defined change-detection method, we now want to build a model that will characterize the normal operation of the system. To do so, we will start by applying this method to clean data, without any attack, such as the example shown in Fig. 4.9, which is a clean signal recorded during normal operation of the system. Each yielded segment should then approximately correspond to one phase of the signal. The idea is to compute some relevant features on each segment, then use these features in a clustering algorithm in order to find clusters that represent each normal state of the signal.

Based on the segmentation performed on this signal, we can extract some features of interest. Here for simplicity, we will be using only a single feature, calculated as the mean of each segment, which is the

estimated slope on the original signal. Other features of interest that could be considered on a real signal are the length of each phase, the estimated noise variance, the start and end value of each phase since we are working on linear segments, etc. For the clustering algorithm, we will use the simple K-means method, as we think it should perform well enough for a first proof of concept. While it requires knowing the number of clusters ahead of time, this is not necessarily an unreasonable assumption to make depending on the system of interest. However we could consider other algorithms if we want to discover the number of clusters, such as DBSCAN [48] or OPTICS [49].

Figure 4.14 shows the result of K-means clustering on the signal LIT101 after using the segmentation shown in Fig. 4.9, assuming that we know there are 4 classes. For this, we used the segmentation performed on the entire signal, which yielded a total of 469 segments.

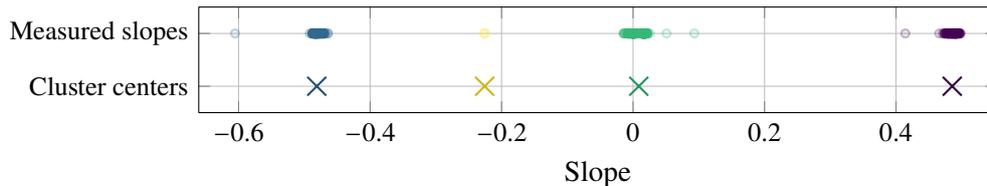


Figure 4.14: K-means output on normal LIT101 signal

At first sight, it may seem like we detected 4 reasonable classes. However, the yellow class centered around -0.2 is actually composed of two data points only, and therefore does not constitute a class of interest. Indeed, from what we could see on the signal, we would expect to have four classes comprised of similar numbers of points, since we can observe all four phases regularly in a cycle. These two yellow points, as well as other the other visible outliers, are actually segmentation artifacts caused by several very short segments, which tend to occur during fast transitions, and therefore do not match any actual phase of the signal. An example of such behavior is shown on Fig. 4.15, where we can see five very short segments detected during some brief peaks of the signal.

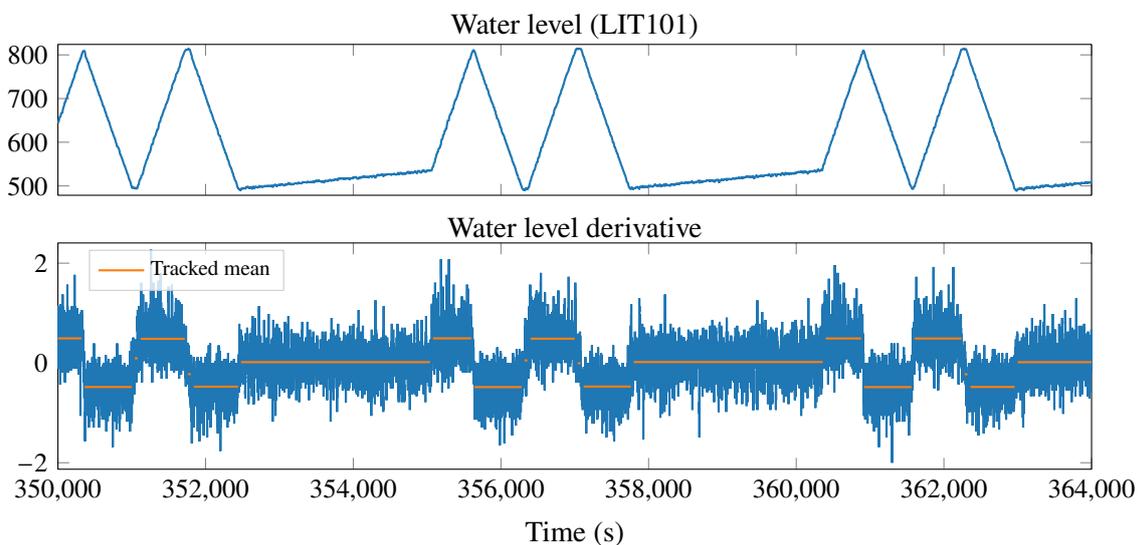


Figure 4.15: Example of segmentation artifacts caused by short segments

The class detected around 0 is actually composed of two classes that were grouped together:

- A class centered on 0 that corresponds to the phases during which the water level is constant, whenever the inflow valve is closed and the outflow pump is off.

- A class centered around 0.15 that corresponds to the slow water level increase, which occurs whenever the inflow valve is open and the outflow pump is on.

These two classes can be seen more clearly on Fig. 4.16.

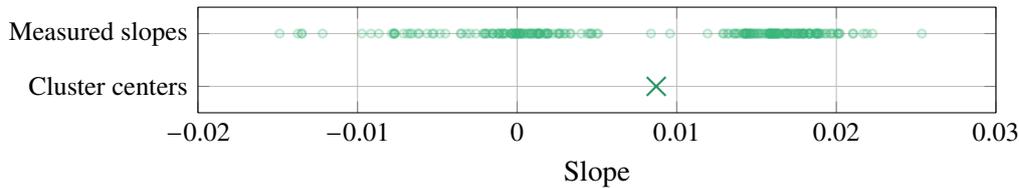


Figure 4.16: K-means output on normal LIT101 signal (zoom around 0)

Of course, the K-means clustering result can vary significantly based on the initialization method used. The result shown here was obtained using K-means++ seeding [50, 51], and consistently yields this result with repeated runs. Through manual seeding with carefully selected starting points, we were able to recover the clusters we were expecting in the first place. We can see this on Fig. 4.17. Here we get a reasonable classification of our points, barring the few outliers that remain.

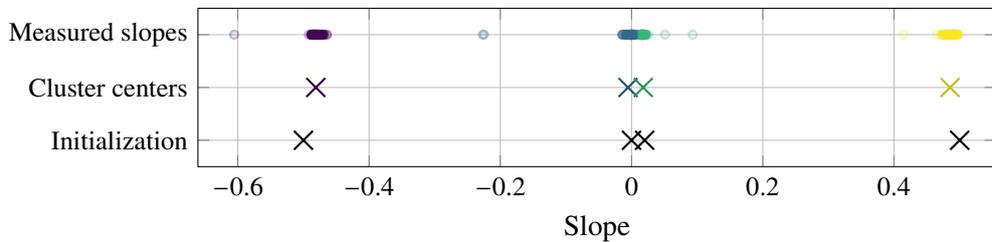


Figure 4.17: K-means output on normal LIT101 signal with manual seeding

Since our intention is to use this result as a basis to build a model of the normal operation conditions of the system, we can also take a look at the distribution of the slopes in presence of attacks. For this, we show on Fig. 4.18 the distribution of the slopes measured both in normal conditions and with attacks. For the normal data, we reused the clustering result shown in Fig. 4.17 as a reference point. For the attack data, we used the segmentation performed in Fig. 4.12, and colored these points according to which of the normal clusters they would belong to.

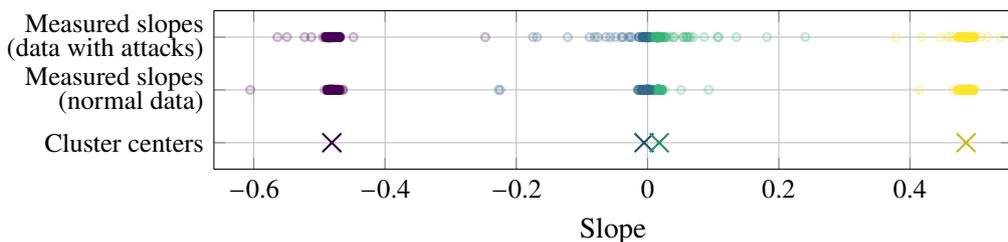


Figure 4.18: Comparison of slopes obtained on signal LIT101 both in normal and attack conditions

How do these two segmentations compare? First, we can see that most of the slopes in the attack dataset are concentrated around the same centers as the normal data. This is normal, because the attack dataset is mainly composed of normal data using the same slopes as the normal dataset. Outside of the normal clusters, we can see many additional slopes, similar to the artifacts we mentioned earlier. From this figure, it is difficult to say whether some of these slopes also match some attacks or if these slopes are all due to segmentation artifacts. This shows that there are still improvements required to get usable data out of our change-detection method.

Assuming that we are able to improve the method to get rid of these artifacts, let us discuss briefly what we could do using the clustering output to define the normal ranges for the slopes. We know that under normal conditions, the measured slopes should remain close to the centers that we discovered. Of course, we need to define what “close” means here. An idea we can consider is estimating the probability density function of these slopes, and then we could use these densities to define the normal regions, for example by using a threshold. We could estimate these densities by using a kernel-based approach for example, or we could also assume that these clusters are Gaussian and model the overall density of these slopes as a Gaussian mixture. Of course many other possibilities could also be considered here.

Once we have a definition of the normal operation ranges, we could then attempt to detect anomalies on unknown signals. To do this, what we envision is performing the same segmentation as previously and then for each segment verify if the slope falls within any of the previously defined normal regions. This process can be performed on streaming data for real time detection, with only a slight delay due to the block-based processing of the segmentation method.

4.3 Conclusion

This chapter has presented a first “proof-of-concept” for a completely autonomous machine learning system aimed at learning the nominal behavior of a signal and detecting anomalies with respect to this learned model. The calculation of the normal behavior is performed via a segmentation tool based on Asymptotic RDT; the segments are then used to extract the slopes representing the various possible phases in the nominal behavior; the slopes are clustered by K-means.

Empirically, the segmentation of the signal before calculating the various slopes pertaining to the nominal behavior seems relevant in that, visually, the segments exhibited by the statistical detection of changes match those that the observers could pinpoint themselves. However, a quantitative assessment of this segmentation is hardly feasible for the following reasons:

1. We do not have any ground truth, provided by either a model of the behavior of the system (for instance, an automaton or some software describing the various actions performed during the running of the system), or an expert that would have observed, segmented and labeled the different segments constituting the dataset.
2. We cannot provide ourselves this ground truth by either modeling the whole system as an automaton (this requires expertise beyond the scope of this thesis) or segmenting the dataset ourselves. In fact, a manual segmentation of this data set would be dubious because of our lack of expertise in this field, and could also be biased if not carried out before running our algorithms. In addition, it can be difficult to decide visually whether a change actually occurred or not in some time intervals, or even possibly multiple changes or only a single one.

Regarding the clustering, the results raise the same kind of questions, in addition with the usual limitations of K-means. Even when the initialization is correct so as to find the four expected clusters, the segmentation artifacts that we highlighted in the analysis make it difficult to properly evaluate the performance of this approach. This also means that we cannot develop, and even less assess, any anomaly/attack detection method at this point, as we saw when comparing the clustering results between the normal and attack datasets.

It is also worth emphasizing that even if we had the means to properly and quantitatively evaluate the approach, as a whole but also at each step, this evaluation would have remained very limited because of the very small number of attacks (9 discontinuities and 36 attacks in total in this dataset).

