

# Nonparametric Estimation via Convex Programming

Anatoli Juditsky\*

Arkadi Nemirovski†

## Abstract

In the paper, we focus primarily on the problem of recovering a linear form  $g^T x$  of unknown “signal”  $x$  known to belong to a given convex compact set  $X \subset \mathbb{R}^n$  from  $N$  independent realizations of a random variable  $\iota$  taking values in a finite set, the distribution  $p$  of  $\iota$  being affinely parameterized by  $x$ :  $p = Ax + b$ . With no additional assumptions on  $X$  and  $A$ , we develop minimax optimal, within an absolute constant factor, and computationally efficient estimation routine. We then apply this routine to recovering  $x$  itself in the Euclidean norm.

## 1 Introduction

In the sequel, we mainly focus on the estimation problem as follows:

**Problem I:** We observe  $N$  independent realizations  $i^1, \dots, i^N$  of a random variable  $\iota$  taking values in a finite set, say, the set  $I = \{1, \dots, M\}$ . The distribution of  $\iota$  (which is identified with a vector  $p$  from the standard simplex  $\mathcal{P}_M = \{y \in \mathbb{R}^M : y \geq 0, \sum_i y_i = 1\}$  by setting  $p_i = \text{Prob}\{\iota = i\}$ ,  $1 \leq i \leq M$ ) is affinely parameterized by  $n$ -dimensional “signal” – vector of unknown parameters  $x$  known to belong to a given convex compact set  $X \subset \mathbb{R}^n$ :  $p = A(x) = [A_1(x); \dots; A_M(x)]$ , where  $A(\cdot)$  is a given affine mapping with  $A(X) \subset \mathcal{P}_M$ .

Our goal is to infer from the observations certain information on  $x$ , primarily, to estimate a given linear form  $g^T z$  of  $z \in \mathbb{R}^n$  at the point  $x$  underlying our observations.

While the unknown  $x$  is assumed to be finite-dimensional, we allow the dimension to be arbitrarily large, thus addressing, essentially, a nonparametric estimation problem. In Nonparametric Statistics, there exists an immense literature on various versions of Problem I, including numerous papers on estimating linear forms, see, e.g., [2–4, 6–8, 10–13, 15–20, 22, 24–26, 28, 30, 32–39] and references therein. To the best of our knowledge, the majority of papers on the subject focus on “concrete” domains  $X$  (e.g., distributions on fine grids obtained by discretization of densities from Sobolev balls), and investigate lower and/or upper bounds on the worst-case, w.r.t.  $x \in X$ , quality to which the problem of interest can be solved. These bounds depend on the number of observations  $N$ , and the question of primary interest is the behaviour of the bounds as  $N \rightarrow \infty$ . When the lower and the upper bounds coincide within a constant factor (or, ideally, within factor  $(1 + o(1))$ ) as  $N \rightarrow \infty$ , the estimation problem is considered as being essentially solved, and the estimation methods underlying the upper bounds are treated as optimal.

---

\*LJK, Université J. Fourier, Grenoble, France, [Anatoli.Juditsky@imag.fr](mailto:Anatoli.Juditsky@imag.fr)

†School of ISyE, Georgia Institute of Technology, Atlanta, USA, [nemirovs@isye.gatech.edu](mailto:nemirovs@isye.gatech.edu)

Research of the second author was partly supported by the NSF grant # 0619977

The approach we take in this paper is of a different spirit. Except for the concluding Section 4, we make no “structural assumptions” on  $X$ , aside of crucial for us assumptions of convexity and compactness, and we make no assumptions on the affine mapping  $A(x)$ . Clearly, with no structural assumptions on  $X$  and  $A(\cdot)$ , explicit bounds on the risks of our estimates, same as bounds on the minimax optimal risk, are impossible. What is possible – and this is our major goal in what follows – is to demonstrate that *when estimating linear forms, the worst-case risk of the estimate we develop is within an absolute constant factor of the “ideal”* (i.e., the minimax optimal) risk. It should be added that while the optimal, within an absolute constant factor, worst-case risk of our estimates is not available in a closed analytical form, it is “available algorithmically” – it can be efficiently computed, provided that  $X$  is computationally tractable<sup>1</sup>.

While we are not aware of general results of the outlined spirit for Problem I, results of this type do exist for the “regression counterpart” of Problem I, namely, for

**Problem II:** Given indirect noisy observations

$$y = Ax + \sigma\xi \tag{1}$$

of unknown signal  $x$  known to belong to a given convex compact set  $X \subset \mathbb{R}^n$  ( $A$  is a given  $m \times n$  matrix,  $\xi \sim \mathcal{N}(0, I_m)$ ,  $\sigma > 0$  is given), we want to estimate the value  $g^T z$  of a given linear form  $g^T z$  of  $z \in \mathbb{R}^n$  at the point  $x$  underlying our observations.

As shown by D. Donoho [8], for all commonly used loss functions, the minimax optimal *affine in  $y$*  estimate in Problem II (this estimate can be easily built, provided that  $X$  is computationally tractable) is minimax optimal, within an absolute constant factor, among *all possible* estimates. In a sense, our results establish similar fact for estimating linear forms of the signal in the context of Problem I, since our estimates also are “affine” – they are affine functions of the empirical distribution of the discrete random variable  $\iota$  induced by our observations.

The rest of this paper is organized as follows. In Section 2 we consider the Hypotheses Testing version of Problem I, where one, given two convex subsets  $X_1, X_2$  in  $X$ , is interested to test the hypothesis  $x \in X_1$  vs. the alternative  $x \in X_2$ . The central Section 3 focuses on the version of Problem I where the goal is to estimate via the observations a given linear form of  $x$ . In the concluding Section 4, we discuss briefly how our results from Section 3 related to Problem I and the aforementioned results of Donoho [8] related to Problem II can be used in order to recover the “entire” signal  $x$  underlying our observations, the model of observations being either Problem I, or Problem II; as a loss function, we use the standard Euclidean norm on  $\mathbb{R}^n$ . When passing from recovering linear forms of the unknown signal to recovering the signal itself, we do impose structural assumptions on  $X$ , but still make no structural assumptions on the affine mapping  $A(x)$  (Problem I) and matrix  $A$  (Problem II), and our “optimality results” become weaker – instead of “optimality within an absolute constant factor” we end up with statements like “the worst-case risk of such-and-such estimate is in-between the minimax optimal risk and the latter risk to the power  $\chi$ ”, with  $\chi$  depending on the geometry of  $X$  (and close to 1 when this geometry is “good enough”). The appendix contains an alternative proof (in our opinion, much simpler than the original one) of the aforementioned Donoho’s theorem on minimax “almost-optimality” of affine estimates in the context of Problem II.

---

<sup>1</sup>For details on computational tractability and complexity issues in Convex Optimization, see, e.g., [1, Chapter 4]. A reader not familiar with this area will not lose much when interpreting a computationally tractable convex set as a set given by a finite system of inequalities  $p_i(x) \leq 0$ ,  $i = 1, \dots, m$ , where  $p_i(x)$  are convex polynomials.

## 2 Problem I: Hypotheses testing

In this Section, we focus on the case of Problem I as follows:

Hypotheses Testing (HT): In the situation of Problem I, given two closed convex subsets  $X_i$ ,  $i = 1, 2$ , of  $X$ , test the hypothesis  $x \in X_1$  vs. the alternative  $x \in X_2$ .

### 2.1 The test

Let  $Y_1, Y_2$  be two closed convex subsets in  $\mathcal{P}_M$  (recall that we identify vectors from  $\mathcal{P}_M$  with probability distributions on the  $M$ -element index set  $I = \{1, \dots, M\}$ ). Assume that we are given  $N$  independent realizations  $i^N = [i_1; \dots; i_N]$  of a random variable  $\iota$  distributed according to  $y \in \mathcal{P}_M$  and want to distinguish between two hypotheses,  $\Pi_1$  and  $\Pi_2$ , stating, respectively, that  $y \in Y_1$  and that  $y \in Y_2$ . A candidate decision rule in this problem is a function  $\psi(i^N)$  taking values 1 and 2; for such a decision rule, its error probabilities  $\epsilon_\kappa(\psi)$ ,  $\kappa = 1, 2$ , are defined as

$$\epsilon_\kappa(\psi) = \sup_{y \in Y_\kappa} \text{Prob}_{i^N \sim y \times \dots \times y} \{ \psi(i^N) \neq \kappa \}.$$

The test we intend to use is as follows.

Test  $\psi_{\phi, c}$ : We choose a “weight vector”  $\phi \in \mathbb{R}^M$  and a threshold  $c \in \mathbb{R}$  and accept hypothesis  $\Pi_1$  when  $\sum_{t=1}^N \phi_{i_t} \geq c$ , otherwise we accept  $\Pi_2$ .

