



# An extended Generalised Variance, with Applications

Luc Pronzato, Henry Wynn, Anatoly Zhigljavsky

► **To cite this version:**

Luc Pronzato, Henry Wynn, Anatoly Zhigljavsky. An extended Generalised Variance, with Applications. 2014. <hal-01086442>

**HAL Id: hal-01086442**

**<https://hal.archives-ouvertes.fr/hal-01086442>**

Submitted on 24 Nov 2014

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

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

# An extended Generalised Variance, with Applications

Luc Pronzato\*, Henry P. Wynn<sup>†</sup> and Anatoly A. Zhigljavsky<sup>‡</sup>

CNRS, London School of Economics and Cardiff University

November 24, 2014

## Abstract

We consider a measure  $\psi_k$  of dispersion which extends the notion of Wilk's generalised variance, or entropy, for a  $d$ -dimensional distribution, and is based on the mean squared volume of simplices of dimension  $k \leq d$  formed by  $k + 1$  independent copies. We show how  $\psi_k$  can be expressed in terms of the eigenvalues of the covariance matrix of the distribution, also when a  $n$ -point sample is used for its estimation, and prove its concavity when raised at a suitable power. Some properties of entropy-maximising distributions are derived, including a necessary and sufficient condition for optimality. Finally, we show how this measure of dispersion can be used for the design of optimal experiments, with equivalence to  $A$  and  $D$ -optimal design for  $k = 1$  and  $k = d$  respectively. Simple illustrative examples are presented.

**MSC:** Primary 94A17; secondary 62K05

**keyword:** dispersion, generalised variance, ,quadratic entropy, ,maximum-entropy measure, design of experiments, optimal design

## 1 Introduction

The idea of dispersion is fundamental to statistics and with different terminology, such as potential, entropy, information and capacity, stretches over a wide area. The variance and standard deviation are the most prevalent for a univariate distribution, and Wilks generalised variance is the term usually reserved for the determinant of the covariance matrix,  $V$ , of a multivariate distribution. Many other measures of dispersion have been introduced and a rich area comprises those that are order-preserving with respect to a dispersion ordering; see [21, 12, 5]. These are sometimes referred to as *measures of peakness* and *peakness ordering*, and are related to the large literature on dispersion measures which grew out of the Gini coefficient, used to measure income inequality [4] and diversity in biology, see [15], which we will discuss briefly below.

In the definitions there are typically two kinds of dispersion, those measuring some kind of mean distance, or squared distance, from a central value, such as in the usual definition of variance, and those based on the expected distance, or squared distance, between two independent copies from the same distribution, such

---

\*pronzato@i3s.unice.fr, Laboratoire I3S, UMR 7172, UNS/CNRS, 2000 route des Lucioles, Les Algorithmes, bât. Euclide B, 06900 Sophia Antipolis, France

<sup>†</sup>h.wynn@lse.ac.uk, London School of Economics, Houghton Street, London, WC2A 2AE, UK

<sup>‡</sup>zhigljavskyAA@cf.ac.uk, School of Mathematics, Cardiff University, Senghennydd Road, Cardiff, CF24 4YH, UK

as the Gini coefficient. It is this second type that will concern us here and we will generalise the idea in several ways by replacing distance by volumes of simplices formed by  $k$  independent copies and by transforming the distance, both inside the expectation and outside.

The area of optimal experimental design is another which has provided a range of dispersion measures. Good designs, it is suggested, are those whose parameter estimates have low dispersion. Typically, this means that the design measure, the spread of the observation sites, *maximises* a measure of dispersion and we shall study this problem.

We think of a dispersion measure as a functional directly on the distribution. The basic functional is an integral, such as a moment. The property we shall stress for such functionals most is concavity: that a functional does not decrease under mixing of the distributions. A fundamental theorem in Bayesian learning is that we expect concave functionals to decrease through taking of observations, see Section 2.2 below.

Our central result (Section 3) is an identity for the mean squared volume of simplices of dimension  $k$ , formed by  $k+1$  independent copies, in terms of the eigenvalues of the covariance matrices or equivalently in terms of sums of the determinants of  $k$ -marginal covariance matrices. Second, we note that after an appropriate (exterior) power transformation the functional becomes concave. We can thus (i) derive properties of measures that maximise this functional (Section 4.1), (ii) use this functional to measure the dispersion of parameter estimates in regression problems, and hence design optimal experiments which minimise this measure of dispersion (Section 4.2).

## 2 Dispersion measures

### 2.1 Concave and homogeneous functionals

Let  $\mathcal{X}$  be a compact subset of  $\mathbb{R}^d$ ,  $\mathcal{M}$  be the set of all probability measures on the Borel subsets of  $\mathcal{X}$  and  $\phi : \mathcal{M} \rightarrow \mathbb{R}^+$  be a functional defined on  $\mathcal{M}$ . We will be interested in the functionals  $\phi(\cdot)$  that are (see Appendix for precise definitions)

- (a) shift-invariant,
- (b) positively homogeneous of a given degree  $q$ , and
- (c) concave:  $\phi[(1 - \alpha)\mu_1 + \alpha\mu_2] \geq (1 - \alpha)\phi(\mu_1) + \alpha\phi(\mu_2)$  for any  $\alpha \in (0, 1)$  and any two measures  $\mu_1, \mu_2$  in  $\mathcal{M}$ .

For  $d = 1$ , a common example of a functional satisfying the above properties, with  $q = 2$  in (b), is the variance

$$\sigma^2(\mu) = E_\mu^{(2)} - E_\mu^2 = \frac{1}{2} \int \int (x_1 - x_2)^2 \mu(dx_1) \mu(dx_2),$$

where  $E_\mu = E_\mu(x) = \int x \mu(dx)$  and  $E_\mu^{(2)} = \int x^2 \mu(dx)$ . Concavity follows from linearity of  $E_\mu^{(2)}$ , that is,  $E_{(1-\alpha)\mu_1 + \alpha\mu_2}^{(2)} = (1-\alpha)E_{\mu_1}^{(2)} + \alpha E_{\mu_2}^{(2)}$ , and Jensen's inequality which implies  $E_{(1-\alpha)\mu_1 + \alpha\mu_2}^2 \leq (1-\alpha)E_{\mu_1}^2 + \alpha E_{\mu_2}^2$ .

Any moment of  $\mu \in \mathcal{M}$  is a homogeneous functional of a suitable degree. However, the variance is the only moment which satisfies (a) and (c). Indeed, the shift-invariance implies that the moment should be central, but the variance is the only concave functional among the central moments, see Appendix. In this sense, one of the aims of this paper is a generalisation of the concept of variance.

In the general case  $d \geq 1$ , the double variance  $2\sigma^2(\mu)$  generalises to

$$\phi(\mu) = \int \int \|x_1 - x_2\|^2 \mu(dx_1) \mu(dx_2) = 2 \int \|x - E_\mu\|^2 \mu(dx) = 2 \text{trace}(V_\mu), \quad (2.1)$$

where  $\|\cdot\|$  is the  $L_2$ -norm in  $\mathbb{R}^d$  and  $V_\mu$  is the covariance matrix of  $\mu$ . This functional, like the variance, satisfies the conditions (a)-(c) with  $q = 2$ .

The functional (2.1) is the double integral of the squared distance between two random points distributed according to the measure  $\mu$ . Our main interest will be concentrated around the general class of functionals defined by

$$\phi(\mu) = \phi_{[k],\delta,\tau}(\mu) = \left( \int \dots \int \mathcal{V}_k^\delta(x_1, \dots, x_{k+1}) \mu(dx_1) \dots \mu(dx_{k+1}) \right)^\tau, \quad k \geq 2 \quad (2.2)$$

for some  $\delta$  and  $\tau$  in  $\mathbb{R}^+$ , where  $\mathcal{V}_k(x_1, \dots, x_{k+1})$  is the volume of the  $k$ -dimensional simplex (its area when  $k = 2$ ) formed by the  $k + 1$  vertices  $x_1, \dots, x_{k+1}$  in  $\mathbb{R}^d$ , with  $k = d$  as a special case. Property (a) for the functionals (2.2) is then a straightforward consequence of the shift-invariance of  $\mathcal{V}_k$ , and positive homogeneity of degree  $q = k \delta \tau$  directly follows from the positive homogeneity of  $\mathcal{V}_k$  with degree  $k$ . Concavity will be proved to hold for  $\delta = 2$  and  $\tau \leq 1/k$  in Section 3. There, we show that this case can be considered as a natural extension of (2.1) (which corresponds to  $k = 1$ ), with  $\phi_{[k],2,\tau}(\mu)$  being expressed as a function of  $V_\mu$ , the covariance matrix of  $\mu$ . The concavity for  $k = \tau = 1$  and all  $0 < \delta \leq 2$ , follows from the Schoenberg theory, which will be discussed briefly below. The functionals (2.2) with  $\delta = 2$  and  $\tau > 0$ ,  $1 \leq k \leq d$ , can be used to define a family of criteria for optimal experimental design, concave for  $\tau \leq 1/k$ , for which an equivalence theorem can be formulated.

## 2.2 Quadratic entropy and learning

In a series of papers [15, 16, 17] C.R. Rao and co-workers have introduced a *quadratic entropy* which is a generalised version of the  $k = 2$  functional of this section but with a general kernel  $K(x_1, x_2)$  in  $\mathbb{R}^d$ :

$$Q_R = \int \int K(x_1, x_2) \mu(dx_1) \mu(dx_2). \quad (2.3)$$

For the discrete version

$$Q_R = \sum_{i=1}^N \sum_{j=1}^N K(x_i, x_j) p_i p_j,$$

Rao and co-workers developed a version of the Analysis of Variance (ANOVA), which they called Analysis of Quadratic Entropy (ANOQE). The Gini coefficient, also used in the continuous and discrete form is a special case with  $d = 1$  and  $K(x_1, x_2) = |x_1 - x_2|$ .

As pointed in [17, Chap. 3], a necessary and sufficient condition for the functional  $Q_R$  to be concave is

$$\int \int K(x_1, x_2) \nu(dx_1) \nu(dx_2) \leq 0 \quad (2.4)$$

for all measures  $\nu$  with  $\int \nu(dx) = 0$ . The discrete version of this is

$$\sum_{i=1}^N \sum_{j=1}^N K(x_i, x_j) q_i q_j \leq 0$$

for any choice of real numbers  $q_1, \dots, q_N$  such that  $\sum_{i=1}^N q_i = 0$ . Schoenberg [19] solves the general problem of finding for what class of functions  $B(\cdot)$  on  $\|x_1 - x_2\|^2$  does the kernel  $K(x_1, x_2) = B(\|x_1 - x_2\|^2)$  satisfy (2.4). The solution is that  $B$  must be a so-called Bernstein function, see [18]. We do not develop these ideas here, but note that  $B(\lambda) = \lambda^\alpha$  is a Bernstein function for all  $0 < \alpha \leq 1$ . This is the reason that, above, we can claim concavity for  $k = 1$  and all  $0 < \delta \leq 2$  in (2.2).