Although the segmentation is promising, especially because it relies on a mathematical framework that has been duly analysed in the foregoing chapters and even it is tempting to design an autonomous learning system, the discussion above emphasizes two issues in machine learning for cybersecurity: one directly connected to the signals considered above and the other from a more general stance.

We have highlighted the lack of expertise and labeled segments to assess our algorithms. In fact, the information that experts could provide is somewhat contained in the signals provided by the actuators. The direct use of these signals to assess our method can hardly be considered without a preliminary assessment of these same signals by an expert. Indeed, for instance, possible delays pinpointed by an expert should be taken into account to achieve an accurate segmentation. However, the signals provided by the actuators could perhaps be harnessed to detect anomalies and attacks by comparison with the outcome of the segmentation, to find incoherencies or even unacceptable deviations. For instance, could a simple correlation between these signals and the segmentation results be used to realign both in time and exhibit abnormal behaviors if the correlation coefficient is not large enough?

From a more general point of view, data is crucial in such studies. It turns out that obtaining relevant data in cybersecurity is very difficult. This is also the case for other very sensitive fields, such as financial data for example. Relevant and realistic data is very hard to produce synthetically, without compromising security and privacy of the systems/people involved, and thus can hardly be released publicly. Depending on the system of interest, it may be possible to record real data and then purge sensitive information (for example financial data can be anonymised to some extent), but this is very demanding and requires extreme thoroughness in order to not risk compromising any critical information. Also, depending on the data that was purged, some important information could be lost in the process, which can limit the scope of the research as well as the relevancy of the results. In addition, it is unlikely that any given dataset is sufficient for any given research work, unless this research work is specifically dedicated to this database, which would probably lead to limited results. Indeed, the outcome of research can be unpredictable and may yield solutions, algorithms or even principles that rely on information not available in the database. The best approach to conduct such research is to have one's own testbed, or to work closely with an industrial partner that can provide relevant data on demand, either from real systems or simulated systems (especially when real data is sensitive) that represent the real systems accurately enough.

Conclusion and perspectives

We started off this thesis with the objective of developing methods for cybersecurity that would allow learning behaviors and detect anomalies and attacks based off these learned behaviors. For this, we focused mainly on the RDT framework, which gave us a good basis to build methods that are applicable to a large number of signals by requiring very little knowledge. A large part of this thesis consisted in extending this framework, allowing us to work on one hand with non-Gaussian noise (Chapter 2), and on the other hand with parameter estimates that can then be incorporated in the RDT model (Chapter 3). Using these tools, we then introduced two methods (Chapter 4) and demonstrated their use on some real signals: first a fairly simple wavelet-based discontinuity detection methods, and then a change-in-mean detection approach, that allows segmenting time series in order to find the phases that compose a given signal. A noteworthy aspect of these methods is that they require few parameters to set up, and these parameters are easily understandable, and therefore fairly easily adjustable in the sense that one can tell how modifying one parameter will affect the performances.

Of course, these methods still remain fairly limited and require more work to build a proper detection system, but we think that they are nonetheless promising proofs of concept, which may lead to usable methods with further work. We notably need to implement a complete anomaly detection system using our change-in-mean detection method, since we only presented some ideas on how to use it for that purpose.

One critical aspect which absolutely requires more work is performance evaluation. We only demonstrated some interesting results visually, but we could never offer proper metrics that would allow comparing our approaches with others. As mentioned several times in Chapter 4, a large part of that is due to the inherent problems that come with working on real data, notably regarding labeling the phenomenons of interest. For instance, the SWaT dataset used throughout this thesis only labels timestamps during which attacks were conducted. However, our methods looked at other aspects of the signals (e.g.: phase changes) which are not necessarily labeled, and therefore we do not have a ground truth that we could use to evaluate any performance metrics. One possibility to circumvent this problem is evaluating it using artificially generated signals, which we could more easily control and therefore giving us the required information.

Another important aspect is moving from anomaly detection to attack detection. So far, we only presented ways to detect anomalies, which may or may not be attacks. For example, incidents such as a natural sensor failure would not be distinguishable from intentional damage caused by an attacker. Such classification would be helpful to allow better responses to incidents and cut down on potential costs due to misattribution of the source of an alarm.

In addition to all this, we also need to keep in mind that the methods developed here only apply to the physical signals of the system. For cyberphysical systems, it is also important to take the cyber part into account, notably through network traffic monitoring. However, the nature of network data is very different from physical signals, and attempting to apply similar methods to what we presented seems very difficult, any may very well not be possible. However, there are already existing intrusion detection systems for network data (e.g.: NIDS). One option to associate both of these detection aspects could be to correlate alarms from different systems in order to improve their detection performances, as increasing the diversity of detection systems may lead to a more resilient system overall, as long as no extra vulnerabilities are introduced.

Regarding the theoretical aspects presented in Chapters 2 and 3, there are also several points that need addressing. For the Generalized RDT (Chapter 2), we were able to prove that for any noise distribution

that presents some invariance, we are able to exhibit a test with a given level γ . We were then able to find a sufficient property for this test to be optimal. However we were unable to find any distribution satisfying this property. While this test might not be optimal depending on the distribution at hand, this result is nonetheless interesting, as ensuring a given level is an important property (see for example the GLRT test, which is not necessarily optimal, but still widely used). Our search for noise distributions for which we might have optimality was also focused on distributions invariant by rotation and using the L_2 -norm as the maximal invariant of choice. It would be worth looking at other groups, and seeing if we can get interesting results for other types of noise distributions.

The computation of the distribution of a maximal invariant [24, Chapter 7, Proposition 7.15] gives some hints to extend the theoretical framework of Chapter 2. Indeed, a non-degenerated normal distribution, which is invariant with respect to the orthogonal group, is generally defined via its density with respect to the Lebesgue measure. The Lebesgue measure is nothing else but the Haar measure associated with the translation group. Mathematically speaking — and that is exactly what Eaton does in [24, Chapter 7, Proposition 7.15] — nothing prevents us from considering random vectors that have normal densities with respect to a Haar measure other than the Lebesgue measure. This gives a hint to go beyond the standard Gaussian case and we can thus wonder to what extent the framework of Chapter 2 can be adapted to Haar measures other than Lebesgue's when defining the noise probability distribution. The questions that raise at this stage are: how the probability distribution of a random vector that has normal density with respect to a given Haar measure actually looks like; how the invariance properties of both the Gaussian density and the Haar measure combine; how relevant in practice are the invariance properties of such a random vector.

For the Asymptotic RDT (Chapter 3), we were able to take into account the estimation of parameters of the RDT model, namely the noise variance and the reference model, which is something that is commonly done in practice. We then presented some simulation results demonstrating how the estimation of the noise variance affects the performance of this test on a detection problem, and also suggest a method to adjust the tolerance to account for this estimation. A notable improvement that could be made is considering the case where the noise is no longer white, and can instead have any covariance matrix C , which is then estimated. We could also go further in our comparison simulations. The comparison that we chose between the RDT and NP tests is not necessarily the most informative one, due to their vastly different natures. It might be more interesting to replace the Neyman-Pearson with one that requires less a priori knowledge. The GLRT could be an interesting one to consider with an adapted problem with composite hypotheses, in which the RDT test might still be applicable. For example, we could consider testing the presence of a signal with unknown phase. Finally, the method that we showed to adjust the tolerance of the RDT test to take the noise estimation into account could use some refinement. As it stands, it is currently quite computationally heavy since it requires a large number of simulations, and it could be worth properly studying in order to find a more usable approach.

As already emphasized in Section 4.3, data is a crux in such a research. It is questionable whether having data beforehand might be detrimental, as it may narrow the research focus or bias the overall approach. Furthermore, in some applications, industrials may be reluctant to provide data, thence the importance of having direct access to the means to generate data easily and rapidly, according to the advances achieved during the research. It is not merely a question of being data-rich, it is also about having the relevant data that we may need at any given time.

In short, regarding the initial objectives of this research, we have a rather advanced theoretical framework to support statistical detection. The application of this framework to learning and anomaly detection cannot be fully assessed because of the lack of labeled data. Qualitatively, the experimental results show that our segmentation approach performs well. Clustering and anomaly/attack detection, are more intricate to evaluate, even empirically, beyond the intrinsic limitations of K-means.

In addition to the previous suggestions for a full detection method, we could also consider using the resulting segmentation in correlation with the state of the actuators present on the system, in order to verify whether they are coherent. For example, if we observe that a pump connected to a tank is active and

emptying that tank, the associated water level sensor is expected to observe a decrease in the water level in that tank. If we do not observe this decrease, this means there is an incoherence in the state of the system and we should trigger an alarm. In order to perform such detection, we would first need to discover the cause-and-effect relations that exist between sensors and actuators.

Overall, there is still work required to achieve a complete anomaly detection method based on this approach, but the work presented in this thesis provides a fairly solid foundation for robust detection methods, with a theoretical basis and encouraging results in our experiments so far.

Appendices

A | Deterministic Distortion Testing

We have seen in Chapter 1 how the RDT problem differs from usual testing problems. But this does not mean that they are completely unrelated. Indeed, with a simple change to the RDT problem, we can fall back on a more classic testing problem, that defines two families of probability distributions and that consists in deciding which of these families contains the distribution that generated the observation. This modified RDT problem is what we call the DDT (Deterministic Distortion Testing) problem, and is defined in Definition A.0.1 hereafter.

Definition A.0.1: DDT problem statement

$$\left\{ \begin{array}{l}
 \textbf{Data model:} \\
 \exists Y \in \mathcal{M}(\Omega, \mathbb{R}^d), \exists \theta \in \mathbb{R}^d, \exists X \in \mathcal{M}(\Omega, \mathbb{R}^d), \\
 \left\{ \begin{array}{l}
 (X \sim \mathcal{N}(0, C)) \\
 \wedge (Y = \theta + X) \\
 \wedge (\forall y \in \mathbb{R}^d, \exists \omega \in \Omega, y = Y(\omega))
 \end{array} \right. \\
 \textbf{Testing problem:} \\
 \text{Given one realization } y = Y(\omega) = \theta + X(\omega), \text{ determine whether:} \\
 \left\{ \begin{array}{l}
 \mathcal{H}_0: \nu_C(\theta - \theta_0) \leq \tau \\
 \text{or} \\
 \mathcal{H}_1: \nu_C(\theta - \theta_0) > \tau
 \end{array} \right. \\
 \text{with } \tau > 0 \text{ and } \theta_0 \in \mathbb{R}^d
 \end{array} \right. \quad (\text{A.1})$$

As a reminder, here is the RDT problem defined in Chapter 1:

Definition 1.2.2: RDT problem statement

$$\left\{ \begin{array}{l}
 \textbf{Data model:} \\
 \exists Y \in \mathcal{M}(\Omega, \mathbb{R}^d), \exists \Theta \in \mathcal{M}(\Omega, \mathbb{R}^d), \exists X \in \mathcal{M}(\Omega, \mathbb{R}^d), \\
 \left\{ \begin{array}{l}
 (X \sim \mathcal{N}(0, C) \text{ with } C \text{ a } d \times d \text{ definite positive matrix}) \\
 \wedge (\Theta \text{ and } X \text{ are independent}) \\
 \wedge (Y = \Theta + X) \\
 \wedge (\forall y \in \mathbb{R}^d, \exists \omega \in \Omega, y = Y(\omega))
 \end{array} \right. \\
 \textbf{Testing problem:} \\
 \text{Given one realization } y = Y(\omega) = \Theta(\omega) + X(\omega), \text{ determine whether:} \\
 \left\{ \begin{array}{l}
 \mathcal{H}_0: \nu_C(\Theta(\omega) - \theta_0) \leq \tau \\
 \text{or} \\
 \mathcal{H}_1: \nu_C(\Theta(\omega) - \theta_0) > \tau
 \end{array} \right. \\
 \text{with } \tau > 0 \text{ and } \theta_0 \in \mathbb{R}^d
 \end{array} \right. \quad (\text{1.12})$$

This problem can be seen as a particular case of the RDT problem, obtained by restricting the set of random variables that we consider to the random variables Θ such that $\Theta = \theta$ almost surely, for some deterministic $\theta \in \mathbb{R}^d$.

