

Hypothesis Testing via Euclidean Separation

Vincent Guigues *

Anatoli Juditsky †

Arkadi Nemirovski ‡

Abstract

We discuss an “operational” approach to testing convex composite hypotheses when the underlying distributions are heavy-tailed. It relies upon Euclidean separation of convex sets and can be seen as an extension of the approach to testing by convex optimization developed in [9, 13]. In particular, we show how one can construct quasi-optimal testing procedures for families of distributions which are majorated, in a certain precise sense, by a sub-spherical symmetric one and study the relationship between tests based on Euclidean separation and “potential-based tests.” We apply the promoted methodology in the problem of sequential detection and illustrate its practical implementation in an application to sequential detection of changes in the input of a dynamic system.

Résumé

Nous proposons une méthode “opérationnelle” pour le problème de tests d’hypothèses composites convexes lorsque les distributions sous-jacentes possèdent des queues lourdes. Elle s’appuie sur la séparation euclidienne des ensembles convexes et peut être vue comme une extension de la méthode développée dans [9, 13] pour l’étude des tests d’hypothèses reposant sur des techniques d’optimisation convexe. En particulier, nous montrons comment construire des tests quasi-optimaux pour des familles de distributions qui sont majorées, dans un sens précis, par une distribution symétrique quasi-sphérique et étudions la relation entre les tests basés sur la séparation euclidienne et les tests utilisant des potentiels. Nous appliquons la méthodologie proposée au problème de la détection séquentielle et décrivons sa mise en oeuvre pour la détection séquentielle de ruptures dans l’entrée d’un système dynamique.

Keywords: Hypothesis testing, nonparametric testing, composite hypothesis testing, statistical applications of convex optimization.

1 Introduction

The following important observation, attributed to H. Chernoff [6] (see also [3, 4]), was the starting point of our research.

Let X_1 and X_2 be two nonempty closed and convex sets, one of them being bounded, in \mathbf{R}^n . Suppose that, given a noisy observation

$$\omega = x + \xi \tag{1}$$

of the unknown signal $x \in X_1 \cup X_2$, where $\xi \sim \mathcal{N}(0, I_n)$ – the standard n -dimensional Gaussian vector, one wants to decide on the hypotheses $H_1 : x \in X_1$ vs. $H_2 : x \in X_2$. Then, assuming that X_1 and X_2 do not intersect (the decision problem is clearly unsolvable otherwise), optimal tests (with respect to different definitions of maximal risks) can be obtained using the following simple

*School of Applied Mathematics FGV/EMAp, 22 250-900 Rio de Janeiro, Brazil vincent.guigues@fgv.br

†LJK, Université Grenoble Alpes, 700 Avenue Centrale 38041 Domaine Universitaire de Saint-Martin-d’Hères, France, anatoli.juditsky@imag.fr

‡Georgia Institute of Technology, Atlanta, Georgia 30332, USA, nemirovs@isye.gatech.edu

Research of the first author was partially supported by an FGV grant, CNPq grants 307287/2013-0, and 401371/2014-0, FAPERJ grant E-26/201.599/2014. Research of the second author was partially supported by the LabEx PERSYVAL-Lab (ANR-11-LABX-0025) and CNPq grant 401371/2014-0. Research of the third author was partially supported by NSF grants CCF-1523768, CCF-1415498, and CNPq grant 401371/2014-0.

construction:

1) solve the following (convex optimization) Euclidean separation problem

$$\text{Opt} = \min_{x^1 \in X_1, x^2 \in X_2} \frac{1}{2} \|x^1 - x^2\|_2 \quad (2)$$

where $\|\cdot\|_2$ is the Euclidean distance.

2) Given an optimal solution (x_*^1, x_*^2) to (2), compute

$$h_* = \frac{x_*^1 - x_*^2}{\|x_*^1 - x_*^2\|_2} \text{ and } c_* = \frac{1}{2} h_*^T (x_*^1 + x_*^2).$$

Then the test \mathcal{T}_* which accepts H_1 if $h_*^T \omega - c_* \geq 0$ and accepts H_2 otherwise minimizes the maximal risk of testing – the maximal over $x \in X_1 \cup X_2$ probability of rejecting the true hypothesis – over the class of all (deterministic or randomized) tests.¹ Furthermore, the risk of \mathcal{T}_* is easily computable: the maximal probability of wrongly rejecting the true hypothesis is

$$\frac{1}{\sqrt{2\pi}} \int_{\text{Opt}}^{\infty} e^{-t^2/2} dt = \text{Erf}(\text{Opt}),$$

where $\text{Erf}(\cdot)$ is the standard error function.

This simple observation had important theoretical consequences (see [11, 15]). Surprisingly, its “practical implications” have been largely overseen. Indeed, when the problem (2) can be solved efficiently,² one can assemble pair-wise tests into multiple-testing procedures to build provably (nearly) optimal tests for a wide class of Gaussian decision problems (e.g., various detection problems [3, 4, 12, 7, 5]). Then in [9] the corresponding framework was extended to “good,” in a certain precise sense, parametric families of distributions, which include, aside from the Gaussian family, families of Poisson and discrete distributions, the “common denominator” of these developments being the fact that for these families the *near-optimal* (plain optimal, in the case of Gaussian family) tests can be built upon using *affine detectors* which can be found by convex optimization. Later, *affine and quadratic detectors* were studied in a more general setting in [13].

In this paper we extend [9], [13] studying the application of tests based on “Euclidean separation” to the problems where the distribution of the observation noise ξ has “heavy tails.” In particular, in Sections 2.2 and 2.3 we discuss the problem of testing convex hypotheses for *sub-spherical families* of distributions, which include Gaussian and Gaussian mixture distributions such as multivariate Student [14] and multivariate Laplace [8] distributions. We study the relationship of sub-spherical families with detector-based tests, and show how “good” detectors can be built for sub-spherical families, as well as for some other (e.g., sub-Gaussian) families of distributions in Section 2.4. Then in Section 3 we explain how tests based on Euclidean separation of pairs of convex hypotheses can be used to construct sequential change detection procedures. Finally, in Section 4 we present a numerical illustration of the proposed techniques: we implement sequential decision rules, developed in Sections 3.3 and 3.4 for a toy problem of detecting changes in the input of a dynamical system – changes in the trend of a simple time series.

2 Basic theory

2.1 Pairwise hypothesis testing: Situation and goal

The basic problem we intend to consider *in a nutshell* is as follows: we are given observation

$$\omega = x + \xi, \quad (3)$$

where $x \in \mathbf{R}^n$ is an unknown signal, and ξ is a random noise with probability density $p(\cdot)$, taken with respect to the Lebesgue measure, known to belong to some given family \mathcal{P} . Our basic goal is, given two nonempty

¹A test which accepts H_1 if $h_*^T \omega - c_* \geq \alpha$ and accepts H_2 otherwise, with a properly chosen α , is also minimax optimal in several other settings, e.g., Neyman-Pearson problem, etc.

²what is the case when sets X_1 and X_2 allow for a “computationally efficient description.” We refer the reader to [2] for precise definitions and details on efficient implementability. For the time being, it is sufficient to assume that (2) can be solved using CVX [10].

closed convex sets $X_\chi \subset \mathbf{R}^n$, $\chi = 1, 2$, with one of these sets bounded, to decide, via observation (3), on the hypotheses H_χ , $\chi = 1, 2$, with H_χ stating that the signal x underlying observation belongs to X_χ . In other words, we are to decide upon the families \mathcal{P}_1 and \mathcal{P}_2 of distributions, where \mathcal{P}_χ is the family of distributions of random vectors $x + \xi$, $\xi \sim p$ with $p \in \mathcal{P}$ and $x \in X_\chi$. For the sake of brevity, we shall simplify the description of the hypotheses H_χ to “ $x \in X_\chi$.”

Stationary repeated observations. We “embed” the just defined inference problem in the family of inference problems (\mathcal{S}_K) , $K = 1, 2, \dots$, where (\mathcal{S}_K) is the problem of deciding whether $x \in X_1$ or $x \in X_2$ via sample $\omega^K = (\omega_1, \dots, \omega_K)$ of K independent observations

$$\omega_k = x + \xi_k, k = 1, \dots, K, \quad (4)$$

with independent noises $\xi_k \sim p(\cdot) \in \mathcal{P}$; we refer to observations (4) with independent across k noises $\xi_k \sim p \in \mathcal{P}$ as to *stationary K -repeated observations*.

Semi-stationary repeated observations. Inference problem (\mathcal{S}_K) can be viewed as a special case of a more general inference problem $(\overline{\mathcal{S}}_K)$ as follows. Suppose that $\{X_1^k, X_2^k : 1 \leq k \leq K\}$ is a collection of nonempty convex and closed sets such that at least one of the set in every pair (X_1^k, X_2^k) is bounded. Let also $\mathcal{P}^1, \dots, \mathcal{P}^K$ be K given families of probability densities with respect to the Lebesgue measure. Suppose that the observation $\omega^K = (\omega_1, \dots, \omega_K)$ is given by

$$\omega_k = x_k + \xi_k, k = 1, \dots, K, \quad (5)$$

where $\{x_k\}_{k=1}^K$ is a deterministic sequence, $\xi_k \sim p_k$ are independent across k noises, and $\{p_k \in \mathcal{P}^k\}_{k=1}^K$ is a deterministic sequence.³ In the sequel, we refer to observations (5) satisfying the just imposed restrictions on $\{x_k\}_{k=1}^K$ and $\{\xi_k\}_{k=1}^K$, as to *semi-stationary K -repeated observations*. Our objective in problem $(\overline{\mathcal{S}}_K)$ is to decide, via the observations (5), on the hypotheses H_χ , $\chi = 1, 2$, with H_χ stating that $x_k \in X_\chi^k$ for all $k = 1, 2, \dots, K$.

A simple test for (\mathcal{S}_K) or $(\overline{\mathcal{S}}_K)$ is, by definition, a function $\mathcal{T}_K(\omega^K)$, $\omega^K = (\omega_1, \dots, \omega_K)$, taking values $\{1, 2\}$, with $\mathcal{T}_K(\omega^K) = \chi$ interpreted as “given observation ω^K , the test accepts H_χ and rejects the alternative.” We define the *partial risks* $\text{Risk}_{\chi\mathcal{S}}(\mathcal{T}_K|\mathcal{P}, X_1, X_2)$, $\chi = 1, 2$ of a test \mathcal{T}_K on the inference problem (\mathcal{S}_K) as

$$\begin{aligned} \text{Risk}_{1\mathcal{S}}(\mathcal{T}_K|\mathcal{P}, X_1, X_2) &= \sup_{x \in X_1} \sup_{p(\cdot) \in \mathcal{P}} \text{Prob}_{x,p}\{\mathcal{T}_K(\omega^K) = 2\}, \\ \text{Risk}_{2\mathcal{S}}(\mathcal{T}_K|\mathcal{P}, X_1, X_2) &= \sup_{x \in X_2} \sup_{p(\cdot) \in \mathcal{P}} \text{Prob}_{x,p}\{\mathcal{T}_K(\omega^K) = 1\}, \end{aligned}$$

where $\text{Prob}_{x,p}$ stands for the probability with respect to the distribution of observations (4). In other words, $\text{Risk}_{\chi\mathcal{S}}(\mathcal{T}_K|\mathcal{P}, X_1, X_2)$, is the worst-case probability for \mathcal{T}_K to reject H_χ when the hypothesis is true. The partial risks $\text{Risk}_{\chi\overline{\mathcal{S}}}(\mathcal{T}_K|[\mathcal{P}^k, X_1^k, X_2^k]_{k=1}^K)$, $\chi = 1, 2$, of a test \mathcal{T}_K on the inference problem $(\overline{\mathcal{S}}_K)$ are defined similarly, but now the supremum is taken with respect to *all* deterministic sequences $\{x_k \in X_\chi^k\}_{k=1}^K$ and $\{p_k \in \mathcal{P}^k\}_{k=1}^K$ participating in (5). Finally, the risks of \mathcal{T}_K on (\mathcal{S}_K) and $(\overline{\mathcal{S}}_K)$ are defined as

$$\begin{aligned} \text{Risk}_{\mathcal{S}}(\mathcal{T}_K|\mathcal{P}, X_1, X_2) &= \max_{\chi=1,2} \text{Risk}_{\chi\mathcal{S}}(\mathcal{T}_K|\mathcal{P}, X_1, X_2), \\ \text{Risk}_{\overline{\mathcal{S}}}(\mathcal{T}_K|[\mathcal{P}^k, X_1^k, X_2^k]_{k=1}^K) &= \max_{\chi=1,2} \text{Risk}_{\chi\overline{\mathcal{S}}}(\mathcal{T}_K|[\mathcal{P}^k, X_1^k, X_2^k]_{k=1}^K). \end{aligned}$$

Our intention is to investigate the performance of specific tests stemming from “Euclidean Separation” of X_1 and X_2 (or X_1^k and X_2^k) to be described in a while.

2.2 Sub-spherical families of distributions

We shall be primarily interested in *sub-spherical families* \mathcal{P} .

³One can easily verify (cf. [9, Section 3.1.2]) that the constructions and results to follow remain intact when the assumptions on observations (5) are weakened to the assumption that x_k and ξ_k are random, and the conditional distribution of ξ_k , given x_1, \dots, x_k and ξ_1, \dots, ξ_{k-1} , always belongs to \mathcal{P}^k .

2.2.1 Sub-spherical families of distributions: definition and basic examples

Definition 2.1 A. A sub-spherical family of distributions $\mathcal{P} = \mathcal{P}_\gamma^n$ on \mathbf{R}^n is specified by an even probability density $\gamma(\cdot)$ on the axis such that γ is positive in a neighbourhood of the origin. \mathcal{P}_γ^n is comprised of all probability densities $p(\cdot)$ on \mathbf{R}^n such that $p(\cdot)$ is even, and

$$\forall(e \in \mathbf{R}^n, \|e\|_2 = 1, \delta \geq 0) : \int_{e^T \xi \geq \delta} p(\xi) d\xi \leq P_\gamma(\delta) := \int_\delta^\infty \gamma(s) ds. \quad (6)$$

that is, p -probability mass of a half-space not containing a neighbourhood of the origin is upper-bounded by $P_\gamma(\delta)$, where δ is the distance from the origin to the half-space, and

$$P_\gamma(r) = \int_r^\infty \gamma(s) ds : \mathbf{R} \rightarrow [0, 1]. \quad (7)$$

B. A sub-spherical family $\mathcal{P} = \mathcal{P}_\gamma^n$ is called monotone, if it contains a cap – a spherically symmetric density $q(\xi) = f(\|\xi\|_2)$ where f is nonincreasing on the nonnegative axis and such that the induced by q density of the distribution of $e^T \xi$, $\|e\|_2 = 1$, is exactly $\gamma(\cdot)$. Note that whenever this is the case, $\gamma(\cdot)$ is nonincreasing on the nonnegative ray.

C. We call a function γ on the real axis nice, if γ is an even probability density which is continuous and is nonincreasing on the nonnegative ray. A sub-family \mathcal{P}^n of a sub-spherical family \mathcal{P}_γ^n is called completely monotone, if γ is nice, and for every $p(\cdot) \in \mathcal{P}^n$ and every $e \in \mathbf{R}^n$, $\|e\|_2 = 1$, the random scalar variable $e^T \xi$, $\xi \sim p$, has probability density $\gamma_{e,p}(\cdot)$, and this density is nice. Note that due to $\mathcal{P}^n \subset \mathcal{P}_\gamma^n$, it holds

$$\int_\delta^\infty \gamma_{e,p}(s) ds \leq \int_\delta^\infty \gamma(s) ds \quad \forall \delta \geq 0. \quad (8)$$

In the sequel, we simplify the notation \mathcal{P}_γ^n to \mathcal{P}_γ when the value of n is clear from the context.

Example: Gaussian scale mixtures. Consider the situation where

$$\xi \sim \sqrt{Z} \eta, \quad (9)$$

where Z is a scalar a.s. positive random variable with given probability distribution $P_Z(t) = \text{Prob}\{Z \leq t\}$, $t \geq 0$ such that $P_Z(0) = 0$, and $\eta \sim \mathcal{N}(0, \Theta)$ is a zero mean Gaussian n -dimensional random vector independent of Z with unknown a priori positive definite covariance matrix Θ which is known to be $\leq I_n$. We refer to Θ as to matrix parameter of the distribution of ξ . Given a unit vector e and $\delta \geq 0$, we have

$$\begin{aligned} \text{Prob}\{e^T \xi \geq \delta\} &= \int_{t>0} [\text{Prob}\{e^T \eta \geq t^{-1/2} \delta\}] dP_Z(t) \\ &= \int_{t>0} [\text{Prob}_{\zeta \sim \mathcal{N}(0, I_n)}\{e^T \Theta^{1/2} \zeta \geq t^{-1/2} \delta\}] dP_Z(t) \\ &= \int_{t>0} \text{Erf}\left(t^{-1/2} \delta / \sqrt{e^T \Theta e}\right) dP_Z(t) \leq \int_{t>0} \text{Erf}(t^{-1/2} \delta) dP_Z(t), \end{aligned}$$

(here $\text{Erf}(\cdot)$ is the standard error function), implying that

$$\text{Prob}\{e^T \xi \geq \delta\} = \int_\delta^\infty \gamma_{e,\Theta,P_Z}(s) ds \leq \int_\delta^\infty \gamma_{P_Z}(s) ds \quad \forall \delta \geq 0,$$

where

$$\begin{aligned} \gamma_{e,\Theta,P_Z}(s) &= \int_{t>0} \frac{1}{\sqrt{2\pi t e^T \Theta e}} \exp\left\{-\frac{s^2}{2t e^T \Theta e}\right\} dP_Z(t), \\ \gamma_{P_Z}(s) &= \int_{t>0} \frac{1}{\sqrt{2\pi t}} \exp\left\{-\frac{s^2}{2t}\right\} dP_Z(t), \quad -\infty < s < \infty. \end{aligned}$$

Clearly, $\gamma_{e,\Theta,P_Z}(s)$ and $\gamma_{P_Z}(s)$ are positive even probability densities nonincreasing on the nonnegative ray; these densities are continuous, provided

$$\int_{t>0} t^{-1/2} dP_Z(t) < \infty. \quad (10)$$

Thus, the family of distributions of random vectors (9) with Z and η as explained above (we refer to these distributions as to *Gaussian scale mixtures*) is contained in the sub-spherical family $\mathcal{P}_{\gamma_{P_Z}}$. The latter family clearly is monotone, the cap being the probability density of $\sqrt{Z}\eta$ with independent $Z \sim P_Z$ and $\eta \sim \mathcal{N}(0, I_n)$. Besides this, in the case of (10) the family $\mathcal{P}_{\gamma_{P_Z}}$ is completely monotone.