We start with a construction of “test parameters”  $(\phi, c)$  based on “Bernstein approximation” [31]. Let us fix  $\nu \in Y_2$ . With  $i^N \sim \nu^N = \nu \times \dots \times \nu$ , the probability for a  $(\phi, c)$ -test to accept hypothesis  $\Pi_1$  is  $\text{Prob}_{i^N \sim \nu^N} \{ \sum_t \phi_{i_t} \geq c \}$ . For every  $\beta > 0$ , this probability does not exceed the quantity

$$\mathbf{E}_{i^N \sim \nu^N} \left\{ \exp \left\{ \sum_t \beta^{-1} \phi_{i_t} \right\} \right\} \exp \{ -\beta^{-1} c \} = \exp \{ -\beta^{-1} c \} \left( \sum_{i \in I} \nu_i \exp \{ \phi_i / \beta \} \right)^N.$$

We conclude that if  $\omega \in (0, 1)$ , then the condition

$$\exists \beta > 0 : N \ln \left( \sum_{i \in I} \nu_i \exp \{ \phi_i / \beta \} \right) - \beta^{-1} c \leq \ln(\omega)$$

is a sufficient condition for the  $\nu^N$ -probability to accept  $\Pi_1$  with test  $\psi_{\phi, c}$  to be  $\leq \omega$ . We rewrite this sufficient condition equivalently as

$$\exists \beta > 0 : \underbrace{\beta N \ln \left( \sum_{i \in I} \nu_i \exp \{ \phi_i / \beta \} \right) - c + \beta \ln(1/\omega)}_{\Psi(\phi, \beta, c; \nu)} \leq 0, \quad (2)$$

the benefit being the fact that *the function  $\Psi(\phi, \beta, c; \nu)$  is convex in  $(\phi, \beta, c)$  in the domain  $\beta > 0$  and is concave in  $\nu \in \mathcal{P}_M$* . Indeed, the concavity in  $\nu$  is evident; to verify the convexity, note that the function  $H(\phi, c; \nu) = N \ln \left( \sum_{i \in I} \nu_i \exp \{ \phi_i \} \right) - c + \ln(1/\omega)$  clearly is convex in  $(\phi, c)$ , and the

“projective transformation”  $F(u) \mapsto \beta F(\beta^{-1}u)$  is known to convert a convex function of  $u$  into a convex in the domain  $\beta > 0$  function of  $(u, \beta)$ .

By similar argument, the condition

$$\exists \alpha > 0 : \alpha N \ln \left( \sum_{i \in I} \mu_i \exp\{-\phi_i/\alpha\} \right) + c + \alpha \ln(1/\omega) \leq 0 \quad (3)$$

guarantees that the  $\mu^N$ -probability to accept  $\Pi_2$  with the test  $\psi_{\phi,c}$  is  $\leq \omega$ . We have arrived at the following

**Proposition 2.1** *Assume that  $\phi \in \mathbb{R}^M$  and  $\alpha, \beta$  are such that*

$$\alpha N \max_{\mu \in Y_1} \ln \left( \sum_{i \in I} \mu_i \exp\{-\phi_i/\alpha\} \right) + \beta N \max_{\nu \in Y_2} \ln \left( \sum_{i \in I} \nu_i \exp\{\phi_i/\beta\} \right) + (\alpha + \beta) \ln(1/\omega) \leq 0, \quad (4)$$

$$\alpha > 0, \beta > 0$$

*Setting*

$$c = \frac{1}{2} \left[ \beta N \max_{\nu \in Y_2} \ln \left( \sum_{i \in I} \nu_i \exp\{\phi_i/\beta\} \right) - \alpha N \max_{\mu \in Y_1} \ln \left( \sum_{i \in I} \mu_i \exp\{-\phi_i/\alpha\} \right) + (\beta - \alpha) \ln(1/\omega) \right], \quad (5)$$

*we ensure that*

$$\epsilon_1(\psi_{\phi,c}) \leq \omega \text{ and } \epsilon_2(\psi_{\phi,c}) \leq \omega.$$

**Proof.** We have

$$\begin{aligned} & \alpha N \max_{\mu \in Y_1} \ln \left( \sum_{i \in I} \mu_i \exp\{-\phi_i/\alpha\} \right) + c + \alpha \ln(1/\omega) \\ &= \frac{1}{2} \left[ \alpha N \max_{\mu \in Y_1} \ln \left( \sum_{i \in I} \mu_i \exp\{-\phi_i/\alpha\} \right) + \beta N \max_{\nu \in Y_2} \ln \left( \sum_{i \in I} \nu_i \exp\{\phi_i/\beta\} \right) + (\alpha + \beta) \ln(1/\omega) \right] \leq 0 \end{aligned}$$

and similarly

$$\begin{aligned} & \beta N \max_{\nu \in Y_2} \ln \left( \sum_{i \in I} \nu_i \exp\{\phi_i/\beta\} \right) - c + \beta \ln(1/\omega) \\ &= \frac{1}{2} \left[ \alpha N \max_{\mu \in Y_1} \ln \left( \sum_{i \in I} \mu_i \exp\{-\phi_i/\alpha\} \right) + \beta N \max_{\nu \in Y_2} \ln \left( \sum_{i \in I} \nu_i \exp\{\phi_i/\beta\} \right) + (\alpha + \beta) \ln(1/\omega) \right] \leq 0. \end{aligned}$$

We see that for every  $\nu \in Y_2$ ,  $\beta, \phi, c$  satisfy (2), so that the  $\nu^N$ -probability to accept  $\Pi_1$  with the test  $\psi_{\phi,c}$  is  $\leq \omega$ , and for every  $\mu \in Y_1$ ,  $\alpha, \phi, c$  satisfy (3), so that the  $\mu^N$ -probability to accept  $\Pi_2$  with the same test also is  $\leq \omega$ . ■

**Remark 2.1** *Note that  $(\phi, \alpha, \beta)$  is a solution to (4) if and only if  $(\phi/(\alpha + \beta), \alpha/(\alpha + \beta), \beta/(\alpha + \beta))$  solves the “normalized” system of constraints*

$$\alpha N \max_{\mu \in Y_1} \ln \left( \sum_{i \in I} \mu_i \exp\{-\phi_i/\alpha\} \right) + \beta N \max_{\nu \in Y_2} \ln \left( \sum_{i \in I} \nu_i \exp\{\phi_i/\beta\} \right) + \ln(1/\omega) \leq 0, \quad (6)$$

$$\alpha > 0, \beta > 0, \alpha + \beta = 1,$$

*and of course every solution to the latter system solves (4) as well. Thus, we lose nothing when working with (6) instead of (4).*

The advantage of the conditions (4), (6) is that they are represented by convex inequalities in variables  $\phi, \alpha, \beta$  with efficiently computable (provided that  $Y_i, i = 1, 2$ , are computationally tractable) left hand sides and as such are “computationally tractable” – we can find efficiently a solution to (4) or (6), provided that one exists.

## 2.2 Efficiency

We are about to demonstrate that our sufficient condition, even in somehow strengthened form, is “nearly necessary” for the existence of a good test for distinguishing between  $\Pi_1$  and  $\Pi_2$ . What follows is a variation of well-known relations between Hellinger affinity and hypotheses testing (cf., e.g., [6, 7]).

**Proposition 2.2** *Let  $\omega \in (0, 1)$ .*

(i) *The relation*

$$\sup_{\mu \in Y_1, \nu \in Y_2} \left[ N \ln \left( \sum_{i \in I} \sqrt{\mu_i \nu_i} \right) \right] + \ln(1/\omega) < 0 \quad (7)$$

*is sufficient for the feasibility of (6). When all vectors  $y \in Y_1 \cup Y_2$  are strictly positive, so that*

$$\kappa \equiv \min_{\mu \in Y_1, \nu \in Y_2, i} \min[\mu_i, \nu_i] > 0,$$

*the same conclusion is true for the non-strict version*

$$\max_{\mu \in Y_1, \nu \in Y_2} \left[ N \ln \left( \sum_{i \in I} \sqrt{\mu_i \nu_i} \right) \right] + \ln(1/\omega) \leq 0 \quad (8)$$

*of (7), and in this case (6) admits a feasible solution  $\phi, \alpha, \beta$  with  $\|\phi\|_\infty \leq S = \frac{1}{4} \ln(1/\kappa)$  and  $\alpha = \beta = 1/2$ .*

(ii) *Assume that (7) is not satisfied. Then one can find  $\bar{\mu} \in Y_1$  and  $\bar{\nu} \in Y_2$  such that for every test  $\psi(i^N)$  one has either*

$$\text{Prob}_{i^N \sim \bar{\mu}^N} \{\psi(i^N) = 2\} \geq \omega_+ \equiv \omega^2/4,$$

*or*

$$\text{Prob}_{i^N \sim \bar{\nu}^N} \{\psi(i^N) = 1\} \geq \omega_+.$$

**Proof.** Consider the function

$$F(\phi; \mu, \nu) = N \left[ \ln \left( \sum_{i \in I} \mu_i \exp\{-\phi_i\} \right) + \ln \left( \sum_{i \in I} \nu_i \exp\{\phi_i\} \right) \right].$$

This function clearly is continuous, concave in  $(\mu, \nu) \in Y_1 \times Y_2$  and convex in  $\phi \in \mathbb{R}^M$ . Since  $Y_1 \times Y_2$  is a convex compact set, we have

$$\inf_{\phi} \max_{\mu \in Y_1, \nu \in Y_2} F(\phi; \mu, \nu) = \max_{\mu \in Y_1, \nu \in Y_2} \inf_{\phi} F(\phi; \mu, \nu). \quad (9)$$

It is immediately seen that for  $(\mu, \nu) \in \mathcal{P}_M \times \mathcal{P}_M$  we have

$$\inf_{\phi} F(\phi; \mu, \nu) = 2N \ln \left( \sum_i \sqrt{\mu_i \nu_i} \right); \quad (10)$$

when both  $\mu, \nu$  are positive, a minimizer of the left hand side is given by

$$\phi_i = \frac{1}{2} \ln(\mu_i/\nu_i), \quad i \in I. \quad (11)$$

From (9), (10) it follows that in the case of (7) there exists  $\phi$  such that

$$\max_{\mu \in Y_1, \nu \in Y_2} N \left[ \ln \left( \sum_{i \in I} \mu_i \exp\{-\phi_i\} \right) + \ln \left( \sum_{i \in I} \nu_i \exp\{\phi_i\} \right) \right] + 2 \ln(1/\omega) \leq 0,$$

meaning that  $(\phi/2, \alpha = \beta = 1/2)$  satisfy (4). In the case when  $\mu_i \geq \kappa > 0$  and  $\nu_i \geq \kappa > 0$  for all  $\mu \in Y_1, \nu \in Y_2$ , the left hand side in (10) admits a minimizer satisfying (11), so that

$$\sup_{\mu \in Y_1, \nu \in Y_2} \inf_{\phi} F(\phi; \mu, \nu) = \max_{\mu \in Y_1, \nu \in Y_2} \min_{\|\phi\|_{\infty} \leq S} F(\phi; \mu, \nu).$$