We can also see the DDT problem as a regular testing problem, where the family of distributions considered is the family of Gaussian distributions with mean $\theta \in \mathbb{R}^d$ and covariance matrix C .

The DDT problem here is shown in comparison to the RDT problem presented in Chapter 1, but the following is also applicable to the GRDT problem studied in Chapter 2.

For this DDT problem, we have an appropriate optimality criterion, the γ -MCP criterion, which is an adaptation of the γ -MCCP criterion presented in Definition 1.2.7 for the RDT problem and in Definition 2.2.5 for the GRDT problem.

Definition A.0.2: γ -MCP (Maximum Constant Power) test

Let $\mathcal{T}^*: \mathbb{R}^d \rightarrow \{0, 1\}$ be a test, and let $\gamma \in (0, 1)$.

The test \mathcal{T}^* is said to have level γ and Maximum Constant Power over \mathfrak{F} — and we simply say that \mathcal{T}^* is γ -MCP — if:

- (i) \mathcal{T}^* has level γ
- (ii) For every $\rho > \tau$, \mathcal{T}^* has constant power function on Υ_ρ
- (iii) For every $\rho > \tau$ and for any test \mathcal{T} with level γ and constant power function on Υ_ρ we have:

$$\forall \theta \in \Upsilon_\rho, \beta_{\mathcal{T}^*}(\theta) \geq \beta_{\mathcal{T}}(\theta)$$

We have the following result which links the γ -MCP and γ -MCCP criteria.

Lemma A.0.3: γ -MCP tests and γ -MCCP tests

A γ -MCCP test is γ -MCP. Conversely, a γ -MCP test with constant power function on every orbit $\Upsilon_\rho \in \mathfrak{F}$ is γ -MCCP.

Proof.

- Let \mathcal{T}^* be a γ -MCCP test. Since \mathcal{T}^* has constant conditional power function given any $\Theta \in \Upsilon_\rho$, we can deduce from Lemma 2.2.10 that \mathcal{T}^* has constant power function on every Υ_ρ for $\rho > \tau$. Let $\rho > \tau$ and let $\mathcal{T} \in \mathcal{K}_\gamma$ be a test with constant power function on Υ_ρ and let Θ be a random variable such that $\Theta \in \Upsilon_\rho$ almost surely. We have:

$$\begin{aligned} \forall \theta \in \Upsilon_\rho, \beta_{\mathcal{T}}(\theta) &= \mathbb{P}[\mathcal{T}(\Theta + X) = 1] && \text{from Lemma 2.2.7} \\ &= \mathbb{P}[\mathcal{T}(\Theta + X) = 1 \mid \Theta \in \Upsilon_\rho] && \text{because } \Theta \in \Upsilon_\rho \text{ almost surely} \\ &\leq \mathbb{P}[\mathcal{T}^*(\Theta + X) = 1 \mid \Theta \in \Upsilon_\rho] && \text{because } \mathcal{T}^* \text{ is } \gamma\text{-MCCP} \\ &\leq \beta_{\mathcal{T}^*}(\theta) \end{aligned}$$

Therefore \mathcal{T}^* is γ -MCP.

- Conversely, let \mathcal{T}^* be a γ -MCP test with constant power function on every orbit $\Upsilon_\rho \in \mathfrak{F}$. From Lemma 2.2.11, \mathcal{T}^* has constant conditional power function given any $\Theta \in \mathcal{M}(\Omega, \mathbb{R}^d)$ on Υ_ρ for $\mathbb{P}M(\Theta)^{-1}$ -almost every ρ . Let $\Theta \in \mathcal{M}(\Omega, \mathbb{R}^d)$, let $\rho > \tau$ and let $\mathcal{T} \in \mathcal{K}_\gamma$ be a test with constant conditional power function given $\Theta \in \Upsilon_\rho$. We have:

$$\begin{aligned} \forall \theta \in \Upsilon_\rho, \mathbb{P}[\mathcal{T}(\Theta + X) = 1 \mid \Theta \in \Upsilon_\rho] &= \beta_{\mathcal{T}}(\theta) \\ \text{and } \forall \theta \in \Upsilon_\rho, \mathbb{P}[\mathcal{T}^*(\Theta + X) = 1 \mid \Theta \in \Upsilon_\rho] &= \beta_{\mathcal{T}^*}(\theta) \end{aligned}$$

Since \mathcal{T}^* is γ -MCP, we have $\beta_{\mathcal{T}^*}(\theta) \geq \beta_{\mathcal{T}}(\theta)$ for any $\theta \in \Upsilon_\rho$, which in turn means that we have:

$$\mathbb{P}[\mathcal{T}^*(\Theta + X) = 1 \mid \Theta \in \Upsilon_\rho] \geq \mathbb{P}[\mathcal{T}(\Theta + X) = 1 \mid \Theta \in \Upsilon_\rho]$$

Therefore \mathcal{T}^* is γ -MCCP. □

As a consequence of this lemma, since the thresholding test $\mathcal{T}_{\lambda_\gamma(\tau)}$ is γ -MCCP (Theorem 1.2.10), it is also γ -MCP.

B Finding suitable families of TP-2 distributions for the GRDT problem

Contents

B.1 Notations and preliminary results.....	121
B.2 Calculations.....	122

In Chapter 2, we described a set of sufficient conditions under which we know that we are able to prove that the test $\mathcal{J}_{\lambda_\gamma(\tau)}$ is γ -MCCP. We need to use the euclidean norm $\|\cdot\|_2$ as the maximal invariant, and the probability density function of the noise X must be such that the family of densities $\{f_\rho, \rho \geq 0\}$ has monotone likelihood ratio (see Definition 1.1.5) where f_ρ is the density of $\|\theta + X\|_2$ with $\|\theta\|_2 = \rho$. In this appendix, we will describe the Calculations to show that this is the case when X follows a Gaussian distribution.

B.1 Notations and preliminary results

- S^{d-1} designates the unit $(d - 1)$ -sphere in \mathbb{R}^d centered on the origin:

$$S^{d-1} = \{x \in \mathbb{R}^d : \|x\|_2 = 1\}$$

Its surface is denoted $m(S^{d-1})$ (m denotes the surface Lebesgue measure on S^{d-1}) and we have:

$$m(S^{d-1}) = \frac{2\pi^{d/2}}{\Gamma(d/2)}$$

- σ denotes the normalized spherical measure on the $(d - 1)$ -sphere:

$$\forall A \in \mathcal{B}(S^{d-1}), \sigma(A) = \frac{m(A)}{m(S^{d-1})}$$

- For any $\rho > 0$, we denote by ρS^{d-1} the sphere centered on the origin with radius ρ , and we have:

$$\int_{\rho S^{d-1}} m(dx) = \rho^{d-1} m(S^{d-1})$$

$$\int_{\rho S^{d-1}} \sigma(dx) = \rho^{d-1}$$

Lemma B.1.1: Spherically invariant random vectors [52, Eq. (12)]

For any spherically invariant random vector $X \in \mathcal{M}(\Omega, \mathbb{R}^d)$ that has a probability density function f_X , the random variable $\|X\|_2$ also has a probability density function $f_{\|X\|_2}$ and we have:

$$\forall x \in \mathbb{R}^d, f_{\|X\|_2}(\|x\|_2) = m(S^{d-1})\|x\|_2^{d-1}f_X(x)$$

where $m(S^{d-1})$ is the surface of the unit $d - 1$ -sphere S^{d-1} :

$$m(S^{d-1}) = \frac{2\pi^{d/2}}{\Gamma(d/2)}$$

Lemma B.1.2: [53, Ch. 1, Eq (1.2)]

For any continuous function f defined on \mathbb{R} , we have:

$$\forall x \in \mathbb{R}^d, \int_{S^{d-1}} f(\langle x, y \rangle) m(dy) = m(S^{d-2}) \int_{-1}^{+1} f(t\|x\|_2) (1-t^2)^{\frac{d-3}{2}} dt$$

B.2 Calculations

Let $X \in \mathcal{M}(\Omega, \mathbb{R}^d)$ be a spherically invariant random vector with density f_X . Its density f_X can be written $\forall x \in \mathbb{R}^d, f_X(x) = h(\|x\|_2^2)$ with $h: \mathbb{R}^+ \rightarrow \mathbb{R}^+$ [25, Theorem 6.2.1]. Let $\rho > 0$ and let $\Xi_\rho \sim U(\rho S^{d-1})$ independent from X uniformly distributed on ρS^{d-1} :

$$\mathbb{P}\Xi_\rho^{-1} = \frac{1}{\rho^{d-1}}\sigma$$

Let $Y = \Xi_\rho + X$. Since X admits a density and is independent from Ξ_ρ , Y also admits a density $f_Y(\cdot, \rho) = f_X * \mathbb{P}\Xi_\rho^{-1}$. Let $y \in \mathbb{R}^d$. We have:

$$\begin{aligned} f_Y(y, \rho) &= \int f_X(y - \xi) \mathbb{P}\Xi_\rho^{-1}(d\xi) \\ &= \int_{\rho S^{d-1}} h(\|y - \xi\|_2^2) \frac{1}{\rho^{d-1}} \sigma(d\xi) \\ &= \int_{S^{d-1}} h(\|y - \rho\xi\|_2^2) \sigma(d\xi) \\ &= \int_{S^{d-1}} h(\|y\|_2^2 + \rho^2 - 2\rho\langle y, \xi \rangle) \sigma(d\xi) \\ &= \frac{1}{m(S^{d-1})} \int_{S^{d-1}} h(\|y\|_2^2 + \rho^2 - 2\rho\langle y, \xi \rangle) m(d\xi) \\ &= \frac{m(S^{d-2})}{m(S^{d-1})} \int_{-1}^{+1} h(\|y\|_2^2 + \rho^2 - 2t\rho\|y\|_2) (1-t^2)^{\frac{d-3}{2}} dt \\ &= \frac{\Gamma(\frac{d}{2})}{\Gamma(\frac{d-1}{2})\sqrt{\pi}} \int_{-1}^{+1} h(\|y\|_2^2 + \rho^2 - 2t\rho\|y\|_2) (1-t^2)^{\frac{d-3}{2}} dt \end{aligned}$$

Assume that h is infinitely differentiable and equal to its Taylor series expansion:

$$\forall a \in \mathbb{R}^+, \forall x \in \mathbb{R}^+, h(x) = \sum_{n=0}^{\infty} \frac{h^{(n)}(a)}{n!} (x-a)^n$$