A standard example of a Gaussian scale mixture is given by *n-variate t-distributions* $t_n(q, \Theta)$, $\Theta \preceq I_n$ (multivariate Student distributions with q degrees of freedom, see [14] and references therein). Here $t_n(q, \Theta)$ is, by definition, the distribution of the random vector $\xi = \sqrt{Z}\eta$ with $Z = q/\zeta$, where ζ is the independent of $\eta \sim \mathcal{N}(0, \Theta)$ random variable following χ^2 -distribution with q degrees of freedom. One can easily see that all one-dimensional projections $e^T \xi$, $\|e\|_2 = 1$, of $\xi \sim t_n(q, I_n)$ are random variables with univariate t_q -distribution, implying that the multidimensional densities in question form a completely monotone sub-family of the sub-spherical family \mathcal{P}_{γ_S} where γ_S is the density of Student's t_q distribution with q degrees of freedom.

Another example of scheme (9) is the *n-variate Laplace distributions* $\mathcal{L}_n(\lambda, \Theta)$, $\Theta \preceq I_n$, where Z is exponentially distributed with parameter λ . In this case all one-dimensional projections $e^T \xi$, $\|e\|_2 = 1$, of $\xi \sim \mathcal{L}_n(\lambda, I_n)$, obey the Laplace distribution with parameter λ , whence the distributions in question form a completely monotone sub-family of the sub-spherical family $\mathcal{P}_{\gamma_{\mathcal{L}}}$, where $\gamma_{\mathcal{L}}$ is the Laplace density

$$\gamma_{\mathcal{L}}(s) = (2\lambda)^{-1} e^{-|s|/\lambda}, \quad s \in \mathbf{R}. \quad (11)$$

Finally, with Z taking value 1 with probability 1, scheme (9) describes Gaussian distributions with zero mean and covariance matrices $\preceq I_n$; all these distributions form a completely monotone sub-family of the sub-spherical family \mathcal{P}_{γ_G} , where γ_G is the standard univariate Gaussian density:

$$\gamma_G(s) = \frac{1}{\sqrt{2\pi}} e^{-s^2/2}.$$

2.2.2 “Calculus” of sub-spherical families of distributions

Sub-spherical families of distributions and their completely monotone subfamilies admit a kind of “calculus” with the basic rules which follow.

The following two facts are immediate:

Proposition 2.1 *Let \mathcal{P}_γ^n be a sub-spherical family of distributions, and let $x \mapsto Qx : \mathbf{R}^n \rightarrow \mathbf{R}^m$ be an onto mapping satisfying $QQ^T \preceq I_m$. Whenever $p(\cdot) \in \mathcal{P}_\gamma^n$, the distribution of the random vector $Q\xi$, $\xi \sim p$, belongs to \mathcal{P}_γ^m . Moreover, if $QQ^T = I_m$ and \mathcal{P}_γ^n has a cap q , then \mathcal{P}_γ^m has a cap as well; this cap is the density of the random vector $Q\xi$, $\xi \sim q$.*

Proposition 2.2 *A sub-spherical family of distributions is closed with respect to taking convex combinations of its members. Besides this, the union $\mathcal{P}_{\gamma_1} \cup \mathcal{P}_{\gamma_2}$ of two sub-spherical families of distributions is contained in the sub-spherical family \mathcal{P}_γ with*

$$\gamma(-s) = \gamma(s) = -\frac{d}{ds} \max_{i=1,2} \int_s^\infty \gamma_i(r) dr, \quad s \geq 0.$$

Complete monotonicity is preserved by taking sums. The precise statement is as follows:

Proposition 2.3 *Let μ and ν be nice, let \mathcal{P}^n be a subfamily of the sub-spherical family \mathcal{P}_μ^n , and let \mathcal{P}^m be a completely monotone subfamily of the sub-spherical family \mathcal{P}_ν^m . Given $r \times n$ matrix A , $r \times m$ matrix B and positive definite matrix Θ such that*

$$\Theta^2 \succeq AA^T, \quad \Theta^2 \succeq BB^T, \quad AA^T + BB^T \succ 0, \quad (12)$$

consider random vectors of the form

$$\xi = \Theta^{-1}[A\eta + B\zeta], \quad (13)$$

where $\eta \sim p(\cdot) \in \mathcal{P}^n$ and $\zeta \sim q(\cdot) \in \mathcal{P}^m$ are independent. Let also

$$\gamma(s) = \int_{-\infty}^{\infty} \mu(s-r)\nu(r)dr. \quad (14)$$

Then

(i) $\gamma(\cdot)$ is nice, and for every $e \in \mathbf{R}^r$, $\|e\|_2 = 1$, and every $\delta \geq 0$ one has

$$\text{Prob}\{e^T \xi \geq \delta\} \leq \int_{\delta}^{\infty} \gamma(s)ds. \quad (15)$$

Besides this, the scalar random variable $e^T \xi$ possesses symmetric density, which combines with (15) to imply that the distribution of ξ belongs to the sub-spherical family \mathcal{P}_{γ}^r .

(ii) If, in addition to the above assumptions, \mathcal{P}^n is completely monotone, then the family of distributions of random variables (13) induced by $\eta \sim p \in \mathcal{P}^n$ and $\zeta \sim q \in \mathcal{P}^m$ is a completely monotone subfamily of the sub-spherical family \mathcal{P}_{γ}^r .

For a proof, see Section A.1.

As an immediate consequence of Proposition 2.3, we get the following

Corollary 2.1 For $1 \leq i \leq I < \infty$, let μ_i be nice functions on the axis, and let \mathcal{P}^{n_i} be subfamilies of the sub-spherical families $\mathcal{P}_{\mu_i}^{n_i}$ such that at least $I - 1$ of these subfamilies are completely monotone. Given $r \times n_i$ matrices A_i such that $A_i A_i^T + A_j A_j^T \succ 0$ whenever $i \neq j$, let $\Theta \succ 0$ be such that

$$\Theta^2 \succeq A_i A_i^T, \quad 1 \leq i \leq I.$$

Let \mathcal{P}^r be the family of probability distributions of random vectors of the form

$$\xi = \Theta^{-1} \sum_{i=1}^I A_i \eta_i$$

where η_1, \dots, η_I are independent of each other and such that $\eta_i \sim p_i$ for some $p_i \in \mathcal{P}^{n_i}$. Then \mathcal{P}^r is contained in the sub-spherical family \mathcal{P}_{γ}^r , where

$$\gamma = \mu_1 \star \mu_2 \star \dots \star \mu_I$$

is nice. If all \mathcal{P}^{n_i} are completely monotone, then so is \mathcal{P}^r .

2.3 Euclidean separation and associated tests

In the sequel, we fix the entities \mathcal{P} , X_1 , X_2 , H_1 , H_2 introduced in the beginning of Section 2.1 and assume that $X_1 \cap X_2 = \emptyset$ (otherwise no test can decide on H_1 vs. H_2 with risk $< 1/2$).

2.3.1 Single-observation Euclidean separation test

Consider the optimization problem

$$\text{Opt} = \min_{x^1 \in X_1, x^2 \in X_2} \|x^1 - x^2\|_2, \quad (16)$$

and let x_*^1, x_*^2 form an optimal solution to the problem (since both X_1, X_2 are nonempty, closed and convex, and one of the sets is bounded, an optimal solution does exist). We set

$$s_*(\omega) = h_*^T \omega - c_*, \quad h_* = \frac{[x_*^1 - x_*^2]}{\|x_*^1 - x_*^2\|_2}, \quad c_* = \frac{1}{2} \left[\min_{x \in X_1} h_*^T x + \max_{x \in X_2} h_*^T x \right] = \frac{1}{2} h_*^T (x_*^1 + x_*^2). \quad (17)$$

Note that while (16) may have many optimal solutions, the vector h_* and the real c_* are uniquely defined by X_χ , $\chi = 1, 2$. The affine function $s_*(\cdot)$ possesses the following properties:

$$\begin{aligned} x \in X_1 &\Rightarrow s_*(x) \geq s_*(x_*^1) = \delta, \quad \delta := \frac{1}{2}\|x_*^1 - x_*^2\|_2, \\ x \in X_2 &\Rightarrow s_*(x) \leq s_*(x_*^2) = -\delta. \end{aligned} \quad (18)$$

We can associate with h_* and c_* the *Euclidean separation test* \mathcal{T}_1 which, given observation ω (3) with x known to belong to $X_1 \cup X_2$, accepts the hypothesis $H_1 : x \in X_1$ and rejects the hypothesis $H_2 : x \in X_2$ when $s_*(\omega) \geq 0$, and accepts H_2 and rejects H_1 otherwise.

Let us make the following immediate observation:

Proposition 2.4 *In the situation of Section 2.1 and in the notation from (17), (18), let $\mathcal{P} = \mathcal{P}_\gamma$ be a sub-spherical family of distributions, and let $\alpha_1 \geq 0$ and $\alpha_2 \geq 0$ be such that $\alpha_1 + \alpha_2 \leq 2\delta$. Whenever $x \in X_1$ and $p \in \mathcal{P}$, the p -probability of the event $\{\xi : s_*(x + \xi) < \frac{1}{2}(\alpha_2 - \alpha_1)\}$ is at most $P_\gamma(\alpha_1)$, and when $x \in X_2$, the p -probability of the event $\{\xi : s_*(x + \xi) \geq \frac{1}{2}(\alpha_2 - \alpha_1)\}$ is at most $P_\gamma(\alpha_2)$. As a result, risks of the test \mathcal{T}_1 satisfy*

$$\begin{aligned} \text{Risk}_{1S}(\mathcal{T}_1|\mathcal{P}, X_1, X_2) &\leq P_\gamma(\alpha_1), \\ \text{Risk}_{2S}(\mathcal{T}_1|\mathcal{P}, X_1, X_2) &\leq P_\gamma(\alpha_2). \end{aligned}$$

In particular, when $\alpha_1 = \alpha_2 = \delta$, the risk $\text{Risk}_S(\mathcal{T}_1|\mathcal{P}, X_1, X_2)$ of the test is at most

$$\varepsilon_* = \varepsilon_*(\delta|\gamma) := P_\gamma(\delta) = \int_{\delta}^{\infty} \gamma(s) ds. \quad (19)$$

For a proof, see Section A.2.

Remark 2.1 *In the situation of Proposition 2.4, let the sub-spherical family \mathcal{P}_γ be monotone. Then the test \mathcal{T}_1 described in the proposition has the minimal risk among all single-observation tests deciding on H_1 vs. H_2 .*

Indeed, denoting by $q(\cdot)$ the cap of \mathcal{P}_γ , and by $p_1(\cdot)$ and $p_2(\cdot)$ the densities of random vectors $x_*^1 + \xi$, $x_*^2 + \xi$, $\xi \sim q(\cdot)$, we clearly have

$$\int \min[p_1(\omega), p_2(\omega)] d\omega = 2\varepsilon_*,$$

implying by the Neyman-Pearson lemma that the risk of any test deciding on two simple hypotheses $x = x_*^1$, $x = x_*^2$ via a single observation (3) is at least ε_* .

2.3.2 Majority tests based on Euclidean separation

Let K be a positive integer, and let $\mathcal{P} = \mathcal{P}_\gamma$ be a sub-spherical family of distributions. The K -observation majority test $\mathcal{T}_K^{\text{maj}}$ for problem (\mathcal{S}_K) works as follows: given observations ω_k , $k = 1, \dots, K$, see (4), the test accepts H_1 and rejects H_2 when $s_*(\omega_k) \geq 0$ for at least $K/2$ values of k , and accepts H_2 and rejects H_1 otherwise. From Proposition 2.4 it follows that the risk of $\mathcal{T}_K^{\text{maj}}$ satisfies the bound

$$\text{Risk}_S(\mathcal{T}_K^{\text{maj}}|\mathcal{P}, X_1, X_2) \leq \sum_{K \geq k \geq K/2} \binom{K}{k} \varepsilon_*^k (1 - \varepsilon_*)^{K-k}. \quad (20)$$

Let us assume that in the problem $(\overline{\mathcal{S}}_K)$ the families $\mathcal{P}^1, \dots, \mathcal{P}^K$ in the definition of $(\overline{\mathcal{S}}_K)$ are sub-spherical families of distributions, and that the sets X_1^k and X_2^k do not intersect for $k = 1, \dots, K$. The majority test $\mathcal{T}_K^{\text{maj}}$ can be easily modified to become applicable to problem $(\overline{\mathcal{S}}_K)$. Let the affine function $s_k(\cdot)$ and positive real δ_k be the entities $s_*(\cdot), \delta$ associated, via (16) – (18), with \mathcal{P}^k in the role of \mathcal{P} and X_χ^k in the role of X_χ , $\chi = 1, 2$. The K -observation majority test $\mathcal{T}_K^{\text{maj}}$ for the problem $(\overline{\mathcal{S}}_K)$ given observations ω_k , $k = 1, \dots, K$, see (5), accepts H_1 and rejects H_2 when $s_k(\omega_k) \geq 0$ for at least $K/2$ values of k , and accepts H_2 and rejects H_1 otherwise. Let now γ_k be the density underlying the sub-spherical family \mathcal{P}^k . When applying Proposition 2.4 to $\mathcal{P} = \mathcal{P}^k$, $X_1 = X_1^k$, and $X_2 = X_2^k$, we conclude that the risk $\text{Risk}_S(\mathcal{T}_1|\mathcal{P}^k, X_1^k, X_2^k)$ of the test \mathcal{T}_1 in the problem of deciding, given a single observation ω_k , upon the hypotheses $H_1^k : x_k \in X_1^k$ vs $H_2^k : x_k \in X_2^k$, does

not exceed $\epsilon_k(\delta_k) = \int_{\delta_k}^{\infty} \gamma_k(s) ds < \frac{1}{2}$. We conclude that the risk of the majority test $\mathcal{T}_K^{\text{maj}}$ in the problem $(\overline{\mathcal{S}}_K)$ satisfies the bound ¹

$$\text{Risk}_{\overline{\mathcal{S}}}(\mathcal{T}_K^{\text{maj}} | [\mathcal{P}^k, X_1^k, X_2^k]_{k=1}^K) \leq \sum_{K/2 \leq k \leq K} p_{k|K}, \quad (21)$$

where $p_{k|K} = p_{k|K}(\epsilon_1(\delta_1), \dots, \epsilon_K(\delta_K))$ is the probability of k successes in the first K non-stationary independent Bernoulli trials with the probability $\epsilon_k(\delta_k)$ of success in k -th trial. Observe that $p_{j|k}$, $0 \leq j \leq k$, $1 \leq k \leq K$, satisfy the recursion:

$$p_{j|k} = (1 - \epsilon_k)p_{j|k-1} + \epsilon_k p_{j-1|k-1}, \text{ with } p_{0|0} = 1, \quad (22)$$

and where, by convention, $p_{-1|k} = 0$, $k = 0, \dots, K-1$. Recursion (22) allows to compute all the quantities $p_{j|k}$, $0 \leq j \leq k \leq K$, and thus the right hand side of (21), in $O(K^2)$ arithmetic operations.

2.4 Potential-based tests with Euclidean separation

2.4.1 Potentials and potential-based tests

Definition 2.2 A *potential* is an odd and nondecreasing Borel real-valued function $\eta(\cdot)$ on the axis. Given a family \mathcal{P} of probability densities on \mathbf{R}^n , a nonnegative δ and a potential $\eta(\cdot)$, we define the δ -risk $\text{risk}_{\delta}(\eta|\mathcal{P})$ of the potential on \mathcal{P} as the smallest ϵ such that

$$\begin{aligned} (a) \quad & \int e^{-\eta(\delta + e^T \xi)} p(\xi) d\xi \leq \epsilon \quad \forall (e \in \mathbf{R}^n : \|e\|_2 = 1, p \in \mathcal{P}), \\ (b) \quad & \int e^{\eta(e^T \xi - \delta)} p(\xi) d\xi \leq \epsilon \quad \forall (e \in \mathbf{R}^n : \|e\|_2 = 1, p \in \mathcal{P}). \end{aligned} \quad (23)$$

Let us make the following immediate observation:

Proposition 2.5 For $k = 1, \dots, K$, let \mathcal{P}^k be a family of probability densities on \mathbf{R}^n , X_1^k and X_2^k be closed nonempty non-intersecting convex sets in \mathbf{R}^n , one of the sets being bounded, and let h_k, c_k and δ_k be associated with $\mathcal{P} = \mathcal{P}^k$, $X_{\chi} = X_{\chi}^k$, $\chi = 1, 2$, via (17) and (18). Given potentials η_1, \dots, η_K , let us define the Euclidean detector induced by the potential η_k and h_k, c_k as the function

$$\phi_k(\omega) = \eta_k(h_k^T \omega - c_k) : \mathbf{R}^n \rightarrow \mathbf{R}, \quad (24)$$

and let

$$\phi^{(K)}(\omega_1, \dots, \omega_K) = \sum_{k=1}^K \phi_k(\omega_k) : \mathbf{R}^{nK} \rightarrow \mathbf{R}, \quad K = 1, 2, \dots \quad (25)$$

Finally, let $\{x_k\}_{k=1}^K$ and $\{p_k \in \mathcal{P}^k\}_{k=1}^K$ be deterministic sequences. Then

$$\begin{aligned} (a) \quad & \int e^{-\phi^{(K)}(x_1 + \xi_1, \dots, x_K + \xi_K)} \prod_{k=1}^K [p_k(\xi_k) d\xi_k] \leq \prod_{k=1}^K \text{risk}_{\delta_k}(\eta_k | \mathcal{P}^k) \quad \forall (x_k \in X_1^k, p_k \in \mathcal{P}^k)_{k=1}^K, \\ (b) \quad & \int e^{\phi^{(K)}(x_1 + \xi_1, \dots, x_K + \xi_K)} \prod_{k=1}^K [p_k(\xi_k) d\xi_k] \leq \prod_{k=1}^K \text{risk}_{\delta_k}(\eta_k | \mathcal{P}^k) \quad \forall (x_k \in X_2^k, p_k \in \mathcal{P}^k)_{k=1}^K. \end{aligned}$$

For a proof, see Section A.3.