Hainy *et al* [6] discuss the link to embedding and review some basic results related to Bayesian learning. One asks what is the class of functionals  $\psi$  on a distribution  $\mu(\theta)$  of a parameter in the Bayesian statistical learning such that for all  $\mu(\theta)$  and all sampling distributions  $\pi(x|\theta)$  one expects to learn, in the preposterior sense:  $\psi(\mu(\theta)) \leq \mathbf{E}_\nu \psi(\pi(\theta|X))$ , with  $X \sim \nu$ . The condition is that  $\psi$  is convex, a result which has a history but is usually attributed to De Groot [2]. This learning is enough to justify calling such a functional a generalised information functional, or a general learning functional. Shannon information falls in this class, and earlier versions of the result were for Shannon information. It follows that wherever, in this paper, we have a concave functional then its *negative* is a learning functional.

### 3 Functionals based on squared volume

In the rest of the paper we focus our attention on the functional

$$\mu \in \mathcal{M} \longrightarrow \psi_k(\mu) = \phi_{[k],2,1}(\mu) = \mathbf{E}_\mu \{ \mathcal{V}_k^2(x_1, \dots, x_{k+1}) \},$$

which corresponds to the mean squared volume of simplices of dimension  $k$  formed by  $k + 1$  independent samples from  $\mu$ . For instance,

$$\psi_2(\mu) = \int \int \int \mathcal{V}_2^2(x_1, x_2, x_3) \mu(dx_1) \mu(dx_2) \mu(dx_3), \quad (3.1)$$

with  $\mathcal{V}_2(x_1, x_2, x_3)$  the area of the triangle formed by the three points with coordinates  $x_1, x_2$  and  $x_3$  in  $\mathbb{R}^d$ ,  $d \geq 2$ . Functionals  $\phi_{[k],\delta,\tau}(\mu)$  for  $\delta \neq 2$  will be considered elsewhere, including the case of negative  $\delta$  and  $\tau$  in connection with space-filling design for computer experiments.

Theorem 3.1 of Section 3.1 indicates how  $\psi_k(\mu)$  can be expressed as a function of  $V_\mu$ , the covariance matrix of  $\mu$ , and shows that  $\phi_{[k],2,1/k}(\cdot)$  satisfies properties (a), (b) and (c) of Section 2.1. The special case of  $k = d$  was known to Wilks [25, 26] in his introduction of generalised variance, see also [24]. The connection with U-statistics is exploited in Section 3.2, where an unbiased minimum-variance estimator of  $\psi_k(\mu)$  based on a sample  $x_1, \dots, x_n$  is expressed in terms of the empirical covariance matrix of the sample.

#### 3.1 Expected squared $k$ -simplex volume

**Theorem 3.1.** *Let the  $x_i$  be i.i.d. with the probability measure  $\mu \in \mathcal{M}$ . Then, for any  $k \in \{1, \dots, d\}$ , we have*

$$\psi_k(\mu) = \frac{k+1}{k!} \sum_{i_1 < i_2 < \dots < i_k} \det\{[V_\mu]_{(i_1, \dots, i_k) \times (i_1, \dots, i_k)}\} \quad (3.2)$$

$$= \frac{k+1}{k!} \sum_{i_1 < i_2 < \dots < i_k} \lambda_{i_1}[V_\mu] \times \dots \times \lambda_{i_k}[V_\mu], \quad (3.3)$$

where  $\lambda_i[V_\mu]$  is the  $i$ -th eigenvalue of the covariance matrix  $V_\mu$  and all  $i_j$  belong to  $\{1, \dots, d\}$ . Moreover, the functional  $\psi_k^{1/k}(\cdot)$  is shift-invariant, homogeneous of degree 2 and concave on  $\mathcal{M}$ .

The proof uses the following two lemmas.

**Lemma 3.1.** *Let the  $k + 1$  vectors  $x_1, \dots, x_{k+1}$  of  $\mathbb{R}^k$  be i.i.d. with the probability measure  $\mu$ ,  $k \geq 2$ . For  $i = 1, \dots, k + 1$ , denote  $z_i = (x_i^\top \ 1)^\top$ . Then*

$$\mathbb{E}_\mu \left\{ \det \left[ \sum_{i=1}^{k+1} z_i z_i^\top \right] \right\} = (k + 1)! \det[V_\mu].$$

*Proof.* We have

$$\mathbb{E}_\mu \left\{ \det \left[ \sum_{i=1}^{k+1} z_i z_i^\top \right] \right\} = (k + 1)! \det \begin{bmatrix} \mathbb{E}_\mu(x_1 x_1^\top) & E_\mu \\ E_\mu^\top & 1 \end{bmatrix} = (k + 1)! \det[V_\mu],$$

see for instance [13, Theorem 1].  $\square$

**Lemma 3.2.** *The matrix functional  $\mu \mapsto V_\mu$  is Loewner-concave on  $\mathcal{M}$ , in the sense that, for any  $\mu_1, \mu_2$  in  $\mathcal{M}$  and any  $\alpha \in (0, 1)$ ,*

$$V_{(1-\alpha)\mu_1 + \alpha\mu_2} \succeq (1 - \alpha)V_{\mu_1} + \alpha V_{\mu_2}, \quad (3.4)$$

where  $A \succeq B$  means that  $A - B$  is nonnegative definite.

*Proof.* Take any vector  $z$  of the same dimension as  $x$ . Then  $z^\top V_\mu z = \text{var}_\mu(z^\top x)$ , which is a concave functional of  $\mu$ , see Section 2.1. This implies that  $z^\top V_{(1-\alpha)\mu_1 + \alpha\mu_2} z = \text{var}_{(1-\alpha)\mu_1 + \alpha\mu_2}(z^\top x) \geq (1 - \alpha)\text{var}_{\mu_1}(z^\top x) + \alpha\text{var}_{\mu_2}(z^\top x) = (1 - \alpha)z^\top V_{\mu_1} z + \alpha z^\top V_{\mu_2} z$ , for any  $\mu_1, \mu_2$  in  $\mathcal{M}$  and any  $\alpha \in (0, 1)$  (see Section 2.1 for the concavity of  $\text{var}_\mu$ ). Since  $z$  is arbitrary, this implies (3.4).  $\square$

*Proof of Theorem 3.1.* When  $k = 1$ , the results follow from  $\psi_1(\mu) = 2 \text{trace}(V_\mu)$ , see (2.1). Using Binet-Cauchy formula, see, e.g., [3, vol. 1, p. 9], we obtain

$$\begin{aligned} \mathcal{V}_k^2(x_1, \dots, x_{k+1}) &= \\ &= \frac{1}{(k!)^2} \det \left( \begin{bmatrix} (x_2 - x_1)^\top \\ (x_3 - x_1)^\top \\ \vdots \\ (x_{k+1} - x_1)^\top \end{bmatrix} \begin{bmatrix} (x_2 - x_1) & (x_3 - x_1) & \cdots & (x_{k+1} - x_1) \end{bmatrix} \right) \\ &= \frac{1}{(k!)^2} \sum_{i_1 < i_2 < \cdots < i_k} \det^2 \begin{bmatrix} \{x_2 - x_1\}_{i_1} & \cdots & \{x_{k+1} - x_1\}_{i_1} \\ \vdots & \vdots & \vdots \\ \{x_2 - x_1\}_{i_k} & \cdots & \{x_{k+1} - x_1\}_{i_k} \end{bmatrix} \\ &= \frac{1}{(k!)^2} \sum_{i_1 < i_2 < \cdots < i_k} \det^2 \begin{bmatrix} \{x_1\}_{i_1} & \cdots & \{x_{k+1}\}_{i_1} \\ \vdots & \vdots & \vdots \\ \{x_1\}_{i_k} & \cdots & \{x_{k+1}\}_{i_k} \\ 1 & \cdots & 1 \end{bmatrix}, \end{aligned}$$

where  $\{x\}_i$  denotes the  $i$ -th component of vector  $x$ . Also, for all  $i_1 < i_2 < \cdots < i_k$ ,

$$\det^2 \begin{bmatrix} \{x_1\}_{i_1} & \cdots & \{x_{k+1}\}_{i_1} \\ \vdots & \vdots & \vdots \\ \{x_1\}_{i_k} & \cdots & \{x_{k+1}\}_{i_k} \\ 1 & \cdots & 1 \end{bmatrix} = \det \left( \sum_{j=1}^{k+1} z_j z_j^\top \right)$$

where we have denoted by  $z_j$  the  $k+1$ -dimensional vector with components  $\{x_j\}_{i_\ell}$ ,  $\ell = 1, \dots, k$ , and 1. When the  $x_i$  are i.i.d. with the probability measure  $\mu$ , using Lemma 3.1 we obtain (3.2), (3.3). Therefore

$$\psi_k(\mu) = \Psi_k[V_\mu] = \frac{k + 1}{k!} \mathcal{E}_k\{\lambda_1[V_\mu], \dots, \lambda_d[V_\mu]\}, \quad (3.5)$$

with  $\mathcal{E}_k\{\lambda_1[V_\mu], \dots, \lambda_d[V_\mu]\}$  the elementary symmetric function of degree  $k$  of the  $d$  eigenvalues of  $V_\mu$ , see, e.g., [11, p. 10]. Note that

$$\mathcal{E}_k[V_\mu] = \mathcal{E}_k\{\lambda_1[V_\mu], \dots, \lambda_d[V_\mu]\} = (-1)^k a_{d-k},$$

with  $a_{d-k}$  the coefficient of the monomial of degree  $d - k$  of the characteristic polynomial of  $V_\mu$ ; see, e.g., [11, p. 21]. We have in particular  $\mathcal{E}_1[V_\mu] = \text{trace}[V_\mu]$  and  $\mathcal{E}_d[V_\mu] = \det[V_\mu]$ . The shift-invariance and homogeneity of degree 2 of  $\psi_k^{1/k}(\cdot)$  follow from the shift-invariance and positive homogeneity of  $\mathcal{V}_k$  with degree  $k$ . Concavity of  $\Psi_k^{1/k}(\cdot)$  follows from [11, p. 116] (take  $p = k$  in eq. (10), with  $\mathcal{E}_0 = 1$ ). From [9], the  $\Psi_k^{1/k}(\cdot)$  are also Loewner-increasing, so that from Lemma 3.2, for any  $\mu_1, \mu_2$  in  $\mathcal{M}$  and any  $\alpha \in (0, 1)$ ,

$$\begin{aligned} \psi_k^{1/k}[(1 - \alpha)\mu_1 + \alpha\mu_2] &= \Psi_k^{1/k}\{V_{(1-\alpha)\mu_1 + \alpha\mu_2}\} \\ &\geq \Psi_k^{1/k}[(1 - \alpha)V_{\mu_1} + \alpha V_{\mu_2}] \\ &\geq (1 - \alpha)\Psi_k^{1/k}[V_{\mu_1}] + \alpha\Psi_k^{1/k}[V_{\mu_2}] \\ &= (1 - \alpha)\psi_k^{1/k}(\mu_1) + \alpha\psi_k^{1/k}(\mu_2). \end{aligned}$$

□

The functionals  $\mu \rightarrow \phi_{[k],2,\tau}(\mu) = \psi_k^\tau(\mu)$  are thus concave for  $0 < \tau \leq 1/k$ , with  $\tau = 1/k$  yielding positive homogeneity of degree 2. The functional  $\psi_1(\mu)$  is a quadratic entropy  $Q_R$ , see (2.3),  $\psi_d(\mu)$  is proportional to Wilks generalised variance, and  $\psi_2^{1/2}(\mu)$ , see (3.1), and  $\psi_3^{1/3}(\mu)$  can be respectively considered as particular versions of cubic and quartic entropies.