We have for any $a \in \mathbb{R}^+$:

$$\begin{aligned} f_Y(y, \rho) &= \frac{\Gamma(\frac{d}{2})}{\Gamma(\frac{d-1}{2})\sqrt{\pi}} \int_{-1}^{+1} \sum_{n=0}^{\infty} \frac{h^{(n)}(a)}{n!} (\|y\|_2^2 + \rho^2 - 2t\rho\|y\|_2 - a)^n (1-t^2)^{\frac{d-3}{2}} dt \\ &= \frac{\Gamma(\frac{d}{2})}{\Gamma(\frac{d-1}{2})\sqrt{\pi}} \sum_{n=0}^{\infty} \frac{h^{(n)}(a)}{n!} \int_{-1}^{+1} (\|y\|_2^2 + \rho^2 - 2t\rho\|y\|_2 - a)^n (1-t^2)^{\frac{d-3}{2}} dt \end{aligned}$$

For $a = \|y\|_2^2 + \rho^2$, we get:

$$\begin{aligned} f_Y(y, \rho) &= \frac{\Gamma(\frac{d}{2})}{\Gamma(\frac{d-1}{2})\sqrt{\pi}} \sum_{n=0}^{\infty} \frac{h^{(n)}(\|y\|_2^2 + \rho^2)}{n!} \int_{-1}^{+1} (-2t\rho\|y\|_2)^n (1-t^2)^{\frac{d-3}{2}} dt \\ &= \frac{\Gamma(\frac{d}{2})}{\Gamma(\frac{d-1}{2})\sqrt{\pi}} \sum_{n=0}^{\infty} \frac{h^{(n)}(\|y\|_2^2 + \rho^2)}{n!} (-2\rho\|y\|_2)^n \int_{-1}^{+1} t^n (1-t^2)^{\frac{d-3}{2}} dt \end{aligned}$$

For every $n \in \mathbb{N}$, we have to compute the integral:

$$I_n = \int_{-1}^{+1} t^n (1-t^2)^{\frac{d-3}{2}} dt$$

If n is odd ($n = 2k + 1, k \in \mathbb{N}$), the function $t \mapsto t^n (1-t^2)^{\frac{d-3}{2}}$ is odd, which means that:

$$\forall k \in \mathbb{N}, I_{2k+1} = 0$$

If n is even ($n = 2k, k \in \mathbb{N}$), the function $t \mapsto t^n (1-t^2)^{\frac{d-3}{2}}$ is even and we have:

$$\forall k \in \mathbb{N}, I_{2k} = 2 \int_0^1 t^{2k} (1-t^2)^{\frac{d-3}{2}} dt$$

By applying the change of variable $u = t^2$ ($du = 2t dt$), we get:

$$\begin{aligned} \forall k \in \mathbb{N}, I_{2k} &= \int_0^1 u^{k-\frac{1}{2}} (1-u)^{\frac{d-3}{2}} du \\ &= \mathbf{B}\left(k + \frac{1}{2}, \frac{d-1}{2}\right) \end{aligned}$$

where \mathbf{B} is the Beta function defined by [54, Eq 6.2.1]:

$$\forall x \geq 0, \forall y \geq 0, \mathbf{B}(x, y) = \int_0^1 t^{x-1} (1-t)^{y-1} dt$$

A property of the Beta function is [54, Eq 6.2.2]:

$$\forall x \geq 0, \forall y \geq 0, \mathbf{B}(x, y) = \frac{\Gamma(x)\Gamma(y)}{\Gamma(x+y)}$$

Therefore:

$$\forall k \in \mathbb{N}, I_{2k} = \frac{\Gamma\left(k + \frac{1}{2}\right)\Gamma\left(\frac{d-1}{2}\right)}{\Gamma\left(k + \frac{d}{2}\right)}$$

$$\begin{aligned}
 f_Y(y, \rho) &= \frac{\Gamma\left(\frac{d}{2}\right)}{\sqrt{\pi}} \sum_{n=0}^{\infty} \frac{h^{(2n)}(\|y\|_2^2 + \rho^2)}{(2n)!} (2\rho\|y\|_2)^{2n} \frac{\Gamma\left(n + \frac{1}{2}\right)}{\Gamma\left(n + \frac{d}{2}\right)} \\
 &= \frac{\Gamma\left(\frac{d}{2}\right)}{\sqrt{\pi}} \sum_{n=0}^{\infty} \frac{h^{(2n)}(\|y\|_2^2 + \rho^2)}{(2n)!} \frac{(2\rho\|y\|_2)^{2n}}{\Gamma\left(n + \frac{d}{2}\right)} \frac{(2n)!}{4^n n!} \sqrt{\pi} \\
 &= \Gamma\left(\frac{d}{2}\right) \sum_{n=0}^{\infty} \frac{h^{(2n)}(\|y\|_2^2 + \rho^2)}{n!} \frac{(\rho\|y\|_2)^{2n}}{\Gamma\left(n + \frac{d}{2}\right)}
 \end{aligned}$$

Applying Lemma B.1.1, we get:

$$\begin{aligned}
 f_{\|Y\|_2}(\|y\|_2, \rho) &= m(S^{d-1})\|y\|_2^{d-1} f_Y(y, \rho) \\
 &= 2\pi^{d/2}\|y\|_2^{d-1} \sum_{n=0}^{\infty} \frac{h^{(2n)}(\|y\|_2^2 + \rho^2)}{n!} \frac{(\rho\|y\|_2)^{2n}}{\Gamma\left(n + \frac{d}{2}\right)}
 \end{aligned} \tag{B.1}$$

From here, it is difficult to get anything out of this expression without considering a specific probability distribution.

To give an example where we know we can get a result out of this calculation, we will consider the case where the random vector X follows a Gaussian distribution with zero mean and independent components: $X \sim \mathcal{N}(0, \sigma_0^2 I_d)$. We then have for any $x \in \mathbb{R}^d$ $f_X(x) = h(\|x\|_2^2)$ with h being defined by:

$$\forall x \geq 0, h(x) = (2\pi\sigma_0^2)^{-d/2} \exp\left(-\frac{x}{2\sigma_0^2}\right)$$

We also have:

$$\forall n \in \mathbb{N}, \forall x \geq 0, h^{(2n)}(x) = \pi^{-d/2} (2\sigma_0^2)^{-2n-d/2} \exp\left(-\frac{x}{2\sigma_0^2}\right)$$

In the gaussian case, we would expect $\|Y\|_2^2$ to follow a non-central χ^2 distribution (potentially up to a scale factor) with a certain non-centrality parameter depending on ρ . As a reminder, the probability density function for the non-central χ^2 distribution with $\nu \in \mathbb{N}$ degrees of freedom and non-centrality parameter $\lambda \geq 0$ is [54, Eq 26.4.25]:

$$\forall x \geq 0, f_{\chi^2_{\nu, \lambda}}(x) = \sum_{n=0}^{\infty} \frac{e^{-\lambda/2}}{n!} \left(\frac{\lambda}{2}\right)^n \frac{x^{n+\frac{\nu}{2}-1} e^{-x/2}}{2^{n+\frac{\nu}{2}} \Gamma\left(n + \frac{\nu}{2}\right)} \tag{B.2}$$

We have:

$$\begin{aligned}
 \forall y \in \mathbb{R}^d, f_{\|Y\|_2}(\|y\|_2, \rho) &= 2\pi^{d/2}\|y\|_2^{d-1} \sum_{n=0}^{\infty} \frac{h^{(2n)}(\|y\|_2^2 + \rho^2)}{n!} \frac{(\rho\|y\|_2)^{2n}}{\Gamma\left(n + \frac{d}{2}\right)} \\
 &= 2\pi^{d/2}\|y\|_2^{d-1} \sum_{n=0}^{\infty} \frac{\pi^{-d/2} (2\sigma_0^2)^{-2n-d/2}}{n!} \exp\left(-\frac{\|y\|_2^2 + \rho^2}{2\sigma_0^2}\right) \frac{(\rho\|y\|_2)^{2n}}{\Gamma\left(n + \frac{d}{2}\right)} \\
 &= \frac{2\|y\|_2^{d-1}}{2^{d/2}\sigma_0^d} \exp\left(-\frac{\|y\|_2^2 + \rho^2}{2\sigma_0^2}\right) \sum_{n=0}^{\infty} \frac{\rho^{2n}\|y\|_2^{2n}}{2^{2n}\sigma_0^{4n}\Gamma\left(n + \frac{d}{2}\right)n!}
 \end{aligned} \tag{B.3}$$

Let $Z = \frac{\|Y\|_2^2}{\sigma_0^2}$, and let $f_Z(\cdot, \rho)$ be its probability density function. Let $\rho_0 = \rho/\sigma_0$. We have:

$$\begin{aligned}
 \forall z \geq 0, f_Z(z, \rho) &= \frac{\sigma_0}{2\sqrt{z}} f_{\|Y\|_2}(\sigma_0\sqrt{z}, \rho) \\
 &= \frac{\sqrt{z}^{d-2}}{2^{d/2}} \exp\left(-\frac{z}{2} - \frac{\rho^2}{2\sigma_0^2}\right) \sum_{n=0}^{\infty} \frac{\rho^{2n} z^n}{2^{2n} \sigma_0^{2n} \Gamma\left(n + \frac{d}{2}\right) n!} \\
 &= \sum_{n=0}^{\infty} \frac{e^{-\rho_0^2/2}}{n!} \left(\frac{\rho_0^2}{2}\right)^n \frac{z^{n+\frac{d}{2}-1} e^{-z/2}}{2^{n+\frac{d}{2}} \Gamma\left(n + \frac{d}{2}\right)}
 \end{aligned} \tag{B.4}$$

Matching this expression with Eq. (B.2), we can see that Z follows a non-central χ^2 distribution with d degrees of freedom and non-centrality parameter $\rho_0^2 = (\rho/\sigma_0)^2$.

The fact that this family of densities has a monotone likelihood ratio is shown by Eaton in [24, Example A.1, p.468].

C | Uniform continuity of the Generalized Marcum function $Q_{d/2}$

This appendix is dedicated to proving that the Generalized Marcum function $Q_{d/2}$ is uniformly continuous, which we need in the proof of Theorem 3.2.2 to use the Portmanteau theorem.

As a reminder, the Generalized Marcum function $Q_{d/2}$ is defined by:

$$Q_{d/2}: [0, \infty) \times [0, \infty) \rightarrow \mathbb{R} \quad (C.1)$$

$$(\rho, \eta) \mapsto 1 - \mathbb{F}_{\chi_d^2(\rho^2)}(\eta^2)$$

where $\mathbb{F}_{\chi_d^2(\rho^2)}$ is the cumulative distribution function of the non-central χ^2 distribution with d degrees of freedom and non-centrality parameter ρ^2 . This means that for any Gaussian random variable $X \sim \mathcal{N}(\theta, I_d)$ with $\theta \in \mathbb{R}^d$, and for any $\lambda \geq 0$, we have $Q_{d/2}(\rho, \lambda) = \mathbb{P}[\|X\|_2^2 \geq \lambda^2]$, where $\rho = \|\theta\|_2$.

One expression of $Q_{d/2}$ is given by:

$$\forall \rho \geq 0, \forall \lambda \geq 0, Q_{d/2}(\rho, \lambda) = \frac{e^{-\rho^2/2}}{2^{\frac{d}{2}-1} \Gamma(d/2)} \int_{\lambda}^{\infty} e^{-t^2/2} t^{d-1} {}_0F_1\left(\frac{d}{2}; \frac{\rho^2 t^2}{4}\right) dt \quad (C.2)$$

where ${}_0F_1$ is a generalized hypergeometric function.