Under the premise of Proposition 2.5, consider a test T_K^{η} which, given K observations (5) with independent of each other $\xi_k \sim p_k \in \mathcal{P}^k$, accepts the hypothesis $H_1 : \{x_k \in X_1^k, k \leq K\}$ whenever $\phi^{(K)}(\omega_1, \dots, \omega_K) \geq 0$, and accepts $H_2 : \{x_k \in X_2^k, k \leq K\}$ otherwise. An immediate corollary of Proposition 2.5 is as follows:

Corollary 2.2 In the notation and under the assumptions from the premise of Proposition 2.5, combining Proposition 2.5 with Markov's inequality we obtain that the risk $\text{Risk}_{\overline{\mathcal{S}}}(\mathcal{T}_K^{\eta} | [\mathcal{P}^k, X_1^k, X_2^k]_{k=1}^K)$ of \mathcal{T}_K^{η} does not exceed the quantity $\prod_{k=1}^K \text{risk}_{\delta_k}(\eta_k | \mathcal{P}^k)$.

¹We use the following fact (if absolutely evident facts indeed exist, this is one of them): given a vector $p = [p_1; \dots; p_n]$ with entries $p_i \in [0, 1]$, and a positive integer $k \leq n$, let $P_{\geq k|n}(p)$ be the probability to get $\geq k$ heads in n independent flips of a coin, with probability p_i to get a head in the i -th trial. Then $P_{\geq k|n}(p)$ is a nondecreasing function of p . Whatever "evident," this fact needs a proof, and here it is. Let $B = \{\omega \in \mathbf{R}^n : 0 \leq \omega_i \leq 1, i \leq n\}$ be the n -dimensional cube equipped with the Lebesgue measure μ . We have $P_{\geq k|n}(p) = \mu(B_p)$, where $B_p = \{\omega \in B : \text{Card}\{i : \omega_i \leq p_i\} \geq k\}$. When $0 \leq p_i \leq p'_i \leq 1$, $1 \leq i \leq n$, we clearly have $B_p \subset B_{p'}$, and consequently $P_{\geq k|n}(p) = \mu(B_p) \leq \mu(B_{p'}) = P_{\geq k|n}(p')$, as claimed.

2.4.2 Potentials for sub-spherical families of distributions

Definition 2.3 A. Given $\delta > 0$, we call a potential $\eta(s) : \mathbf{R} \rightarrow \mathbf{R}$ δ -regular, if the function

$$H_{\delta\eta}(s) = e^{-\eta(\delta-s)} + e^{-\eta(\delta+s)}$$

is nondecreasing on the ray $s \geq 0$.

B. The δ -index of the δ -regular potential η on a sub-spherical family \mathcal{P}_γ is the quantity

$$\epsilon_\delta(\eta|\gamma) := \int_{-\infty}^{\infty} e^{-\eta(r)} \gamma(r - \delta) dr = \int_{-\infty}^{\infty} e^{-\eta(\delta+s)} \gamma(s) ds = \int_0^{\infty} H_{\delta\eta}(s) \gamma(s) ds, \quad (26)$$

where the concluding equality is due to the fact that $\gamma(\cdot)$ is even.

Proposition 2.6 Let $\mathcal{P} = \mathcal{P}_\gamma$ be a sub-spherical family of probability densities on \mathbf{R}^n , let $\delta \geq 0$, and let η be a δ -regular potential. Then

$$\text{risk}_\delta(\eta|\mathcal{P}_\gamma) \leq \epsilon_\delta(\eta|\gamma). \quad (27)$$

For a proof, see Section A.4.

Example 1: step potential. Assume that $\mathcal{P} = \mathcal{P}_\gamma$ is a sub-spherical family of distributions, and let $\delta > 0$, implying that $\epsilon_\star = \epsilon_\star(\delta|\gamma)$, as given by (19), belongs to $[0, 1/2)$ (recall that γ is an even probability density positive in a neighbourhood of the origin). We define the *step potential* as⁴

$$\eta(s) = \frac{1}{2} \ln \left(\frac{1 - \epsilon_\star}{\epsilon_\star} \right) \text{sign}(s). \quad (28)$$

Taking into account that $0 \leq \epsilon_\star < 1/2$, it is immediately seen that η is a δ -regular potential. The δ -index of the step potential satisfies

$$\epsilon_\delta(\eta|\gamma) = 2\sqrt{\epsilon_\star(1 - \epsilon_\star)}, \quad (29)$$

as is shown by the following computation:

$$\begin{aligned} \epsilon_\delta(\eta|\gamma) &= \int_{-\infty}^{\infty} e^{-\eta(s)} \gamma(s - \delta) ds = \sqrt{\frac{1 - \epsilon_\star}{\epsilon_\star}} \int_{-\infty}^{-\delta} \gamma(s) ds + \sqrt{\frac{\epsilon_\star}{1 - \epsilon_\star}} \int_{-\delta}^{\infty} \gamma(s) ds \\ &= \sqrt{\frac{1 - \epsilon_\star}{\epsilon_\star}} \epsilon_\star + \sqrt{\frac{\epsilon_\star}{1 - \epsilon_\star}} (1 - \epsilon_\star) = 2\sqrt{\epsilon_\star(1 - \epsilon_\star)}. \end{aligned}$$

Note that in the situation of Section 2.1 with stationary K -repeated observations (4), a sub-spherical family $\mathcal{P} = \mathcal{P}_\gamma$ and non-intersecting X_1, X_2 , the test \mathcal{T}_K^η associated with the step potential η is exactly the majority test $\mathcal{T}_K^{\text{maj}}$ defined in Section 2.3.2. Because the index $\epsilon_\delta(\eta|\gamma)$ of the step potential is < 1 due to $\epsilon_\star < 1/2$, its δ -risk $\text{risk}_\delta(\eta|\mathcal{P}_\gamma)$ is < 1 as well.

Example 2: ramp potential. Assume that \mathcal{P} is the sub-spherical family $\mathcal{P}_{\gamma_\mathcal{L}}$ with $\gamma_\mathcal{L}(s) = \frac{1}{2\lambda} e^{-|s|/\lambda}$, $s \in \mathbf{R}$ (see the ‘‘Laplace’’ example in Section 2.2). Given $\delta > 0$, let us consider the *ramp potential* η_δ which minimizes the risk $\epsilon_\delta(\eta|\gamma_\mathcal{L})$:

$$\eta_\delta(s) = \frac{1}{2} \ln \left(\frac{\gamma_\mathcal{L}(s - \delta)}{\gamma_\mathcal{L}(s + \delta)} \right) = \begin{cases} s/\lambda & \text{for } |s| \leq \delta, \\ \frac{\delta}{\lambda} \text{sign}(s), & \text{for } |s| > \delta. \end{cases}$$

One can easily verify that η_δ is δ -regular for any $\delta > 0$, and the corresponding δ -index satisfies

$$\epsilon_\delta(\eta_\delta|\gamma_\mathcal{L}) = \int_{-\infty}^{\infty} \sqrt{\gamma_\mathcal{L}(s - \delta)\gamma_\mathcal{L}(s + \delta)} ds = e^{-\delta/\lambda} \left(\frac{\delta}{\lambda} + 1 \right),$$

which is < 1 for all $\delta > 0$.

⁴Hereafter we put $\text{sign}(s) = s/|s|$ if $s \neq 0$ and $\text{sign}(0) = 0$.

Examples above show that in the situation of Section 2.1, assuming that $\mathcal{P} = \mathcal{P}_\gamma$ is a sub-spherical family, there exist potentials η with δ -indexes $\epsilon_\delta(\eta|\gamma) < 1$ and therefore, using Proposition 2.6, with risks $\text{risk}_\delta(\eta|\mathcal{P}) < 1$, provided $\delta > 0$. In this situation, when solving problem $(\overline{\mathcal{S}}_K)$ with $X_1^k \equiv X_1, X_2^k \equiv X_2$, X_1 and X_2 not intersecting, and specifying $\delta > 0$ according to (18), we can decide on the hypotheses H_1 and H_2 with any desired risk $\epsilon \in (0, 1)$ via semi-stationary K -repeated observations (5), provided that

$$K \geq K(\epsilon) := \left\lceil \frac{\ln(\epsilon^{-1})}{\ln(\text{risk}_\delta(\eta|\mathcal{P})^{-1})} \right\rceil, \quad (30)$$

see Corollary 2.2 (here $\lfloor x \rfloor$ stands for the smallest integer larger or equal to x).

2.4.3 Potentials for the sub-Gaussian family

Aside from sub-spherical families, there are other families of probability distributions allowing, in the case of $X_1 \cap X_2 = \emptyset$, for potentials with risks < 1 and thus for tests with an arbitrarily low risk, provided stationary or semi-stationary K -repeated observations with properly selected K are available. The simplest family of this type is the family \mathcal{P}_{sG}^n of sub-Gaussian probability densities $p(\cdot)$, with parameters 0 and I_n , that is, probability densities p on \mathbf{R}^n such that

$$\int e^{h^T \xi} p(\xi) d\xi \leq e^{\frac{1}{2} h^T h} \quad \forall h \in \mathbf{R}^n. \quad (31)$$

Assuming that $\mathcal{P} = \mathcal{P}_{sG}$ and given $\delta > 0$, let us put

$$\eta_{sG, \delta}(s) = \delta s. \quad (32)$$

Proposition 2.7 *Whenever $\delta \geq 0$, we have*

$$\text{risk}_\delta(\eta_{sG, \delta} | \mathcal{P}_{sG}) \leq e^{-\delta^2/2}. \quad (33)$$

Proposition 2.7 is a special case of Proposition 3.3 in [13]; to make the presentation self-contained, we provide its proof in Section A.5.

Note that if \mathcal{P} is the sub-spherical family $\mathcal{P}_{\gamma G}$ with $\gamma G(s) = \frac{1}{\sqrt{2\pi}} e^{-s^2/2}$, $s \in \mathbf{R}$ (recall that this family contains, for instance, Gaussian distributions with zero mean and covariance matrix $\preceq I_n$, cf. Section 2.2.1) the potential $\eta_{sG, \delta}(s) = \frac{1}{2} \ln \left(\frac{\gamma G(s-\delta)}{\gamma G(s+\delta)} \right) = \delta s$ minimizes the δ -index over $\mathcal{P}_{\gamma G}$ with $\epsilon_\delta(\eta_{sG, \delta} | \gamma G) = e^{-\delta^2/2}$.

2.4.4 “Majority of means” tests

Let now $X_1 \subset \mathbf{R}^n$ and $X_2 \subset \mathbf{R}^n$ be closed convex nonempty sets, one of the sets being bounded, and such that $X_1 \cap X_2 = \emptyset$. In this situation, given $\epsilon \in (0, 1)$ and K -repeated stationary observations ω^K , the potential-based tests developed in Sections 2.4.1 – 2.4.3 allow to decide on H_1, H_2 with risk $\leq \epsilon$, where K grows logarithmically with ϵ^{-1} . For instance, in the case of a sub-Gaussian family of distributions, the corresponding test attains the risk ϵ provided that $K \geq K(\epsilon) = O(\ln(\epsilon^{-1})/\delta^2)$ (cf. (33)), where $2\delta > 0$ is the Euclidean distance between X_1 and X_2 , see (16) – (18). On the other hand, these tests rely upon the “Cramer-type” Definition 2.2 of the risk of the potential, and thus assume control of the exponential moment of $\eta(\xi)$. In this section we present a different testing procedure, which only uses second order characteristics of the potential (mean and variance), and yet allows to achieve arbitrarily low risks of testing for essentially the same sizes of observation sample. We describe this modification in the simplest case of stationary repeated observations, generalizations to more general settings (e.g., that of Proposition 2.5) being straightforward.

Now, let x_*^1, x_*^2 and h_* be associated with X_1 and X_2 via (16) and (17), and let \mathcal{P} be a family of probability distributions on \mathbf{R}^n . We suppose that a potential $\eta(\cdot)$ and $c \in \mathbf{R}$ are such that for some $\varrho > 0$ and all $p \in \mathcal{P}$,

$$\begin{aligned} (a) \quad & \mathbf{E}_{\xi \sim p} \{ \eta(h_*^T(x_*^1 + \xi) + c) \} - \mathbf{E}_{\xi \sim p} \{ \eta(h_*^T(x_*^2 + \xi) + c) \} \geq \varrho, \\ (b) \quad & \mathbf{Var}_{\xi \sim p} \{ \eta(h_*^T(x + \xi) + c) \} \leq 1 \end{aligned} \quad (34)$$

for all $x \in X_1 \cup X_2$.⁵

⁵Here and below \mathbf{Var} denotes the “usual” variance: for a probability density p on \mathbf{R}^n and $f : \mathbf{R}^n \rightarrow \mathbf{R}$, $\mathbf{Var}_{\xi \sim p} \{ f(\xi) \} = \mathbf{E}_{\xi \sim p} \{ f^2(\xi) \} - [\mathbf{E}_{\xi \sim p} \{ f(\xi) \}]^2$.

Example. Let \mathcal{P} be a family of zero-mean distributions on \mathbf{R}^n with covariance matrix $\preceq \sigma^2 I_n$:

$$\int \xi p(\xi) d\xi = 0, \quad \int \xi \xi^T p(\xi) d\xi \preceq I_n, \quad \forall p \in \mathcal{P}.$$

For the linear potential $\eta(t) = t$ with $c = 0$ we clearly have

$$\begin{aligned} \mathbf{E}_{\xi \sim p} \{ \eta(h_*^T(x_*^1 + \xi) + c) \} - \mathbf{E}_{\xi \sim p} \{ \eta(h_*^T(x_*^2 + \xi) + c) \} &= h_*^T(x_*^1 - x_*^2) = 2\delta =: \varrho, \quad \forall p \in \mathcal{P} \\ \mathbf{Var}_{\xi \sim p} \{ \eta(h_*^T(x + \xi) + c) \} &= \mathbf{E}_{\xi \sim p} \{ (h_*^T \xi)^2 \} \leq 1 \quad \forall (x \in X_1 \cup X_2, p \in \mathcal{P}) \end{aligned}$$

(recall that $\delta = \frac{1}{2} \|x_*^1 - x_*^2\|_2$, and $\|h_*\|_2 = 1$).

Now, given K -repeated stationary observations $\omega^K = [x + \xi_1, \dots, x + \xi_K]$, and $\epsilon_1, \epsilon_2 \in (0, 1)$, we consider the inference problem (\mathcal{S}_K) of deciding via the observation ω^K whether $x \in X_1$ (hypothesis H_1) or $x \in X_2$ (hypothesis H_2). Our objective is to build the test \mathcal{T}_K for \mathcal{S}^K such that, uniformly over $x \in X_1 \cup X_2$ and $p \in \mathcal{P}$, the probability of wrongly rejecting H_1 (accepting H_2) is $\leq \epsilon_1$, and the probability of wrongly rejecting H_2 (accepting H_1) is $\leq \epsilon_2$. and we want to attain this goal using the smallest possible size K of observation sample ω^K .

For the sake of definiteness, assume that $\epsilon_1 \leq \epsilon_2$. We denote $\kappa = \frac{1}{2} \left(1 - \frac{\ln(\epsilon_2^{-1})}{\ln(\epsilon_1^{-1})} \right)$ (note that $0 \leq \kappa < \frac{1}{2}$). Let now $m = \lfloor 4e(e^{-\kappa} + 1)^2 \varrho^{-2} \rfloor$, and let

$$\psi_j(\omega^K) = \frac{1}{m} \sum_{i=(j-1)m+1}^{mj} \eta(h_*^T \omega_i + c), \quad j = 1, 2, \dots$$

Denote

$$c_* = \mathbf{E}_{\xi \sim p} \{ \eta(h_*^T(x_*^1 + \xi) + c) \} - \frac{e^\kappa \varrho}{1 + e^\kappa} \geq \mathbf{E}_{\xi \sim p} \{ \eta(h_*^T(x_*^2 + \xi) + c) \} + \frac{\varrho}{1 + e^\kappa}.$$

The K -observation majority of means test $\mathcal{T}_K^{\text{mm}}$ for \mathcal{S}_K is as follows: given

$$K := Jm \geq \lfloor 2 \ln(\epsilon_1^{-1}) \rfloor m \tag{35}$$

observations ω_k , $k = 1, \dots, K$, $\mathcal{T}_K^{\text{mm}}$ accepts the hypothesis H_1 when $\psi_j(\omega^K) \geq c_*$ for at least $J/2$ values of j , and accepts the hypothesis H_2 otherwise.

Proposition 2.8 *In the just described situation, the risks of the test $\mathcal{T}_K^{\text{mm}}$ meet the problem specifications, namely,*

$$\text{Risk}_{1\mathcal{S}}(\mathcal{T}_K | \mathcal{P}, X_1, X_2) \leq \epsilon_1, \quad \text{Risk}_{2\mathcal{S}}(\mathcal{T}_K | \mathcal{P}, X_1, X_2) \leq \epsilon_2.$$

For a proof, see Section A.6.

Let us consider the test $\mathcal{T}_K^{\text{mm}}$ using the linear potential with $c = 0$. Under the premise of the proposition, i.e., in the situation where the noise covariance matrix is $\preceq I_n$ and the distance between the sets X_1 and X_2 is $\geq 2\delta$, the size of the stationary K -repeated observation sufficient for the test to satisfy the risk specifications is

$$\lfloor e(e^{-\kappa} + 1)^2 \delta^{-2} \rfloor \lfloor 2 \ln(\epsilon_1^{-1}) \rfloor = O(\delta^{-2} \ln(\epsilon_1^{-1})).$$

Note that when $\epsilon_1 = \epsilon_2 = \epsilon$, the above number of observations sufficient to decide on H_1, H_2 with risk ϵ is within absolute constant factor of the number of observations, as given by (30) and (33), needed for the same purpose in the case when $\mathcal{P} = \mathcal{P}_{\text{SG}}$.

3 Sequential detection via Euclidean separation

3.1 Motivating example

We start introducing a motivating example which will help understand the general setting discussed in Sections 3.2, 3.3, 3.4 and which will be used in Section 4.3 to test our methodology.