From the well-known expression of the coefficients of the characteristic polynomial of a matrix  $V$ , we have

$$\begin{aligned} \Psi_k(V) &= \frac{k+1}{k!} \mathcal{E}_k(V) \\ &= \frac{k+1}{(k!)^2} \det \begin{bmatrix} \text{trace}(V) & k-1 & 0 & \dots \\ \text{trace}(V^2) & \text{trace}(V) & k-2 & \dots \\ \dots & \dots & \dots & \dots \\ \text{trace}(V^{k-1}) & \text{trace}(V^{k-2}) & \dots & 1 \\ \text{trace}(V^k) & \text{trace}(V^{k-1}) & \dots & \text{trace}(V) \end{bmatrix}, \end{aligned} \quad (3.6)$$

see, e.g., [10, p. 28], and the  $\mathcal{E}_k(V)$  satisfy the recurrence relations (Newton identities):

$$\mathcal{E}_k(V) = \frac{1}{k} \sum_{i=1}^k (-1)^{i-1} \mathcal{E}_{k-i}(V) \mathcal{E}_1(V^i), \quad (3.7)$$

see, e.g., [3, Vol. 1, p. 88] and [9]. Particular forms of  $\psi_k(\cdot)$  are

$$\begin{aligned} k = 1 : & \quad \psi_1(\mu) = 2 \text{trace}(V_\mu), \\ k = 2 : & \quad \psi_2(\mu) = \frac{3}{4} [\text{trace}^2(V_\mu) - \text{trace}(V_\mu^2)], \\ k = 3 : & \quad \psi_3(\mu) = \frac{1}{9} [\text{trace}^3(V_\mu) - 3 \text{trace}(V_\mu^2) \text{trace}(V_\mu) + 2 \text{trace}(V_\mu^3)], \\ k = d : & \quad \psi_d(\mu) = \frac{d+1}{d!} \det(V_\mu). \end{aligned}$$

### 3.2 The empirical version and unbiased estimates

Let  $x_1, \dots, x_n$  be a sample of  $n$  vectors of  $\mathbb{R}^d$ , i.i.d. with the measure  $\mu$ . This sample can be used to obtain an empirical estimate  $(\widehat{\psi}_1)_n$  of  $\psi_k(\mu)$ , through the consideration of the  $\binom{n}{k+1}$   $k$ -dimensional simplices that can be constructed with the  $x_i$ . Below we show how a much simpler (and still unbiased) estimation of  $\psi_k(\mu)$  can be obtained through the empirical variance-covariance matrix of the sample. See also [25, 26].

Denote

$$\begin{aligned}\widehat{x}_n &= \frac{1}{n} \sum_{i=1}^n x_i, \\ \widehat{V}_n &= \frac{1}{n-1} \sum_{i=1}^n (x_i - \widehat{x}_n)(x_i - \widehat{x}_n)^\top = \frac{1}{n(n-1)} \sum_{i < j} (x_i - x_j)(x_i - x_j)^\top,\end{aligned}$$

respectively the empirical mean and variance-covariance matrix of  $x_1$ . Note that both are unbiased. We thus have

$$(\widehat{\psi}_1)_n = \frac{2}{n(n-1)} \sum_{i < j} \|x_i - x_j\|^2 = 2 \operatorname{trace}[\widehat{V}_n] = \Psi_1(\widehat{V}_n) = \psi_1(\mu_n),$$

with  $\mu_n$  the empirical measure of the sample, and the estimator  $(\widehat{\psi}_1)_n$  is unbiased. More generally, for  $k \geq 1$  we have the following.

**Theorem 3.2.** *For  $x_1, \dots, x_n$  a sample of  $n$  vectors of  $\mathbb{R}^d$ , i.i.d. with the measure  $\mu$ , and for any  $k \in \{1, \dots, d\}$ , the quantity*

$$(\widehat{\psi}_k)_n = \frac{(n-k-1)!(n-1)^k}{(n-1)!} \Psi_k(\widehat{V}_n) = \frac{(n-k-1)!(n-1)^k}{(n-1)!} \psi_k(\mu_n), \quad (3.8)$$

*forms an unbiased estimator of  $\psi_k(\mu)$  and has minimum variance among all unbiased estimators.*

*Proof.* Denote

$$(\widehat{\psi}_k)_n = \binom{n}{k+1}^{-1} \sum_{j_1 < j_2 < \dots < j_{k+1}} \mathcal{V}_k^2(x_{j_1}, \dots, x_{j_{k+1}}). \quad (3.9)$$

It forms a U-statistics for the estimation of  $\psi_k(\mu)$  and is thus unbiased and has minimum variance, see, e.g., [20, Chap. 5]. We only need to show that it can be written as (3.8).

We can write

$$\begin{aligned}(\widehat{\psi}_k)_n &= \binom{n}{k+1}^{-1} \\ &\times \sum_{j_1 < j_2 < \dots < j_{k+1}} \frac{1}{(k!)^2} \sum_{i_1 < i_2 < \dots < i_k} \det^2 \begin{bmatrix} \{x_{j_1}\}_{i_1} & \cdots & \{x_{j_{k+1}}\}_{i_1} \\ \vdots & \vdots & \vdots \\ \{x_{j_1}\}_{i_k} & \cdots & \{x_{j_{k+1}}\}_{i_k} \\ 1 & \cdots & 1 \end{bmatrix}, \\ &= \binom{n}{k+1}^{-1} \frac{1}{(k!)^2} \sum_{i_1 < i_2 < \dots < i_k} \det \left( \sum_{j=1}^n \{z_j\}_{i_1, \dots, i_k} \{z_j\}_{i_1, \dots, i_k}^\top \right),\end{aligned}$$



where we have used Binet-Cauchy formula and where  $\{z_j\}_{i_1, \dots, i_k}$  denotes the  $k+1$  dimensional vector with components  $\{x_j\}_{i_\ell}$ ,  $\ell = 1, \dots, k$ , and 1. This gives

$$\begin{aligned}
(\widehat{\psi}_k)_n &= \binom{n}{k+1}^{-1} \frac{n^{k+1}}{(k!)^2} \sum_{i_1 < i_2 < \dots < i_k} \det \left( \frac{1}{n} \sum_{j=1}^n \{z_j\}_{i_1, \dots, i_k} \{z_j\}_{i_1, \dots, i_k}^\top \right), \\
&= \binom{n}{k+1}^{-1} \frac{n^{k+1}}{(k!)^2} \\
&\quad \times \sum_{i_1 < i_2 < \dots < i_k} \det \begin{bmatrix} (1/n) \{ \sum_{j=1}^n x_j x_j^\top \}_{(i_1, \dots, i_k) \times (i_1, \dots, i_k)} & \{ \widehat{x}_n \}_{i_1, \dots, i_k} \\ \{ \widehat{x}_n \}_{i_1, \dots, i_k}^\top & 1 \end{bmatrix}, \\
&= \binom{n}{k+1}^{-1} \frac{n^{k+1}}{(k!)^2} \sum_{i_1 < i_2 < \dots < i_k} \det \left[ \frac{n-1}{n} \{ \widehat{V}_n \}_{(i_1, \dots, i_k) \times (i_1, \dots, i_k)} \right],
\end{aligned}$$

and thus (3.8).  $\square$

Using the notation of Theorem 3.1, since  $\mathcal{E}_k(V) = (-1)^k a_{d-k}(V)$ , with  $a_{d-k}(V)$  the coefficient of the monomial of degree  $d-k$  of the characteristic polynomial of  $V$ , for a nonsingular  $V$  we obtain

$$\mathcal{E}_k(V) = \det(V) \mathcal{E}_{d-k}(V^{-1}), \quad (3.10)$$

see also [9, Eq. 4.2]. Therefore, we also have

$$(\widehat{\psi}_{d-k})_n = \frac{(n-d+k-1)!(n-1)^{d-k}}{(n-1)!} \frac{(d-k+1)k!}{(k+1)(d-k)!} \det(\widehat{V}_n) \Psi_k(\widehat{V}_n^{-1}), \quad (3.11)$$

which forms an unbiased and minimum-variance estimator of  $\psi_{d-k}(\mu)$ . Note that the estimation of  $\psi_k(\mu)$  is much simpler through (3.8) or (3.11) than using the direct construction (3.9).

One may notice that  $\psi_1(\mu_n)$  is clearly unbiased due to the linearity of  $\Psi_1(\cdot)$ , but it is remarkable that  $\psi_k(\mu_n)$  becomes unbiased after a suitable scaling, see (3.8). Since  $\Psi_k(\cdot)$  is highly nonlinear for  $k > 1$ , this property would not hold if  $\widehat{V}_n$  were replaced by another unbiased estimator of  $V_\mu$ .

The value of  $(\widehat{\psi}_k)_n$  only depend on  $\widehat{V}_n$ , with  $\mathbb{E}\{(\widehat{\psi}_k)_n\} = \psi_k(V_\mu)$ , but its variance depends on the distribution itself. From [20, Lemma A, p. 183], the variance of  $(\widehat{\psi}_k)_n$  satisfies

$$\text{var}[(\widehat{\psi}_k)_n] = \frac{(k+1)^2}{n} \omega + O(n^{-2}),$$

where  $\omega = \text{var}[h(x)]$ , with  $h(x) = \mathbb{E}\{\mathcal{Y}_k^2(x_1, x_2, \dots, x_{k+1}) | x_1 = x\}$ . Obviously,  $\mathbb{E}[h(x)] = \psi_k(\mu)$  and calculations similar to those in the proof of Theorem 3.1 give

$$\begin{aligned}
\omega &= \frac{1}{(k!)^2} \sum_{I, J} \det[\{V_\mu\}_{I \times I}] \det[\{V_\mu\}_{J \times J}] \\
&\quad \times \left[ \mathbb{E} \left\{ (E_\mu - x)_I^\top \{V_\mu\}_{I \times I}^{-1} (E_\mu - x)_I (E_\mu - x)_J^\top \{V_\mu\}_{J \times J}^{-1} (E_\mu - x)_J \right\} - k^2 \right],
\end{aligned} \quad (3.12)$$

where  $I$  and  $J$  respectively denote two sets of indices  $i_1 < i_2 < \dots < i_k$  and  $j_1 < j_2 < \dots < j_k$  in  $\{1, \dots, k+1\}$ , the summation being over all possible such sets. Simplifications occur in some particular cases. For instance, when  $\mu$  is a normal measure, then

$$\begin{aligned}
\omega &= \frac{2}{(k!)^2} \sum_{I, J} \det[\{V_\mu\}_{I \times I}] \det[\{V_\mu\}_{J \times J}] \\
&\quad \times \text{trace} \left[ \{V_\mu\}_{J \times J}^{-1} \{V_\mu\}_{J \times I} \{V_\mu\}_{I \times I}^{-1} \{V_\mu\}_{I \times J} \right].
\end{aligned}$$