Proving that $Q_{d/2}$ is uniformly continuous means that we have to prove the following:

$$\forall \varepsilon > 0, \exists \alpha > 0, \forall (\rho_1, \lambda_1, \rho_2, \lambda_2) \in [0, \infty)^4,$$

$$\sqrt{(\rho_1 - \rho_2)^2 + (\lambda_1 - \lambda_2)^2} < \alpha \Rightarrow |Q_{d/2}(\rho_1, \lambda_1) - Q_{d/2}(\rho_2, \lambda_2)| \leq \varepsilon \quad (C.3)$$

In order to slightly simplify the following calculations, we will work on the function $R_{d/2} = 1 - Q_{d/2}$, which is completely equivalent. We have:

$$\forall \rho \geq 0, \forall \lambda \geq 0, R_{d/2}(\rho, \lambda) = \frac{e^{-\rho^2/2}}{2^{\frac{d}{2}-1} \Gamma(d/2)} \int_0^{\lambda} e^{-t^2/2} t^{d-1} {}_0F_1\left(\frac{d}{2}; \frac{\rho^2 t^2}{4}\right) dt \quad (C.4)$$

and for any Gaussian random variable $X \sim \mathcal{N}(\theta, I_d)$ with $\theta \in \mathbb{R}^d$ and $\rho = \|\theta\|_2$:

$$\forall \lambda \geq 0, R_{d/2}(\rho, \lambda) = \mathbb{P}[\|X\|_2^2 \leq \lambda^2] = \mathbb{F}_{\chi_d^2(\rho^2)}(\lambda^2) \quad (C.5)$$

To prove that $R_{d/2}$ is uniformly continuous on $[0, \infty) \times [0, \infty)$, we will prove that its partial derivatives are bounded.

Lemma C.0.1

$\frac{\partial R_{d/2}}{\partial \rho}$ is bounded on $[0, \infty) \times [0, \infty)$.

Proof. From Eq. (C.5), for any $\lambda \geq 0$, $\rho \geq 0$, and $\theta \in \mathbb{R}^d$ such that $\|\theta\|_2 = \rho$, we have:

$$R_{d/2}(\rho, \lambda) = \int_{B(0, \lambda)} \frac{1}{(2\pi)^{d/2}} \exp\left(-\frac{\|x - \theta\|_2^2}{2}\right) dx \quad (\text{C.6})$$

We will start by considering the case $d = 1$. Taking $\theta = \rho$, we get:

$$R_{1/2}(\rho, \lambda) = \frac{1}{\sqrt{2\pi}} \int_{-\lambda}^{\lambda} \exp\left(-\frac{(x - \rho)^2}{2}\right) dx \quad (\text{C.7})$$

The function $f: x, \rho \mapsto \exp\left(-\frac{(x - \rho)^2}{2}\right)$ defined on $[-\lambda, \lambda] \times [0, +\infty)$ is differentiable with respect to ρ and we have:

$$\frac{\partial f}{\partial \rho}(x, \rho) = (x - \rho) \exp\left(-\frac{(x - \rho)^2}{2}\right) \quad (\text{C.8})$$

$\frac{\partial f}{\partial \rho}$ is continuous on $[-\lambda, \lambda] \times [0, +\infty)$. Therefore $R_{1/2}$ is differentiable with respect to ρ and we have:

$$\frac{\partial R_{1/2}}{\partial \rho}(\rho, \lambda) = \frac{1}{\sqrt{2\pi}} \int_{-\lambda}^{\lambda} (x - \rho) \exp\left(-\frac{(x - \rho)^2}{2}\right) dx \quad (\text{C.9})$$

We can now bound this expression:

$$\begin{aligned} \left| \frac{\partial R_{1/2}}{\partial \rho}(\rho, \lambda) \right| &\leq \frac{1}{\sqrt{2\pi}} \int_{-\lambda}^{\lambda} |x - \rho| \exp\left(-\frac{(x - \rho)^2}{2}\right) dx \\ &\leq \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} |t| \exp\left(-\frac{t^2}{2}\right) dt \\ &\leq \sqrt{\frac{2}{\pi}} \end{aligned} \quad (\text{C.10})$$

Therefore $\frac{\partial R_{1/2}}{\partial \rho}$ is bounded.

We now consider the case when $d \geq 2$. For any $\rho \geq 0$, we can consider the vector $\theta = (\rho, 0, \dots, 0) \in \mathbb{R}^d$, which yields:

$$\begin{aligned} R_{d/2}(\rho, \lambda) &= \int \mathbb{1}_{B_d(0, \lambda)}(x_1, \dots, x_d) \frac{1}{(2\pi)^{d/2}} \exp\left(-\frac{(x_1 - \rho)^2}{2}\right) \exp\left(-\frac{1}{2} \sum_{i=2}^d x_i^2\right) dx_1 \dots dx_d \\ &= \frac{1}{(2\pi)^{d/2}} \int_{x_1 = -\lambda}^{\lambda} \exp\left(-\frac{(x_1 - \rho)^2}{2}\right) \left(\int_{x_2 \dots x_d} \mathbb{1}_{B_d(0, \lambda)}(x_1, \dots, x_d) \exp\left(-\frac{1}{2} \sum_{i=2}^d x_i^2\right) dx_2 \dots dx_d \right) dx_1 \end{aligned} \quad (\text{C.11})$$

For any $x_1 \in [-\lambda, \lambda]$, we have:

$$\begin{aligned} &\int_{x_2 \dots x_d} \mathbb{1}_{B_d(0, \lambda)}(x_1, \dots, x_d) \exp\left(-\frac{1}{2} \sum_{i=2}^d x_i^2\right) dx_2 \dots dx_d \\ &= \int_{x_2 \dots x_d} \mathbb{1}_{B_{d-1}(0, \sqrt{\lambda^2 - x_1^2})}(x_2, \dots, x_d) \exp\left(-\frac{1}{2} \sum_{i=2}^d x_i^2\right) dx_2 \dots dx_d \\ &= \int_{B_{d-1}(0, \sqrt{\lambda^2 - x_1^2})} \exp\left(-\frac{\|y\|_2^2}{2}\right) dy \\ &= (2\pi)^{\frac{d-1}{2}} \mathbb{F}_{\chi_{d-1}^2}(\lambda^2 - x_1^2) \end{aligned} \quad (\text{C.12})$$

Therefore:

$$R_{d/2}(\rho, \lambda) = \frac{1}{\sqrt{2\pi}} \int_{-\lambda}^{\lambda} \exp\left(-\frac{(x_1 - \rho)^2}{2}\right) \mathbb{F}_{\chi_{d-1}^2}(\lambda^2 - x_1^2) dx_1 \quad (\text{C.13})$$

Let f be the function defined for every $x_1 \in [-\lambda, \lambda]$ and $\rho \in [0, +\infty)$ by:

$$f(x_1, \rho) = \exp\left(-\frac{(x_1 - \rho)^2}{2}\right) \mathbb{F}_{\chi_{d-1}^2}(\lambda^2 - x_1^2) \quad (\text{C.14})$$

The function f is differentiable with respect to ρ and we have:

$$\frac{\partial f}{\partial \rho}(x_1, \rho) = (x_1 - \rho) \exp\left(-\frac{(x_1 - \rho)^2}{2}\right) \mathbb{F}_{\chi_{d-1}^2}(\lambda^2 - x_1^2) \quad (\text{C.15})$$

From this expression, we can see that $\frac{\partial f}{\partial \rho}$ is continuous on $[-\lambda, \lambda] \times [0, +\infty)$. Therefore $R_{d/2}$ is differentiable with respect to ρ , and for any $\lambda \geq 0$ and $\rho \geq 0$:

$$\frac{\partial R_{d/2}}{\partial \rho}(\rho, \lambda) = \frac{1}{\sqrt{2\pi}} \int_{-\lambda}^{\lambda} (x_1 - \rho) \exp\left(-\frac{(x_1 - \rho)^2}{2}\right) \mathbb{F}_{\chi_{d-1}^2}(\lambda^2 - x_1^2) dx_1 \quad (\text{C.16})$$

We can now conclude that $\frac{\partial R_{d/2}}{\partial \rho}(\rho, \lambda)$ is bounded:

$$\begin{aligned} \left| \frac{\partial R_{d/2}}{\partial \rho}(\rho, \lambda) \right| &\leq \frac{1}{\sqrt{2\pi}} \int_{-\lambda}^{\lambda} |x_1 - \rho| \exp\left(-\frac{(x_1 - \rho)^2}{2}\right) \mathbb{F}_{\chi_{d-1}^2}(\lambda^2 - x_1^2) dx_1 \\ &\leq \frac{1}{\sqrt{2\pi}} \int_{-\lambda}^{\lambda} |x_1 - \rho| \exp\left(-\frac{(x_1 - \rho)^2}{2}\right) dx_1 \\ &\leq \sqrt{\frac{2}{\pi}} \end{aligned} \quad (\text{C.17})$$

Therefore $\frac{\partial R_{d/2}}{\partial \rho}(\rho, \lambda)$ is bounded for any integer $d \geq 1$. \square

Before proving that the second partial derivative $\frac{\partial R_{d/2}}{\partial \lambda}$ is bounded, we introduce several useful inequalities. In the following, $I_\nu(x)$ is the modified Bessel function of the first kind.

Lemma C.0.2

We have the following inequalities:

$$\forall d \geq 2, \forall x \geq 0, I_{d/2-1}(x) \leq \frac{\sqrt{\pi}}{\Gamma\left(\frac{d-1}{2}\right)} \left(\frac{x}{2}\right)^{\frac{d}{2}-1} I_0(x) \quad (\text{C.18})$$

$$\forall x \geq 0, I_0(x) \leq \sqrt{\frac{\pi}{8}} \frac{e^x}{\sqrt{x}} \quad (\text{found in [55]}) \quad (\text{C.19})$$

Proof of Eq. (C.18). From [56, Eq 8.43.3], for any $\nu > -1/2$ and $x \geq 0$, we have:

$$I_\nu(x) = \frac{1}{\Gamma\left(\nu + \frac{1}{2}\right) \sqrt{\pi}} \left(\frac{x}{2}\right)^\nu \int_0^\pi \exp(x \cos \theta) \sin^{2\nu} \theta d\theta \quad (\text{C.20})$$

Thus, for any $d \geq 2$, taking $\nu = d/2 - 1$ yields:

$$I_{d/2-1}(x) = \frac{1}{\Gamma\left(\frac{d-1}{2}\right)\sqrt{\pi}} \left(\frac{x}{2}\right)^{d/2-1} \int_0^\pi \exp(x \cos \theta) \sin^{d-2} \theta d\theta \quad (\text{C.21})$$

For any $d \geq 2$ and any $\theta \in [0, \pi]$, we have $\sin^{d-2} \theta \leq 1$, hence:

$$\int_0^\pi \exp(x \cos \theta) \sin^{d-2} \theta d\theta \leq \pi I_0(x) \quad (\text{C.22})$$

Combining Eqs. (C.21) and (C.22) gives the desired result. \square

Proof of Eq. (C.19). We have:

$$\begin{aligned} I_0(x) &= \frac{1}{\pi} \int_0^\pi \exp(x \cos \theta) d\theta \\ &= \frac{2}{\pi} \int_0^{\pi/2} \exp(x \cos 2\theta) d\theta \\ &= \frac{2e^x}{\pi} \int_0^{\pi/2} \exp(-2x \sin^2 \theta) d\theta \end{aligned} \quad (\text{C.23})$$

For any $\theta \in [0, \pi/2]$, we have $\sin \theta \geq \frac{2\theta}{\pi}$. Therefore:

$$\begin{aligned} I_0(x) &\leq \frac{2e^x}{\pi} \int_0^{\pi/2} \exp\left(-\frac{8x}{\pi^2} \theta^2\right) d\theta \\ &\leq \frac{2e^x}{\pi} \int_0^{+\infty} \exp\left(-\frac{8x}{\pi^2} \theta^2\right) d\theta \\ &\leq \sqrt{\frac{\pi}{8}} \frac{e^x}{\sqrt{x}} \end{aligned} \quad (\text{C.24}) \quad \square$$

In addition to Eqs. (C.18) and (C.19), we will also make use of the following inequalities:

$$\forall x > 0, \forall \mu > \nu \geq 0, I_\nu(x) > I_\mu(x) \quad (\text{C.25})$$

$$\forall x \geq 0, I_0(x) \leq e^x \quad (\text{C.26})$$

The first inequality can be found in [57], and the second can be deduced immediately from [58, Eq 6.25].