Consider the simple time series model

$$\begin{aligned} y_k &= \alpha_k + \zeta_k \\ \alpha_k &= \alpha_{k-1} + \eta_k + u_k \end{aligned}, \quad k = 1, 2, \dots, d, \tag{36}$$

where

- y_k is the observation at time k , and $u = [u_1; \dots; u_d] \in \mathbf{R}^d$ is the deterministic input,
- $\zeta = [\zeta_1; \dots; \zeta_d]$ is zero mean d -dimensional Gaussian random vector with unknown covariance matrix known to be $\preceq \sigma^2 I_d$;
- $\eta = [\eta_1; \dots; \eta_d]$ is independent of ζ d -dimensional random vector obeying multivariate Student distribution with ν degrees of freedom and matrix parameter $\preceq I_d$. Here ν is a positive integer or $+\infty$, with $\nu = \infty$ interpreted as the fact that η is a zero mean Gaussian random vector with covariance matrix $\preceq I_d$.

We intend to decide from observations y_1, \dots, y_d on the nuisance hypothesis $u = 0$ vs. signal alternative

$$u \in \bigcup_{\substack{1 \leq i \leq d, \\ 0 < \rho < R}} [U_i^+(\rho) \cup U_i^-(\rho)], \quad (37)$$

where $R \in (0, \infty)$ is a parameter, $U_i^-(\rho) = -U_i^+(\rho)$, and $[U_i^+(\rho) \cup U_i^-(\rho)]$ is the set of “signal inputs of shape i and magnitude $\geq \rho$.” We consider two cases:

- *pulse signal inputs*: $U_i^+(\rho) = \{u \in \mathbf{R}^d : u_k = 0, k \neq i, \rho \leq u_i \leq R\}$, $1 \leq i \leq d$;
- *step signal inputs*: $U_i^+(\rho) = \{u \in \mathbf{R}^d : u_k = 0, k < i, \rho \leq u_i = u_{i+1} = \dots = u_d \leq R\}$, $1 \leq i \leq d$.

To make the notation consistent with the one used in the general setting described in Section 3.2, we set $U_{2i-1}(\rho) = U_i^+(\rho)$ and $U_{2i}(\rho) = U_i^-(\rho)$, thus getting $N = 2d$ parametric families of signal inputs and we refer to signal input $u \in U_j(\rho)$ as to signal input of shape j and magnitude $\geq \rho$.

Control parameters are $R > 0$ appearing in (37) and tolerance $\epsilon \in (0, 1/2)$ responsible for the risks of our decision rules.

Our goal is to find decision rules \mathcal{T}_k and positive reals $\rho_{kj} \in (0, R]$, $1 \leq k \leq d$, $1 \leq j \leq N$, such that

- \mathcal{T}_k makes a decision given observation $y^k = [y_1; \dots; y_k]$, and this decision, depending on y^k , is either “signal conclusion,” or “nuisance conclusion” (exactly one of them). In the case of signal conclusion at step k , the inference procedure is terminated at this step, otherwise we pass to time $k + 1$ (when $k < d$) or terminate (when $k = d$).
- The following risk specifications are met:
 - When the nuisance hypothesis is true (i.e., $u = 0$), the probability of terminating with signal conclusion (false alarm) somewhere on the time horizon $1, \dots, d$ is $\leq \epsilon$.
 - For every $k \leq d$ and every $j \leq N$, if the input u underlying our observation belongs to $U_j(\rho)$ with $\rho_{kj} \leq \rho < R$, the probability of signal conclusion somewhere on the time horizon $1, \dots, k$ should be at least $1 - \epsilon$.

We intend to achieve this goal with as small ρ_{kj} as possible.

We now explain how to cast the problem we have just described as the more general detection problem discussed in Sections 3.2, 3.3, 3.4. To this end, we need to build a new observation scheme. More precisely, we have

$$y^k = \alpha_0 \underbrace{[1; \dots; 1]}_k + B_k \eta + C_k \zeta + B_k u$$

for some known matrices B_k and C_k . Let F_k be $d \times (k - 1)$ matrix with columns forming an orthonormal basis of the $(k - 1)$ -dimensional subspace of \mathbf{R}^k comprised of vectors with zero mean. Let us set

$$z^k := F_k^T y^k = D_k \eta + E_k \zeta + D_k u \quad (38)$$

where $(k - 1) \times d$ matrices D_k and E_k are given by $D_k = F_k^T B_k$, $E_k = F_k^T C_k$, and, as is immediately seen, have rank $k - 1$. We treat z^k , rather than y^k , as an intermediate observation at time k .

We then find a positive definite matrix Θ_k such that

$$\Theta_k^2 \succeq D_k D_k^T \ \& \ \Theta_k^2 \succeq E_k E_k^T.$$

Finally, the observation ω^k at time k , is the vector (cf. (54))

$$\omega^k = \Theta_k^{-1} z^k = \underbrace{[\Theta_k^{-1} D_k]}_{A_k} u + \underbrace{\Theta_k^{-1} [D_k \eta + E_k \zeta]}_{\xi_k}. \quad (39)$$

We have come to observation scheme (39) which can be handled by the change detection procedure we are about to describe in Sections 3.2, 3.3, 3.4. In particular, by Proposition 2.3, ξ_k has a symmetric density $p_k(\cdot) \in \mathcal{P}_\gamma^{k-1}$ with $\gamma = \gamma_S \star \gamma_\sigma$, where γ_S is the density of the standard univariate Student's t_ν distribution with ν degrees of freedom, and γ_σ is the density of $\mathcal{N}(0, \sigma^2)$.

3.2 Situation and goal

In this Section our objective is to make decisions about an unknown vector $u \in \mathbf{R}^{n_u}$ representing inputs of a linear system in the situation where the information about u is acquired sequentially, and the goal is to decide whether the input observed so far is a nuisance or is ‘‘meaningful.’’ Our modeling methodology goes back to [5]. Specifically, suppose that the observation $\omega^k \in \mathbf{R}^{m_k}$ available at step $k = 1, \dots, K$, is

$$\omega^k = A_k u + \xi_k, \quad (40)$$

where u is the unknown system input, $\xi_k \in \mathbf{R}^{m_k}$ are random noises, and $A_k \in \mathbf{R}^{m_k \times n_u}$ are known matrices. Throughout this section we suppose that we are given

1. a family \mathcal{P} of distributions on \mathbf{R}^{n_ξ} and matrices $Z_k \in \mathbf{R}^{m_k \times n_\xi}$ such that Z_k is of rank m_k and

$$\xi_k = Z_k \xi, \quad k = 1, \dots, K, \quad (41)$$

where ξ obeys some (possibly unknown) distribution $p_\xi \in \mathcal{P}$;

2. a convex compact set $U_{\text{inp}} \subset \mathbf{R}^{n_u}$ of *admissible inputs*;
3. a convex compact *nuisance set* $V_{\text{nuis}} \subset \text{int } U_{\text{inp}}$ such that $0 \in V_{\text{nuis}}$;
4. N convex and closed ‘‘activation’’ sets $W_j \subset \mathbf{R}^{n_u}$ such that $0 \notin W_j$ and

$$w \in W_j, \rho \geq 1 \Rightarrow \rho w \in W_j; \quad (42)$$

5. N nonempty convex compact ‘‘drag sets’’ V_{drag}^j such that $0 \in V_{\text{drag}}^j$.

For the example of Section 3.1, we have $\xi = [\eta; \zeta]$, $Z_k[\eta; \zeta] = \Theta_k^{-1} [D_k \eta + E_k \zeta]$, $U_{\text{inp}} = \{u \in \mathbf{R}^d : 0 \leq u_i \leq R\}$, $V_{\text{nuis}} = V_{\text{drag}}^j = \{0\}$, and

- for *pulse signal inputs*: $W_{2i-1} = \{u \in \mathbf{R}^d : u_k = 0, k \neq i, 1 \leq u_i\}$, $W_{2i} = -W_{2i-1}$, $1 \leq i \leq d$;
- for *step signal inputs*: $W_{2i-1} = \{u \in \mathbf{R}^d : u_k = 0, k < i, 1 \leq u_i = u_{i+1} = \dots = u_d\}$, $W_{2i} = -W_{2i-1}$, $1 \leq i \leq d$.

We call an input $u \in V_{\text{nuis}}$ a *nuisance*, and a vector z of the form $z = \rho w$, with $w \in W_j$ and $\rho > 0$, an *activation of shape j and magnitude $\geq \rho$* . Note that if $0 < \rho' \leq \rho$ and z is an activation of shape j and magnitude $\geq \rho$, then, as it should be, z is an activation of shape j and magnitude $\geq \rho'$, due to $z = \rho' w'$, $w' = (\rho/\rho')w$, and $w' \in W_j$ due to $0 < \rho' \leq \rho$ and (42).

We call input u a *signal of shape j and magnitude $> \rho$* (or $\geq \rho > 0$), if $u \in U_{\text{inp}}$ and

$$u = v + \rho' w$$

with $v \in V_{\text{drag}}^j$, $w \in W_j$, and $\rho' > \rho$ (resp. $\rho' \geq \rho$), and denote by

$$U_j(\rho) = \{u = v + \rho w : u \in U_{\text{inp}}, v \in V_{\text{drag}}^j, w \in W_j\}$$

the set of all signal inputs of shape j and magnitude $\geq \rho$.

Our goal is to decide via observations ω^k , $k \leq K$, on the nuisance hypothesis “the input u to (40) is a nuisance” (i.e., $u \in V_{\text{nuil}}$) vs. the signal alternative “the input u to (40) is a signal of some shape and (positive) magnitude.” Note that in this case, under the signal alternative the input u belongs to a nonconvex set expressed as a union of convex sets. More precisely, we assume that we are given tolerances

$$\{\epsilon_k \in (0, 1/2), 1 \leq k \leq K\}, \{\epsilon_{kj} \in (0, 1/2), 1 \leq k \leq K, 1 \leq j \leq N\}$$

and want to design a sequence of decision rules $\{\mathcal{T}_k : 1 \leq k \leq K\}$ along with thresholds $\rho_{kj} > 0$, $1 \leq k \leq K$, $1 \leq j \leq N$ with the following properties: for every $k \leq K$,

- rule \mathcal{T}_k makes a decision based on observation ω^k , and this decision is either to accept the null hypothesis or to accept the signal one (but not both);
- if the input is a nuisance, the probability for \mathcal{T}_k to accept the signal alternative is $\leq \epsilon_k$;
- if, for some j , the input is a signal of shape j and magnitude $> \rho_{kj}$, the probability for \mathcal{T}_k to accept the nuisance hypothesis is at most ϵ_{kj} .

Given tolerances $\epsilon_k, \epsilon_{kj}$, we would like to achieve the outlined goal with thresholds ρ_{kj} as small as possible.

3.3 Tests \mathcal{T}_k , Scheme I.

3.3.1 Assumptions on the distribution of noise ξ

Throughout Section 3.3, we make the following assumption on the family \mathcal{P} of probability densities of random disturbance ξ in (41):

Assumption A1 For every $k \leq K$ we can point out a parametric family $\mathcal{H}_k = \{\eta_\delta^k(\cdot) : \delta \geq 0\}$ of potentials and a continuous nonincreasing function $\mathcal{R}_k(\delta) : \mathbf{R}_+ \rightarrow (0, 1]$ such that $\mathcal{R}_k(0) = 1$ and

$$\text{risk}_\delta(\eta_\delta^k | \mathcal{P}^k) \leq \mathcal{R}_k(\delta) \quad \forall \delta \geq 0, \quad (43)$$

where \mathcal{P}^k is the family of probability densities of random vectors $\xi_k = Z_k \xi$ with $\xi \sim p_\xi \in \mathcal{P}$.

Note that Assumption A1 indeed holds in the situations of our primary interest. Specifically, assume that $Z_k Z_k^T \preceq I_{m_k}$.⁶ Then

- when $\mathcal{P} = \mathcal{P}_{\text{sG}}^{n_\xi}$ is the family of sub-Gaussian distributions on \mathbf{R}^{n_ξ} , with parameters $0, I_{n_\xi}$, then \mathcal{P}^k belongs to the family $\mathcal{P}_{\text{sG}}^{m_k}$ of sub-Gaussian, with parameters $0, I_{m_k}$, distributions on \mathbf{R}^{m_k} . By Proposition 2.7, (43) is ensured by the choice

$$\mathcal{H}_k = \{\eta_\delta^k(s) = \delta s, \delta \geq 0\}, \quad \mathcal{R}_k(\delta) = e^{-\delta^2/2};$$

- when $\mathcal{P} = \mathcal{P}_\gamma$ is a sub-spherical family of distributions on \mathbf{R}^{n_ξ} , (41) combines with Proposition 2.1 to ensure that the family \mathcal{P}^k is contained in the sub-spherical family of distributions $\mathcal{P}_\gamma^{m_k}$. Invoking the example of the step potential from Section 2.4.2, (43) is ensured by setting

$$\mathcal{H}_k = \left\{ \eta_\delta^k(s) = \frac{1}{2} \ln \left(\frac{1 - \varepsilon_\star(\delta)}{\varepsilon_\star(\delta)} \right) \text{sign}(s), \delta \geq 0 \right\}, \quad \mathcal{R}_k(\delta) = 2\sqrt{\varepsilon_\star(\delta)(1 - \varepsilon_\star(\delta))},$$

$$\varepsilon_\star(\delta) = \int_\delta^\infty \gamma(s) ds.$$

3.3.2 Building decision rules

For every $k \leq K$ and every $j \leq N$, consider the parametric convex optimization problem

$$\text{Opt}_{kj}(\rho) = \min_{v', v, w} \left\{ \frac{1}{2} \|A_k(v' - [v + \rho w])\|_2 : v' \in V_{\text{nuil}}, v \in V_{\text{drg}}^j, w \in W_j, v + \rho w \in U_{\text{inp}} \right\} \quad (P_{kj}[\rho])$$

where the parameter ρ is positive.

From our assumptions on $V_{\text{nuil}}, U_{\text{inp}}, V_{\text{drg}}^j$ and W_j it immediately follows that

⁶this always can be achieved by appropriate scaling of observations (40).

1. the set Δ_j of those $\rho > 0$ for which $(P_{kj}[\rho])$ is feasible, is a half-open segment $(0, R_j]$ with $0 < R_j < \infty$, and
2. when $\rho \in \Delta_j$, problem $(P_{kj}[\rho])$ is solvable, and $\text{Opt}_{kj}(\rho)$ is a real-valued continuous nondecreasing function on Δ_j such that $\lim_{\rho \rightarrow +0} \text{Opt}_{kj}(\rho) = 0$.

Given $k \leq K$, we specify the test \mathcal{T}_k as follows.

1. We compute the quantities R_j .⁷
2. We select somehow a (perhaps, empty) set $J_k \subset \{1, 2, \dots, N\}$ and reals $\rho_{kj} \in (0, R_j]$, $j \in J_k$, in such a way that
 - (a) $\rho_{kj} = R_j$ when $j \notin J_k$,
 - (b) we have

$$\sum_{j \in J_k} \epsilon_{kj}^{-1} \mathcal{R}_k^2(\text{Opt}_{kj}(\rho_{kj})) \leq \epsilon_k. \quad (44)$$

Note that (44) implies that $\text{Opt}_{kj}(\rho_{kj}) > 0$ for $j \in J_k$, otherwise the left hand side in (44) is at least 1 due to $\mathcal{R}_k(0) = 1$, and we have assumed that $\epsilon_k < 1/2$.

3. If $J_k = \emptyset$, \mathcal{T}_k accepts the nuisance hypothesis. When $J_k \neq \emptyset$, we act as follows:
 - (a) for $j \in J_k$, we solve the optimization problem $(P_{kj}[\rho_{kj}])$ and denote by $(v'_{kj}; v_{kj}, w_{kj})$ an optimal solution to the problem. We set

$$\begin{aligned} \delta_{kj} &= \text{Opt}_{kj}(\rho_{kj}), \quad u_{kj} = v_{kj} + \rho_{kj} w_{kj}, \quad h_{kj} = \frac{A_k[v'_{kj} - u_{kj}]}{\|A_k[v'_{kj} - u_{kj}]\|_2}, \\ c_{kj} &= \frac{1}{2} h_{kj}^T A_k[v'_{kj} + u_{kj}], \quad \phi_{kj}(\omega^k) = \eta_{\delta_{kj}}^k (h_{kj}^T \omega^k - c_{kj}), \end{aligned} \quad (45)$$

with $\eta_{\delta}^k(\cdot)$ given by Assumption A1, and define α_{kj} from the relation

$$e^{\alpha_{kj}} = \epsilon_{kj} / \mathcal{R}_k(\delta_{kj}). \quad (46)$$

- (b) Finally, given observation ω^k , the rule \mathcal{T}_k accepts the nuisance hypothesis if $\phi_{kj}(\omega^k) + \alpha_{kj} \geq 0$ for all $j \in J_k$, and accepts the signal hypothesis otherwise.

3.3.3 Performance analysis

We now check that the just defined decision rules meet the goal stated in Section 3.2. Let us fix $k \leq K$ and $j \in J_k$. Given feasible input u , let P_u^k be the distribution of observation ω^k , the input being u . By construction and due to Proposition 2.5, applied to observations $\omega^k = A_k u + \xi_k$ and $X_1 = A_k V_{\text{nuis}}$, $X_2 = A_k U_j(\rho_{kj})$, we have

$$\begin{aligned} \int e^{-\phi_{kj}(\omega^k)} P_u^k(d\omega^k) &\leq \mathcal{R}_k(\delta_{kj}), \quad \text{when } u \text{ is a nuisance,} \\ \int e^{\phi_{kj}(\omega^k)} P_u^k(d\omega^k) &\leq \mathcal{R}_k(\delta_{kj}), \quad \text{when } u \text{ is a signal of shape } j \text{ and magnitude } \geq \rho_{kj}. \end{aligned} \quad (47)$$

Assume first that the input u is a nuisance, and let us upper-bound the P_u^k -probability of rejecting the nuisance hypothesis at step k . This takes place only when $\phi_{kj}(\omega^k) + \alpha_{kj} < 0$ for some $j \in J_k$, and, by the first relation in (47), for a given $j \in J_k$ the P_u^k -probability of this event is at most $e^{-\alpha_{kj}} \mathcal{R}_k(\delta_{kj}) = \mathcal{R}_k^2(\delta_{kj}) / \epsilon_{kj}$ where the last equality is due to (46). As a result, the P_u^k -probability to reject the nuisance hypothesis at step k is at most

$$\sum_{j \in J_k} \epsilon_{kj}^{-1} \mathcal{R}_k^2(\delta_{kj}) = \sum_{j \in J_k} \epsilon_{kj}^{-1} \mathcal{R}_k^2(\text{Opt}_{kj}(\rho_{kj})) \leq \epsilon_k,$$

where the concluding inequality is due to (44).