By the standard Saddle Point Theorem, we can interchange max and min in the right hand side, thus arriving at

$$\begin{aligned} \min_{\|\phi\|_{\infty} \leq S} \max_{\mu \in Y_1, \nu \in Y_2} F(\phi; \mu, \nu) &= \max_{\mu \in Y_1, \nu \in Y_2} \min_{\|\phi\|_{\infty} \leq S} F(\phi; \mu, \nu) \\ &= \max_{\mu \in Y_1, \nu \in Y_2} \inf_{\phi} F(\phi; \mu, \nu) = \max_{\mu \in Y_1, \nu \in Y_2} N \ln \left( \sum_i \sqrt{\mu_i \nu_i} \right). \end{aligned}$$

From this relation, same as above, it follows that if  $\phi$  is a minimizer in  $\min_{\|\phi\|_{\infty} \leq S} \max_{\mu \in Y_1, \nu \in Y_2} F(\phi; \mu, \nu)$  and (8) takes place, then  $(\phi/2, \alpha = \beta = 1/2)$  satisfies (4). (i) is proved.

(ii): The function under sup in the left hand side of (7) is concave and upper semi-continuous on  $Y_1 \times Y_2$ ; thus, it either equals to  $-\infty$  on this set (this is the case when  $\mu_i \nu_i = 0$  for all  $\mu \in Y_1, \nu \in Y_2$  and all  $i$ ), or the sup is achieved. In the first case, (7) clearly is satisfied, which is not so under the premise of (ii); thus, under this premise there exist  $\bar{\mu} \in Y_1$  and  $\bar{\nu} \in Y_2$  such that

$$N \ln \left( \sum_i \sqrt{\bar{\mu}_i \bar{\nu}_i} \right) + \ln(1/\omega) \geq 0. \quad (12)$$

Now assume, on the contrary to what should be proved, that there exists a test  $\psi$  such that

$$\text{Prob}_{i^N \sim \bar{\nu}^N} \{\psi(i^N) = 1\} < \omega_+, \quad \text{Prob}_{i^N \sim \bar{\mu}^N} \{\psi(i^N) = 2\} < \omega_+. \quad (13)$$

Setting  $A = \{i^N : \psi(i^N) = 1\}$  and  $B = \{i^N : \psi(i^N) = 2\}$ , we have

$$\sum_{i^N \in A} \bar{\nu}_{i^N}^N < \omega_+, \quad \sum_{i^N \in B} \bar{\mu}_{i^N}^N < \omega_+,$$

whence

$$\begin{aligned} \sum_{i^N} \sqrt{\bar{\mu}_{i^N}^N \bar{\nu}_{i^N}^N} &= \sum_{i^N \in A} \sqrt{\bar{\mu}_{i^N}^N \bar{\nu}_{i^N}^N} + \sum_{i^N \in B} \sqrt{\bar{\mu}_{i^N}^N \bar{\nu}_{i^N}^N} \\ &\leq \left( \sum_{i^N \in A} \bar{\mu}_{i^N}^N \right)^{1/2} \left( \sum_{i^N \in A} \bar{\nu}_{i^N}^N \right)^{1/2} + \left( \sum_{i^N \in B} \bar{\mu}_{i^N}^N \right)^{1/2} \left( \sum_{i^N \in B} \bar{\nu}_{i^N}^N \right)^{1/2} < 2\sqrt{\omega_+}. \end{aligned}$$

On the other hand, we have

$$\sum_{i^N} \sqrt{\bar{\mu}_{i^N}^N \bar{\nu}_{i^N}^N} = \sum_{i_1, \dots, i_N \in I} \sqrt{\bar{\mu}_{i_1} \dots \bar{\mu}_{i_N} \bar{\nu}_{i_1} \dots \bar{\nu}_{i_N}} = \left( \sum_{i \in I} \sqrt{\bar{\mu}_i \bar{\nu}_i} \right)^N,$$

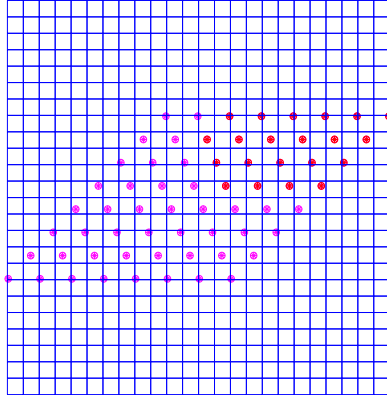


Figure 1: The support of  $\eta$  (magenta and red):  $8 \times 8$  point grid (grid steps 2) and the partitioning of the range of  $\eta + \xi$  (blue) into  $24 \times 24$  squares (grid steps 1).

and we arrive at the inequality

$$\left( \sum_{i \in I} \sqrt{\mu_i \bar{v}_i} \right)^N < 2\sqrt{\omega_+} = \omega,$$

which is impossible by (12). ■

### 2.3 Numerical illustration

Consider the situation where we observe the value of the sum  $\eta + \xi$  of two independent 2D random vectors, with  $\xi \sim \mathcal{N}(0, I_2)$  and  $\eta$  being a finite-valued random variable taking values from the 64-point grid  $G$  shown on Fig. 1. We denote by  $x$  the (unknown in advance) distribution of  $\eta$ ;  $X$  is the space of all probability distributions on  $G$  (that is,  $X$  is the standard 64-dimensional simplex). We denote by  $G_+$  the part of  $G$  belonging to the non-negative quadrant (red points on Fig. 1). Our goal is to distinguish between two hypotheses on  $x$ :  $x \in X_1 = \{x \in X : \sum_{g \in G_+} x_g \leq 0.4\}$  and  $x \in X_2 = \{x \in X : \sum_{g \in G_+} x_g \geq 0.6\}$ .

We reduce the situation to the one described in Problem I by partitioning the range  $\mathbb{R}^2$  of  $\eta + \xi$  into 577 parts: the  $576 = 24 \times 24$  squares shown on Fig. 1 and the complement to the union of these squares;  $\iota = \iota(\eta + \xi)$  is the index of the part to which  $\eta + \xi$  belongs. With  $N = 1$ , we can easily find the smallest  $\omega = \omega_{\min}$  for which the inequality (6) associated with the problem at hand has a solution  $(\alpha_*, \beta_*, \phi_*)$ ; computations result in  $\omega_{\min} = 1 - 0.004657$ . The solution  $(\alpha_*, \beta_*, \phi_*)$  implies the optimal, with our approach, decision rule for distinguishing between our two hypotheses  $x \in X_i$ ,  $i = 1, 2$ , for every number  $N$  of observations; the upper bound on error probabilities with  $N$  observations is  $\omega_{\min}^N \approx \exp\{-0.0047N\}$ . E.g., with  $N = 1000$ , the error probabilities are guaranteed to be  $\leq 0.0094$ . In an experiment with 20,000 sample of randomly generated probability distributions on  $G$ , with 10,000 distributions belonging to  $X_1$  and the remaining 10,000 belonging to  $X_2$ , the empirical probabilities to reject the hypothesis  $x \in X_i$  when it is true was 0.001 for  $i = 1$  and 0 for  $i = 2$ .

### 3 Problem I: Estimating linear form

Here we focus on the problem

Estimating Linear Form (ELF): In the situation of Problem I, given a linear form  $g^T z$  on  $\mathbb{R}^n$ , estimate  $g^T x$ .

#### 3.1 Simple estimates

An estimate – candidate solution to our problem – is a function  $\widehat{g}(i^N) : I \times \dots \times I \rightarrow \mathbb{R}$ . Given tolerance  $\epsilon \in (0, 1)$ , we define the  $\epsilon$ -risk of such an estimate on  $X$  as

$$\text{Risk}(\widehat{g}; \epsilon) = \inf \left[ \delta : \sup_{x \in X} \text{Prob}_{i^N \sim A^N(x)=A(x) \times \dots \times A(x)} \{ |\widehat{g}(i^N) - g^T x| > \delta \} < \epsilon \right],$$

and the minimax optimal  $\epsilon$ -risk as

$$\text{Risk}_*(\epsilon) = \inf_{\widehat{g}(\cdot)} \text{Risk}(\widehat{g}; \epsilon).$$

We call an estimate  $\widehat{g}(i^N)$  *simple*, if it is of the form

$$\widehat{g}(i^N) = \sum_{t=1}^N h_{i_t} + c,$$

where  $h \in \mathbb{R}^M$  and  $c \in \mathbb{R}$ . In other words, a simple estimate is an estimate which is affine in the empirical distribution

$$\pi(i^N) = \frac{1}{N} [\#\{t : i_t = 1\}; \dots; \#\{t : i_t = M\}]$$

of  $\iota$  associated with observations  $i^N$ . We denote by  $\text{RiskS}(\epsilon)$  the best  $\epsilon$ -risk achievable with simple estimates.