Lemma C.0.3

$\frac{\partial R_{d/2}}{\partial \lambda}$ is bounded on $[0, +\infty) \times [0, +\infty)$.

Proof. We start by considering the case when $d = 1$. In this case, from Eq. (C.7), $R_{1/2}$ is differentiable with respect to λ and we have:

$$\forall \rho \geq 0, \forall \lambda \geq 0, \frac{\partial R_{1/2}}{\partial \lambda}(\rho, \lambda) = \frac{1}{\sqrt{2\pi}} \left(\exp\left(-\frac{(\lambda - \rho)^2}{2}\right) - \exp\left(-\frac{(\lambda + \rho)^2}{2}\right) \right) \quad (\text{C.27})$$

Therefore $\frac{\partial R_{1/2}}{\partial \lambda}$ is bounded by $\frac{1}{\sqrt{2\pi}}$.

We now consider the general case when $d \geq 2$.

From Eq. (C.4), we have:

$$\forall \rho \geq 0, \forall \lambda \geq 0, R_{d/2}(\rho, \lambda) = \frac{e^{-\rho^2/2}}{2^{\frac{d}{2}-1} \Gamma(d/2)} \int_0^\lambda e^{-t^2/2} t^{d-1} {}_0F_1\left(\frac{d}{2}; \frac{\rho^2 t^2}{4}\right) dt \quad (C.28)$$

$R_{d/2}$ is differentiable with respect to λ and we have:

$$\forall \rho \geq 0, \forall \lambda \geq 0, \frac{\partial R_{d/2}}{\partial \lambda}(\rho, \lambda) = \frac{e^{-\rho^2/2} e^{-\lambda^2/2}}{2^{\frac{d}{2}-1} \Gamma(d/2)} \lambda^{d-1} {}_0F_1\left(\frac{d}{2}; \frac{\rho^2 \lambda^2}{4}\right) \quad (C.29)$$

From [54, Eq. 9.6.47], we have:

$$\forall x \geq 0, \forall \nu \geq 0, I_\nu(x) = \frac{(x/2)^\nu}{\Gamma(\nu+1)} {}_0F_1\left(\nu+1; \frac{x^2}{4}\right) \quad (C.30)$$

Therefore for any $\rho > 0$ and $\lambda > 0$, and taking $\nu = d/2 - 1$, we get:

$$\frac{\partial R_{d/2}}{\partial \lambda}(\rho, \lambda) = e^{-\rho^2/2} e^{-\lambda^2/2} \frac{\lambda^{d-1}}{(\rho \lambda)^{d/2-1}} I_{d/2-1}(\rho \lambda) \quad (C.31)$$

We can already note here that we have $\frac{\partial R_{d/2}}{\partial \lambda} \geq 0$

Using Eqs. (C.18) and (C.26), we get:

$$\frac{\partial R_{d/2}}{\partial \lambda}(\rho, \lambda) \leq e^{-\frac{1}{2}(\rho-\lambda)^2} \lambda^{d-1} \frac{\sqrt{\pi}}{\Gamma\left(\frac{d-1}{2}\right) 2^{\frac{d}{2}-1}} \quad (C.32)$$

Similarly, using Eqs. (C.19) and (C.25), we get:

$$\frac{\partial R_{d/2}}{\partial \lambda}(\rho, \lambda) \leq e^{-\frac{1}{2}(\rho-\lambda)^2} \left(\frac{\lambda}{\rho}\right)^{\frac{d-1}{2}} \sqrt{\frac{\pi}{8}} \quad (C.33)$$

For any $\rho > 0$ and any $\lambda > 0$, let:

$$f_1(\rho, \lambda) = e^{-\frac{1}{2}(\rho-\lambda)^2} \lambda^{d-1} \frac{\sqrt{\pi}}{\Gamma\left(\frac{d-1}{2}\right) 2^{\frac{d}{2}-1}} \quad (C.34)$$

$$f_2(\rho, \lambda) = e^{-\frac{1}{2}(\rho-\lambda)^2} \left(\frac{\lambda}{\rho}\right)^{\frac{d-1}{2}} \sqrt{\frac{\pi}{8}} \quad (C.35)$$

and for any $a > 0$ and any $\lambda > 0$, let:

$$g_{1,a}(\lambda) = f_1\left(\frac{a}{\lambda}, \lambda\right) = e^a e^{-\frac{1}{2}\left(\frac{a^2}{\lambda^2} + \lambda^2\right)} \lambda^{d-1} \frac{\sqrt{\pi}}{\Gamma\left(\frac{d-1}{2}\right) 2^{\frac{d}{2}-1}} \quad (C.36)$$

$$g_{2,a}(\lambda) = f_2\left(\frac{a}{\lambda}, \lambda\right) = e^a e^{-\frac{1}{2}\left(\frac{a^2}{\lambda^2} + \lambda^2\right)} \lambda^{d-1} \sqrt{\frac{\pi}{8}} \frac{1}{a^{\frac{d-1}{2}}} \quad (C.37)$$

With these notations, we have:

$$\forall \rho > 0, \forall \lambda > 0, \frac{\partial R_{d/2}}{\partial \lambda}(\rho, \lambda) \leq g_{1,\rho\lambda}(\lambda) \quad \text{and} \quad \frac{\partial R_{d/2}}{\partial \lambda}(\rho, \lambda) \leq g_{2,\rho\lambda}(\lambda) \quad (C.38)$$

Studying the functions $g_{1,a}$ and $g_{2,a}$ shows that, for every $a > 0$, both functions present a maximum at $\lambda_a = \sqrt{\frac{d-1+\sqrt{(d-1)^2+4a^2}}{2}}$. Thus, for any $a > 0$, we have:

$$\begin{aligned} g_{1,a}(\lambda_a) &= e^{ae^{-\frac{4a^2+(d-1)^2+(d-1)\sqrt{(d-1)^2+4a^2}}{2(d-1+\sqrt{(d-1)^2+4a^2})}}} \left(\frac{d-1+\sqrt{(d-1)^2+4a^2}}{2} \right)^{\frac{d-1}{2}} \frac{\sqrt{\pi}}{\Gamma\left(\frac{d-1}{2}\right)2^{\frac{d}{2}-1}} \\ g_{2,a}(\lambda_a) &= e^{ae^{-\frac{4a^2+(d-1)^2+(d-1)\sqrt{(d-1)^2+4a^2}}{2(d-1+\sqrt{(d-1)^2+4a^2})}}} \left(\frac{d-1+\sqrt{(d-1)^2+4a^2}}{2} \right)^{\frac{d-1}{2}} \sqrt{\frac{\pi}{8}} \frac{1}{a^{\frac{d-1}{2}}} \end{aligned} \quad (\text{C.39})$$

These two expressions present the following limits:

$$\begin{aligned} g_{1,a}(\lambda_a) &\xrightarrow{a \rightarrow 0} e^{-\frac{d-1}{2}} (d-1)^{\frac{d-1}{2}} \frac{\sqrt{\pi}}{\Gamma\left(\frac{d-1}{2}\right)2^{\frac{d}{2}-1}} \\ g_{2,a}(\lambda_a) &\xrightarrow{a \rightarrow +\infty} \sqrt{\frac{\pi}{8}} \end{aligned} \quad (\text{C.40})$$

Let $a_0 > 0$.

1. On $(0, a_0]$, the function $a \mapsto g_{1,a}(\lambda_a)$ is continuous and admits a finite limit in 0. It is therefore bounded on that interval: there exists $C_1 \in \mathbb{R}$ such that, for any $a \in (0, a_0]$, we have $g_{1,a}(\lambda_a) < C_1$.
2. Similarly, on $[a_0, +\infty)$, the function $a \mapsto g_{2,a}(\lambda_a)$ is continuous and admits a finite limit in $+\infty$. It is therefore bounded on that interval: there exists $C_2 \in \mathbb{R}$ such that, for any $a \in [a_0, +\infty)$, we have $g_{2,a}(\lambda_a) < C_2$.

Therefore:

$$\begin{aligned} \forall a \in (0, a_0], \forall \lambda > 0, g_{1,a}(\lambda) &\leq C_1 \\ \forall a \in [a_0, +\infty), \forall \lambda > 0, g_{2,a}(\lambda) &\leq C_2 \end{aligned} \quad (\text{C.41})$$

From Eq. (C.38), we can then deduce that:

$$\forall \rho > 0, \forall \lambda > 0, \frac{\partial R_{d/2}}{\partial \lambda}(\rho, \lambda) \leq \max(C_1, C_2) \quad (\text{C.42})$$

We now have to study the cases $\rho = 0$ and $\lambda = 0$. For any $d \geq 2$, we have ${}_0F_1\left(\frac{d}{2}; 0\right) = 1$. Using Eq. (C.29), we get:

$$\forall \rho \geq 0, \frac{\partial R_{d/2}}{\partial \lambda}(\rho, 0) = 0 \quad (\text{C.43})$$

$$\forall \lambda \geq 0, \frac{\partial R_{d/2}}{\partial \lambda}(0, \lambda) = \frac{e^{-\lambda^2/2}}{2^{\frac{d}{2}-1}\Gamma(d/2)} \lambda^{d-1} \quad (\text{C.44})$$

The first expression is clearly bounded, and the second one is continuous on $[0, +\infty)$ and has a finite limit as λ grows to $+\infty$, thus it is also bounded.

Therefore $\frac{\partial R_{d/2}}{\partial \lambda}$ is bounded on $[0, +\infty) \times [0, +\infty)$. \square

Theorem C.0.4

The Generalized Marcum function $Q_{d/2}$ is uniformly continuous on $[0, +\infty) \times [0, +\infty)$.