Now, let the input u be a signal of shape j and magnitude $> \rho_{kj}$. Due to $\rho_{kj} = R_j$ for $j \notin J_k$, we have $j \in J_k$. Let us upper-bound the P_u^k -probability of the nuisance conclusion at step k . By construction, the latter

⁷The simplest way to identify R_j is to run bisection in ρ on a large initial range of ρ in order to find the largest ρ for which $(P_{kj}[\rho])$ is feasible.

occurs only when $\phi_{kj}(\omega^k) + \alpha_{kj} \geq 0$, and, by the second relation in (47), this can happen with P_u^k -probability at most $e^{\alpha_{kj}} \mathcal{R}_k(\delta_{kj}) = \epsilon_{kj}$, where the concluding equality is due to (46).

The bottom line is that the probability of false alarm at step k (a signal conclusion when the input is nuisance) is $\leq \epsilon_k$, and the probability of ρ_{kj} -miss (the probability to make a nuisance conclusion when the input is a signal of shape j and magnitude $> \rho_{kj}$) is at most ϵ_{kj} .

3.4 Decision rules \mathcal{T}_k , Scheme II

Throughout Section 3.4, we make the following

Assumption A: The family \mathcal{P} of probability densities of the disturbance ξ in (41) is such that for all $k \leq K$, probability densities of noises ξ_k belong to sub-spherical families $\mathcal{P}_\gamma^{m_k}$ with some common γ .

Note this assumption takes place, e.g., when \mathcal{P} is a sub-spherical family $\mathcal{P}_\gamma^{n_\xi}$ and $Z_k Z_k^T \preceq I_{m_k}$, $k \leq K$, see Proposition 2.1.

We put (cf. (19))

$$\varepsilon_*(\delta|\gamma) = \int_{\delta}^{\infty} \gamma(s) ds, \quad (48)$$

3.4.1 Building the decision rules

Given $k \leq K$, we specify \mathcal{T}_k as follows.

1. Exactly as in Scheme I, for every $j \leq N$, we consider the parametric convex optimization problem $(P_{kj}[\rho])$ with positive ρ , and compute R_j , the largest ρ for which the problem is feasible.
2. We select somehow a (perhaps, empty) set $J_k \subset \{1, 2, \dots, N\}$ and reals $\rho_{kj} \in (0, R_j]$, $1 \leq j \leq N$, and $\alpha_{\chi kj}$, $\chi = 1, 2$, $j \in J_k$ in such a way that $\rho_{kj} = R_j$ when $j \notin J_k$, and we have

$$\begin{aligned} (a) \quad & \varepsilon_*(\alpha_{2kj}|\gamma) \leq \epsilon_{kj}, \quad \forall j \in J_k, \\ (b) \quad & \sum_{j \in J_k} \varepsilon_*(\alpha_{1kj}|\gamma) \leq \epsilon_k, \\ (c) \quad & \alpha_{1kj} + \alpha_{2kj} \leq 2\delta_{kj}, \quad \delta_{kj} = \text{Opt}_{kj}(\rho_{kj}), \quad \forall j \in J_k. \end{aligned} \quad (49)$$

Note that for $j \in J_k$ we have $\alpha_{1kj} > 0$ and $\alpha_{2kj} > 0$ (by (49.a-b) combined with $\epsilon_{kj} < 1/2$, $\epsilon_k < 1/2$; recall that $\varepsilon_*(s|\gamma) \geq 1/2$ when $s \leq 0$). As a result, (49.c) implies that $\delta_{kj} > 0$ whenever $j \in J_k$.

3. If $J_k = \emptyset$, \mathcal{T}_k accepts the nuisance hypothesis. When $J_k \neq \emptyset$, we act as follows:

- (a) for $j \in J_k$, we solve the optimization problem $(P_{kj}[\rho_{kj}])$ and denote by $(v'_{kj}; v_{kj}, w_{kj})$ an optimal solution to the problem. Similarly to (45), we set

$$\begin{aligned} u_{kj} &= v_{kj} + \rho_{kj} w_{kj}, \quad h_{kj} = \frac{A_k[v'_{kj} - u_{kj}]}{\|A_k[v'_{kj} - u_{kj}]\|_2}, \\ c_{kj} &= \frac{1}{2} h_{kj}^T A_k[v'_{kj} + u_{kj}], \quad \phi_{kj}(\omega^k) = h_{kj}^T \omega^k - c_{kj}. \end{aligned} \quad (50)$$

- (b) Finally, given observation ω^k , the rule \mathcal{T}_k accepts the nuisance hypothesis if $\phi_{kj}(\omega^k) \geq \frac{1}{2}(\alpha_{2kj} - \alpha_{1kj})$ for all $j \in J_k$, and accepts the signal hypothesis otherwise.

3.4.2 Performance analysis

Let $k \leq K$ be fixed, and let P_u^k be the probability distribution of observation ω^k , the input being u . Taking into account (49.c) and applying Proposition 2.4 with $\alpha_\chi = \alpha_{\chi kj}$, $\delta = \delta_{kj}$ and $X_1 = A_k V_{\text{nuis}}$, $X_2 = A_k U_j(\rho_{kj})$, we obtain

$$\begin{aligned} (a) \quad & P_u^k \left\{ \phi_{kj}(\omega^k) < \frac{1}{2}(\alpha_{2kj} - \alpha_{1kj}) \right\} \leq \varepsilon_*(\alpha_{1kj}|\gamma), \quad \text{if } u \in V_{\text{nuis}}, \\ (b) \quad & P_u^k \left\{ \phi_{kj}(\omega^k) \geq \frac{1}{2}(\alpha_{2kj} - \alpha_{1kj}) \right\} \leq \varepsilon_*(\alpha_{2kj}|\gamma), \quad \text{if } u \in U_j(\rho_{kj}). \end{aligned} \quad (51)$$

Assume first that the input u is a nuisance, and let us upper-bound the P_u^k -probability of rejecting the nuisance hypothesis at step k . This rejection implies that $\phi_{kj}(\omega^k) < \frac{1}{2}(\alpha_{2kj} - \alpha_{1kj})$ for some $j \in J_k$, and by (51.a) the P_u^k -probability of this event for a given $j \in J_k$ is at most $\varepsilon_\star(\alpha_{1kj}|\gamma)$. As a result, the P_u^k -probability of signal conclusion at step k when u is a nuisance is at most

$$\sum_{j \in J_k} \varepsilon_\star(\alpha_{1kj}|\gamma) \leq \epsilon_k,$$

where the inequality is due to (49.b).

Now let the input u be a signal of shape j and magnitude $> \rho_{kj}$, i.e., $u \in U_j(\rho)$ with $\rho > \rho_{kj}$. Due to $\rho_{kj} = R_j$ for $j \notin J_k$ this means that $j \in J_k$. Let us upper-bound the P_u^k -probability of a nuisance conclusion at step k . The nuisance hypothesis is not rejected only when $\phi_{kj}(\omega^k) \geq \frac{1}{2}(\alpha_{2kj} - \alpha_{1kj})$, and by (51.b) the P_u^k -probability of the latter event does not exceed $\varepsilon_\star(\alpha_{2kj}|\gamma) \leq \epsilon_{kj}$ (recall that by definition of sets $U_j(\rho)$, $u \in U_j(\rho)$ with $\rho > \rho_{kj}$ implies that $u \in U_j(\rho_{kj})$), where the concluding inequality is due to (49.a).

The bottom line is that the P_u^k -probability of false alarm at step k (rejecting the nuisance hypothesis when it is true) is $\leq \epsilon_k$, and the P_u^k -probability of ρ_{kj} -miss (making a nuisance conclusion when the input is a signal of shape j and magnitude $> \rho_{kj}$) is at most ϵ_{kj} .

4 Application: change detection in linear dynamical system

4.1 Problem statement

We consider the change detection problem as follows.

1. We are given a discrete time linear time invariant system

$$\begin{aligned} x_t &= P_t x_{t-1} + Q_t u + R_t \xi, \\ y_t &= C_t x + D_t u + S_t \xi, \quad t = 1, \dots, d, \end{aligned} \tag{52}$$

where

- $x = [x_0; x_1; \dots; x_d]$, $x_t \in \mathbf{R}^{n_x}$ is the state trajectory, $u \in \mathbf{R}^{n_u}$ is the input, $y_t \in \mathbf{R}^{n_y}$ is the output at time t ,
- $\xi \in \mathbf{R}^{n_\xi}$ is a random disturbance with (unknown) distribution P ,
- P_t, \dots, S_t , $1 \leq t \leq d$, are known matrices of appropriate sizes.

Given $\tau \leq d$, we set $y^\tau = [y_1; \dots; y_\tau]$. We denote $E_x = \underbrace{\mathbf{R}^{n_x} \times \dots \times \mathbf{R}^{n_x}}_d$ and similarly for E_u .

2. We are given sets $U_{\text{inp}} \subset E_u$ (admissible inputs), $V_{\text{nuis}} \subset E_u$ (nuisances), V_{drag}^j (drags), $W_j \subset E_u$ (activations of shape j and magnitude ≥ 1), $1 \leq j \leq N$, meeting the requirements of Section 3.2. These sets, exactly as in Section 3.2, give rise to the notions of admissible and nuisance inputs to (52), same as signal inputs of shape j and magnitude $\geq \rho$.
3. The distribution P of disturbance ξ is known to belong to a given family \mathcal{P} of probability densities on \mathbf{R}^{n_ξ} . We are also given tolerances $\epsilon_t \in (0, 1/2)$, $\epsilon_{tj} \in (0, 1/2)$, $1 \leq t \leq d$, $1 \leq j \leq N$.

We acquire observations y_τ one by one, so that at time t we have at our disposal the observation $y^t = [y_1; \dots; y_t]$. Our objective is to design tests \mathcal{T}_t and thresholds $\rho_{tj} > 0$, $1 \leq t \leq d$, $1 \leq j \leq N$, meeting requirements completely similar to those from Section 3.2:

- given observation y^t , the test \mathcal{T}_t should make either a nuisance, or a signal conclusion (but not both);
- if the input to (52) is a nuisance, the probability of the non-nuisance conclusion at time t should be at most ϵ_t , and if, for some j , the input is a signal of shape $j \leq N$ and magnitude $> \rho_{tj}$, the probability of the nuisance conclusion at time t should be at most ϵ_{tj} .

Both these requirements should be satisfied for every $t \leq d$, and we would like to meet them with as small thresholds ρ_{tj} as possible.

We are about to demonstrate that the just outlined change detection problem can be handled via the techniques developed in Sections 3.3 and 3.4. Basically all we need to this end is to convert our observation scheme into the one considered in Section 3, and this is what we are about to do next.

4.2 Building the observation scheme

Given an input u , a noise ξ and $t \leq d$, the observation y^t is not uniquely determined by the input and the noise; it is also affected by the initial state x_0 of the system. To get rid of the influence of the initial condition we act as follows.

1. We denote by F^t the linear subspace of $E_y^t = \{[y_1; \dots; y_t] \in \mathbf{R}^{tn_y}\}$ comprised of all outputs $y^t = [y_1; \dots; y_t]$ which in the noiseless case $\xi = 0$ stem from zero input and some initial state of the system, build an orthonormal basis of the orthogonal complement of F^t in E_y^t and make the vectors of this basis the rows of a matrix, thus arriving at a $\mu_t \times (tn_y)$ matrix M_t . Note that μ_t is a nondecreasing function of t .
2. We set $z^t = M_t[y_1; \dots; y_t]$, where y_τ is given by (52).

It may happen that $\mu_t = 0$ for some t . In this case, our decision rule \mathcal{T}_t by construction accepts the nuisance hypothesis, so that nontrivial decision rules will be associated only with those time instants t for which $\mu_t \geq 1$. Let $t = \kappa + 1$ be the first instant such that the corresponding $\mu_t \geq 1$. Note that time instants t with $\mu_t \geq 1$ form the final segment $\{\kappa + 1, \kappa + 2, \dots, \kappa + K = d\}$ of $1, \dots, d$ (recall that μ_t is nondecreasing in t). We set $m_k = \mu_{\kappa+k}$, $1 \leq k \leq K$. By construction,

$$z^{\kappa+k} = \bar{A}_k u + \bar{B}_k \xi, \quad (53)$$

with some $m_k \times n_u$ matrix \bar{A}_k and $m_k \times n_\xi$ matrix \bar{B}_k readily given by our data.

From now on, we assume that⁸ $K > 0$ and that \bar{B}_k are of full row rank, i.e. $\text{rank}(\bar{B}_k) = m_k$, $1 \leq k \leq K$. Finally, we select somehow invertible $m_k \times m_k$ matrices L_k and pass from observations (53) to observations

$$\omega^k = A_k u + \xi_k, \text{ where } \omega^k = L_k z^{\kappa+k}, A_k = L_k \bar{A}_k, \xi_k = Z_k \xi, Z_k = L_k \bar{B}_k. \quad (54)$$

For example, we can set $L_k = (\bar{B}_k \bar{B}_k^T)^{-1/2}$, thus ensuring that $Z_k Z_k^T = I_{m_k}$.

Notice that observations (54) meet the requirements imposed in Section 3.2. As a result, we find ourselves in the situation considered in Section 3.2 and therefore can apply to the change detection problem in question the machinery developed in Sections 3.3 and 3.4.

4.3 Illustration: detecting changes in the trend of a simple time series

We consider the example of Section 3.1 and apply, with minor modifications, the construction outlined in Section 3.4.

4.3.1 Constructing decision rules

To attain our goal (see Section 3.1) we act as follows. The rule \mathcal{T}_1 is trivial – it always accepts the nuisance hypothesis. To describe how \mathcal{T}_k , $k > 1$, is built, let us fix $k \in \{2, \dots, d\}$.

1. *Building observation ω^k .* Setting $\xi = [\eta; \zeta]$ and $Z_k[\eta; \zeta] = \Theta_k^{-1} [D_k \eta + E_k \zeta]$, our observations (39) are as required in (40), and we meet Assumption A.

In our implementation, we use $\Theta_k = \Omega_k^{1/2}$, where Ω_k is the minimum trace matrix satisfying $\Omega_k \succeq D_k D_k^T$, $\Omega_k \succeq E_k E_k^T$.

We now apply the decision rules of Scheme II from Section 3.4 to observation (39). To define parameters $J_k, \rho_{kj}, \alpha_{1kj}$ and α_{2kj} we proceed as follows.

⁸Observe that when $K = 0$ our approach results in trivial tests always accepting the nuisance hypothesis – in this case the input-related component in the observations is fully masked by the influence of initial condition x_0 .

2. We set

$$\widehat{J}_k = \{1, \dots, 2k\}; \widehat{\epsilon} = \frac{\epsilon}{d(d+1)-2}; \epsilon_k = 2k\widehat{\epsilon}; \epsilon_{kj} = \epsilon, j \in \widehat{J}_k; \rho_{kj} = R, j \notin \widehat{J}_k. \quad (55)$$

For $j \in \widehat{J}_k$, we specify $\alpha_{1kj}, \alpha_{2kj}, \delta_{kj}$ by the relations

$$\int_{\alpha_{1kj}}^{\infty} \gamma(s) ds = \widehat{\epsilon}, \int_{\alpha_{2kj}}^{\infty} \gamma(s) ds = \epsilon_{kj} = \epsilon, \delta_{kj} = \frac{1}{2}[\alpha_{1kj} + \alpha_{2kj}].$$

Note that, by construction, setting $\epsilon_1 = 0$ we have

$$\forall k : \sum_{j \in \widehat{J}_k} \int_{\alpha_{1kj}}^{\infty} \gamma(s) ds \leq \epsilon_k, \sum_{k=1}^d \epsilon_k = \epsilon \text{ and } \alpha_{1kj} + \alpha_{2kj} = 2\delta_{kj}, j \in \widehat{J}_k,$$

cf. (49).

3. For $j \in \widehat{J}_k$, we consider the convex optimization problem (cf. $(P_{kj}[\rho])$). $\text{Opt}_{kj}(\rho)$ clearly is continuous and nonincreasing in $\rho > 0$ and $\lim_{\rho \rightarrow +0} \text{Opt}(\rho) = 0$. When $\text{Opt}_{kj}(R) \leq \delta_{kj}$ we set $\rho_{kj} = R$. Otherwise, we find the smallest $\rho = \rho_{kj}$ such that $\text{Opt}_{kj}(\rho) \geq \delta_{kj}$; note that in the latter case we have $\rho_{kj} \in (0, R)$ and $\text{Opt}_{kj}(\rho_{kj}) = \delta_{kj}$. We have specified ρ_{kj} for all $j \in \widehat{J}_k$.
4. We set $J_k = \{j \in \widehat{J}_k : \rho_{kj} < R\}$ thus ensuring that $\rho_{kj} < R$ if and only if $j \in J_k$.⁹
5. Same as in Section 3.4, for $j \in J_k$ we denote by u_{kj} an optimal solution to problem $(P_{kj}[\rho])$ with $\rho = \rho_{kj}$ and set (cf. (50))

$$h_{kj} = -\frac{A_k u_{kj}}{\|A_k u_{kj}\|_2} = -\frac{A_k u_{kj}}{2\delta_{kj}}, \quad c_{kj} = \frac{1}{2} h_{kj}^T A_k u_{kj}, \quad \phi_{kj}(\omega) = h_{kj}^T \omega - c_{kj}.$$

6. Finally, given observation ω^k , our rule \mathcal{T}_k makes the nuisance conclusion if and only if $\phi_{kj}(\omega_k) \geq \frac{1}{2}[\alpha_{2kj} - \alpha_{1kj}]$ for all $j \in J_k$, and makes signal conclusion otherwise, cf. Section 3.4. Invoking the results of Section 3.4.2, the decision rules we have built do satisfy risk specifications of Section 3.1.