If, moreover,  $V_\mu$  is the diagonal matrix  $\text{diag}\{\lambda_1, \dots, \lambda_d\}$ , then

$$\omega = \frac{2}{(k!)^2} \sum_{I,J} \beta(I,J) \prod_I \lambda_i \prod_J \lambda_j,$$

with  $\beta(I,J)$  denoting the number of coincident indices between  $I$  and  $J$  (i.e., the size of  $I \cap J$ ). When  $\mu$  is such that the components of  $x$  are i.d.d. with variance  $\sigma^2$ , then  $V_\mu = \sigma^2 I_d$ , with  $I_d$  the  $d$ -dimensional identity matrix, and

$$\begin{aligned} \mathbb{E} \left\{ (E_\mu - x)_I^\top \{V_\mu\}_{I \times I}^{-1} (E_\mu - x)_I (E_\mu - x)_J^\top \{V_\mu\}_{J \times J}^{-1} (E_\mu - x)_J \right\} = \\ \mathbb{E} \left\{ \left( \sum_{i \in I} z_i^2 \right) \left( \sum_{j \in J} z_j^2 \right) \right\}, \end{aligned}$$

where the  $z_i = \{x - E_\mu\}_i / \sigma$  are i.i.d. with mean 0 and variance 1. We then obtain

$$\omega = \frac{\sigma^{4k}}{(k!)^2} (\mathbb{E}\{z_i^4\} - 1) \beta_{d,k},$$

where

$$\begin{aligned} \beta_{d,k} = \sum_{I,J} \beta(I,J) &= \sum_{i=1}^k i \binom{d}{i} \binom{d-i}{k-i} \binom{d-i-(k-i)}{k-i} \\ &= \frac{(d-k+1)^2}{d} \binom{d}{k-1}^2. \end{aligned}$$

**Example 1** We generate 1,000 independent samples of  $n$  points for different measures  $\mu$ . Figure 1 presents a box-plot of the ratios  $(\hat{\psi}_k)_n / \psi_k(\mu)$  for various values of  $k$  and  $n = 100$  (left),  $n = 1,000$  (right), when  $\mu = \mu_1$  uniform in  $[0, 1]^{10}$ . Figure 2 presents the same information when  $\mu = \mu_2$  which corresponds to the normal distribution  $\mathcal{N}(0, I_{10}/12)$  in  $\mathbb{R}^{10}$ . Note that  $V_{\mu_1} = V_{\mu_2}$  but the dispersions are different in the two figures. Table 1 gives  $(10^3 \times)$  the values of  $\widehat{\mathbb{E}}\{(\hat{\psi}_k)_n\} / \psi_k(\mu) - 1$  for  $\mu = \mu_1, \mu_2$  and the same series of values for  $k$ , with  $\widehat{\mathbb{E}}\{(\hat{\psi}_k)_n\}$  denoting the empirical mean over the 1,000 independent repetitions.

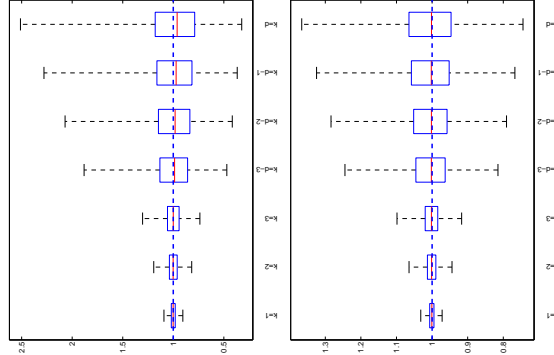


Figure 1: Box-plot of  $(\hat{\psi}_k)_n / \psi_k(\mu)$  for different values of  $k$ :  $\mu$  is uniform in  $[0, 1]^{10}$ ,  $n = 100$  (Left) and  $n = 1,000$  (Right) — 1,000 repetitions; minimum, median and maximum values are indicated, together with 25% and 75% quantiles.

Other properties of U-statistics apply to the estimator  $(\hat{\psi}_k)_n$ , including almost-sure consistency and the classical law of the iterated logarithm, see [20, Section 5.4].

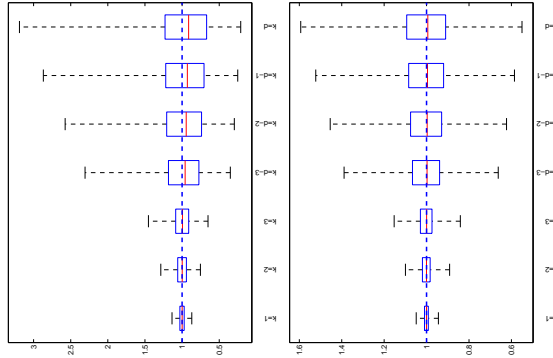


Figure 2: Same as in Figure 1 but for  $\mu$  normal  $\mathcal{N}(0, I_{10}/12)$ .

Table 1:  $\delta(n, k) = 10^3 \times [\widehat{\mathbb{E}}\{(\widehat{\psi}_k)_n\} / \psi_k(\mu) - 1]$  for different values of  $n$  and  $k$ ;  $\mu_1$  is uniform in  $[0, 1]^{10}$ ,  $\mu_2$  is normal  $\mathcal{N}(0, I_{10}/12)$ .

	$k$	1	2	3	$d-3$	$d-2$	$d-1$	$d$
$\mu_1$	$n = 100$	0.9	1.9	2.9	7.7	9.0	10.3	11.6
	$n = 1,000$	0.7	1.4	2.1	4.9	5.6	6.3	7.0
$\mu_2$	$n = 100$	1.2	2.3	3.2	4.3	3.9	3.0	1.9
	$n = 1,000$	0.5	1.0	1.5	3.6	4.1	4.7	5.2

In particular,  $(\widehat{\psi}_k)_n$  is asymptotically normal  $\mathcal{N}(\psi_k(\mu), (k+1)^2\omega/n)$ , with  $\omega$  given by (3.12). This is illustrated in Figure 3-left below for  $\mu$  uniform in  $[0, 1]^{10}$ , with  $n = 1,000$  and  $k = 3$ . The distribution is already reasonably close to normality for small values of  $n$ , see Figure 3-right for which  $n = 20$ .

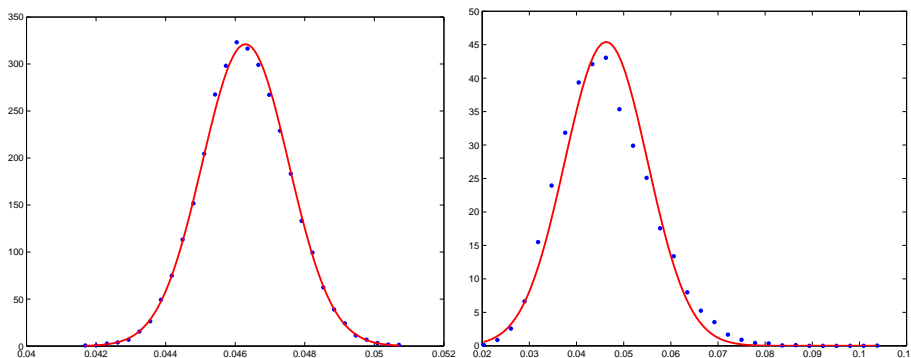


Figure 3: Dots: empirical distribution of  $(\widehat{\psi}_k)_n$  (histogram for 10,000 independent repetitions); solid line: asymptotic normal distribution  $\mathcal{N}(\psi_k(\mu), (k+1)^2\omega/n)$ ;  $\mu$  is uniform in  $[0, 1]^{10}$  and  $k = 3$ ; left:  $n = 1,000$ ; right:  $n = 20$ .

## 4 Maximum-entropy measures and optimal designs

In this section we consider two types of optimisation problems on  $\mathcal{M}$  related to the functions  $\Psi_k(\cdot)$  introduced in Theorem 3.1. First, in Section 4.1, we are interested in the characterisation and construction of maximum-entropy measures; that is, measures  $\mu_k^* \in \mathcal{M}$  which maximize  $\psi_k(\mu) = \Psi_k(V_\mu)$ . The existence of an optimal measure follows from the compactness of  $\mathcal{X}$  and continuity of  $\mathcal{V}_k(x_1, \dots, x_{k+1})$  in

each  $x_i$ , see [1, Th. 1]; the concavity and differentiability of the functional  $\psi_k^{1/k}(\cdot)$  allow us to derive a necessary and sufficient condition for optimality.

In Section 4.2 we consider the problem of optimal design of experiments, where the covariance matrix  $V$  is the inverse of the information matrix  $M(\xi)$  for some regression model.

## 4.1 Maximum-entropy measures

### 4.1.1 Necessary and sufficient condition

Since the functionals  $\psi_k^{1/k}(\cdot)$  are concave and differentiable, for all  $k = 1, \dots, d$ , we can easily derive a necessary and sufficient condition for a probability measure  $\mu_k^*$  on  $\mathcal{X}$  to maximise  $\psi_k(\mu)$ , in the spirit of the celebrated Equivalence Theorem of Kiefer and Wolfowitz [8].

Denote by  $\nabla_{\Psi_k}[V]$  the gradient of  $\Psi_k(\cdot)$  at matrix  $V$  (a matrix of the same size as  $V$ ) and by  $F_{\psi_k}(\mu; \nu)$  the directional derivative of  $\psi_k(\cdot)$  at  $\mu$  in the direction  $\nu$ ;

$$F_{\psi_k}(\mu; \nu) = \lim_{\alpha \rightarrow 0^+} \frac{\psi_k[(1-\alpha)\mu + \alpha\nu] - \psi_k(\mu)}{\alpha}.$$

From the expression (3.6) of  $\Psi_k(V)$ , we have

$$\nabla_{\Psi_k}[V] = \frac{k+1}{k!} \nabla_{\mathcal{E}_k}[V],$$

where  $\nabla_{\mathcal{E}_k}[V]$  denotes the gradient of  $\mathcal{E}_k(\cdot)$  at  $V$ , which, using (3.7), can be shown by induction to satisfy

$$\nabla_{\mathcal{E}_k}[V] = \sum_{i=0}^{k-1} (-1)^i \mathcal{E}_{k-i-1}(V) V^i, \quad (4.1)$$

see [9]. We thus obtain in particular

$$\begin{aligned} k=1: \quad & \nabla_{\Psi_1}[V] = 2I_d, \\ k=2: \quad & \nabla_{\Psi_2}[V] = \frac{3}{2} [\text{trace}(V)I_d - V], \\ k=3: \quad & \nabla_{\Psi_3}[V] = \frac{1}{3} [\text{trace}^2(V) - \text{trace}(V^2)]I_d - \frac{2}{3} \text{trace}(V)V + \frac{2}{3} V^2, \\ k=d: \quad & \nabla_{\Psi_d}[V] = \frac{d+1}{d!} \det(V) V^{-1}. \end{aligned}$$

Using the differentiability of  $\Psi_k(\cdot)$ , direct calculation gives

$$F_{\psi_k}(\mu; \nu) = \text{trace} \left\{ \nabla_{\Psi_k}[V_\mu] \frac{dV_{(1-\alpha)\mu + \alpha\nu}}{d\alpha} \Big|_{\alpha=0} \right\},$$

with

$$\frac{dV_{(1-\alpha)\mu + \alpha\nu}}{d\alpha} \Big|_{\alpha=0} = \int [xx^\top - (E_\mu x^\top + x E_\mu^\top)] \nu(dx) - \int xx^\top \mu(dx) + 2E_\mu E_\mu^\top. \quad (4.2)$$

Notice that  $dV_{(1-\alpha)\mu + \alpha\nu}/d\alpha|_{\alpha=0}$  is linear in  $\nu$ .

Then, from the concavity of  $\psi_k^{1/k}(\cdot)$ ,  $\mu_k^*$  maximises  $\psi_k(\mu)$  with respect to  $\mu \in \mathcal{M}$  if and only if  $\psi_k(\mu_k^*) > 0$  and  $F_{\psi_k}(\mu_k^*; \nu) \leq 0$  for all  $\nu \in \mathcal{M}$ , that is

$$\text{trace} \left\{ \nabla_{\Psi_k}[V_{\mu_k^*}] \frac{dV_{(1-\alpha)\mu_k^* + \alpha\nu}}{d\alpha} \Big|_{\alpha=0} \right\} \leq 0, \quad \forall \nu \in \mathcal{M}. \quad (4.3)$$

We obtain the following.