Proof. From the previous lemmas, $\left| \frac{\partial R_{d/2}}{\partial \rho} \right|$ and $\left| \frac{\partial R_{d/2}}{\partial \lambda} \right|$ are bounded. Let $C \geq 0$ such that $\left| \frac{\partial R_{d/2}}{\partial \rho} \right| \leq C$ and $\left| \frac{\partial R_{d/2}}{\partial \lambda} \right| \leq C$. Let $\varepsilon > 0$, and let $\rho_1 \geq 0$, $\lambda_1 \geq 0$, $\rho_2 \geq 0$ and $\lambda_2 \geq 0$. We have:

$$|R_{d/2}(\rho_1, \lambda_1) - R_{d/2}(\rho_2, \lambda_2)| \leq |R_{d/2}(\rho_1, \lambda_1) - R_{d/2}(\rho_1, \lambda_2)| + |R_{d/2}(\rho_1, \lambda_2) - R_{d/2}(\rho_2, \lambda_2)| \quad (\text{C.45})$$

From the mean value theorem, there exist $\rho' \geq 0$ and $\lambda' \geq 0$ such that:

$$|R_{d/2}(\rho_1, \lambda_1) - R_{d/2}(\rho_1, \lambda_2)| = \left| \frac{\partial R_{d/2}}{\partial \lambda}(\rho_1, \lambda') \right| |\lambda_1 - \lambda_2| \quad (\text{C.46})$$

$$|R_{d/2}(\rho_1, \lambda_2) - R_{d/2}(\rho_2, \lambda_2)| = \left| \frac{\partial R_{d/2}}{\partial \rho}(\rho', \lambda_2) \right| |\rho_1 - \rho_2| \quad (\text{C.47})$$

The partial derivatives of $R_{d/2}$ are bounded by C , therefore:

$$\begin{aligned} |R_{d/2}(\rho_1, \lambda_1) - R_{d/2}(\rho_2, \lambda_2)| &\leq C(|\lambda_1 - \lambda_2| + |\rho_1 - \rho_2|) \\ &\leq C\sqrt{2}\sqrt{(\lambda_1 - \lambda_2)^2 + (\rho_1 - \rho_2)^2} \end{aligned} \quad (\text{C.48})$$

Thus choosing $\alpha = \frac{\varepsilon}{C\sqrt{2}}$ lets us conclude that $R_{d/2}$ is uniformly continuous on $[0, +\infty) \times [0, +\infty)$. \square

Acronyms

CUSUM	Cumulative Sum
DCS	Distributed Control System (Système numérique de contrôle-commande)
DDT	Deterministic Distortion Testing
DNR	Maximum-Distortion-to-Noise ratio
FIR	Finite Impulse Response
GLRT	Generalized Likelihood Ratio Test
ICS	Industrial Control System (Système de contrôle industriel)
IDS	Intrusion Detection System (Système de détection d'intrusions)
	HIDS Host-based Intrusion Detection System
	NIDS Network-based Intrusion Detection System
IoT	Internet of Things (Internet des objets)
MAD	Median Absolute Deviation
MCCP	Maximum Constant Conditional Power
MCP	Maximum Constant Power
NP	Neyman-Pearson
PLC	Programmable Logic Controller (Automate programmable industriel)
RDT	Random Distortion Testing
	ARDT Asymptotic Random Distortion Testing
	GRDT Generalized Random Distortion Testing
ROC	Receiver Operating Characteristic
SCADA	Supervisory Control And Data Acquisition (Système de contrôle et d'acquisition de données)
SNR	Signal-to-Noise ratio
SPRT	Sequential Probability Ratio Test
SWaT	Secure Water Treatment
UMP	Uniformly Most Powerful
UMPI	Uniformly Most Powerful Invariant

List of Publications

- [1] Guillaume Ansel et al. “Asymptotic Random Distortion Testing and Application to Change-in-Mean Detection”. In: ISIVC’2020. 2021-04.
- [2] Dominique Pastor and Guillaume Ansel. “Asymptotic Random Distortion Testing for Anomaly Detection”. In: ARCI’2021. 2021-02.

Bibliography

- [1] Nicolas Falliere, Liam O Murchu, and Eric Chien. *W32.Stuxnet Dossier*. 2011-02. URL: https://www.wired.com/images_blogs/threatlevel/2011/02/Symantec-Stuxnet-Update-Feb-2011.pdf (visited on 2020-05-25).
- [2] Jonathan Goh et al. “A Dataset to Support Research in the Design of Secure Water Treatment Systems”. In: *Critical Information Infrastructures Security*. Ed. by Grigore Havarneanu et al. Lecture Notes in Computer Science. Cham: Springer International Publishing, 2017. ISBN: 978-3-319-71368-7. DOI: 10.1007/978-3-319-71368-7_8.
- [3] Varun Chandola, Arindam Banerjee, and Vipin Kumar. “Anomaly Detection: A Survey”. In: *ACM Computing Surveys* 41.3 (2009-07-01). ISSN: 03600300. DOI: 10.1145/1541880.1541882. URL: <http://portal.acm.org/citation.cfm?doid=1541880.1541882> (visited on 2017-09-04).
- [4] *1999 DARPA Intrusion Detection Evaluation Dataset | MIT Lincoln Laboratory*. URL: <https://www.ll.mit.edu/r-d/datasets/1999-darpa-intrusion-detection-evaluation-dataset> (visited on 2020-06-16).
- [5] Mahbod Tavallaei et al. “A Detailed Analysis of the KDD CUP 99 Data Set”. In: *Computational Intelligence for Security and Defense Applications, 2009. CISDA 2009. IEEE Symposium On*. IEEE, 2009. URL: <http://ieeexplore.ieee.org/abstract/document/5356528/> (visited on 2017-09-19).
- [6] Tommy Morris. *Industrial Control System (ICS) Cyber Attack Datasets*. URL: <https://sites.google.com/a/uah.edu/tommy-morris-uah/ics-data-sets> (visited on 2020-06-18).
- [7] David I. Urbina et al. *Survey and New Directions for Physics-Based Attack Detection in Control Systems*. 2016-11-21. DOI: 10.6028/NIST.GCR.16-010.
- [8] K. J. Åström and P. Eykhoff. “System Identification—A Survey”. In: *Automatica* 7.2 (1971-03-01). ISSN: 0005-1098. DOI: 10.1016/0005-1098(71)90059-8. URL: <http://www.sciencedirect.com/science/article/pii/0005109871900598> (visited on 2018-09-19).
- [9] Lennart Ljung. *System Identification: Theory for the User*. 2nd ed. Prentice Hall Information and System Sciences Series. Upper Saddle River, NJ: Prentice Hall PTR, 1999. ISBN: 978-0-13-656695-3.
- [10] A. Nanduri and L. Sherry. “Anomaly Detection in Aircraft Data Using Recurrent Neural Networks (RNN)”. In: *2016 Integrated Communications Navigation and Surveillance (ICNS)*. 2016 Integrated Communications Navigation and Surveillance (ICNS). 2016-04. DOI: 10.1109/ICNSURV.2016.7486356.
- [11] Pavel Filonov, Andrey Lavrentyev, and Artem Vorontsov. *Multivariate Industrial Time Series with Cyber-Attack Simulation: Fault Detection Using an LSTM-Based Predictive Data Model*. 2016-12-20. arXiv: 1612.06676 [cs, stat]. URL: <http://arxiv.org/abs/1612.06676> (visited on 2018-07-16).
- [12] Pavel Filonov, Fedor Kitashov, and Andrey Lavrentyev. *RNN-Based Early Cyber-Attack Detection for the Tennessee Eastman Process*. 2017-09-07. arXiv: 1709.02232 [cs]. URL: <http://arxiv.org/abs/1709.02232> (visited on 2018-06-21).

- [13] Dmitry Shalyga, Pavel Filonov, and Andrey Lavrentyev. *Anomaly Detection for Water Treatment System Based on Neural Network with Automatic Architecture Optimization*. 2018-07-19. arXiv: 1807.07282 [cs, stat]. URL: <http://arxiv.org/abs/1807.07282> (visited on 2018-09-04).
- [14] Zachary C. Lipton. *The Mythos of Model Interpretability*. 2017-03-06. arXiv: 1606.03490 [cs, stat]. URL: <http://arxiv.org/abs/1606.03490> (visited on 2020-11-24).
- [15] Alejandro Barredo Arrieta et al. *Explainable Artificial Intelligence (XAI): Concepts, Taxonomies, Opportunities and Challenges toward Responsible AI*. Version 1. 2019-10-22. arXiv: 1910.10045 [cs]. URL: <http://arxiv.org/abs/1910.10045> (visited on 2020-11-24).
- [16] T. Kailath and H. V. Poor. "Detection of Stochastic Processes". In: *IEEE Transactions on Information Theory* 44.6 (1998-10). ISSN: 1557-9654. DOI: 10.1109/18.720538.
- [17] Dominique Pastor and Quang-Thang Nguyen. "Random Distortion Testing and Optimality of Thresholding Tests". In: *IEEE Transactions on Signal Processing* 61.16 (2013-08). ISSN: 1053-587X. DOI: 10.1109/TSP.2013.2265680.
- [18] Dominique Pastor and Francois-Xavier Socheleau. "Random Distortion Testing with Linear Measurements". In: *Signal Processing* 145 (2018-04). ISSN: 0165-1684. DOI: 10.1016/j.sigpro.2017.11.017. URL: <https://www.sciencedirect.com/science/article/pii/S0165168417304127> (visited on 2017-12-15).
- [19] Prashant Khanduri et al. "Sequential Random Distortion Testing of Non-Stationary Processes". In: *IEEE Transactions on Signal Processing* 67.21 (2019-11). ISSN: 1053-587X, 1941-0476. DOI: 10.1109/TSP.2019.2940124.
- [20] Prashant Khanduri et al. "Truncated Sequential Non-Parametric Hypothesis Testing Based on Random Distortion Testing". In: *IEEE Transactions on Signal Processing* 67.15 (2019-08). ISSN: 1941-0476. DOI: 10.1109/TSP.2019.2923140.
- [21] Bernard C. Levy. *Principles of Signal Detection and Parameter Estimation*. Boston, MA: Springer US, 2008. ISBN: 978-0-387-76544-0. DOI: 10.1007/978-0-387-76544-0. URL: <http://link.springer.com/10.1007/978-0-387-76544-0> (visited on 2020-07-01).
- [22] Jerzy Neyman and Egon S. Pearson. "On the Problem of the Most Efficient Tests of Statistical Hypotheses". In: *Philosophical Transactions of the Royal Society of London. Series A, Containing Papers of a Mathematical or Physical Character* 231 (1933). ISSN: 0264-3952. JSTOR: 91247.
- [23] Samuel Karlin and Herman Rubin. "The Theory of Decision Procedures for Distributions with Monotone Likelihood Ratio". In: *The Annals of Mathematical Statistics* 27.2 (1956-06). ISSN: 0003-4851. DOI: 10.1214/aoms/1177728259. URL: <http://projecteuclid.org/euclid.aoms/1177728259> (visited on 2019-09-26).
- [24] Morris L. Eaton. *Multivariate Statistics: A Vector Space Approach*. 2007-01-01.
- [25] Erich L. Lehmann and Joseph P. Romano. *Testing Statistical Hypotheses*. 3. ed. Springer Texts in Statistics. New York, NY: Springer, 2005. ISBN: 978-0-387-98864-1.
- [26] Abraham Wald. "Tests of Statistical Hypotheses Concerning Several Parameters When the Number of Observations Is Large". In: *Transactions of the American Mathematical Society* 54.3 (1943). ISSN: 0002-9947. DOI: 10.2307/1990256. JSTOR: 1990256.
- [27] C. Radhakrishna Rao. "Large Sample Tests of Statistical Hypotheses Concerning Several Parameters with Applications to Problems of Estimation". In: *Mathematical Proceedings of the Cambridge Philosophical Society* 44.1 (1948-01). ISSN: 1469-8064, 0305-0041. DOI: 10.1017/S0305004100023987. URL: <https://www.cambridge.org/core/journals/mathematical-proceedings-of-the-cambridge-philosophical-society/article/large-sample-tests-of-statistical-hypotheses-concerning-several-parameters-with-applications-to-problems-of-estimation/B83FAA6838A7E7D933EA3582C784ED06> (visited on 2020-10-13).