4.3.2 Quantifying conservatism: performance indexes

Let us pass from intermediate observations (38) to observations

$$w^k = [D_k D_k^T]^{-1/2} z^k = \underbrace{Q_k \eta + S_k \zeta}_{\lambda^k} + Q_k u, \quad (56)$$

where $Q_k = [D_k D_k^T]^{-1/2} D_k$ satisfies $Q_k Q_k^T = I_{k-1}$ and $S_k = [D_k D_k^T]^{-1/2} E_k$. Since Q_k has orthonormal rows, specifying the distribution of η to be multivariate Student's $t_d(\nu, I_d)$ distribution on \mathbf{R}^d with ν degrees of freedom and unit matrix parameter, the distribution $p(\cdot)$ of the random variable $Q_k \eta$ will be multivariate Student's $t_{k-1}(\nu, I_{k-1})$ distribution. Now, let θ be the smallest nonvanishing singular value of S_k (or, equivalently, θ^2 is the smallest eigenvalue of $S_k S_k^T$). Clearly, we can specify the covariance matrix Σ of zero mean d -dimensional Gaussian random vector ζ to satisfy $\Sigma \preceq \sigma^2 I_d$ and to be such that the covariance matrix of $S_k \zeta$ is $\theta^2 \sigma^2 I_{k-1}$. We conclude that we can point out distributions of η and ζ , satisfying specifications of the model (36), and such that the random noise λ^k in (56) will be the sum of two independent zero-mean random vectors, one with $(k-1)$ -dimensional Student distribution $t_{k-1}(\nu, I_{k-1})$, and the other – Gaussian, with covariance matrix $(\theta\sigma)^2 I_{k-1}$. As it is immediately seen, λ^k has a probability density $p(\cdot)$ of the form $f(\|\cdot\|_2)$ with nonincreasing f and whenever $e \in \mathbf{R}^{k-1}$ is a unit vector, the probability density $\gamma_k(\cdot)$ of the scalar random variable $e^T \lambda^k$ is

$$\gamma_k = \gamma_S \star \gamma_{\theta\sigma}.$$

⁹Note that, typically, $|J_k| < |\widehat{J}_k|$. Thus, one can easily improve the estimation procedure by better accounting for the “remaining at step k ” part of false alarm probability. A simple “dynamic” management of false alarm probabilities of tests is implemented in the numerical experiments described in the next section. The detailed construction is presented in the Online complement of the paper available at <http://arxiv.org/abs/1705.07196>.

Here, as above, γ_S is the density of the univariate Student's t_ν distribution, and $\gamma_{\theta_k\sigma}$ is the density of $\mathcal{N}(0, (\theta\sigma)^2)$.

Now let V_1, V_2 be two closed convex sets in the space \mathbf{R}^d of inputs such that the sets $Q_k V_\chi$, $\chi = 1, 2$, are closed, and one of these two sets is bounded, and let (u_*^1, u_*^2) be an optimal solution to the convex optimization problem

$$\delta = \min_{u^1 \in V_1, u^2 \in V_2} \frac{1}{2} \|Q_k[u^1 - u^2]\|_2. \quad (57)$$

By Remark 2.1, no test based on observation (56) can decide on two simple hypotheses $u = u_*^1$, $u = u_*^2$ with risk $< \int_{\frac{\delta}{2}}^{\infty} \gamma_k(s) ds$, implying that the same lower risk bound also holds true for all tests utilizing observations y^k rather than w^k .

Now let $V_1 = \{0\}$ and $V_2 = U_j(\rho)$, so that δ as defined in (57) becomes a (clearly, continuous and non-increasing when $\rho > 0$) function $\delta_{k,j}(\rho)$ of ρ . Let us define $\rho_{k,j}^*$ as follows: if $\int_{\delta_{k,j}(R)}^{\infty} \gamma_k(s) ds > \epsilon$, we set $\rho_{k,j}^* = R$, otherwise $\rho_{k,j}^* \in (0, R]$ is the smallest $\rho > 0$ such that $\int_{\delta_{k,j}(\rho)}^{\infty} \gamma_k(s) ds \leq \epsilon$. By construction, for every $\rho \in (0, \rho_{k,j}^*)$ there is no test which, given an observation y^k , would decide with risk $\leq \epsilon$ on the hypothesis “ $u = 0$ ” vs. the alternative “ u is a signal from $U_j(\rho)$ with $\rho > 0$.” It is natural to quantify the conservatism of our decision rules \mathcal{T}_k by the *performance indexes* $\rho_{k,j}/\rho_{k,j}^*$; the less are these indexes, the less is the conservatism.

4.3.3 Numerical results

We operate on time horizon $d = 8$ and deal with $\sigma = 1$ and with 5 values of the number ν of degrees of freedom of the Student distribution of η , specifically, the values 1, 2, 3, 6, ∞ . In the experiments of this section we use parameter values $R = 10^4$, $\epsilon = 0.01$.¹⁰ The range of parameters $\rho_{k,j}$ and of ratios $\rho_{k,j}/\rho_{k,i}^*$ are presented in Figure 1. Some comments are in order.

1. Relation $U_{2i-1}(\rho) = -U_{2i}(\rho)$ implies that $\rho_{k,2i-1} = \rho_{k,2i}$, $\rho_{k,2i-1}^* = \rho_{k,2i}^*$, $1 \leq i \leq d$.
2. We display the range of quantities $\rho_{k,j}$ and $\rho_{k,j}/\rho_{k,j}^*$ only for those values of k, j for which $\rho_{k,j}^* < R$, that is, ignore pairs k, j for which already an optimistic lower bounds $\rho_{k,j}^*$ on the magnitude of signal inputs of shape j which can be detected, with the required risk, at time k should be $\geq R$, which is forbidden by (37). On a closest inspection, the ignored pairs k, j are the pairs of the form $k, j = 2i - 1$ and $k, j = 2i$ where
 - (a) $i > k$, or
 - (b) $k = 1$, or
 - (c) [only for pulses!] $i = 1 < k$.

The reasons are clear: (a) stems from the fact that at time k it is impossible to detect a whatever large signal input which starts at time $i > k$. (b) reflects the fact that the contribution of a whatever large signal input, if any, to the very first observation is fully masked by the initial condition α_0 , and in our model we do not impose any restrictions on this initial condition. (c) is of a similar origin: when the signal inputs are pulses, of a whatever magnitude, at time 1, are fully masked by the initial conditions and thus cannot be detected at all.

3. As it could be guessed, when the signal inputs are steps, the quantities $\rho_{k,j}$, for j fixed, decrease as $k \geq j$ grows, since influence of step-change on our observations accumulates with time. In contrast, no such phenomenon is observed for pulse signal inputs where there is “nothing to accumulate.”¹¹
4. The conservatism of our decision rules, as presented in the tables, while unpleasant, seems to be not too high when $\nu \geq 3$, and becomes really arresting when $\nu = 1$. The origin of this phenomenon is quite transparent. The conservatism seems to stem primarily from systematic use in our constructions, for absence of something better, of the union bounds for probabilities. For example, when computing $\rho_{k,j}^*$, we allow for the probability of false alarm at time k to be as large as ϵ , while in our decision rules, we “distribute” this probability between d instants where we make our decisions. Similarly, when computing $\rho_{k,j}^*$, we act as if the only alternative to the nuisance hypothesis were a particular signal hypothesis $u \in U_j(\rho)$, while in fact we have several signal hypotheses to consider and should take into account the

¹⁰We use **Mosek** and YALMIP Matlab toolbox [1, 16] to solve corresponding optimization problems.

¹¹Or, rather, that the noise in the states α_k of the model accumulates at the same rate, thus cancelling the effect of the growing observation sample.

resulting “accumulation of risk.” As a result, we require from pairwise tests participating in \mathcal{T}_k to have risk essentially smaller than ϵ , which in the case of a “heavy tail” noise distribution allowed by our model requires an essentially larger magnitudes of detectable signal inputs than those allowing for detection when the shape of signal input is known in advance. And indeed, we see that the ratios ρ_{kj}/ρ_{kj}^* rapidly increase as ν decreases.

We report on Figure 2 the evolution of parameters ρ_{kj} for $k = 8$ for step and pulse signals as a function of the risk of the test ϵ and the standard deviation σ of the Gaussian component ζ_t of the noise. As expected, ρ_{kj} increase with σ and when ϵ decreases.

4.3.4 Change detection in linear dynamic system revisited

The methodology developed in Section 3.4 allows for a straightforward refinement which hopefully improves the resulting inference performance. Note that the inference rules, as given in Sections 3.3 and 3.4 use very conservative bound for the probability of false alarm – for multiple tests this probability is simply the sum of probabilities of false rejections of the nuisance hypothesis for each test. This results in the increase of the testing thresholds ρ_{kj} , which amounts to a “logarithmic factor” in the case of Gaussian observation noise, but becomes much more severe in the case of a heavy-tail noise distribution. One way to make the decision less cautious is to reduce the number of hypotheses to test by aggregating the alternatives. Here we illustrate the idea of the proposed modification on the simple numerical example in section 4.3. Specifically, when building the detection procedure, at time k we act as follows:

- We compute the quantities ρ_{kj}^* , $j = 1, \dots, 2d$, and denote by J_k^* the set of those j for which $\rho_{kj}^* < R$. As it was explained, there is no reason to bother to detect at time k signal inputs of shape $j \notin J_k^*$.
- Assume we have somehow associated thresholds $\rho_{kj} \geq \rho_{kj}^*$ to indexes $j \in J_k^*$; our goal, same as before, is to build a decision rule \mathcal{T}_k which, given ω^k ,
 - with probability at least $1 - \epsilon$, makes signal conclusion at time k , provided the input belongs to $U_j(\rho_{kj})$ with some $j \in J_k^*$;
 - has false alarm probability at time k (the probability to make signal conclusion when the input is a nuisance) $\leq \epsilon_k$, $\sum_k \epsilon_k = \epsilon$.

Next, let us color the sets $U_j(\rho_{kj})$, $j \in J_k^*$ (and their indexes) in a number L of colors; let I_ℓ be the set of indexes $j \in J_k^*$ colored by color ℓ . We associate with each ℓ a convex alternative U^ℓ – the convex hull of “alternatives of color ℓ ”:

$$U^\ell = \text{Conv} \left(\bigcup_{j \in I_\ell} U_j(\rho_{kj}) \right), \ell = 1, \dots, L,$$

and replace the original detection problem with the following one: given observation ω^k we want to decide on the null hypothesis H_0 “the input is nuisance” (in our case, zero) vs. the alternative H_1 “the input belongs to $\bigcup_{\ell=1}^L U^\ell$.” Same as before, our goal is to ensure probability of false alarm $\leq \epsilon_k$ and probability of miss $\leq \epsilon$.

Note that if we are able to do so, we meet our initial design specifications – with input from $U_j(\rho_{kj})$ for some $j \in J_k^*$, the probability of signal conclusion at time k will be at least $1 - \epsilon$. To build the decision rule, let us use pairwise tests given by the construction from Section 3.4: we need L tests, ℓ -th of them deciding on H_0 vs. the alternative $H_1^\ell : u \in U^\ell$, with risks ϵ_k^ℓ of false alarm and ϵ of miss. If we can ensure $\sum_{\ell=1}^L \epsilon_k^\ell \leq \epsilon_k$, we are done – the decision rule which makes nuisance conclusion when all our L tests “vote” for H_0 , and makes the signal conclusion otherwise, is what we are looking for. Note that the original construction in Section 4.3 is of exactly this structure, with $L = \text{Card}(J_k^*)$ (i.e., every $U_j(\rho_{kj})$, $j \in J_k^*$, has its own color). It could make sense, however, to use aggregated alternatives, thus reducing the number of colors. When doing so,

— on one hand, it is more difficult to decide on H_0 vs. the alternatives, because the image $A_k U^\ell$ of the “aggregated” set of signal inputs is closer to the image $\{0\}$ of the nuisance set than the images $A_k U_j(\rho_{kj})$, $j \in I_\ell$, of individual sets of signal inputs participating in the aggregation. Therefore, to ensure the same risk, we now need a somewhat larger separation of the images of the nuisance and the signal inputs in the observation space;

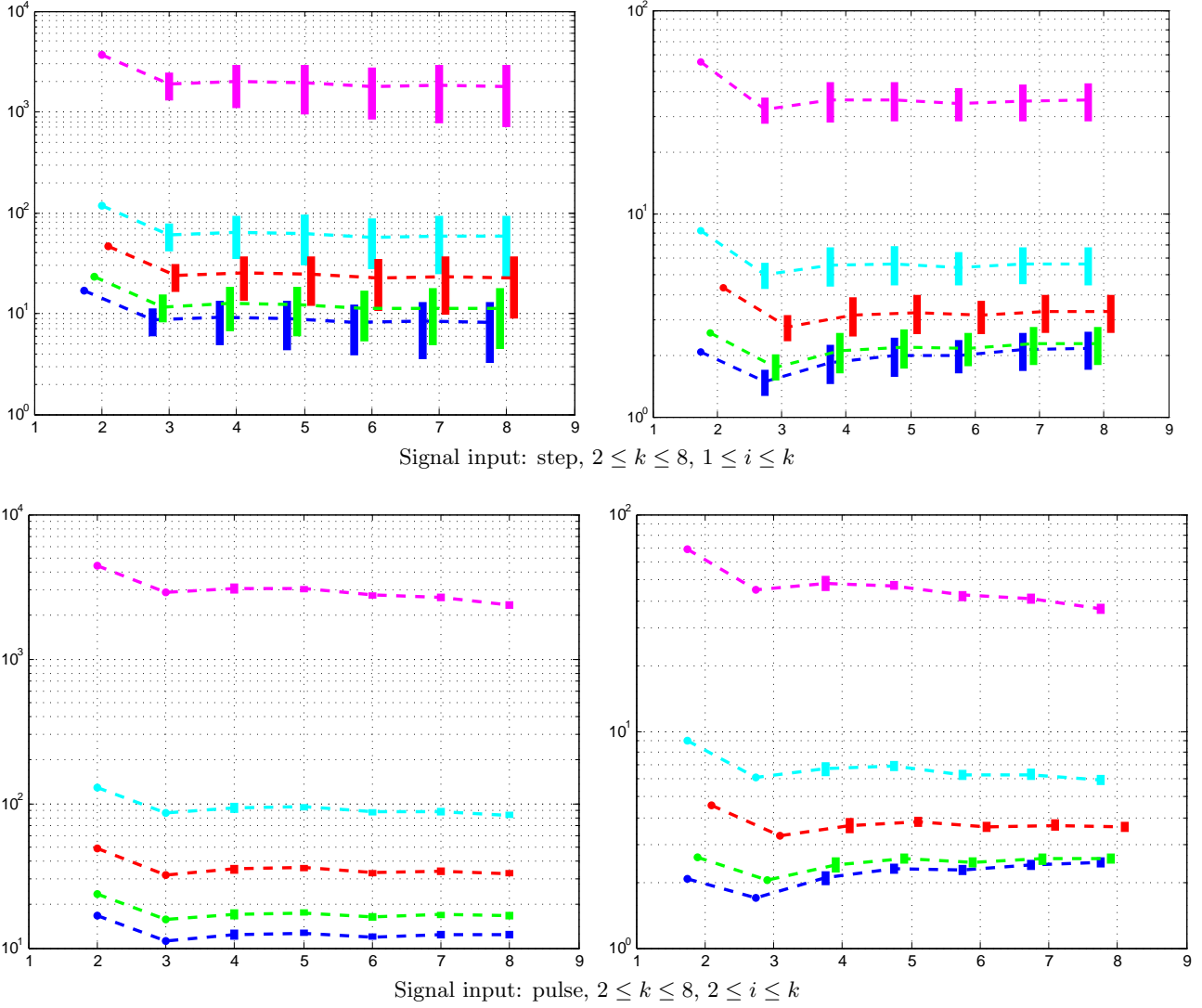


Figure 1: Detecting changes in the trend of a simple time series. Blue/green/red/cyan/magenta: $\nu = \infty/6/3/2/1$. Left plots: ranges (vertical segments) of $\rho_{k,2i-1} = \rho_{k,2i}$, vs. k . Right plots: ranges (vertical segments) of performance indexes $\rho_{k,2i-1}/\rho_{k,2i-1}^* = \rho_{k,2i}/\rho_{k,2i}^*$ vs. k . Ranges of i and k cover the domain where $\rho_{k,2i}^* = \rho_{k,2i-1}^* < R = 10^4$. Charts are shifted horizontally to improve the plot readability. On these plots both $\rho_{k,2i-1}$ and $\rho_{k,2i-1}/\rho_{k,2i-1}^*$ decrease with ν .