**Theorem 4.1.** *The probability measure  $\mu_k^*$  such that  $\psi_k(\mu_k^*) > 0$  is  $\psi_k$ -optimal, that is, maximises  $\psi_k(\mu)$  with respect to  $\mu \in \mathcal{M}$ ,  $k \in \{1, \dots, d\}$ , if and only if*

$$\max_{x \in \mathcal{X}} (x - E_{\mu_k^*})^\top \frac{\nabla_{\Psi_k}[V_{\mu_k^*}]}{\Psi_k(V_{\mu_k^*})} (x - E_{\mu_k^*}) \leq k. \quad (4.4)$$

Moreover,

$$(x - E_{\mu_k^*})^\top \frac{\nabla_{\Psi_k}[V_{\mu_k^*}]}{\Psi_k(V_{\mu_k^*})} (x - E_{\mu_k^*}) = k \quad (4.5)$$

for all  $x$  in the support of  $\mu_k^*$ .

*Proof.* First note that the Newton equations (3.7) and the recurrence (4.1) for  $\nabla_{\mathcal{E}_k}[\cdot]$  imply that  $\text{trace}(V \nabla_{\Psi_k}[V]) = k \Psi_k(V)$  for all  $k = 1, \dots, d$ .

The condition (4.4) is sufficient. Indeed, suppose that  $\mu_k^*$  such that  $\psi_k(\mu_k^*) > 0$  satisfies (4.4). We obtain

$$\int (x - E_{\mu_k^*})^\top \nabla_{\Psi_k}[V_{\mu_k^*}](x - E_{\mu_k^*}) \nu(dx) \leq \text{trace} \{V_{\mu_k^*} \nabla_{\Psi_k}[V_{\mu_k^*}]\}$$

for any  $\nu \in \mathcal{M}$ , which gives (4.3) when we use (4.2). The condition is also necessary since (4.3) must be true in particular for  $\delta_x$ , the delta measure at any  $x \in \mathcal{X}$ , which gives (4.4). The property (4.5) on the support of  $\mu_k^*$  follows from the observation that  $\int (x - E_{\mu_k^*})^\top \nabla_{\Psi_k}[V_{\mu_k^*}](x - E_{\mu_k^*}) \mu_k^*(dx) = \text{trace} \{V_{\mu_k^*} \nabla_{\Psi_k}[V_{\mu_k^*}]\}$ .  $\square$

Note that for  $k < d$ , the covariance matrix  $V_{\mu_k^*}$  of a  $\psi_k$ -optimal measure  $\mu_k^*$  is not necessarily unique and may be singular; see, e.g., Example 1 below. Also,  $\psi_k(\mu) > 0$  implies that  $\psi_{k-1}(\mu) > 0$ ,  $k = 2, \dots, d$ .

**Remark 4.1.** *As a natural extension of the concept of potential in case of order-two interactions ( $k = 1$ ), we call  $P_{k,\mu}(x) = \psi_k(\mu, \dots, \mu, \delta_x)$  the potential of  $\mu$  at  $x$ , where*

$$\psi_k(\mu_1, \dots, \mu_{k+1}) = \int \dots \int \mathcal{Y}_k^2(x_1, \dots, x_{k+1}) \mu_1(dx_1) \dots \mu_{k+1}(dx_{k+1}).$$

*This yields  $F_{\psi_k}(\mu; \nu) = (k+1)[\psi_k(\mu, \dots, \mu, \nu) - \psi_k(\mu)]$ , where  $\mu$  appears  $k$  times in  $\psi_k(\mu, \dots, \mu, \nu)$ . Therefore, Theorem 4.4 states that  $\mu_k^*$  is  $\psi_k$ -optimal if and only if  $\psi_k(\mu_k^*, \dots, \mu_k^*, \nu) \leq \psi_k(\mu_k^*)$  for any  $\nu \in \mathcal{M}$ , or equivalently  $P_{k,\mu_k^*}(x) \leq \psi_k(\mu_k^*)$  for all  $x \in \mathcal{X}$ .*

**Remark 4.2.** *Consider Kiefer's  $\Phi_p$ -class of orthogonally invariant criteria and their associated functional  $\varphi_p(\cdot)$ , defined by*

$$\varphi_p(\mu) = \Phi_p(V_\mu) = \begin{cases} \lambda_{\max}(V_\mu) & \text{for } p = \infty, \\ \{\frac{1}{d} \text{trace}(V_\mu^p)\}^{1/d} & \text{for } p \neq 0, \pm\infty, \\ \det^{1/d}(V_\mu) & \text{for } p = 0, \\ \lambda_{\min}(V_\mu) & \text{for } p = -\infty, \end{cases}$$

*where  $V_\mu$  is a  $d \times d$  matrix; see, e.g., [14, Chap. 6]. From a result in [7], if a measure  $\mu_p$  optimal for some  $\varphi_p(\cdot)$  with  $p \in (-\infty, 1]$  is such that  $V_{\mu_p}$  is proportional to the identity matrix  $I_d$ , then  $\mu_p$  is simultaneously optimal for all orthogonally invariant criteria. A measure  $\mu_p$  having this property is therefore  $\psi_k$ -optimal for all  $k = 1, \dots, d$ . Notice that  $\psi_1(\cdot)$  and  $\psi_d^{1/d}(\cdot)$  respectively coincide with  $\varphi_1(\cdot)$  and  $\varphi_0(\cdot)$  (up to a multiplicative scalar).*

**Remark 4.3.** Using (3.10), when  $V$  is nonsingular we obtain the property

$$\Psi_k(V) = \frac{(k+1)(d-k)!}{(d-k+1)k!} \det(V) \Psi_{d-k}(V^{-1})$$

which implies that maximising  $\Psi_k(V)$  is equivalent to maximising  $\log \det(V) + \log \Psi_{d-k}(V^{-1})$ . Therefore, Theorem 4.4 can be reformulated as:  $\mu_k^*$  maximises  $\psi_k(\mu)$  if and only if

$$\max_{x \in \mathcal{X}} (x - E_{\mu_k^*})^\top \left[ V_{\mu_k^*}^{-1} - V_{\mu_k^*}^{-1} \frac{\nabla_{\Psi_{d-k}}[V_{\mu_k^*}^{-1}]}{\Psi_{d-k}(V_{\mu_k^*}^{-1})} V_{\mu_k^*}^{-1} \right] (x - E_{\mu_k^*}) \leq d - k,$$

with equality for  $x$  in the support of  $\mu_k^*$ . When  $k$  is large (and  $d - k$  is small), one may thus check the optimality of  $\mu_k^*$  without using the complicated expressions of  $\Psi_k(V)$  and  $\nabla_{\Psi_k}[V]$ .

#### 4.1.2 A duality property

The characterisation of maximum-entropy measures can also be approached from the point of view of duality theory.

When  $k = 1$ , the determination of a  $\psi_1$ -optimal measure  $\mu_1^*$  is equivalent to the dual problem of constructing the minimum-volume ball  $\mathcal{B}_d^*$  containing  $\mathcal{X}$ . If this ball has radius  $\rho$ , then  $\psi_1(\mu_1^*) = 2\rho^2$ , and the support points of  $\mu_1^*$  are the points of contact between  $\mathcal{X}$  and  $\mathcal{B}_d^*$ ; see [1, Th. 6]. Moreover, there exists an optimal measure with no more than  $d + 1$  points.

The determination of an optimal measure  $\mu_d^*$  is also dual to a simple geometrical problem: it corresponds to the determination of the minimum-volume ellipsoid  $\mathcal{E}_d^*$  containing  $\mathcal{X}$ . This is equivalent to a  $D$ -optimal design problem in  $\mathbb{R}^{d+1}$  for the estimation of  $\beta = (\beta_0, \beta_1^\top)^\top$ ,  $\beta_1 \in \mathbb{R}^d$ , in the linear regression model with intercept  $\beta_0 + \beta_1^\top x$ ,  $x \in \mathcal{X}$ , see [23]. Indeed, denote

$$W_\mu = \int_{\mathcal{X}} (1 \ x^\top)^\top (1 \ x^\top) \mu(dx).$$

Then  $\mathcal{E}_{d+1}^* = \{z \in \mathbb{R}^{d+1} : z^\top W_{\mu_d^*}^{-1} z \leq d + 1\}$ , with  $\mu_d^*$  maximising  $\det(W_\mu)$ , is the minimum-volume ellipsoid centered at the origin and containing the set  $\{z \in \mathbb{R}^{d+1} : z = (1 \ x^\top)^\top, x \in \mathcal{X}\}$ . Moreover,  $\mathcal{E}_d^*$  corresponds to the intersection between  $\mathcal{E}_{d+1}^*$  and the hyperplane  $\{z\}_1 = 1$ ; see, e.g., [22]. This gives  $\psi_d(\mu_d^*) = (d+1)/d! \det(W_{\mu_d^*})$ . The support points of  $\mu_d^*$  are the points of contact between  $\mathcal{X}$  and  $\mathcal{E}_d^*$ , there exists an optimal measure with no more than  $d(d+3)/2 + 1$  points, see [23].

The property below generalises this duality property to any  $k \in \{1, \dots, d\}$ .

#### Theorem 4.2.

$$\max_{\mu \in \mathcal{M}} \Psi_k^{1/k}(V_\mu) = \min_{M, c: \mathcal{X} \subset \mathcal{E}(M, c)} \frac{1}{\phi_k^\infty(M)},$$

where  $\mathcal{E}(M, c)$  denotes the ellipsoid  $\mathcal{E}(M, c) = \{x \in \mathbb{R}^d : (x - c)^\top M (x - c) \leq 1\}$  and  $\phi_k^\infty(M)$  is the polar function

$$\phi_k^\infty(M) = \inf_{V \succeq 0: \text{trace}(MV)=1} \frac{1}{\Psi_k^{1/k}(V)}. \quad (4.6)$$

The proof is given in Appendix. The polar function  $\phi_k^\infty(\cdot)$  possesses the properties of what is called an information function in [14, Chap. 5]; in particular, it is concave on the set of symmetric non-negative definite matrices. This duality property has the following consequence.