#### 3.2 Main result

Our main result is as follows:

**Theorem 3.1** *Let  $\epsilon \in (0, 1/4)$ . Then, for every  $\delta > 0$ , we can point out a simple estimate  $\widehat{g}_{\epsilon, \delta}(\cdot)$  satisfying the relation*

$$\text{Risk}(\widehat{g}_{\epsilon, \delta}; \epsilon) \leq \vartheta(\epsilon) \text{Risk}_*(\epsilon) + \delta, \quad \vartheta(\epsilon) = \frac{2 \ln \left( \frac{2}{\epsilon} \right)}{\ln \left( \frac{1}{4\epsilon} \right)} \quad (14)$$

(note that  $\vartheta(\epsilon) \rightarrow 2$  as  $\epsilon \rightarrow +0$ ); in particular,

$$\text{RiskS}(\epsilon) \leq \vartheta(\epsilon) \text{Risk}_*(\epsilon).$$

In addition, the estimate  $\widehat{g}_{\epsilon, \delta}$  is readily given by a solution to an explicit convex program and as such can be found in a computationally efficient fashion, provided that  $X$  is computationally tractable.



**Proof.** Consider the function

$$\begin{aligned}\Phi(\alpha, \beta, \phi; x, y) &= \alpha N \ln \left( \sum_i A_i(x) \exp\{\alpha^{-1} \phi_i\} \right) + \beta N \ln \left( \sum_i A_i(y) \exp\{-\beta^{-1} \phi_i\} \right) \\ &\quad + g^T y - g^T x : U \times V \rightarrow \mathbb{R}, \\ U &= \{(\alpha, \beta, \phi) : \alpha > 0, \beta > 0, \phi \in \mathbb{R}^M\}, \\ V &= \{(x, y) : x, y \in X\}.\end{aligned}\tag{15}$$

For  $\omega \geq 0$ , let

$$2S(\omega) = \inf_{(\alpha, \beta, \phi) \in U} \max_{(x, y) \in V} [\Phi(\alpha, \beta, \phi; x, y) + (\alpha + \beta)\omega].\tag{16}$$

We are about to demonstrate that

- A. Whenever  $0 < \epsilon < 1/4$  and  $\delta > 0$ , we can point out a simple estimate  $\widehat{g}_{\epsilon, \delta}(\cdot)$  such that  $\text{Risk}(\widehat{g}_{\epsilon, \delta}; \epsilon) \leq S(\ln(2/\epsilon)) + \delta$ , and  $\widehat{g}$  is readily given by a solution to a convex program;
- B. One has  $S(\ln(2/\epsilon)) \leq \vartheta(\epsilon) \text{Risk}_*(\epsilon)$ .

Clearly A and B combine to yield the statement of Theorem.

**A** is given by the following

**Lemma 3.1** *Let  $\epsilon \in (0, 1)$  and  $\delta > 0$ , so that, by definition of  $S(\cdot)$ , there exists  $(\bar{\alpha}, \bar{\beta}, \bar{\phi}) \in U$  such that*

$$\max_{(x, y) \in V} \Phi(\bar{\alpha}, \bar{\beta}, \bar{\phi}; x, y) + (\bar{\alpha} + \bar{\beta}) \ln(2/\epsilon) < 2[S(\ln(2/\epsilon)) + \delta].\tag{17}$$

*Setting*

$$\begin{aligned}c &= \frac{1}{2} \left[ \underbrace{\max_{y \in X} \left\{ \bar{\beta} N \ln \left( \sum_i A_i(y) \exp\{-\bar{\beta}^{-1} \bar{\phi}_i\} \right) + g^T y + \bar{\beta} \ln(2/\epsilon) \right\}}_J \right. \\ &\quad \left. - \underbrace{\max_{x \in X} \left\{ \bar{\alpha} N \ln \left( \sum_i A_i(x) \exp\{\bar{\alpha}^{-1} \bar{\phi}_i\} \right) - g^T x + \bar{\alpha} \ln(2/\epsilon) \right\}}_I \right],\end{aligned}\tag{18}$$

*the simple estimate*

$$\widehat{g}_{\epsilon, \delta}(i^N) = \sum_{t=1}^n \bar{\phi}_{i_t} + c$$

*satisfies*  $\text{Risk}(\widehat{g}_{\epsilon, \delta}; \epsilon) \leq S(\ln(2/\epsilon)) + \delta$ .

**Proof.** Setting  $D = S(\ln(2/\epsilon)) + \delta$ , observe that the left hand side in (17) is nothing but  $I + J$ , that is,  $I + J < 2D$ . We now have

$$\max_{x \in X} \left\{ \bar{\alpha} N \ln \left( \sum_i A_i(x) \exp\{\bar{\alpha}^{-1} \bar{\phi}_i\} \right) + c - g^T x + \bar{\alpha} \ln(2/\epsilon) \right\} = I + \frac{J - I}{2} = \frac{I + J}{2} < D\tag{19}$$

and

$$\max_{y \in X} \left\{ \bar{\beta} N \ln \left( \sum_i A_i(y) \exp\{\bar{\beta}^{-1} \bar{\phi}_i\} \right) - c + g^T y + \bar{\beta} \ln(2/\epsilon) \right\} = J - \frac{J - I}{2} = \frac{I + J}{2} < D.\tag{20}$$

Now, for  $x \in X$  we have

$$\begin{aligned}
& \text{Prob}_{i^N \sim A^N(x)} \left\{ \sum_t \bar{\phi}_{i_t} + c > g^T x + D \right\} \\
&= \text{Prob}_{i^N \sim A^N(x)} \left\{ \exp\left\{ \sum_t \bar{\alpha}^{-1} \bar{\phi}_{i_t} \right\} \exp\{\bar{\alpha}^{-1}(c - g^T x)\} \exp\{-\bar{\alpha}^{-1}D\} > 1 \right\} \\
&\leq \mathbf{E}_{i^N \sim A^N(x)} \left\{ \exp\left\{ \sum_t \bar{\alpha}^{-1} \bar{\phi}_{i_t} \right\} \exp\{\bar{\alpha}^{-1}(c - g^T x)\} \exp\{-\bar{\alpha}^{-1}D\} \right\} \\
&= \left( \sum_i A_i(x) \exp\{\bar{\alpha}^{-1} \bar{\phi}_i\} \right)^N \exp\{\bar{\alpha}^{-1}(c - g^T x)\} \exp\{-\bar{\alpha}^{-1}D\},
\end{aligned}$$

whence

$$\begin{aligned}
& \ln \left( \text{Prob}_{i^N \sim A^N(x)} \left\{ \sum_t \bar{\phi}_{i_t} + c > g^T x + D \right\} \right) \\
&\leq N \ln \left( \sum_i A_i(x) \exp\{\bar{\alpha}^{-1} \bar{\phi}_i\} \right) + \bar{\alpha}^{-1}(c - g^T x) - \bar{\alpha}^{-1} \frac{I+J}{2} - \bar{\alpha}^{-1} \left[ D - \frac{I+J}{2} \right] \\
&\leq \ln(\epsilon/2) - \bar{\alpha}^{-1} \left[ D - \frac{I+J}{2} \right],
\end{aligned} \tag{21}$$

where the concluding inequality is given by (19). Similarly,

$$\begin{aligned}
& \text{Prob}_{i^N \sim A^N(x)} \left\{ \sum_t \bar{\phi}_{i_t} + c < g^T x - D \right\} \\
&= \text{Prob}_{i^N \sim A^N(x)} \left\{ \exp\left\{ -\sum_t \bar{\beta}^{-1} \bar{\phi}_{i_t} \right\} \exp\{\bar{\beta}^{-1}(-c + g^T x)\} \exp\{-\bar{\beta}^{-1}D\} > 1 \right\} \\
&\leq \mathbf{E}_{i^N \sim A^N(x)} \left\{ \exp\left\{ -\sum_t \bar{\beta}^{-1} \bar{\phi}_{i_t} \right\} \exp\{\bar{\beta}^{-1}(-c + g^T x)\} \exp\{-\bar{\beta}^{-1}D\} \right\} \\
&= \left( \sum_i A_i(x) \exp\{-\bar{\beta}^{-1} \bar{\phi}_i\} \right)^N \exp\{\bar{\beta}^{-1}(-c + g^T x)\} \exp\{-\bar{\beta}^{-1}D\},
\end{aligned}$$

whence

$$\begin{aligned}
& \ln \left( \text{Prob}_{i^N \sim A^N(x)} \left\{ \sum_t \bar{\phi}_{i_t} + c < g^T x - D \right\} \right) \\
&\leq N \ln \left( \sum_i A_i(x) \exp\{-\bar{\beta}^{-1} \bar{\phi}_i\} \right) + \bar{\beta}^{-1}(-c + g^T x) - \bar{\beta}^{-1} \frac{I+J}{2} - \bar{\beta}^{-1} \left[ D - \frac{I+J}{2} \right] \\
&\leq \ln(\epsilon/2) - \bar{\beta}^{-1} \left[ D - \frac{I+J}{2} \right],
\end{aligned} \tag{22}$$

where the concluding inequality is given by (20). We see that

$$\sup_{x \in X} \text{Prob}_{i^N \sim A^N(x)} \left\{ \left| \sum_t \bar{\phi}_{i_t} + c - g^T x \right| > D \right\} \leq \epsilon \exp\left\{ -\min[\bar{\alpha}^{-1}, \bar{\beta}^{-1}] \left[ D - \frac{I+J}{2} \right] \right\} < \epsilon,$$

as claimed. ■.

**B** is given by the following fact:

**Lemma 3.2** *When  $\epsilon \in (0, 1/4)$ , one has*

$$\text{Risk}_*(\epsilon) \geq \frac{\ln\left(\frac{1}{4\epsilon}\right)}{2 \ln\left(\frac{2}{\epsilon}\right)} S(\ln(2/\epsilon)). \tag{23}$$

**Proof.**  $\mathbf{1}^0$ . Assume, on the contrary to what should be proved, that

$$\text{Risk}_*(\epsilon) + \kappa \leq \frac{\ln\left(\frac{1}{4\epsilon}\right)}{2 \ln\left(\frac{2}{\epsilon}\right)} S(\ln(2/\epsilon)) \tag{24}$$

with some  $\kappa > 0$ , and let us lead this assumption to a contradiction. Let us make a simple observation as follows:

**Lemma 3.3**  $S(\gamma)$  is a nonnegative concave function of  $\gamma \geq 0$ . In particular,  $S(\theta\gamma) \leq \theta S(\gamma)$  whenever  $\theta \geq 1$  and  $\gamma \geq 0$ .