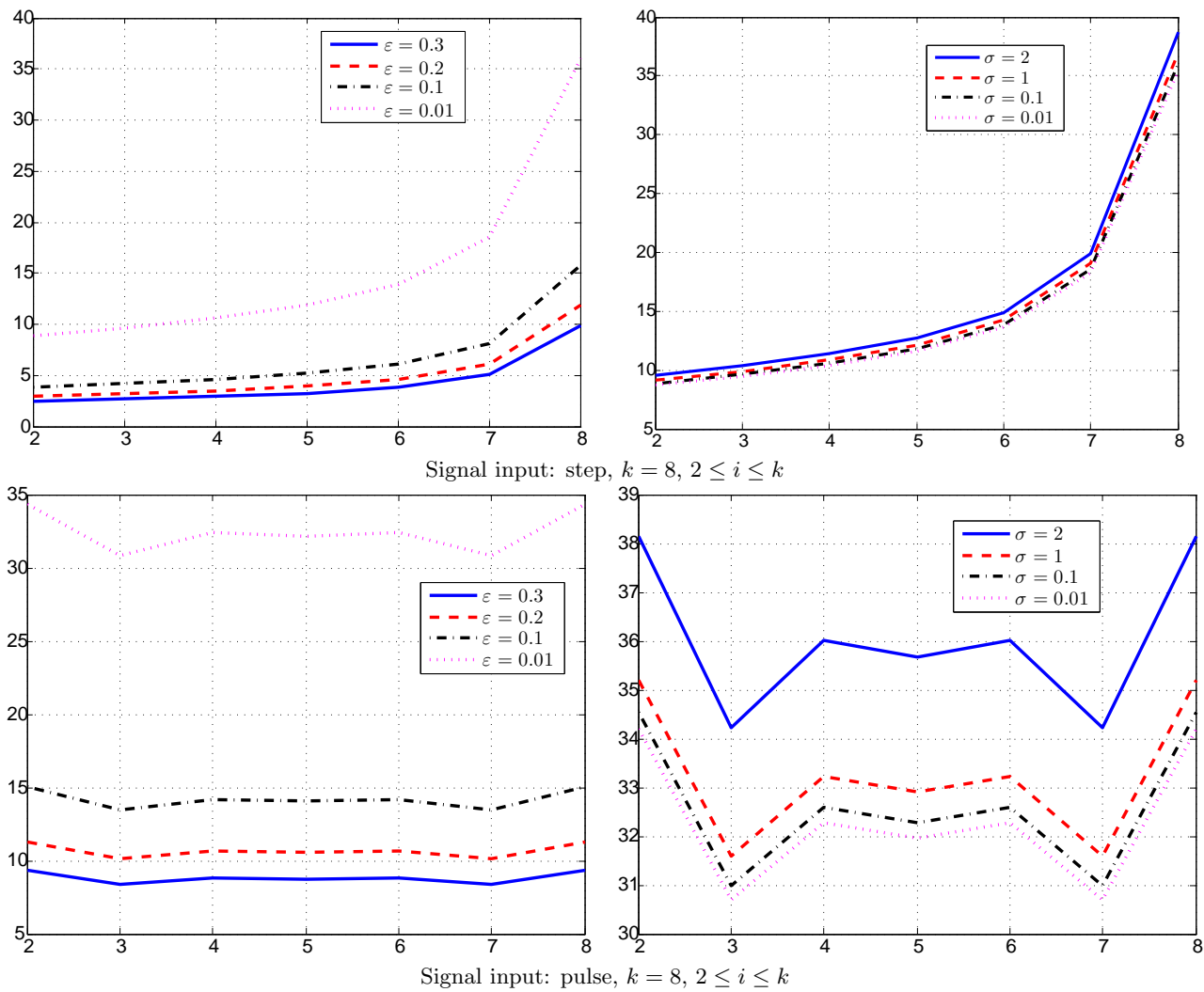


Figure 2: Coefficients $\rho_{k,2i-1} = \rho_{k,2i}$ for $\nu = 3$, $k = 8$, and $i = 2, \dots, 8$, as a function of the desired risk ϵ of the test (left pane) and the standard deviation σ of the Gaussian noise ζ_t (right pane).

— on the other hand, we should now “distribute” ϵ_k among $L < \text{Card}(J_k^*)$ miss probabilities ϵ_k^ℓ . This would allow to operate with larger miss probabilities, thus reducing the necessary separation of the images of the nuisance and the signal inputs in the observation space.

It is hard to tell in advance which of these two opposite effects will prevail; an answer, however, could be provided by computation, and it makes sense to give to the outlined modification a try. To make things as simple as possible, let us act as follows.

- After the colors are assigned and the sets U^ℓ , $\ell \leq L$, are built, we specify ϵ_k^ℓ and the quantities α_{1k} , α_{2k} , δ_k to meet the requirements

$$\epsilon_k^\ell = \frac{\epsilon_k}{L}, \quad 1 \leq \ell \leq L; \quad \int_{\alpha_{1k}}^{\infty} \gamma(s) ds = \frac{\epsilon_k}{L}; \quad \int_{\alpha_{2k}}^{\infty} \gamma(s) ds = \epsilon; \quad \delta_k = \frac{1}{2}[\alpha_{1k} + \alpha_{2k}].$$

Let us suppose that the relation

$$\min_{u \in U^\ell} \frac{1}{2} \|A_k u\|_2 \geq \delta_k, \quad \ell = 1, \dots, L, \quad (58)$$

is satisfied. We set

$$h_{k\ell} = A_k u_{k\ell} / \|A_k u_{k\ell}\|_2, \quad c_{k\ell} = \frac{1}{2} h_{k\ell}^T A_k u_{k\ell}, \quad \phi_{k\ell}(\omega^k) = h_{k\ell}^T \omega_k - c_{k\ell},$$

where $u_{k\ell}$ are optimal solutions to the optimization problems in (58). It is immediately seen that making at time k the nuisance conclusion if and only if $\phi_{kj}(\omega^k) \geq \frac{1}{2}[\alpha_{2k} - \alpha_{1k}]$, we ensure simultaneously the probability of false alarm at time k at most ϵ_k , and the probability of miss when the input belongs to $\bigcup_{j \in J_k^*} U_j(\rho_{kj})$ at most ϵ , thus meeting our design specifications.

Specifying ρ_{kj} . The question we did not address so far is how to choose ρ_{kj} , $j \in J_k^*$. What we expect of these quantities is to ensure the validity of (58), and the simplest way to achieve this goal is as follows. Setting $\rho_{kj} = \theta \rho_{kj}^*$, the left hand side in (58) is a nondecreasing function of θ , and we can find by bisection the smallest $\theta = \theta_k \geq 1$ for which (58) takes place. After θ_k is found, we set $\rho_{kj} = \theta_k \rho_{kj}^*$, $j \in J_k^*$. Clearly, with this approach, the performance indexes ρ_{kj} / ρ_{kj}^* , $j \in J_k^*$, are all equal to θ_k .

How it works. We applied the just outlined construction to the data underlying the numerical experiment described in Section 4.3. Our implementation was the simplest possible: given k , we looked at all j 's such that $\rho_{kj}^* < R$; the set J_k^* of these j 's with our data is nonempty only when $k \geq 2$ and is either $\{1, \dots, 2k\}$ (step signals), or $\{3, 4, \dots, 2k\}$ (pulse signals). A set $U_j(\rho)$ with odd index $j = 2i - 1$ (even index $j = 2i$), $j \in J_k^*$, is comprised of signals which are zero before time i and “jump up/jump down” at time i depending on whether j is even or odd. We color these sets in $L = 2$ colors, depending on whether the corresponding indexes j are odd or even, that is, we use at step k

$$U^\ell = \bigcup_{j \in J_k^*, j \bmod 2 = \ell} U_j(\theta_k \rho_{kj}^*), \quad \ell = 1, 2,$$

with θ_k as explained above.

We present on Figure 3 the comparison of performance indexes ρ_{kj} / ρ_{kj}^* , $j \in J_k^*$, for the original inference routine (these indexes are presented on Figure 1) and the performance indexes of the just described modified inference. For our initial routine, the performance indexes slightly vary with $j \in J_k^*$, and we present their ranges; for the new routine, the performance indexes do not depend on j . We observe that in the considered example the proposed straightforward aggregation of signal inputs typically results in degradation of the performance indexes, and improves these indexes significantly for the “heavy tailed” noise distributions (the case of $\nu = 1$).

References

- [1] E. D. Andersen and K. D. Andersen. *The MOSEK optimization toolbox for MATLAB manual. Version 7.0*, 2013. <http://docs.mosek.com/7.0/toolbox/>.

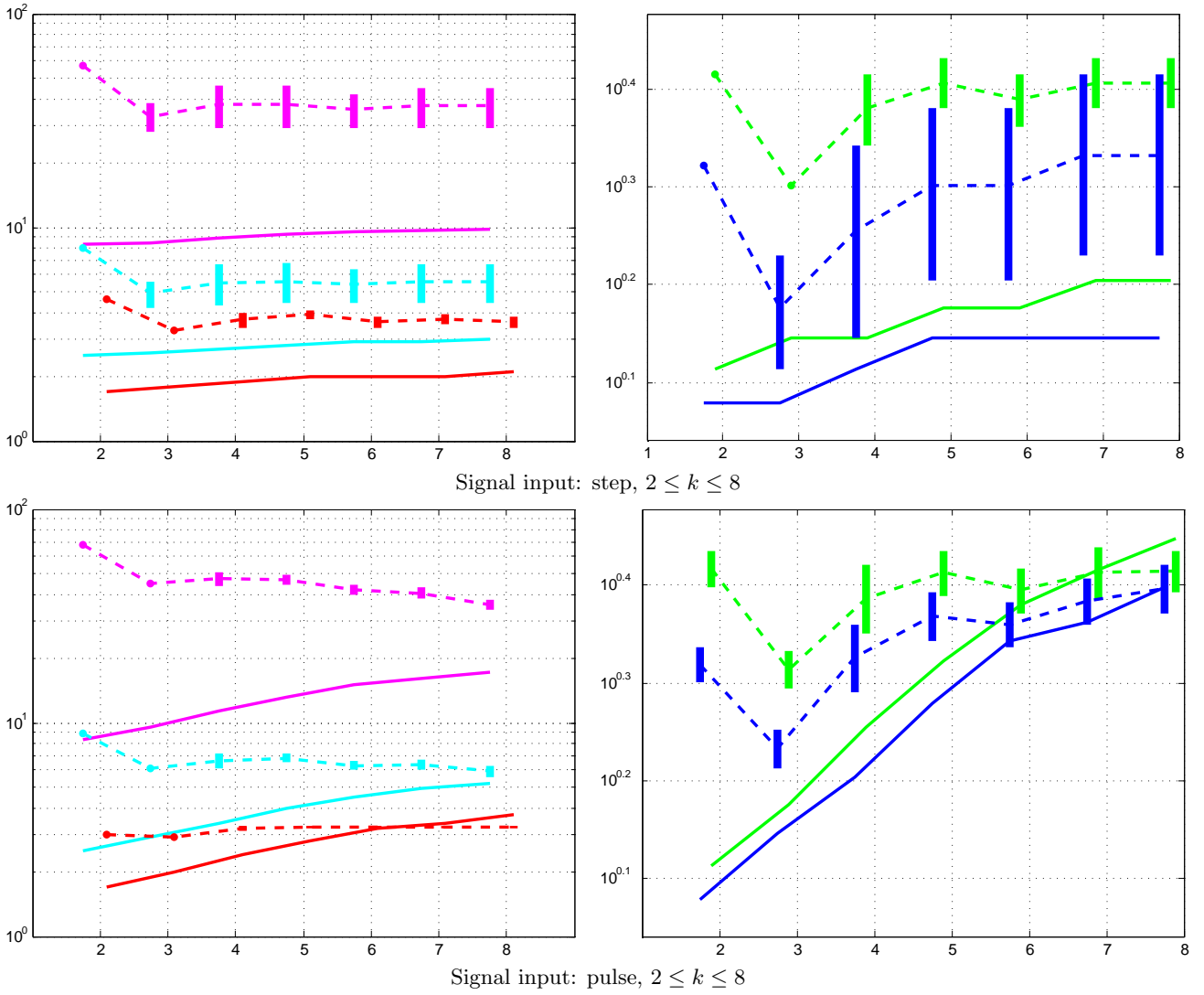


Figure 3: Ranges of performance indexes $\rho_{k,j}/\rho_{k,j}^*$ (cf. right plots of Figure 1) as compared to the performance index for the refined inference (solid lines). Red/cyan/magenta: $\nu = 3/2/1$ (left plots), blue/green: $\nu = \infty/6$ (right plots). Charts are shifted horizontally to improve the plot readability. On these plots $\rho_{k,j}/\rho_{k,j}^*$ decrease with ν .

- [2] A. Ben-Tal and A. Nemirovski. *Lectures on modern convex optimization: analysis, algorithms, and engineering applications*, volume 2. Siam, 2001.
- [3] M. Burnashev. On the minimax detection of an inaccurately known signal in a white noise background. *Theory Probab. Appl.*, 24:107–119, 1979.
- [4] M. Burnashev. Discrimination of hypotheses for gaussian measures and a geometric characterization of the gaussian distribution. *Math. Notes*, 32:757–761, 1982.
- [5] Y. Cao, V. Guigues, A. Juditsky, A. Nemirovski, and Y. Xie. Change detection via affine and quadratic detectors. *arXiv preprint arXiv:1608.00524*, 2016.
- [6] H. Chernoff. A measure of asymptotic efficiency for tests of a hypothesis based on the sum of observations. *The Annals of Mathematical Statistics*, pages 493–507, 1952.
- [7] L. Dumbgen and V. G. Spokoiny. Multiscale testing of qualitative hypotheses. *Annals of Statistics*, pages 124–152, 2001.
- [8] T. Eltoft, T. Kim, and T.-W. Lee. On the multivariate laplace distribution. *IEEE Signal Processing Letters*, 13(5):300–303, 2006.
- [9] A. Goldenshluger, A. Juditsky, and A. Nemirovski. Hypothesis testing by convex optimization. *Electronic Journal of Statistics*, 9(2):1645–1712, 2015.
- [10] M. Grant and S. Boyd. *The CVX Users’ Guide. Release 2.1*, 2014. <http://web.cvxr.com/cvx/doc/CVX.pdf>.
- [11] Y. Ingster and I. A. Suslina. *Nonparametric goodness-of-fit testing under Gaussian models*, volume 169 of *Lecture Notes in Statistics*. Springer, 2002.
- [12] A. Juditski and A. Nemirovski. On sequential hypotheses testing via convex optimization. *Automation and Remote Control*, 76:809–825, 2015.
- [13] A. Juditsky and A. Nemirovski. Hypothesis testing via affine detectors. *Electronic Journal of Statistics*, 10:2204–2242, 2016.
- [14] S. Kotz and S. Nadarajah. *Multivariate t-distributions and their applications*. Cambridge University Press, 2004.
- [15] E. L. Lehmann and J. P. Romano. *Testing statistical hypotheses*. Springer Science & Business Media, 2006.
- [16] J. Löfberg. Yalmip : A toolbox for modeling and optimization in matlab. *In Proceedings of the IEEE CACSD Conference*, 2004.

A Proofs

A.1 Proof of Proposition 2.3

In what follows, for functions $f, g : \mathbf{R} \rightarrow \mathbf{R}$, we say that f dominates g (notation: $f \succeq g$, or, equivalently, $g \preceq f$), if

$$\int_{\delta}^{\infty} f(s)ds \geq \int_{\delta}^{\infty} g(s)ds, \forall \delta \geq 0.$$

Let

- \mathcal{E} be the family of even probability densities on the real axis,
- \mathcal{N} be the family of nice functions on the axis.

1°. Note that \succeq clearly is transitive: if $f \succeq g$ and $g \succeq h$, then $f \succeq h$. Furthermore, $\mathcal{N} \subset \mathcal{E}$ (by definition of \mathcal{N}), and \mathcal{E} is closed with respect to taking convolution (evident). We also need the following technical facts.

1°.a \mathcal{N} is closed with respect to taking convolution.

Indeed, let $f, g \in \mathcal{N}$. The fact that $f \star g$ is even and continuous on the real axis is evident. Therefore in order to show that $f \star g \in \mathcal{N}$ it suffices to verify that $h = f \star g$ is nonincreasing on the nonnegative ray. For $0 < z \leq f(0)$, denoting $\bar{y}(z) = \min \{y \geq 0 : f(y) = z\}$, we have for every $x \geq 0$,

$$h(x) = \int_{-\infty}^{\infty} f(x-y)g(y)dy = \int_{-\infty}^{\infty} \left(\int_0^{f(x-y)} dz \right) g(y)dy = \int_0^{f(0)} H(x, z) dz$$

where $H(x, z) = \int_{x-\bar{y}(z)}^{x+\bar{y}(z)} g(y)dy$ for $0 < z \leq f(0)$ and $H(x, 0) = 0$. To conclude we observe that for every fixed $f(0) \geq z \geq 0$, $H(\cdot, z)$ is a differentiable, nonnegative, and nonincreasing function on the nonnegative ray.

Indeed, for $0 < z \leq f(0)$, the derivative $H'_x(x, z)$ of this function at x is $H'_x(x, z) = g(x + \bar{y}(z)) - g(x - \bar{y}(z))$ and it is clear that this quantity is ≤ 0 since $g \in \mathcal{N}$. We have checked that for every fixed $z \geq 0$, $H(\cdot, z)$ is nonincreasing on the nonnegative ray. It follows that h is nonincreasing on the nonnegative ray.

1°.b Let $\bar{f}, f \in \mathcal{E}$, $g \in \mathcal{N}$, and let $\bar{f} \succeq f$. Then $\bar{f} \star g \succeq f \star g$.

Let us verify that for all $x \geq 0$, $\bar{H}(x) \leq H(x)$, where \bar{H} and H are cumulative distribution functions (c.d.f.) of the densities $\bar{f} \star g$ and $f \star g$, respectively. Observe that

$$H(x) = \int_{-\infty}^{\infty} g(s)F(x-s)ds = \int_{-\infty}^{\infty} g(x-t)F(t)dt,$$

same as

$$\bar{H}(x) = \int_{-\infty}^{\infty} g(x-t)\bar{F}(t)dt,$$

where F and \bar{F} are the c.d.f.'s of the densities f and \bar{f} , respectively. Thus, when setting $\Delta(s) = F(s) - \bar{F}(s)$, we get

$$\begin{aligned} H(x) - \bar{H}(x) &= \int_{-\infty}^{\infty} g(x-t)\Delta(t)dt = \int_0^{\infty} g(x-t)\Delta(t)dt + \int_0^{\infty} g(x+t)\Delta(-t)dt \\ [\text{because } \Delta(-t) = -\Delta(t)] &= \int_0^{\infty} \Delta(t)[g(x-t) - g(x+t)]dt \geq 0, \end{aligned}$$

where the final " \geq " is due to $\Delta(t) \geq 0$ and $g(x-t) \geq g(x+t)$ for $x, t \geq 0$ (recall that $g \in \mathcal{N}$) for $t \geq 0$.

1°.c. Let $f \in \mathcal{E}$, $\rho \in (0, 1)$, and let $f_\rho(s) = \rho^{-1}f(\rho^{-1}s)$. Then $f_\rho \in \mathcal{E}$ and $f_\rho \preceq f$.

Indeed, the inclusion $f_\rho \in \mathcal{E}$ is obvious. On the other hand, if $\xi \sim f$, one has $P(\rho\xi \geq x) = P(\xi \geq x/\rho) \leq P(\xi \geq x)$ for $0 < \rho < 1$ and $x \geq 0$, what is exactly $f_\rho \preceq f$.

1°.d. Let $f \in \mathcal{N}$ and $g \in \mathcal{E}$. Then $f \star g \succeq f$.