**Corollary 4.1.** *The determination of a covariance matrix  $V_k^*$  that maximises  $\Psi_k(V_\mu)$  with respect to  $\mu \in \mathcal{M}$  is equivalent to the determination of an ellipsoid  $\mathcal{E}(M_k^*, c_k^*)$  containing  $\mathcal{X}$ , minimum in the sense that  $M_k^*$  maximizes  $\phi_k^\infty(M)$ . The points of contact between  $\mathcal{E}(M_k^*, c_k^*)$  and  $\mathcal{X}$  form the support of  $\mu_k^*$ .*

For any  $V \succeq 0$ , denote by  $M_*(V)$  the matrix

$$M_*(V) = \frac{\nabla_{\Psi_k}[V]}{k \Psi_k(V)} = \frac{1}{k} \nabla_{\log \Psi_k}[V]. \quad (4.7)$$

Note that  $M_*(V) \succeq 0$ , see [14, Lemma 7.5], and that

$$\text{trace}[VM_*(V)] = 1,$$

see the proof of Theorem 4.4. The matrix  $V \succeq 0$  maximises  $\Psi_k(V)$  under the constraint  $\text{trace}(MV) = 1$  for some  $M \succeq 0$  if and only if  $V[M_*(V) - M] = 0$ . Therefore, if  $M$  is such that there exists  $V_* = V_*(M) \succeq 0$  such that  $M = M_*[V_*(M)]$ , then  $\phi_k^\infty(M) = \Psi_k^{-1/k}[V_*(M)]$ . When  $k < d$ , the existence of such a  $V_*$  is not ensured for all  $M \succeq 0$ , but happens when  $M = M_k^*$  which maximises  $\phi_k^\infty(M)$  under the constraint  $\mathcal{X} \in \mathcal{E}(M, c)$ . Moreover, in that case there exists a  $\mu_k^* \in \mathcal{M}$  such that  $M_k^* = M_*(V_{\mu_k^*})$ , and this  $\mu_k^*$  maximises  $\psi_k(\mu)$  with respect to  $\mu \in \mathcal{M}$ .

Consider in particular the case  $k = 1$ . Then,  $M_*(V) = I_d/\text{trace}(V)$  and  $\phi_1^\infty(M) = \lambda_{\min}(M)/2$ . The matrix  $M_k^*$  of the optimal ellipsoid  $\mathcal{E}(M_k^*, c_k^*)$  is proportional to the identity matrix and  $\mathcal{E}(M_k^*, c_k^*)$  is the ball of minimum-volume that encloses  $\mathcal{X}$ .

When  $k = 2$  and  $I_d \succeq (d-1)M/\text{trace}(M)$ , direct calculations show that  $\phi_2^\infty(M) = \Psi_2^{-1/2}[V_*(M)]$ , with

$$V_*(M) = [I_d \text{trace}(M)/(d-1) - M][\text{trace}^2(M)/(d-1) - \text{trace}(M^2)]^{-1};$$

the optimal ellipsoid is then such that  $\text{trace}^2(M)/(d-1) - \text{trace}(M^2)$  is maximised.

### 4.1.3 Examples

**Example 2** Take  $\mathcal{X} = [0, 1]^d$ ,  $d \geq 1$  and denote by  $v_i$ ,  $i = 1, \dots, 2^d$  the  $2^d$  vertices of  $\mathcal{X}$ . Consider  $\mu^* = (1/2^d) \sum_{i=1}^{2^d} \delta_{v_i}$ , with  $\delta_v$  the Dirac delta measure at  $v$ . Then,  $V_{\mu^*} = I_d/4$  and one can easily check that  $\mu^*$  is  $\psi_1$ -optimal. Indeed,  $E_{\mu^*} = \mathbf{1}_d/2$ , with  $\mathbf{1}_d$  the  $d$ -dimensional vector of ones, and  $\max_{x \in \mathcal{X}} (x - \mathbf{1}_d/2)^\top (2I_d)(x - \mathbf{1}_d/2) = d/2 = \text{trace}\{V_{\mu^*} \nabla_{\Psi_1}[V_{\mu^*}]\}$ . From Remark 4.2, the measure  $\mu^*$  is  $\psi_k$ -optimal for all  $k = 1, \dots, d$ .

Note that the two-point measure  $\mu_1^* = (1/2)[\delta_{\mathbf{0}} + \delta_{\mathbf{1}_d}]$  is such that  $V_{\mu_1^*} = (\mathbf{1}_d \mathbf{1}_d^\top)/4$  and  $\psi_1(\mu_1^*) = d/2 = \psi_1(\mu^*)$ , and is therefore  $\psi_1$ -optimal too. It is not  $\psi_k$ -optimal for  $k > 1$ , since  $\psi_k(\mu_1^*) = 0$ ,  $k > 1$ .

**Example 3** Take  $\mathcal{X} = \mathcal{B}_d(\mathbf{0}, \rho)$ , the closed ball of  $\mathbb{R}^d$  centered at the origin  $\mathbf{0}$  with radius  $\rho$ . Let  $\mu_0$  be the uniform measure on the sphere  $\mathcal{S}_d(\mathbf{0}, \rho)$  (the boundary of  $\mathcal{B}_d(\mathbf{0}, \rho)$ ). Then,  $V_{\mu_0}$  is proportional to the identity matrix  $I_d$ , and  $\text{trace}[V_{\mu_0}] = \rho^2$  implies that  $V_{\mu_0} = \rho^2 I_d/d$ . Take  $k = d$ . We have  $E_{\mu_0} = \mathbf{0}$  and

$$\max_{x \in \mathcal{X}} (x - E_{\mu_0})^\top \nabla_{\Psi_d}[V_{\mu_0}](x - E_{\mu_0}) = \frac{(d+1)\rho^{2d}}{d^{d-1}d!} = \text{trace}\{V_{\mu_0} \nabla_{\Psi_d}[V_{\mu_0}]\},$$

so that  $\mu_0$  is  $\psi_d$ -optimal from (4.4).

Let  $\mu_d$  be the measure that allocates mass  $1/(d+1)$  at each vertex of a  $d$  regular simplex having its  $d+1$  vertices on  $\mathcal{S}_d(\mathbf{0}, \rho)$ , with squared volume  $\rho^{2d}(d+$

$1)^{d+1}/[d^d(d!)^2]$ . We also have  $V_{\mu_d} = \rho^2 I_d/d$ , so that  $\mu_d$  is  $\psi_d$ -optimal too. In view of Remark 4.2,  $\mu_0$  and  $\mu_d$  are  $\psi_k$ -optimal for all  $k$  in  $\{1, \dots, d\}$ .

Let now  $\mu_k$  be the measure that allocates mass  $1/(k+1)$  at each vertex of a  $k$  regular simplex  $\mathcal{P}_k$ , centered at the origin, with its vertices on  $\mathcal{S}_d(\mathbf{0}, \rho)$ . The squared volume of  $\mathcal{P}_k$  equals  $\rho^{2k} (k+1)^{k+1}/[k^k(k!)^2]$ . Without any loss of generality, we can choose the orientation of the space such that  $V_{\mu_k}$  is diagonal, with its first  $k$  diagonal elements equal to  $\rho^2/k$  and the other equal to zero. Note that  $\psi_{k'}(\mu_k) = 0$  for  $k' > k$ . Direct calculations based on (3.6) give

$$\psi_k(\mu_k) = \frac{k+1}{k!} \frac{\rho^{2k}}{k^k} \leq \psi_k(\mu_0) = \frac{k+1}{k!} \binom{d}{k} \frac{\rho^{2k}}{d^k},$$

with equality for  $k = 1$  and  $k = d$ , the inequality being strict otherwise. Figure 4 presents the efficiency  $[\psi_k(\mu_k)/\psi_k(\mu_0)]^{1/k}$  as a function of  $k$  when  $d = 20$ .

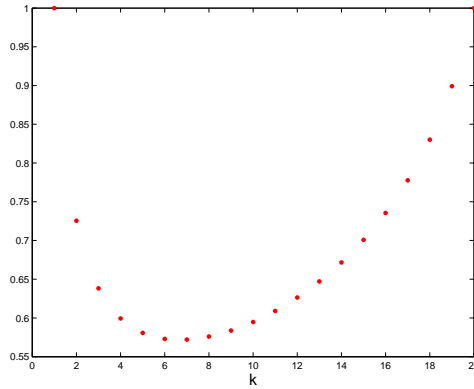


Figure 4: Efficiency  $[\psi_k(\mu_k)/\psi_k(\mu_0)]^{1/k}$  as a function of  $k$  when  $d = 20$  in Example 3.

## 4.2 Optimal design in regression models

In this section we consider the case when  $V = M^{-1}(\xi)$ , where  $M(\xi)$  is the information matrix

$$M(\xi) = \int_{\mathbb{T}} f(t)f^\top(t) \xi(dt)$$

in a regression model  $Y_j = \theta^\top f(t_j) + \varepsilon_j$  with parameters  $\theta \in \mathbb{R}^d$ , for a design measure  $\xi \in \Xi$ . Here  $\Xi$  denotes the set of probability measures on a set  $\mathbb{T}$  such that  $\{f(t) : t \in \mathbb{T}\}$  is compact, and  $M^{-1}(\xi)$  is the (asymptotic) covariance matrix of an estimator  $\hat{\theta}$  of  $\theta$  when the design variables  $t$  are distributed according to  $\xi$ . The value  $\psi_k(\mu)$  of Theorem 3.1 defines a measure of dispersion for  $\hat{\theta}$ , that depends on  $\xi$  through  $V_\mu = M^{-1}(\xi)$ . The design problem we consider consists in choosing  $\xi$  that minimises this dispersion, as measured by  $\Psi_k[M^{-1}(\xi)]$ , or equivalently that maximises  $\Psi_k^{-1}[M^{-1}(\xi)]$ .

### 4.2.1 Properties

It is customary in optimal design theory to maximise a concave and Loewner-increasing function of  $M(\xi)$ , see [14, Chap. 5] for desirable properties of optimal design criteria. Here we have the following.



**Theorem 4.3.** *The functions  $M \rightarrow \Psi_k^{-1/k}(M^{-1})$ ,  $k = 1, \dots, d$ , are Loewner-increasing, concave and differentiable on the set  $\mathbb{M}^+$  of  $d \times d$  symmetric positive-definite matrices. The functions  $\Psi_k(\cdot)$  are also orthogonally invariant.*

*Proof.* The property (3.10) yields

$$\Psi_k^{-1/k}(M^{-1}) = \left( \frac{k+1}{k!} \right)^{-1/k} \frac{\det^{1/k}(M)}{\mathcal{E}_{d-k}^{1/k}(M)} \quad (4.8)$$

which is a concave function of  $M$ , see Eq. (10) of [11, p. 116]. Since  $\Psi_k(\cdot)$  is Loewner-increasing, see [9], the function  $M \rightarrow \Psi_k^{-1/k}(M^{-1})$  is Loewner-increasing too. Its orthogonal invariance follows from the fact that it is defined in terms of the eigenvalues of  $M$ .  $\square$

Note that Theorems 3.1 and 4.3 imply that the functions  $M \rightarrow -\log \Psi_k(M)$  and  $M \rightarrow \log \Psi_k(M^{-1})$  are convex for all  $k = 1, \dots, d$ , a question which was left open in [9].

As a consequence of Theorem 4.3, we can derive a necessary and sufficient condition for a design measure  $\xi_k^*$  to maximise  $\Psi_k^{-1/k}[M^{-1}(\xi)]$  with respect to  $\xi \in \Xi$ , for  $k = 1, \dots, d$ .

**Theorem 4.4.** *The design measure  $\xi_k^*$  such that  $M(\xi_k^*) \in \mathbb{M}^+$  maximises  $\tilde{\psi}_k(\xi) = \Psi_k^{-1/k}[M^{-1}(\xi)]$  with respect to  $\xi \in \Xi$  if and only if*

$$\max_{t \in \mathbb{T}} f^\top(t) M^{-1}(\xi_k^*) \frac{\nabla_{\Psi_k}[M^{-1}(\xi_k^*)]}{\Psi_k[M^{-1}(\xi_k^*)]} M^{-1}(\xi_k^*) f(t) \leq k \quad (4.9)$$

or, equivalently,

$$\max_{t \in \mathbb{T}} \left\{ f^\top(t) M^{-1}(\xi_k^*) f(t) - f^\top(t) \frac{\nabla_{\Psi_{d-k}}[M(\xi_k^*)]}{\Psi_{d-k}[M(\xi_k^*)]} f(t) \right\} \leq d - k. \quad (4.10)$$

Moreover, there is equality in (4.9) and (4.10) for all  $t$  in the support of  $\xi_k^*$ .