- [28] Robert F. Engle. “Chapter 13 Wald, Likelihood Ratio, and Lagrange Multiplier Tests in Econometrics”. In: *Handbook of Econometrics*. Vol. 2. Elsevier, 1984. ISBN: 978-0-444-86186-3. DOI: 10.1016/S1573-4412(84)02005-5. URL: <https://linkinghub.elsevier.com/retrieve/pii/S1573441284020055> (visited on 2020-10-16).
- [29] Yin Sun, Árpád Baricz, and Shidong Zhou. “On the Monotonicity, Log-Concavity, and Tight Bounds of the Generalized Marcum and Nuttall Q-Functions”. In: *IEEE Transactions on Information Theory* 56.3 (2010-03). ISSN: 1557-9654. DOI: 10.1109/TIT.2009.2039048.
- [30] Peter J. Huber and Elvezio Ronchetti. *Robust Statistics*. 2nd ed. Wiley Series in Probability and Statistics. Hoboken, N.J: Wiley, 2009. ISBN: 978-0-470-12990-6.
- [31] Serge Lang. *Algebra*. Vol. 1. New York: Springer, 2005. ISBN: 978-1-4612-6551-1.
- [32] Sabrina Bourmani. “Binary Decision for Observations with Unknown Distribution : An Optimal and Invariance-Based Framework”. These de doctorat. Ecole nationale supérieure Mines-Télécom Atlantique Bretagne Pays de la Loire, 2020-02-20. URL: <https://www.theses.fr/2020IMTA0173> (visited on 2020-06-26).
- [33] Patrick Billingsley. *Probability and Measure*. 3. ed. Wiley Series in Probability and Mathematical Statistics. New York, NY: Wiley, 1995. ISBN: 978-0-471-00710-4.
- [34] Robert J. Serfling. *Approximation Theorems of Mathematical Statistics*. Paperback ed. Wiley Series in Probability and Statistics. New York, NY: Wiley, 2002. ISBN: 978-0-471-21927-9.
- [35] Patrick Billingsley. *Convergence of Probability Measures*. 2. ed. Wiley Series in Probability and Statistics. New York, NY: Wiley, 1999. ISBN: 978-0-471-19745-4.
- [36] Mervin E. Muller. “A Note on a Method for Generating Points Uniformly on N -Dimensional Spheres”. In: *Communications of the ACM* 2.4 (1959-04-01). ISSN: 0001-0782. DOI: 10.1145/377939.377946. URL: <https://doi.org/10.1145/377939.377946> (visited on 2020-11-17).
- [37] Quang-Thang Nguyen. “Contributions to Statistical Signal Processing with Applications in Biomedical Engineering”. PhD thesis. Télécom Bretagne, Université de Bretagne Occidentale, 2012-11-23. URL: <https://tel.archives-ouvertes.fr/tel-00818320/document> (visited on 2017-12-24).
- [38] Ingrid Daubechies. *Ten Lectures on Wavelets*. CBMS-NSF Regional Conference Series in Applied Mathematics. Society for Industrial and Applied Mathematics, 1992-01-01. ISBN: 978-0-89871-274-2. DOI: 10.1137/1.9781611970104. URL: <https://epubs.siam.org/doi/book/10.1137/1.9781611970104> (visited on 2020-09-14).
- [39] T. Pham-Gia and T.L. Hung. “The Mean and Median Absolute Deviations”. In: *Mathematical and Computer Modelling* 34.7-8 (2001-10-01). ISSN: 0895-7177. DOI: 10.1016/S0895-7177(01)00109-1. URL: <https://www.sciencedirect.com/science/article/pii/S0895717701001091> (visited on 2020-09-14).
- [40] Michele Basseville and Igor V Nikiforov. *Detection of Abrupt Changes: Theory and Application*.
- [41] Charles Truong, Laurent Oudre, and Nicolas Vayatis. *A Review of Change Point Detection Methods*. 2018-01-02. arXiv: 1801.00718 [cs, stat]. URL: <http://arxiv.org/abs/1801.00718> (visited on 2018-10-23).
- [42] Samaneh Aminikhanghahi and Diane J. Cook. “A Survey of Methods for Time Series Change Point Detection”. In: *Knowledge and Information Systems* 51.2 (2017-05). ISSN: 0219-1377, 0219-3116. DOI: 10.1007/s10115-016-0987-z. URL: <http://link.springer.com/10.1007/s10115-016-0987-z> (visited on 2020-01-26).
- [43] W. A. Shewhart. *Economic Control of Quality of Manufactured Product*. Oxford, England: Van Nostrand, 1931.

- [44] E. S. Page. “Continuous Inspection Schemes”. In: *Biometrika* 41.1/2 (1954). ISSN: 0006-3444. DOI: 10.2307/2333009. JSTOR: 2333009.
- [45] Abraham Wald. “Sequential Tests of Statistical Hypotheses”. In: *The Annals of Mathematical Statistics* 16.2 (1945-06). ISSN: 0003-4851, 2168-8990. DOI: 10.1214/aoms/1177731118. URL: <https://projecteuclid.org/euclid.aoms/1177731118> (visited on 2018-10-05).
- [46] Dominique Pastor and Quang-Thang Nguyen. “Robust Statistical Process Control in Block-RDT Framework”. In: IEEE, 2015-04. ISBN: 978-1-4673-6997-8. DOI: 10.1109/ICASSP.2015.7178701. URL: <http://ieeexplore.ieee.org/document/7178701/> (visited on 2018-07-19).
- [47] Peter J. Rousseeuw and Christophe Croux. “Alternatives to the Median Absolute Deviation”. In: *Journal of the American Statistical Association* 88.424 (1993). ISSN: 0162-1459. DOI: 10.2307/2291267. JSTOR: 2291267.
- [48] Martin Ester, Hans-Peter Kriegel, and Xiaowei Xu. “A Density-Based Algorithm for Discovering Clusters in Large Spatial Databases with Noise”. In: ().
- [49] Mihael Ankerst et al. “OPTICS: Ordering Points To Identify the Clustering Structure”. In: ACM Press, 1999.
- [50] David Arthur and Sergei Vassilvitskii. “K-Means++: The Advantages of Careful Seeding”. In: *Proceedings of the Eighteenth Annual ACM-SIAM Symposium on Discrete Algorithms*. SODA '07. USA: Society for Industrial and Applied Mathematics, 2007-01-07. ISBN: 978-0-89871-624-5.
- [51] Fabian Pedregosa et al. “Scikit-Learn: Machine Learning in Python”. In: *Journal of Machine Learning Research* 12.85 (2011). URL: <http://jmlr.org/papers/v12/pedregosa11a.html> (visited on 2020-09-25).
- [52] Dominique Pastor, Roger Gay, and Albert Groenenboom. “A Sharp Upper Bound for the Probability of Error of the Likelihood Ratio Test for Detecting Signals in White Gaussian Noise”. In: *IEEE Transactions on Information Theory* 48.1 (2002-01). DOI: 10.1109/18.971751.
- [53] Fritz John. *Plane Waves and Spherical Means: Applied to Partial Differential Equations*. Reprinted. New York: Springer, 1981. ISBN: 978-0-387-90565-5.
- [54] Milton Abramowitz and Irene A. Stegun. *Handbook of Mathematical Functions with Formulas, Graphs, and Mathematical Tables*. 9. Dover print. Dover Books on Mathematics. Dover Publ, 1983. ISBN: 978-0-486-61272-0.
- [55] *Reference Request - Bounding a Modified Bessel Function of the First Kind*. Mathematics Stack Exchange. URL: <https://math.stackexchange.com/questions/893613/bounding-a-modified-bessel-function-of-the-first-kind> (visited on 2020-10-29).
- [56] I. S. Gradshteyn, I. M. Ryzhik, and Alan Jeffrey. *Table of Integrals, Series, and Products*. 7th ed. Amsterdam ; Boston: Academic Press, 2007. ISBN: 978-0-12-373637-6.
- [57] A. L. Jones. “An Extension of an Inequality Involving Modified Bessel Functions”. In: *Journal of Mathematics and Physics* 47.1-4 (1968). ISSN: 1467-9590. DOI: 10.1002/sapm1968471220. URL: <http://onlinelibrary.wiley.com/doi/abs/10.1002/sapm1968471220> (visited on 2020-10-29).
- [58] Yudell L. Luke. “Inequalities for Generalized Hypergeometric Functions”. In: *Journal of Approximation Theory* 5.1 (1972-01-01). ISSN: 0021-9045. DOI: 10.1016/0021-9045(72)90028-7. URL: <https://www.sciencedirect.com/science/article/pii/0021904572900287> (visited on 2020-10-29).

Titre : Contributions aux statistiques et méthodes d'apprentissage automatique pour la détection d'attaques basée sur la physique dans les systèmes industriels

Mots clés : Cybersécurité, détection d'anomalies, détection d'attaques, random distortion testing, apprentissage non-supervisé

Résumé : L'objectif de cette thèse est de développer de nouvelles méthodes de détection de cyberattaques basées sur des techniques d'apprentissage automatique. Ces travaux se sont concentrés sur l'étude de systèmes industriels, et plus spécifiquement la caractérisation du comportement normal des signaux physiques du système. Cette caractérisation permet ensuite de détecter des anomalies du système comme étant des déviations du comportement du système par rapport à ce modèle nominal appris. Les travaux effectués sont basés sur la théorie RDT (Random Distortion Testing), issue de la théorie statistique de la décision et permettant de donner un test optimal pour déterminer si une grandeur est suffisamment proche ou non d'un modèle donné, sans en connaître la distribution de probabilité. Cette théorie a été utilisée comme base afin de développer une méthode de détection de changement et a été appliquée avec succès sur des signaux réels, permettant également de contrôler le taux de fausses alarmes. Une extension de la théorie RDT a été développée afin de prendre en compte l'estimation du modèle et de la variance du bruit, supposés connus dans la théorie initiale. La méthode de détection de changements développée a ensuite été utilisée comme base pour caractériser les différentes phases des signaux du système via des méthodes de clustering. Les premiers essais d'une méthode complète d'apprentissage et de détection d'anomalies effectués sur des signaux réels offrent des résultats encourageants pour la détection d'attaques.

Title: Contributions to Statistics and Machine Learning for Physics-based Attack Detection in Industrial Systems

Keywords: Cybersecurity, anomaly detection, attack detection, random distortion testing, non-supervised learning

Abstract: The objective of this thesis is the development of new cyberattack detection methods based on machine learning techniques. This work focused on industrial systems, and consisted in characterizing the normal behavior of the physical signals of the system. This characterization then allows us to detect anomalies that occur in the system, these anomalies being deviations from the learned nominal behavior. The work presented here is based on the RDT (Random Distortion Testing) theory, which stems from statistical decision theory, and offers an optimal test to determine whether some phenomenon lies close enough to some model, without requiring any knowledge about its probability distribution. We used this theory as a basis to develop a change-detection method, which was then successfully applied to real signals, while also allowing control of the false-alarm rate. We developed an extension of the RDT framework to account for the estimation of the model and of the noise variance, which were assumed to be known in the initial theory. This change-detection method has then been used as a basis to characterize the different phases of the signals in the system using clustering methods. Our first attempts to develop a complete learning and anomaly-detection method have yielded encouraging results for attack detection on real signals.