**Proof.** Let  $\bar{x} \in X$ . Then

$$\Phi(\alpha, \beta, \phi; \bar{x}, \bar{x}) = \alpha N \ln \left( \underbrace{\sum_i A_i(\bar{x}) \exp\{\alpha^{-1}\phi_i\}}_{\geq \exp\{\alpha^{-1} \sum_i A_i(\bar{x})\phi_i\}} \right) + \beta N \ln \left( \underbrace{\sum_i A_i(\bar{x}) \exp\{-\beta^{-1}\phi_i\}}_{\geq \exp\{-\beta^{-1} \sum_i A_i(\bar{x})\phi_i\}} \right) \geq 0,$$

whence  $\max_{x,y \in X} \Phi(\alpha, \beta, \phi; x, y) \geq 0$  and therefore  $S(\gamma)$  is nonnegative when  $\gamma \geq 0$ . By construction,  $S(\gamma)$  is the infimum of a family of affine functions of  $\gamma$  and as such is concave. ■

Let us set

$$\gamma = \frac{1}{2} \ln \left( \frac{1}{4\epsilon} \right). \quad (25)$$

**2<sup>0</sup>.** Since  $S(\cdot)$  is concave and nonnegative on the nonnegative ray, we have

$$2S(\gamma) = \inf_{(\alpha, \beta, \phi) \in U} \max_{x, y \in X} \underbrace{[\Phi(\alpha, \beta, \phi; x, y) + (\alpha + \beta)\gamma]}_{\Psi(\alpha, \beta, \phi; x, y)} \geq (\gamma / \ln(2/\epsilon)) 2S(\ln(2/\epsilon)) \geq 2 \text{Risk}_*(\epsilon) + 2\kappa, \quad (26)$$

the concluding  $\geq$  being given by (24). The function  $\Psi$  is continuous on  $U \times V$ , concave in  $(x, y) \in V$  and convex in  $(\alpha, \beta, \phi) \in U$ . Since both  $U$  and  $V$  are convex sets and  $V$  is compact, we have

$$\inf_{\alpha, \beta, \phi} \max_{x, y \in X} \Psi(\alpha, \beta, \phi; x, y) = \max_{x, y \in X} \underbrace{\inf_{(\alpha, \beta, \phi) \in U} \Psi(\alpha, \beta, \phi; x, y)}_{\psi(x, y)}, \quad (27)$$

and, besides this,  $\psi(x, y)$  is upper semicontinuous on  $V$  and concave on  $V$  function. As such, it attains its maximum over  $x, y \in X$  at certain  $(\bar{x}, \bar{y})$ . Invoking (26), (27), we get

$$\forall (\alpha > 0, \beta > 0, \phi) : \Phi(\alpha, \beta, \phi; \bar{x}, \bar{y}) + (\alpha + \beta)\gamma \geq 2 \text{Risk}_*(\epsilon) + 2\kappa,$$

whence, recalling the definition of  $\Phi$ ,

$$\begin{aligned} & \forall (\alpha > 0, \beta > 0, \phi) : \\ & \alpha N \ln \left( \sum_i A_i(\bar{x}) \exp\{\alpha^{-1}\phi_i\} \right) + \beta N \ln \left( \sum_i A_i(\bar{y}) \exp\{-\beta^{-1}\phi_i\} \right) + (\alpha + \beta)\gamma \\ & \geq 2 \text{Risk}_*(\epsilon) + g^T[\bar{x} - \bar{y}] + 2\kappa. \end{aligned} \quad (28)$$

**3<sup>0</sup>.** We claim that

$$\begin{aligned} (a) \quad & g^T[\bar{y} - \bar{x}] \geq 2 \text{Risk}_*(\epsilon) + 2\kappa, \\ (b) \quad & \forall (\alpha > 0, \beta > 0, \phi) : \\ & \alpha N \ln \left( \sum_i A_i(\bar{x}) \exp\{\alpha^{-1}\phi_i\} \right) + \beta N \ln \left( \sum_i A_i(\bar{y}) \exp\{-\beta^{-1}\phi_i\} \right) + (\alpha + \beta)\gamma \geq 0. \end{aligned} \quad (29)$$

Indeed, the inequality in (28) holds true for  $\phi = 0$  and all positive  $\alpha, \beta$ , that is,  $(\alpha + \beta)\gamma \geq 2 \text{Risk}_*(\epsilon) + g^T[\bar{x} - \bar{y}] + 2\kappa$  for all  $\alpha, \beta \geq 0$ , which implies (29.a). To prove (29.b), assume for a moment that there exist  $\bar{\alpha} > 0$ ,  $\bar{\beta} > 0$  and  $\bar{\phi}$  such that

$$w \equiv \bar{\alpha} N \ln \left( \sum_i A_i(\bar{x}) \exp\{\bar{\alpha}^{-1}\bar{\phi}_i\} \right) + \bar{\beta} N \ln \left( \sum_i A_i(\bar{y}) \exp\{-\bar{\beta}^{-1}\bar{\phi}_i\} \right) + (\bar{\alpha} + \bar{\beta})\gamma < 0.$$

Setting  $\alpha_t = t\bar{\alpha}$ ,  $\beta_t = t\bar{\beta}$  and  $\phi^t = t\bar{\phi}$ , we see that along the sequence  $\{\alpha_t, \beta_t, \phi^t\}_{t=1}^\infty$  the left hand side in the inequality in (28) is  $tw \rightarrow -\infty$ ,  $t \rightarrow \infty$ , which is forbidden by (28). This contradiction completes the proof of (29.b).

4<sup>0</sup>. By the definition of  $\epsilon$ -risk, there exists  $\epsilon' < \epsilon$  and an estimate  $\widehat{g}(i^N)$  such that

$$\sup_{x \in X} \text{Prob}_{i^N \sim \mathcal{A}^N(x)} \{|\widehat{g}(i^N) - g^T x| \geq \text{Risk}_*(\epsilon) + \kappa/2\} \leq \epsilon'.$$

Setting

$$\psi(i^N) = \begin{cases} 1, & \widehat{g}(i^N) \leq \frac{1}{2}g^T[\bar{x} + \bar{y}] \\ -1, & \text{otherwise} \end{cases},$$

we get

$$\begin{aligned} \text{Prob}_{i^N \sim \mathcal{A}^N(\bar{x})} \{\psi(i^N) \neq 1\} &\leq \text{Prob}_{i^N \sim \mathcal{A}^N(\bar{x})} \{\widehat{g}(i^N) \geq \frac{1}{2}g^T[\bar{y} + \bar{x}]\} \\ &\leq \text{Prob}_{i^N \sim \mathcal{A}^N(\bar{x})} \{\widehat{g}(i^N) \geq g^T \bar{x} + \text{Risk}_*(\epsilon) + \kappa\}, \end{aligned}$$

since by (29.a) we have  $g^T \bar{x} \leq \frac{1}{2}g^T[\bar{y} + \bar{x}] - \text{Risk}_*(\epsilon) - \kappa$ . It follows that

$$\text{Prob}_{i^N \sim \mathcal{A}^N(\bar{x})} \{\psi(i^N) \neq 1\} \leq \epsilon'. \quad (30)$$

By similar argument, we have

$$\text{Prob}_{i^N \sim \mathcal{A}^N(\bar{y})} \{\psi(i^N) \neq -1\} \leq \epsilon'. \quad (31)$$

Setting  $\mu_i = A_i(\bar{x})$ ,  $\nu_i = A_i(\bar{y})$  and  $\mathcal{I} = \{i^N : \psi(i^N) = -1\}$ , we conclude that

$$\sum_{i^N \in \mathcal{I}} \prod_{t=1}^N \mu_{i_t} \leq \epsilon' \quad \text{and} \quad \sum_{i^N \notin \mathcal{I}} \prod_{t=1}^N \nu_{i_t} \leq \epsilon',$$

whence

$$\begin{aligned} (\sum_i \sqrt{\mu_i \nu_i})^N &= \sum_{i^N} \sqrt{\mu_{i^N}^N \nu_{i^N}^N} = \sum_{i^N \in \mathcal{I}} \sqrt{\mu_{i^N}^N \nu_{i^N}^N} + \sum_{i^N \notin \mathcal{I}} \sqrt{\mu_{i^N}^N \nu_{i^N}^N} \\ &\leq \sqrt{\sum_{i^N \in \mathcal{I}} \mu_{i^N}^N} \sqrt{\sum_{i^N \in \mathcal{I}} \nu_{i^N}^N} + \sqrt{\sum_{i^N \notin \mathcal{I}} \mu_{i^N}^N} \sqrt{\sum_{i^N \notin \mathcal{I}} \nu_{i^N}^N} \\ &\leq 2\sqrt{\epsilon'}, \end{aligned}$$

whence

$$N \ln \left( \sum_i \sqrt{\mu_i \nu_i} \right) \leq \frac{1}{2} \ln(4\epsilon'). \quad (32)$$