*Proof.* From (4.8), the maximisation of  $\tilde{\psi}_k(\xi)$  is equivalent to the maximisation of  $\tilde{\phi}_k(\xi) = \log \det[M(\xi)] - \log \Psi_{d-k}[M(\xi)]$ . The proof is similar to that of Theorem 4.4 and is based on the following expressions for the directional derivatives of these two functionals at  $\xi$  in the direction  $\nu \in \Xi$ ,

$$F_{\tilde{\psi}_k}(\xi; \nu) = \text{trace} \left( \frac{1}{k} M^{-1}(\xi) \frac{\nabla_{\Psi_k}[M^{-1}(\xi)]}{\Psi_k[M^{-1}(\xi)]} M^{-1}(\xi) [M(\nu) - M(\xi)] \right)$$

and

$$F_{\tilde{\phi}_k}(\xi; \nu) = \text{trace} \left( \left\{ M^{-1}(\xi) - \frac{\nabla_{\Psi_{d-k}}[M(\xi)]}{\Psi_{d-k}[M(\xi)]} \right\} [M(\nu) - M(\xi)] \right),$$

and on the property  $\text{trace}\{M \nabla_{\Psi_j}[M]\} = j \Psi_j(M)$ .  $\square$

In particular, consider the following special cases for  $k$  (note that  $\Psi_0(M) = \mathcal{E}_0(M) = 1$  for any  $M$ ).

$$\begin{aligned} k = d : & \quad \tilde{\psi}_d(\xi) = \log \det[M(\xi)], \\ k = d - 1 : & \quad \tilde{\psi}_{d-1}(\xi) = \log \det[M(\xi)] - \log \text{trace}[M(\xi)] - \log 2, \\ k = d - 2 : & \quad \tilde{\psi}_{d-2}(\xi) = \log \det[M(\xi)] \\ & \quad - \log \{ \text{trace}^2[M(\xi)] - \text{trace}[M^2(\xi)] \} - \log(3/4). \end{aligned}$$

The necessary and sufficient condition (4.10) then takes the following form:

$$\begin{aligned}
k = d : & \quad \max_{t \in \mathbb{T}} f^\top(t) M^{-1}(\xi_k^*) f(t) \leq d, \\
k = d - 1 : & \quad \max_{t \in \mathbb{T}} \left\{ f^\top(t) M^{-1}(\xi_k^*) f(t) - \frac{f^\top(t) f(t)}{\text{trace}[M(\xi_k^*)]} \right\} \leq d - 1, \\
k = d - 2 : & \quad \max_{t \in \mathbb{T}} \left\{ f^\top(t) M^{-1}(\xi_k^*) f(t) \right. \\
& \quad \left. - 2 \frac{\text{trace}[M(\xi_k^*)] f^\top(t) f(t) - f^\top(t) M(\xi_k^*) f(t)}{\text{trace}^2[M(\xi_k^*)] - \text{trace}[M^2(\xi_k^*)]} \right\} \leq d - 2.
\end{aligned}$$

Also, for  $k = 1$  condition (4.9) gives

$$\max_{t \in \mathbb{T}} f^\top(t) \frac{M^{-2}(\xi_1^*)}{\text{trace}[M^{-1}(\xi_1^*)]} f(t) \leq 1$$

(which corresponds to  $A$ -optimal design), and for  $k = 2$

$$\max_{t \in \mathbb{T}} \frac{\text{trace}[M^{-1}(\xi_2^*)] f^\top(t) M^{-2}(\xi_2^*) f(t) - f^\top(t) M^{-3}(\xi_2^*) f(t)}{\text{trace}^2[M^{-1}(\xi_2^*)] - \text{trace}[M^{-2}(\xi_2^*)]} \leq 1.$$

Finally, note that a duality theorem, in the spirit of Theorem 4.2, can be formulated for the maximisation of  $\Psi_k^{-1/k}[M^{-1}(\xi)]$ ; see [14, Th. 7.12] for the general form a such duality properties in optimal experimental design.

#### 4.2.2 Examples

**Example 4** For the linear regression model on  $\theta_0 + \theta_1 x$  on  $[-1, 1]$ , the optimal design for  $\tilde{\psi}_k(\cdot)$  with  $k = d = 2$  or  $k = 1$  is

$$\xi_k^* = \left\{ \begin{array}{cc} -1 & 1 \\ 1/2 & 1/2 \end{array} \right\},$$

where the first line corresponds to support points and the second indicates their respective weights.

**Example 5** For linear regression with the quadratic polynomial model  $\theta_0 + \theta_1 t + \theta_2 t^2$  on  $[-1, 1]$ , the optimal designs for  $\tilde{\psi}_k(\cdot)$  have the form

$$\xi_k^* = \left\{ \begin{array}{ccc} -1 & 0 & 1 \\ w_k & 1 - 2w_k & w_k \end{array} \right\},$$

with  $w_3 = 1/3$ ,  $w_2 = (\sqrt{33} - 1)/16 \simeq 0.2965352$  and  $w_1 = 1/4$ . Define the efficiency  $\text{Eff}_k(\xi)$  of a design  $\xi$  as

$$\text{Eff}_k(\xi) = \frac{\tilde{\psi}_k(\xi)}{\tilde{\psi}_k(\xi_k^*)}.$$

Table 2 gives the efficiencies  $\text{Eff}_k(\xi_j^*)$  for  $j, k = 1, \dots, d = 3$ . The design  $\xi_2^*$ , optimal for  $\tilde{\psi}_2(\cdot)$ , appears to make a good compromise between  $A$ -optimality (which corresponds to  $\tilde{\psi}_1(\cdot)$ ) and  $D$ -optimality (which corresponds to  $\tilde{\psi}_3(\cdot)$ ).

**Example 6** For linear regression with the cubic polynomial model  $\theta_0 + \theta_1 t + \theta_2 t^2 + \theta_3 t^3$  on  $[-1, 1]$ , the optimal designs for  $\tilde{\psi}_k(\cdot)$  have the form

$$\xi_k^* = \left\{ \begin{array}{cccc} -1 & -z_k & z_k & 1 \\ w_k & 1/2 - w_k & 1/2 - w_k & w_k \end{array} \right\},$$

Table 2: Efficiencies  $\text{Eff}_k(\xi_j^*)$  for  $j, k = 1, \dots, d$  in Example 5.

	$\text{Eff}_1$	$\text{Eff}_2$	$\text{Eff}_3$
$\xi_1^*$	1	0.9770	0.9449
$\xi_2^*$	0.9654	1	0.9886
$\xi_3^*$	0.8889	0.9848	1

Table 3: Efficiencies  $\text{Eff}_k(\xi_j^*)$  for  $j, k = 1, \dots, d$  in Example 6.

	$\text{Eff}_1$	$\text{Eff}_2$	$\text{Eff}_3$	$\text{Eff}_4$
$\xi_1^*$	1	0.9785	0.9478	0.9166
$\xi_2^*$	0.9694	1	0.9804	0.9499
$\xi_3^*$	0.9180	0.9753	1	0.9897
$\xi_4^*$	0.8527	0.9213	0.9872	1

where

$$\begin{aligned} z_4 &= 1/\sqrt{5} \simeq 0.4472136, & w_4 &= 0.25, \\ z_3 &\simeq 0.4350486, & w_3 &\simeq 0.2149859, \\ z_2 &\simeq 0.4240013, & w_2 &\simeq 0.1730987, \\ z_1 &= \sqrt{3\sqrt{7} - 6}/3 \simeq 0.4639509, & w_1 &= (4 - \sqrt{7})/9 \simeq 0.1504721, \end{aligned}$$

with  $z_3$  satisfying the equation  $2z^6 - 3z^5 - 45z^4 + 6z^3 - 4z^2 - 15z + 3 = 0$  and

$$w_3 = \frac{5z^6 + 5z^4 + 5z^2 + 1 - \sqrt{z^{12} + 2z^{10} + 3z^8 + 60z^6 + 59z^4 + 58z^2 + 73}}{12(z^6 + z^4 + z^2 - 3)},$$

with  $z = z_3$ . For  $k = d - 2 = 2$ , the numbers  $z_2$  and  $w_2$  are too difficult to express analytically. Table 3 gives the efficiencies  $\text{Eff}_k(\xi_j^*)$  for  $j, k = 1, \dots, d$ . Here again the design  $\xi_2^*$  appears to make a good compromise: it maximises the minimum efficiency  $\min_k \text{Eff}_k(\cdot)$  among the designs considered.

## Appendix

**Shift-invariance and positive homogeneity** Denote by  $\mathcal{M}$  the set of probability measures defined on the Borel subsets of  $\mathcal{X}$ , a compact subset of  $\mathbb{R}^d$ . For any  $\mu \in \mathcal{M}$ , any  $\theta \in \mathbb{R}^d$  and any  $\lambda \in \mathbb{R}^+$ , respectively denote by  $T_{-\theta}[\mu]$  and  $H_{\lambda^{-1}}[\mu]$  the measures defined by:

$$\text{for any } \mu\text{-measurable } \mathcal{A} \subseteq \mathcal{X}, \quad T_{-\theta}[\mu](\mathcal{A} + \theta) = \mu(\mathcal{A}), \quad H_{\lambda^{-1}}[\mu](\lambda\mathcal{A}) = \mu(\mathcal{A}),$$

where  $\mathcal{A} + \theta = \{x + \theta : x \in \mathcal{A}\}$  and  $\lambda\mathcal{A} = \{\lambda x : x \in \mathcal{A}\}$ . The shift-invariance of  $\phi(\cdot)$  then means that  $\phi(T_{-\theta}[\mu]) = \phi(\mu)$  for any  $\mu \in \mathcal{M}$  and any  $\theta \in \mathbb{R}^d$ , positive homogeneity of degree  $q$  means that  $\phi(H_{\lambda^{-1}}[\mu]) = \lambda^q \phi(\mu)$  for any  $\mu \in \mathcal{M}$  and any  $\lambda \in \mathbb{R}^+$ .

**The variance is the only concave central moment** For  $q \neq 2$ , the  $q$ -th central moment  $\Delta_q(\mu) = \int \|x - E_\mu\|^q \mu(dx)$  is shift-invariant and homogeneous of degree  $q$ , but it is not concave on  $\mathcal{M}(\mathcal{X})$ . Indeed, consider for instance the two-point probability measures

$$\mu_1 = \left\{ \begin{array}{cc} 0 & 1 \\ 1/2 & 1/2 \end{array} \right\} \text{ and } \mu_2 = \left\{ \begin{array}{cc} 0 & 101 \\ w & 1-w \end{array} \right\},$$

where the first line denotes the support points and the second one their respective weights. Then, for

$$w = 1 - \frac{1}{404} \frac{201^{q-1} - 202q + 405}{201^{q-1} - 101q + 102}$$

one has  $\partial^2 \Delta_q[(1-\alpha)\mu_1 + \alpha\mu_2]/\partial\alpha^2|_{\alpha=0} \geq 0$  for all  $q \geq 1.84$ , the equality being obtained at  $q = 2$  only. Counterexamples are easily constructed for values of  $q$  smaller than 1.84.