Note that the c.d.f. H of $f \star g$ satisfies

$$H(x) = \int_{-\infty}^{\infty} g(s)F(x-s)ds = \int_0^{\infty} g(s)[F(x-s) + F(x+s)]ds,$$

where F is the c.d.f. of f (recall that $g \in \mathcal{E}$), and therefore for $x \geq 0$ it holds

$$H(x) - F(x) = \int_0^{\infty} g(s)[F(x-s) + F(x+s) - 2F(x)]ds \leq 0$$

due to the concavity of F on \mathbf{R}_+ .

2°. Now let $p(\cdot) \in \mathcal{P}^n$, $q(\cdot) \in \mathcal{P}^m$, and let $e \in \mathbf{R}^r$, $\|e\|_2 = 1$, be given. Let us set $e_1 = A^T \Theta^{-1}e$, $e_2 = B^T \Theta^{-1}e$, let $\eta \sim p$ and $\zeta \sim q$ be independent, and let $\xi = \Theta^{-1}[A\eta + B\zeta]$. We have

$$\omega := e^T \xi = [A^T \Theta^{-1}e]^T \eta + [B^T \Theta^{-1}e]^T \zeta = e_1^T \eta + e_2^T \zeta.$$

Observe that $e_1^T e_1 = e^T \Theta^{-1} A A^T \Theta^{-1} e \leq 1$ due to $\Theta^{-1} A A^T \Theta^{-1} \preceq I_r$, and for similar reasons $e_2^T e_2 \leq 1$. We denote $\rho_1 = \|e_1\|_2$ and $\rho_2 = \|e_2\|_2$, so that $\rho_1, \rho_2 \in [0, 1]$. When $\rho_\chi > 0$, $\chi = 1, 2$, we set $\bar{e}_\chi = \rho_\chi^{-1} e_\chi$. Now, let $f_1 \in \mathcal{E}$ (respectively, $f_2 \in \mathcal{E}$) be the density of the scalar random variable $e_1^T \eta$ (respectively, $e_2^T \zeta$). Note that when $\rho_1 > 0$ ($\rho_2 > 0$) random variable $e_1^T \eta$ ($e_2^T \zeta$) indeed has a density, and this density is even. Let also $\bar{f}_1 \in \mathcal{E}$ ($\bar{f}_2 \in \mathcal{E}$) be the density of $\bar{e}_1^T \eta$ ($\bar{e}_2^T \zeta$).

2⁰.a. Assume for a moment that $\rho_1 > 0$ and $\rho_2 > 0$. Then

- $p \in \mathcal{P}_\mu^n$, whence $f_1 \in \mathcal{E}$, $\bar{f}_1 \in \mathcal{E}$ and $f_1(s) = \rho_1^{-1} \bar{f}_1(\rho_1^{-1}s)$, whence $f_1 \preceq \bar{f}_1$ by **1^o.c**. Since $\bar{f}_1 \preceq \mu$ and \preceq is transitive, we have also $f_1 \preceq \mu$.
- $q \in \mathcal{P}^m$, and \mathcal{P}^m is a completely monotone subfamily of \mathcal{P}_ν^m , whence $f_2 \in \mathcal{N}$, $\bar{f}_2 \in \mathcal{N}$ and, same as above, $f_2 \preceq \bar{f}_2 \preceq \nu$, whence also $f_2 \preceq \nu$.
- The density of ω is $f_1 \star f_2$. We have

$$\begin{aligned} f_1 \star f_2 &\preceq \mu \star f_2 && \text{[by 1^o.b in view of } f_1, \mu \in \mathcal{E}, f_2 \in \mathcal{N} \text{ and } f_1 \preceq \mu], \\ f_2 \star \mu &\preceq \mu \star \nu && \text{[by 1^o.b in view of } f_2, \nu \in \mathcal{E}, \mu \in \mathcal{N} \text{ and } f_2 \preceq \nu]. \end{aligned}$$

Hence, $f_1 \star f_2 \preceq \gamma := \mu \star \nu$, such that $\gamma \in \mathcal{E}$, and (15) follows. Besides this, in the case in question ω has an even density.

2⁰.b. Now let $\rho_1 = 0$. Then $\rho_2 > 0$ due to $AA^T + BB^T \succ 0$, and the probability density of ω is f_2 . We have $f_2 \preceq \bar{f}_2 \preceq \nu \preceq \mu \star \nu$ (the concluding \preceq is due to **1^o.d**), and (15) follows. Similarly, when $\rho_2 = 0$, we have $\rho_1 > 0$, and the probability density of ω is f_1 . Similarly to the case where $\rho_1 = 0$, we have $f_1 \preceq \bar{f}_1 \preceq \mu \preceq \mu \star \nu$, and (15) follows. Furthermore, as we have seen, ω always has an even probability density. Finally, $\gamma = \mu \star \nu$ is nice by **1^o.a**. (i) is proved.

3⁰. To prove (ii), note that in the notation of item 2⁰ and under the premise of (ii), $\rho_\chi > 0$ implies that $f_\chi \in \mathcal{N}$. Hence, due to **1^o.a**, the distribution of ω has density from \mathcal{N} (this density is $f_1 \star f_2$ when $\rho_1, \rho_2 > 0$, f_2 when $\rho_1 = 0$, and f_1 when $\rho_2 = 0$), which combines with (i) to imply (ii). \square

A.2 Proof of Proposition 2.4

Proof. When $x \in X_1$, we have $h_*^T x \geq h_*^T x_*^1 = c_* + \delta$, thus

$$\begin{aligned} \{\xi : s_*(x + \xi) < \tfrac{1}{2}(\alpha_2 - \alpha_1)\} &= \{\xi : h_*^T(x + \xi) < c_* + \tfrac{1}{2}(\alpha_2 - \alpha_1)\} \\ &= \{\xi : h_*^T \xi < c_* - h_*^T x + \tfrac{1}{2}(\alpha_2 - \alpha_1)\} \\ &\subseteq \{\xi : h_*^T \xi < c_* - h_*^T x_*^1 + \tfrac{1}{2}(\alpha_2 - \alpha_1)\} \\ &= \{\xi : h_*^T \xi < -\delta + \tfrac{1}{2}(\alpha_2 - \alpha_1)\} \subseteq \{h_*^T \xi < -\alpha_1\}, \end{aligned}$$

where the last inclusion is due to $\alpha_1 + \alpha_2 \leq 2\delta$. Therefore, for $p \in \mathcal{P}$ it holds

$$\text{Prob}_{\xi \sim p} \{s_*(x + \xi) < \tfrac{1}{2}(\alpha_2 - \alpha_1)\} \leq \underbrace{\int_{h_*^T \xi < -\alpha_1} p(\xi) d\xi}_{(a)} \stackrel{(a)}{=} \underbrace{\int_{h_*^T \xi > \alpha_1} p(\xi) d\xi}_{(b)} \stackrel{(b)}{\leq} P_\gamma(\alpha_1)$$

where (a) is due to the fact that $p(\cdot)$ is even, and (b) is due to $p \in \mathcal{P}_\gamma$ and $\alpha_1 \geq 0$. When $x \in X_2$, we have $h_*^T x \leq h_*^T x_*^2 = c_* - \delta$, and using a completely similar argument we conclude that for $p \in \mathcal{P}$ it holds

$$\text{Prob}_{\xi \sim p} \{s_*(x + \xi) \geq \tfrac{1}{2}(\alpha_2 - \alpha_1)\} \leq P_\gamma(\alpha_2). \quad \square$$

A.3 Proof of Proposition 2.5

When $x_k \in X_1^k$ and $p_k \in \mathcal{P}^k$ for all $k \leq K$, due to the origin of h_k , δ_k and c_k , we have

$$h_k^T(x_k + \xi_k) - c_k \geq \delta_k + h_k^T \xi_k,$$

and, because $\eta_k(\cdot)$ is nondecreasing,

$$\int e^{-\eta_k(h_k^T[x_k + \xi_k] - c_k)} p_k(\xi_k) d\xi_k \leq \int e^{-\eta_k(\delta_k + h_k^T \xi_k)} p_k(\xi_k) d\xi_k \leq \text{risk}_{\delta_k}(\eta_k | \mathcal{P}^k), \quad k = 1, \dots, K,$$

where the concluding “ \leq ” is due to (23.a). Hence,

$$\begin{aligned} \int e^{-\phi^{(K)}(x_1+\xi_1, \dots, x_K+\xi_K)} \prod_{k=1}^K [p_k(\xi_k) d\xi_k] &= \prod_{k=1}^K \left[\int e^{-\phi_k(x_k+\xi_k)} p_k(\xi_k) d\xi_k \right] \\ &= \prod_{k=1}^K \left[\int e^{-\eta_k(h_k^T[x_k+\xi_k]-c_k)} p_k(\xi_k) d\xi_k \right] \leq \prod_{k=1}^K [\text{risk}_{\delta_k}(\eta_k | \mathcal{P}^k)]. \end{aligned}$$

When $x_k \in X_2^k$ and $p_k \in \mathcal{P}^k$ for all $k \leq K$, in a completely similar way we obtain

$$\int e^{\phi^{(K)}(x_1+\xi_1, \dots, x_K+\xi_K)} \prod_{k=1}^K [p_k(\xi_k) d\xi_k] \leq \prod_{k=1}^K [\text{risk}_{\delta_k}(\eta_k | \mathcal{P}^k)]. \quad \square$$

A.4 Proof of Proposition 2.6

Let $p(\cdot) \in \mathcal{P}$, let $e \in \mathbf{R}^n$ be a unit vector, and let $q(\cdot)$ be the probability density of the scalar random variable $e^T \xi$ induced by the density $p(\cdot)$ of ξ . We start with the following well known observation:

Lemma A.1 *Let f, g be two probability densities on \mathbf{R} such that*

$$\begin{aligned} (a) \quad & \int_0^\infty f(s) ds = \int_0^\infty g(s) ds, \\ (b) \quad & \int_r^\infty f(s) ds \geq \int_r^\infty g(s) ds, \quad \forall r \geq 0, \end{aligned} \quad (59)$$

and let $h(s)$ be a nondecreasing real-valued function on the nonnegative ray such that $\int_0^\infty h(s) f(s) ds < \infty$. Then

$$\int_0^\infty h(s) f(s) ds \geq \int_0^\infty h(s) g(s) ds. \quad (60)$$

To make the presentation self-contained, here is the proof of the lemma:

In view of (59.a), we can assume w.l.o.g. that $h(0) = 0$. Let us extend $h(s)$ from the nonnegative ray to the entire real axis by setting $h(s) = 0$, $s < 0$, thus arriving at a monotone on the axis nonnegative function. Let $\eta \sim f$ and $\zeta \sim g$. When denoting H the c.d.f. of $h(\eta)$, under the premise of the lemma we clearly have

$$\mathbf{E}_{\eta \sim f}\{h(\eta)\} = \int_0^\infty t dH(t) = \int_0^\infty (1 - H(t)) dt = \int_0^\infty \text{Prob}\{h(\eta) > t\} dt.$$

On the other hand, for $t \geq 0$ the set $\{s : h(s) > t\}$ is a ray either of the form $[a_t, +\infty)$ or $(a_t, +\infty)$ with $a_t \geq 0$, so that (59.b) implies that

$$\forall t \geq 0 \quad \text{Prob}\{h(\eta) > t\} \geq \text{Prob}\{h(\zeta) > t\}.$$

As a result,

$$\mathbf{E}_{\eta \sim f}\{h(\eta)\} = \int_0^\infty \text{Prob}\{h(\eta) > t\} dt \geq \int_0^\infty \text{Prob}\{h(\zeta) > t\} dt.$$

We conclude that $\mathbf{E}_{\zeta \sim g}\{h(\zeta)\}$ is finite and satisfies

$$\mathbf{E}_{\zeta \sim g}\{h(\zeta)\} = \int_0^\infty \text{Prob}\{h(\zeta) > t\} dt \leq \mathbf{E}_{\eta \sim f}\{h(\eta)\}. \quad \square$$

1⁰. Note that $q(\cdot)$ is an even probability density on the axis and

$$\int_s^\infty [\gamma(r) - q(r)]dr = P_\gamma(s) - \text{Prob}_{\xi \sim p}\{\xi : e^T \xi \geq s\} \begin{cases} \geq 0, & s \geq 0, \\ = 0, & s = 0. \end{cases}$$

We have

$$\int_{\mathbf{R}^n} e^{-\eta(\delta+e^T\xi)} p(\xi) d\xi = \int_{-\infty}^\infty e^{-\eta(\delta+s)} q(s) ds \stackrel{(a)}{=} \int_0^\infty H_{\delta\eta}(s) q(s) ds \stackrel{(b)}{\leq} \int_0^\infty H_{\delta\eta}(s) \gamma(s) ds \stackrel{(c)}{=} \epsilon_\delta(\eta|\gamma),$$

where (a) is due to the fact that q is even, (b) is a result of applying Lemma A.1 to densities γ , q and nondecreasing $H_{\delta\eta}$, and (c) is due to the definition (26) of the δ -index.

2⁰. We have

$$\begin{aligned} & \int_{\mathbf{R}^n} e^{\eta(-\delta+e^T\xi)} p(\xi) d\xi = \int_{\mathbf{R}^n} e^{-\eta(\delta-e^T\xi)} p(\xi) d\xi \text{ [since } \eta(\cdot) \text{ is odd]} \\ & = \int_{\mathbf{R}^n} e^{-\eta(\delta+e^T\xi)} p(\xi) d\xi \text{ [since } p(\cdot) \text{ is even],} \end{aligned}$$

and we have already seen in 1⁰ that the concluding quantity is $\leq \epsilon_\delta(\eta|\gamma)$. The bottom line is that inequalities (23) hold true with $\epsilon = \epsilon_\delta(\eta|\gamma)$, and (27) follows. \square

A.5 Proof of Proposition 2.7

Let $e \in \mathbf{R}^n$ be a unit vector, $\delta \geq 0$, and $p \in \mathcal{P}_{sG}$. We have

$$\int_{\mathbf{R}^n} e^{-\eta(\delta+e^T\xi)} p(\xi) d\xi = \int_{\mathbf{R}^n} e^{-\delta^2 - \delta e^T \xi} p(\xi) d\xi \leq e^{-\delta^2 + \delta^2/2} = e^{-\delta^2/2},$$

where the concluding \leq is due to $p \in \mathcal{P}_{sG}$ and $\|e\|_2 = 1$. Similarly,

$$\int_{\mathbf{R}^n} e^{\eta(e^T\xi - \delta)} p(\xi) d\xi = \int_{\mathbf{R}^n} e^{-\delta^2 + \delta e^T \xi} p(\xi) d\xi \leq e^{-\delta^2 + \delta^2/2} = e^{-\delta^2/2}.$$

The resulting inequalities hold true for all unit vectors e and all $p \in \mathcal{P}_{sG}$, implying (33). \square

A.6 Proof of Proposition 2.8

1^o. Let $p \in \mathcal{P}$ and $x \in X_1$ be fixed. Due to the monotonicity of η , and by the definition of c_* , we get from (34.a):

$$\mathbf{E}_{\xi \sim p}\{\eta(h_*^T(x + \xi) + c)\} \geq \mathbf{E}_{\xi \sim p}\{\eta(h_*^T(x_*^1 + \xi) + c)\} = c_* + \frac{e^\kappa \rho}{1 + e^\kappa},$$

and so

$$\mathbf{E}_{\xi^K \sim p \times \dots \times p} \psi_j([x + \xi_1; \dots; x + \xi_K]) = \mathbf{E}_{\xi \sim p}\{\eta(h^T(x + \xi) + c)\} \geq c_* + \frac{e^\kappa \rho}{1 + e^\kappa}.$$

On the other hand, by (34.b) we have

$$\mathbf{Var}_{\xi^K \sim p \times \dots \times p}(\psi_j(\omega^K)) = m^{-1} \mathbf{Var}_{\xi \sim p}\{\eta(h^T(x + \xi) + c)\} \leq m^{-1}.$$

Now, by the Chebyshev inequality,

$$\begin{aligned} \text{Prob}_{\xi^K \sim p \times \dots \times p} \{\psi_j(\omega^K) < c_*\} & \leq \text{Prob}_{\xi^K \sim p \times \dots \times p} \left\{ \psi_j(\omega^K) - \mathbf{E}_{\xi^K \sim p \times \dots \times p} \{\psi_j(\omega^K)\} < -\frac{e^\kappa \rho}{1 + e^\kappa} \right\} \\ & \leq \frac{(1 + e^\kappa)^2}{m e^{2\kappa} \rho^2} \leq \frac{1}{4e}. \end{aligned}$$

As a result, the risk $\text{Risk}_{1\mathcal{S}}$ of the test $\mathcal{T}_K^{\text{mm}}$ satisfies the bound

$$\text{Risk}_{1\mathcal{S}}(\mathcal{T}_K^{\text{mm}}|\mathcal{P}, X_1, X_2) \leq \sum_{J/2 \leq j \leq J} \binom{J}{j} (4e)^{-j} \left(1 - \frac{1}{4e}\right)^{J-j} \leq 2^J (4e)^{-J/2} = e^{-J/2} \leq \epsilon_1,$$

where the final inequality is due to (35).

2°. The same argument, as applied to $x \in X_2$ results in

$$\mathbf{E}_{\xi^K \sim p \times \dots \times p} \psi_j([x + \xi_1; \dots; x + \xi_K]) = \mathbf{E}_{\xi \sim p} \{\eta(h_*^T(x + \xi) + c)\} \leq c_* - \frac{\varrho}{1 + e^\kappa},$$

and

$$\begin{aligned} \text{Prob}_{\xi^K \sim p \times \dots \times p} \{\psi_j(\omega^K) \geq c_*\} &\leq \text{Prob}_{\xi^K \sim p \times \dots \times p} \left\{ \psi_j(\omega^K) - \mathbf{E}_{\xi^K \sim p \times \dots \times p} \{\psi_j(\omega^K)\} \geq \frac{\varrho}{1 + e^\kappa} \right\} \\ &\leq \frac{(1 + e^\kappa)^2}{m\varrho^2} \leq \frac{1}{4} e^{2\kappa - 1}. \end{aligned}$$

Same as above, we conclude that

$$\text{Risk}_{2\mathcal{S}}(\mathcal{T}_K^{\text{mm}}|\mathcal{P}, X_1, X_2) \leq e^{J(\kappa - \frac{1}{2})} = \exp\left(-\frac{J \ln \epsilon_2^{-1}}{2 \ln \epsilon_1^{-1}}\right) \leq \epsilon_2,$$

where the concluding inequality is due to (35). □