On the other hand, from (29.b) it follows that

$$\inf_{\phi} \left[ N \ln \left( \sum_i \mu_i \exp\{\phi_i\} \right) + N \ln \left( \sum_i \nu_i \exp\{-\phi_i\} \right) + 2\gamma \right] \geq 0,$$

whence, passing to limit along the sequence  $\phi^t$  of values of  $\phi$  given by

$$\phi_i^t \begin{cases} = \frac{1}{2} \ln(\nu_i/\mu_i), & \mu_i > 0 \quad \text{and} \quad \nu_i > 0 \\ \rightarrow +\infty, t \rightarrow \infty, & \mu_i = 0 \quad \text{and} \quad \nu_i > 0 \\ \rightarrow -\infty, t \rightarrow \infty, & \mu_i > 0 \quad \text{and} \quad \nu_i = 0 \\ = 0, & \mu_i = \nu_i = 0. \end{cases},$$

we get  $2N \ln \left( \sum_i \sqrt{\mu_i \nu_i} \right) \geq -2\gamma$ , that is,

$$N \ln \left( \sum_i \sqrt{\mu_i \nu_i} \right) \geq \frac{1}{2} \ln(4\epsilon),$$

which is impossible due to (32) combined with  $\epsilon' < \epsilon$ . We have arrived at a desired contradiction.  $\blacksquare$

### 3.3 Modifications

#### 3.3.1 A simplification

A close inspection reveals that the proof of Theorem 3.1 remains valid when we restrict the parameters  $\alpha, \beta$  to be equal to each other instead of being independent. Thus, the following version of Theorem 3.1 and Lemmas 3.1 – 3.2 takes place:

**Theorem 3.2** *Let  $\epsilon \in (0, 1/4)$ . Then, for every  $\delta > 0$ , we can point out a simple estimate  $\tilde{g}_{\epsilon, \delta}(\cdot)$  satisfying the relation (14). This estimate can be built as follows. Let*

$$\begin{aligned} \tilde{\Phi}(\alpha, \phi; x, y) &= \alpha N \left[ \ln \left( \sum_i A_i(x) \exp\{\alpha^{-1} \phi_i\} \right) + \ln \left( \sum_i A_i(y) \exp\{-\alpha^{-1} \phi_i\} \right) \right. \\ &\quad \left. + g^T y - g^T x : \tilde{U} \times V \rightarrow \mathbb{R}, \right. \\ \tilde{U} &= \{(\alpha, \phi) : \alpha > 0, \phi \in \mathbb{R}^M\}, \\ V &= \{(x, y) : x, y \in X\}. \end{aligned} \quad (33)$$

and

$$2\tilde{S}(\omega) = \inf_{(\alpha, \phi) \in \tilde{U}} \max_{(x, y) \in V} \left[ \tilde{\Phi}(\alpha, \phi; x, y) + 2\alpha\omega \right], \quad (34)$$

so that for a given  $\delta > 0$  there exist  $\tilde{\alpha} > 0$  and  $\tilde{\phi}$  such that

$$\max_{x, y \in X} \left[ \tilde{\Phi}(\tilde{\alpha}, \tilde{\phi}) + 2\tilde{\alpha} \right] \leq 2\tilde{S}(\ln(2/\epsilon)) + \delta. \quad (35)$$

Setting

$$\begin{aligned} c &= \frac{1}{2} \left[ \max_{y \in X} \left\{ \tilde{\alpha} N \ln \left( \sum_i A_i(y) \exp\{-\tilde{\alpha}^{-1} \tilde{\phi}_i\} \right) + g^T y + \tilde{\alpha} \ln(2/\epsilon) \right\} \right. \\ &\quad \left. - \max_{x \in X} \left\{ \tilde{\alpha} N \ln \left( \sum_i A_i(x) \exp\{\tilde{\alpha}^{-1} \tilde{\phi}_i\} \right) - g^T x + \tilde{\alpha} \ln(2/\epsilon) \right\} \right], \end{aligned} \quad (36)$$

the simple estimate

$$\tilde{g}_{\epsilon, \delta}(i^N) = \sum_{t=1}^n \tilde{\phi}_{i_t} + c$$

satisfies  $\text{Risk}(\tilde{g}_{\epsilon, \delta}; \epsilon) \leq \tilde{S}(\ln(2/\epsilon)) + \delta$ .

Besides this,

$$\text{Risk}_*(\epsilon) \geq \vartheta^{-1}(\epsilon) \tilde{S}(\ln(2/\epsilon)), \quad \vartheta(\epsilon) = \frac{2 \ln \left( \frac{2}{\epsilon} \right)}{\ln \left( \frac{1}{4\epsilon} \right)},$$

whence

$$\text{Risk}(\tilde{g}_{\epsilon, \delta}; \epsilon) \leq \vartheta(\epsilon) \text{Risk}_*(\epsilon) + \delta.$$

### 3.3.2 Computational issues

Our local goal is to investigate the structure of function  $\tilde{S}(\cdot)$  and the possibilities to compute associated simple estimates. Same as in the proof of Lemma 3.2, we can rewrite the expression for  $\tilde{S}$  equivalently as

$$\begin{aligned} 2\tilde{S}(\omega) &= \max_{x,y \in X} \Psi_\omega(x,y), \\ \Psi_\omega(x,y) &= \inf_{\alpha > 0, \phi} \left\{ \alpha N \left[ \ln \left( \sum_i A_i(x) \exp\{\alpha^{-1}\phi\} \right) + \ln \left( \sum_i A_i(y) \exp\{-\alpha^{-1}\phi\} \right) \right] + 2\alpha\omega \right. \\ &\quad \left. + g^T y - g^T x \right\} \end{aligned} \quad (37)$$

We now claim that

$$\begin{aligned} \Psi_\omega(x,y) &= \begin{cases} g^T[y-x], & (x,y) \in \tilde{V}_\omega \\ -\infty, & (x,y) \in (X \times X) \setminus \tilde{V}_\omega \end{cases}, \\ \tilde{V}_\omega &= \left\{ (x,y) \in X \times X : N \ln \left( \sum_i \sqrt{A_i(x)A_i(y)} \right) + \omega \geq 0 \right\}. \end{aligned} \quad (38)$$

Indeed, it is immediately seen that

$$\begin{aligned} &\inf_\phi N\alpha \left[ \ln \left( \sum_i A_i(x) \exp\{\alpha^{-1}\phi\} \right) + \ln \left( \sum_i A_i(y) \exp\{-\alpha^{-1}\phi\} \right) \right] \\ &= 2N\alpha \ln \left( \sum_i \sqrt{A_i(x)A_i(y)} \right), \end{aligned}$$

so that (37) reads

$$\Psi_\omega(x,y) = \inf_{\alpha > 0} 2\alpha \left[ N \ln \left( \sum_i \sqrt{A_i(x)A_i(y)} \right) + \omega \right] + g^T[y-x].$$

The right hand side in this expression is  $-\infty$  outside of  $\tilde{V}_\omega$  and is equal to  $g^T[y-x]$  on  $\tilde{V}_\omega$ , and we arrive at (38).

Recalling the first equality in (37), we get

$$\tilde{S}(\omega) = \max_{(x,y) \in \tilde{V}_\omega} \frac{1}{2} g^T[y-x]. \quad (39)$$

To proceed to computational issues, assume that there exists  $\kappa > 0$  such that  $A_i(x) \geq \kappa$  for all  $x \in X$  and all  $i$ . Note that by (39)  $\tilde{S}(\omega)$ ,  $\omega = \ln(2/\epsilon)$ , is the optimal value in the convex optimization problem

$$\max_{x,y} \left\{ \begin{array}{l} \frac{1}{2} g^T[y-x] : \quad x \in X, y \in X \quad (a) \\ f(x,y) \equiv \exp\{-\omega/N\} - \sum_i \sqrt{A_i(x)A_i(y)} \leq 0 \quad (b) \end{array} \right\} \quad (P)$$

((b) is an equivalent form of the constraint cutting  $\tilde{V}_\omega$  off  $X \times X$ ). The problem clearly is solvable. By convex programming theory, an optimal solution  $(\bar{x}, \bar{y})$  to this problem can be augmented

by Lagrange multiplier  $\nu \geq 0$  such that the vectors  $e_x = \frac{\partial}{\partial x} \Big|_{(x,y)=(\bar{x},\bar{y})} [\frac{1}{2}g^T x + \nu f(x,y)]$ ,  $e_y = \frac{\partial}{\partial y} \Big|_{(x,y)=(\bar{x},\bar{y})} [-\frac{1}{2}g^T x + \nu f(x,y)]$  belong to normal cones of  $X$  at the points  $\bar{x}$ ,  $\bar{y}$ , respectively:

$$\forall(x, y \in X) : e_x^T(x - \bar{x}) \geq 0, e_y^T(y - \bar{y}) \geq 0. \quad (40)$$

After  $\bar{x}, \bar{y}, \nu$  are found, we can convert these data straightforwardly into a simple estimate with the  $\epsilon$ -risk  $\tilde{S}(\omega) = \tilde{S}(\ln(2/\epsilon))$  as follows. There are two possible cases: (a)  $\nu = 0$  and (b)  $\nu > 0$ . In the case of (a), (40) implies that  $(\bar{x}, \bar{y})$  is an optimal solution to the problem obtained from (P) by eliminating the constraint (P.b), that is,  $g^T \bar{y}$  is the maximum, and  $g^T \bar{x}$  is the minimum of  $g^T x$  on  $X$ . In this case, an estimate which reproduces  $g^T x$ ,  $x \in X$ , with the risk  $\frac{1}{2}g^T[\bar{y} - \bar{x}] = \tilde{S}(\omega)$  is trivial – this is the constant estimate  $\tilde{g}(i^N) = \frac{1}{2}g^T[\bar{y} + \bar{x}]$ . Note that this case indeed takes place when  $\omega = \ln(2/\epsilon)$  is large enough, that is, given the number of observations, our reliability requirement is too strong to allow for a nontrivial estimation.