**Proof of Theorem 4.2**

(i) The fact that  $\max_{\mu \in \mathcal{M}} \Psi_k^{1/k}(V_\mu) \geq \min_{M,c: \mathcal{X} \subset \mathcal{E}(M,c)} 1/\phi_k^\infty(M)$  is a consequence of Theorem 4.4. Indeed, the measure  $\mu_k^*$  maximises  $\Psi_k^{1/k}(V_\mu)$  if and only if

$$(x - E_{\mu_k^*})^\top M_*(V_{\mu_k^*})(x - E_{\mu_k^*}) \leq 1 \text{ for all } x \text{ in } \mathcal{X}. \quad (4.11)$$

Denote  $M_k^* = M_*(V_{\mu_k^*})$ ,  $c_k^* = E_{\mu_k^*}$ , and consider the Lagrangian  $L(V, \alpha; M)$  for the maximisation of  $(1/k) \log \Psi_k(V)$  with respect to  $V \succeq 0$  under the constraint  $\text{trace}(MV) = 1$ :  $L(V, \alpha; M) = (1/k) \log \Psi_k(V) - \alpha[\text{trace}(MV) - 1]$ . We have

$$\left. \frac{\partial L(V, 1; M_k^*)}{\partial V} \right|_{V=V_{\mu_k^*}} = M_k^* - M_k^* = 0$$

and  $\text{trace}(M_k^* V_{\mu_k^*}) = 1$ , with  $V_{\mu_k^*} \succeq 0$ . Therefore,  $V_{\mu_k^*}$  maximises  $\Psi_k(V)$  under the constraint  $\text{trace}(M_k^* V) = 1$ , and, moreover,  $\mathcal{X} \subset \mathcal{E}(M_k^*, c_k^*)$  from (4.11). This implies

$$\begin{aligned} \Psi_k^{1/k}(V_{\mu_k^*}) &= \max_{V \succeq 0: \text{trace}(M_k^* V)=1} \Psi_k^{1/k}(V) \\ &\geq \min_{M,c: \mathcal{X} \subset \mathcal{E}(M,c)} \max_{V \succeq 0: \text{trace}(MV)=1} \Psi_k^{1/k}(V) = \min_{M,c: \mathcal{X} \subset \mathcal{E}(M,c)} \frac{1}{\phi_k^\infty(M)}. \end{aligned}$$

(ii) We prove now that  $\min_{M,c: \mathcal{X} \subset \mathcal{E}(M,c)} 1/\phi_k^\infty(M) \geq \max_{\mu \in \mathcal{M}} \Psi_k^{1/k}(V_\mu)$ . Note that we do not have an explicit form for  $\phi_k^\infty(M)$  and that the infimum in (4.6) can be attained at a singular  $V$ , not necessarily unique, so that we cannot differentiate  $\phi_k^\infty(M)$ . Also note that compared to the developments in [14, Chap. 7], here we consider covariance matrices instead of moment matrices.

Consider the maximisation of  $\log \phi_k^\infty(M)$  with respect to  $M$  and  $c$  such that  $\mathcal{X} \subset \mathcal{E}(M, c)$ , with Lagrangian

$$L(M, c, \beta) = \log \phi_k^\infty(M) + \sum_{x \in \mathcal{X}} \beta_x [1 - (x - c)^\top M(x - c)], \quad \beta_x \geq 0 \text{ for all } x \text{ in } \mathcal{X}.$$

For the sake of simplicity we consider here  $\mathcal{X}$  to be finite, but  $\beta$  may denote any positive measure on  $\mathcal{X}$  otherwise. Denote the optimum by  $T^* = \max_{M,c: \mathcal{X} \subset \mathcal{E}(M,c)} \log \phi_k^\infty(M)$ . It satisfies  $T^* = \max_{M,c} \min_{\beta \geq 0} L(M, c, \beta) \leq \min_{\beta \geq 0} \max_{M,c} L(M, c, \beta)$ , and  $\max_{M,c} L(M, c, \beta)$  is attained for any  $c$  such that  $Mc = M \sum_{x \in \mathcal{X}} \beta_x x / (\sum_{x \in \mathcal{X}} \beta_x)$ , that is, in particular for

$$c^* = \frac{\sum_{x \in \mathcal{X}} \beta_x x}{\sum_{x \in \mathcal{X}} \beta_x},$$

and for  $M^*$  such that  $0 \in \partial_M L(M, c^*, \beta)|_{M=M^*}$ , the subdifferential of  $L(M, c^*, \beta)$  with respect to  $M$  at  $M^*$ . This condition can be written as

$$\sum_{x \in \mathcal{X}} \beta_x (x - c^*)(x - c^*)^\top = \tilde{V} \in \partial \log \phi_k^\infty(M)|_{M=M^*},$$

with  $\partial \log \phi_k^\infty(M)$  the subdifferential of  $\log \phi_k^\infty(M)$ ,

$$\partial \log \phi_k^\infty(M) = \{V \succeq 0 : \Psi_k^{1/k}(V) \phi_k^\infty(M) = \text{trace}(MV) = 1\},$$

see [14, Th. 7.9]. Since  $\text{trace}(MV) = 1$  for all  $V \in \partial \log \phi_k^\infty(M)$ ,  $\text{trace}(M^* \tilde{V}) = 1$  and thus  $\sum_{x \in \mathcal{X}} \beta_x (x - c^*)^\top M^* (x - c^*) = 1$ . Also,  $\Psi_k^{1/k}(\tilde{V}) = 1/\phi_k^\infty(M^*)$ , which gives

$$L(M^*, c^*, \beta) = -\log \Psi_k^{1/k} \left[ \sum_{x \in \mathcal{X}} \beta_x (x - c^*)(x - c^*)^\top \right] + \sum_{x \in \mathcal{X}} \beta_x - 1.$$

We obtain finally

$$\begin{aligned} & \min_{\beta \geq 0} L(M^*, c^*, \beta) \\ &= \min_{\gamma > 0, \alpha \geq 0} \left\{ -\log \Psi_k^{1/k} \left[ \sum_{x \in \mathcal{X}} \alpha_x (x - c^*)(x - c^*)^\top \right] + \gamma - \log(\gamma) - 1 \right\}, \\ &= \min_{\alpha \geq 0} -\log \Psi_k^{1/k} \left[ \sum_{x \in \mathcal{X}} \alpha_x (x - c^*)(x - c^*)^\top \right] = -\log \Psi_k^{1/k}(V_k^*), \end{aligned}$$

where we have denoted  $\gamma = \sum_{x \in \mathcal{X}} \beta_x$  and  $\alpha_x = \beta_x/\gamma$  for all  $x$ . Therefore  $T^* \leq -\log \Psi_k^{1/k}(V_k^*)$ , that is,  $\log [\min_{M, c: \mathcal{X} \subset \mathcal{E}(M, c)} 1/\phi_k^\infty(M)] \geq \log \Psi_k^{1/k}(V_k^*)$ .

## References

- [1] G. Björck. Distributions of positive mass, which maximize a certain generalized energy integral. *Arkiv för Matematik*, 3(21):255–269, 1956.
- [2] M.H. DeGroot and M.M. Rao. Bayes estimation with convex loss. *Ann. Math. Statist.*, 34(3):839–846, 1963.
- [3] F.R. Gantmacher. *Théorie des Matrices*. Dunod, Paris, 1966.
- [4] C. Gini. Measurement of inequality of incomes. *Economic J.*, 31(121):124–126, 1921.
- [5] A. Giovagnoli and H.P. Wynn. Multivariate dispersion orderings. *Stat. & Prob. Lett.*, 22(4):325–332, 1995.
- [6] M. Hainy, W.G. Müller, and H.P. Wynn. Learning functions and approximate bayesian computation design: ABCD. *Entropy*, 16(8):4353–4374, 2014.
- [7] R. Harman. Lower bounds on efficiency ratios based on  $\phi_p$ -optimal designs. In A. Di Bucchianico, H. Läuter, and H.P. Wynn, editors, *mODa’7 – Advances in Model-Oriented Design and Analysis, Proceedings of the 7th Int. Workshop, Heeze (Netherlands)*, pages 89–96, Heidelberg, 2004. Physica Verlag.
- [8] J. Kiefer and J. Wolfowitz. The equivalence of two extremum problems. *Canad. J. Math.*, 12:363–366, 1960.
- [9] J. López-Fidalgo and J.M. Rodríguez-Díaz. Characteristic polynomial criteria in optimal experimental design. In A.C. Atkinson, L. Pronzato, and H.P. Wynn, editors, *Advances in Model-Oriented Data Analysis and Experimental Design, Proceedings of MODA’5, Marseilles, June 22–26, 1998*, pages 31–38. Physica Verlag, Heidelberg, 1998.
- [10] I.G. Macdonald. *Symmetric functions and Hall polynomials*. Oxford University Press, Oxford, 1995. [2nd ed.].

- [11] M. Marcus and H. Minc. *A Survey of Matrix Theory and Matrix Inequalities*. Dover, New York, 1964.
- [12] H. Oja. Descriptive statistics for multivariate distributions. *Stat. & Prob. Lett.*, 1(6):327–332, 1983.
- [13] L. Pronzato. On a property of the expected value of a determinant. *Stat. & Prob. Lett.*, 39:161–165, 1998.
- [14] F. Pukelsheim. *Optimal Experimental Design*. Wiley, New York, 1993.
- [15] C.R. Rao. Diversity and dissimilarity coefficients: a unified approach. *Theoret. Popn Biol.*, 21(1):24–43, 1982.
- [16] C.R. Rao. Diversity: Its measurement, decomposition, apportionment and analysis. *Sankhyā: Indian J. Statist., Series A*, 44(1):1–22, 1982.
- [17] C.R. Rao. Convexity properties of entropy functions and analysis of diversity. In *Inequalities in Statistics and Probability*, volume 5, pages 68–77. Lecture Notes-Monograph Series, IMS, Hayward, CA, 1984.
- [18] R.L. Schilling, R. Song, and Z. Vondracek. *Bernstein Functions: Theory and Applications*. de Gruyter, Berlin/Boston, 2012.
- [19] I.J. Schoenberg. Metric spaces and positive definite functions. *Trans. Amer. Math. Soc.*, 44(3):522–536, 1938.
- [20] R.J. Serfling. *Approximation Theorems of Mathematical Statistics*. Wiley, New York, 1980.
- [21] M. Shaked. Dispersive ordering of distributions. *J. Appl. Prob.*, pages 310–320, 1982.
- [22] N.Z. Shor and O.A. Berezovski. New algorithms for constructing optimal circumscribed and inscribed ellipsoids. *Optim. Meth. Soft.*, 1:283–299, 1992.
- [23] D.M. Titterington. Optimal design: some geometrical aspects of  $D$ -optimality. *Biometrika*, 62(2):313–320, 1975.
- [24] H.R. van der Vaart. A note on Wilks’ internal scatter. *Ann. Math. Statist.*, 36(4):1308–1312, 1965.
- [25] S.S. Wilks. Certain generalizations in the analysis of variance. *Biometrika*, 24:471–494, 1932.
- [26] S.S. Wilks. Multidimensional statistical scatter. In I. Olkin, S.G. Ghurye, W. Hoeffding, W.G. Madow, and H.B. Mann, editors, *Contributions to Probability and Statistics: Essays in Honor of Harold Hotelling*, pages 486–503. Stanford University Press, Stanford, 1960.