Now assume that  $\nu > 0$ . In this case, when setting

$$W = \sum_i \sqrt{A_i(\bar{x})A_i(\bar{y})}, \alpha = \frac{\nu}{N}W, \phi_i = \alpha \ln \left( \frac{\nu W}{\alpha N} \sqrt{A_i(\bar{y})/A_i(\bar{x})} \right), \quad (41)$$

it is straightforward to verify that the function  $\tilde{\Phi}(\alpha, \phi; x, y) + 2\alpha\omega$  attains its maximum over  $x, y \in X$  at the point  $(\bar{x}, \bar{y})$ , and the maximal value is exactly  $2\tilde{S}(\omega)$ . Recalling that  $\omega = \ln(2/\epsilon)$  and invoking Theorem 3.2, our  $\alpha, \phi$  can be straightforwardly converted into a simple estimate with  $\epsilon$ -risk  $\tilde{S}(\ln(2/\epsilon))$ . The bottom line is that such an estimate is readily given by an optimal solution to the convex optimization problem (P) augmented by the associated Lagrange multiplier  $\nu$ . Note that (high accuracy approximations of) these data are produced, in a polynomial time fashion, by every “intelligent” convex programming algorithm.

### 3.4 Example: Tomography problem

Consider the following somehow simplified model of Positron Emission Tomography (PET). There are  $n$  sources of particles; the particles originating at  $j$ -th source form a Poisson process with parameter  $x_j$ , and these processes are independent of each other. The vector  $x = [x_1; \dots; x_n]$  is known to belong to a given convex compact set  $X \subset \mathbb{R}_+^n$ . There are  $M$  detectors, and a particle emitted by source  $j$  can be either registered at the same instant by exactly one of the detectors, or is missed at all. The probabilities for a particle emitted by source  $j$  at certain time instant  $t$  to be registered at this instant by detector  $i$  are known and form a  $M \times n$  matrix  $P = [p_{ij} \geq 0]$ , with  $\sum_j p_{ij} \leq 1$ ,  $1 \leq i \leq M$ . We observe the numbers  $\mu_i$  of particles registered by every one of the detectors in time interval  $0 \leq t \leq 1$  and intend to infer from these observations an estimate of  $g^T x$ ,  $g \in \mathbb{R}^n$  being given.

In “actual PET”, a patient is injected a radioactive fluid (“tracer”) and is placed in a cylinder with the surface split into small cells. Now, every disintegration act in the tracer produces a positron which annihilates with a near-by electron, producing two  $\gamma$ -quanta running in opposite directions from the annihilation point at the speed of light. The orientation of the resulting “line of response” (LOR) passing through the annihilation point is completely random. If this orientation is such that the line crosses the cylinder’s surface at two points, the corresponding

cells are nearly simultaneously hit by  $\gamma$ -quanta. Such an event (two cells are hit within a very narrow time window) is registered. When modelling this process, the sources of “particles” are voxels – small 3D cubes into which we split the interior of the cylinder, the detectors are pairs of cells (called bins), and  $p_{ij}$  is (easily computable from the geometry of the device) probability for a LOR (“particle”) originating in voxel  $j$  to be registered by pair  $i$  of the cells.

The tracer is chosen in such a way that it concentrates in “areas of interest” (e.g., in the areas of high metabolic activity in tumor diagnosis), and the distribution of tracer between the voxels plays the role of the parameter vector  $x$ .

This problem as it is stated does not fit the framework of Problem I – instead of  $N$  independent observations of a discrete random variable with distribution affinely parameterized by  $x$ , we have a single observation of random variable with distribution nonlinearly depending on  $x$ . Nevertheless, the problem can be reduced to Problem I, namely, as follows. Given integer  $N > \max_{x \in X} \sum_j x_j$ , let us split the time interval  $[0, 1]$  into  $N$  consecutive intervals  $D_1, \dots, D_N$  of duration  $\Delta = 1/N$  each, and consider an “artificial” model as follows: in time interval  $D_t$ , the sources either emit no particles at all, which happens with probability  $1 - \Delta \sum_j x_j$ , or exactly one particle is emitted, and its origin is  $j$ -th source with probability  $x_j \Delta$ . The emitted particle, if any, is registered by detector  $i$  with probability  $p_{ij}$ , where  $j$  is the source of the particle. Finally, emission/detection processes at time slots  $D_1, \dots, D_N$  are independent of each other. In this new problem, we do have  $N$  independent observations of a discrete random variable  $\iota$  which takes value  $i$ ,  $1 \leq i \leq M$ , with probability  $A_i^\Delta(x) = \Delta \sum_j p_{ij} x_j$ , and takes additional value 0 (“no particle in time slot of duration  $\Delta$  is registered”) with probability  $A_0^\Delta(x) = 1 - \sum_{i=1}^M A_i^\Delta(x)$ . Note that a simple estimate, associated with our new observation model, is an affine function of the numbers of particles registered by every one of the detectors in the time interval  $[0, 1]$ . At the same time, as  $N \rightarrow \infty$ , the distribution of the random vector comprised of these numbers approaches the distribution of the observations in the “true” model, so that the larger is  $N$ , the closer is the artificial model to the true one.

Now, problem (39) responsible for a “good” simple estimate for the artificial model reads

$$\tilde{S}^N(\omega) = \max_{x, y \in X} \left\{ \frac{1}{2} g^T [y - x] : \sum_{i=0}^M \sqrt{A_i^\Delta(x) A_i^\Delta(y)} \geq \exp\{-\omega/N\} \right\}, \quad \omega = \ln(2/\epsilon). \quad (P_N)$$

When replacing the constraint with the equivalent constraint

$$F_N(x, y) \equiv N \left[ \sum_{i=0}^M \sqrt{A_i^\Delta(x) A_i^\Delta(y)} - \exp\{-\omega/N\} \right] \geq 0 \quad (42)$$

we observe that as  $N \rightarrow \infty$ , the function  $F_N(x, y)$  converges uniformly on  $X \times X$  to the function  $F(x, y) = \sum_{i=1}^M \sqrt{a_i(x) a_i(y)} - \frac{1}{2} \left[ \sum_{i=1}^M [a_i(x) + a_i(y)] \right] + \omega$ , where  $a_i(x) = \sum_j p_{ij} x_j$ . Thus, the problems  $(P_N)$  “have a limit” as  $N \rightarrow \infty$ , and the limiting problem is

$$\text{Opt} = \max_{x, y \in X} \left\{ \frac{1}{2} g^T [y - x] : h(x, y) \equiv \frac{1}{2} \sum_{i=1}^M \left[ \sqrt{a_i(x)} - \sqrt{a_i(y)} \right]^2 - \omega \leq 0 \right\}. \quad (P_\infty)$$

We can expect that an optimal solution  $(\bar{x}, \bar{y})$  of  $(P_\infty)$  can be converted into an affine in  $\mu = [\mu_1; \dots; \mu_M]$  estimate of  $g^T x$  via the actual observations  $\mu$ , and that the  $\epsilon$ -risk of this estimate on  $X$  is within an absolute constant factor of the minimax optimal  $\epsilon$ -risk associated with the true problem. We are about to demonstrate that this is indeed the case, provided that  $a_i(x)$  are positive on  $X$ ,  $i = 1, \dots, M$ .



**Building the estimate.** Let  $\omega = \ln(2/\epsilon)$ . The associated problem  $(P_\infty)$  clearly is convex and feasible, and  $h$  is smooth on  $X \times X$  due to  $a_i(x) > 0$ ,  $x \in X$ . It follows that the problem is solvable, and an optimal solution  $(\bar{x}, \bar{y})$  to this problem can be augmented by a nonnegative Lagrange multiplier  $\nu$  in such a way that for the vectors

$$\begin{aligned} e_x &\equiv \nabla_x \Big|_{x=\bar{x}, y=\bar{y}} \left[ \frac{1}{2} g^T [x - y] + \nu h(x, y) \right] = \frac{1}{2} \left[ g + \nu p - \nu \sum_i \sqrt{a_i(\bar{y})/a_i(\bar{x})} p_i \right], \\ p_i &= \nabla a_i(x), p = \sum_i p_i, \\ e_y &\equiv \nabla_y \Big|_{x=\bar{x}, y=\bar{y}} \left[ \frac{1}{2} g^T [x - y] + \nu h(x, y) \right] = \frac{1}{2} \left[ -g + \nu p - \nu \sum_i \sqrt{a_i(\bar{x})/a_i(\bar{y})} p_i \right] \end{aligned} \quad (43)$$

it holds

$$\forall x, y \in X : e_x^T (x - \bar{x}) \geq 0, e_y^T (y - \bar{y}) \geq 0, \quad (44)$$

and, in addition,

$$\nu h(\bar{x}, \bar{y}) = 0. \quad (45)$$

It is possible that  $\nu = 0$ . Then (43), (44) imply that  $(\bar{x}, \bar{y})$  maximize  $\frac{1}{2} g^T [y - x]$  in  $x, y \in X$ , so that  $\bar{y}$  is a maximizer,  $\bar{x}$  is a minimizer of  $g^T x$  over  $x \in X$ , and  $\text{Opt} = \frac{1}{2} g^T [\bar{y} - \bar{x}]$ . In this case, the estimate of  $g^T x$  associated with the optimal solution of  $(P_\infty)$  is the constant estimate  $\hat{g} \equiv \frac{1}{2} g^T [\bar{y} + \bar{x}]$ , and the  $\epsilon$ -risk of this estimate clearly is equal to  $\text{Opt}$ .

Now let  $\nu > 0$ . Let us set

$$\phi_i = \frac{1}{2} \ln(a_i(\bar{y})/a_i(\bar{x})),$$

and consider the affine in  $x, y$  function

$$\Psi(x, y) = \nu \sum_i [a_i(x) \exp\{\phi_i\} - a_i(x)] + \nu \sum_i [a_i(y) \exp\{-\phi_i\} - a_i(y)] + g^T [y - x] + 2\nu\omega.$$

We have

$$\nabla_x \Psi(x, y) = \nu \sum_i \sqrt{a_i(\bar{y})/a_i(\bar{x})} p_i - \nu p - g = -2e_x, \quad \nabla_y \Psi(x, y) = \nu \sum_i \sqrt{a_i(\bar{x})/a_i(\bar{y})} p_i - \nu p + g = -2e_y,$$

whence, by (44),  $\Psi$  attains its maximum in  $x, y \in X$  at the point  $(\bar{x}, \bar{y})$ :

$$\begin{aligned} \forall x, y \in X : \Psi(x, y) &\leq \Psi(\bar{x}, \bar{y}) = \nu \sum_i \left[ a_i(\bar{x}) \sqrt{a_i(\bar{y})/a_i(\bar{x})} - a_i(\bar{x}) \right] + \nu \sum_i \left[ a_i(\bar{y}) \sqrt{a_i(\bar{x})/a_i(\bar{y})} - a_i(\bar{y}) \right] \\ &+ g^T [\bar{y} - \bar{x}] + 2\nu\omega = 2\nu \left[ \sum_i \sqrt{a_i(\bar{x})a_i(\bar{y})} - \frac{1}{2} \sum_i [a_i(\bar{x}) + a_i(\bar{y})] + \omega \right] + g^T [\bar{y} - \bar{x}] = 2\text{Opt}, \end{aligned} \quad (46)$$

where the concluding inequality is given by (45) combined with  $\nu > 0$ . Recalling the definition of  $\Psi$ , (46) reads

$$\underbrace{\max_{x \in X} \left[ \nu \sum_i [a_i(x) \exp\{\phi_i\} - a_i(x)] + \nu\omega - g^T x \right]}_I + \underbrace{\max_{y \in X} \left[ \nu \sum_i [a_i(y) \exp\{-\phi_i\} - a_i(y)] + \nu\omega + g^T y \right]}_J \leq 2\text{Opt}. \quad (47)$$

Now let us set

$$c = \frac{J - I}{2};$$

Then (47) reads

$$\begin{aligned} (a) \quad \max_{x \in X} \left[ \nu \sum_i [a_i(x) \exp\{\phi_i\} - a_i(x)] + \nu\omega - g^T x + c \right] &= I + \frac{J - I}{2} = \frac{I + J}{2} \leq \text{Opt}, \\ (a) \quad \max_{x \in X} \left[ \nu \sum_i [a_i(x) \exp\{-\phi_i\} - a_i(x)] + \nu\omega + g^T x - c \right] &= J - \frac{J - I}{2} = \frac{I + J}{2} \leq \text{Opt}. \end{aligned} \quad (48)$$

Consider the affine in  $\mu$  estimate

$$\widehat{g}(\mu) = \nu \sum_i \phi_i \mu_i + c \quad (49)$$

and let us prove that its  $\epsilon$ -risk on  $X$  does not exceed  $\text{Opt}$ . Indeed, let  $\delta > 0$  and  $x \in X$ . Using the standard properties of Poisson processes we conclude that the components  $\mu_i$ ,  $i = 1, \dots, M$ , of the observation  $\mu$  in the true problem, the underlying signal being  $x$ , are *independent of each other* Poisson random variables with the parameters of the Poisson distributions  $a_1(x), \dots, a_M(x)$ , respectively. Denoting by  $P(x)$  the distribution of observations,  $x$  being the underlying signal, and setting  $a(x) = \sum_i a_i(x)$ , we have

$$\begin{aligned} & \text{Prob}_{\mu \sim P(x)} \{ \widehat{g}(\mu) - g^T x > \text{Opt} + \delta \} \leq \mathbf{E}_{\mu \sim P(x)} \{ \exp\{ \nu^{-1} [\widehat{g}(\mu) - g^T x - \text{Opt} - \delta] \} \} \\ & = \exp\{ \nu^{-1} [c - g^T x - \text{Opt} - \delta] \} \prod_{i=1}^M \mathbf{E}_{\mu_i \sim \text{Poisson}(a_i(x))} \{ \exp\{ \mu_i \phi_i \} \} \\ & = \exp\{ \nu^{-1} [c - g^T x - \text{Opt} - \delta] \} \prod_{i=1}^M \left[ \sum_{k=0}^{\infty} \frac{a_i^k(x) \exp\{-a_i(x)\}}{k!} \exp\{k \phi_i\} \right] \\ & = \exp\{ \nu^{-1} [c - g^T x - \text{Opt} - \delta] \} \prod_{i=1}^M \exp\{ a_i(x) \exp\{ \phi_i \} - a_i(x) \} \\ & = \exp\{ \nu^{-1} [ \nu [\sum_i a_i(x) \exp\{ \phi_i \} - a(x)] + c - g^T x - \text{Opt} - \delta] \} \leq \exp\{-\omega - \nu^{-1} \delta\}, \end{aligned}$$

where the concluding inequality is given by (48.a). Recalling the definition of  $\omega$ , we arrive at

$$\text{Prob}_{\mu \sim P(x)} \{ \widehat{g}(\mu) - g^T x > \text{Opt} + \delta \} \leq \epsilon'/2 < \epsilon/2,$$

with  $\epsilon'$  independent of  $x$ . Similar computation, with (48.b) substituted for (48.a), results in

$$\text{Prob}_{\mu \sim P(x)} \{ \widehat{g}(\mu) - g^T x < -\text{Opt} - \delta \} \leq \epsilon'/2 < \epsilon/2;$$

since  $x \in X$  is arbitrary, we see that the  $\epsilon$ -risk of the estimate  $\widehat{g}$  on  $X$  is  $\leq \text{Opt} + \delta$  for every  $\delta \geq 0$ , that is, this risk is  $\leq \text{Opt}$ , as claimed.

**Near-optimality of the estimate.** We claim that the optimal value  $\text{Opt}$  in  $(P_\infty)$  (which, as we have just seen, is the  $\epsilon$ -risk of the affine in  $\mu$  estimate we have built) is within the factor  $\vartheta = \frac{2 \ln(\frac{2}{\epsilon})}{\ln(\frac{1}{4\epsilon})}$  of the minimax optimal  $\epsilon$ -risk  $\text{Risk}_*(\epsilon)$  in the Tomography problem.

Indeed, assume that our claim is not valid, that is, that there exist an estimate  $g_*(\mu)$ ,  $R < \text{Opt}/\vartheta$  and  $\epsilon' < \epsilon$  such that

$$\sup_{x \in X} \text{Prob}_{\mu \sim P(x)} \{ |g_*(\mu) - g^T x| > R \} \leq \epsilon', \quad (50)$$

where, as above,  $P(x)$  is the distribution of observations  $\mu$  in the Tomography problem, the underlying signal being  $x$ . Let us lead this assumption to a contradiction. Consider  $N$ -th approximating model of observations, and let  $P_N(\cdot)$  be the distribution of the ‘‘particle count’’ in this model, that is, the random vector with  $i$ -th coordinate,  $i = 0, 1, \dots, M$ , being the total number of particles registered by  $i$ -th detector in the time period  $0 \leq t \leq 1$ . Both  $P_N(x)$  and  $P(x)$  are probability distributions on a common discrete set (the set  $\mathbb{Z}_+^M$  of nonnegative  $M$ -dimensional integral vectors); from the compactness of  $X$  and the construction of our approximating models it follows that the distributions  $P_N(x)$  uniformly in  $x \in X$  converge in probability, as  $N \rightarrow \infty$ , to  $P(x)$ ; since we are speaking about distributions on a common countable set, it follows that for every  $\delta > 0$  there exists  $N_\delta$  such that

$$\forall (x \in X, A \subset \mathbb{Z}_+^M, N \geq N_\delta) : |P_N(x)\{A\} - P(x)\{A\}| \leq \delta.$$

Now let  $\epsilon'' \in (\epsilon', \epsilon)$ . From the just described convergence of distributions  $P_N(\cdot)$  to  $P(\cdot)$  it follows immediately that with properly chosen  $N_0$  we have

$$N \geq N_0 \Rightarrow \sup_{x \in X} \text{Prob}_{\mu \sim P_N(x)} \{ |g_*(\mu) - g^T x| > R \} < \epsilon'',$$

whence also

$$\forall(N \geq N_0) : \text{Risk}_{N,*}(\epsilon'') \leq R, \quad (51)$$

where  $\text{Risk}_{N,*}(\epsilon'')$  is the minimax optimal  $\epsilon''$ -risk in  $N$ -th approximating problem. Now, recall that  $\text{Opt}$  is the optimal value in problem  $(P_\infty)$  with  $\omega = \ln(2/\epsilon)$ , and  $(\bar{x}, \bar{y})$  is an optimal solution to this problem. Setting  $\omega'' = \ln(2/\epsilon'') > \omega$ , we conclude that

$$\frac{1}{2} \sum_{i=1}^M \left[ \sqrt{a_i(\bar{x})} - \sqrt{a_i(\bar{y})} \right]^2 - \omega'' < 0,$$

whence

$$N \left[ \sum_{i=0}^M \sqrt{A_i^\Delta(\bar{x}) A_i^\Delta(\bar{y})} - \exp\{-\omega''/N\} \right] < 0$$

for all large enough values of  $N$  (cf. the derivation of  $(P_\infty)$ ). In other words, for all large enough values of  $N$ ,  $(\bar{x}, \bar{y})$  is a feasible solution to the problem  $(P_N)$  with  $\omega$  replaced with  $\omega''$ , whence  $\tilde{S}^N(\omega'') \equiv \tilde{S}^N(\ln(2/\epsilon'')) \geq \text{Opt}$  for these  $N$ . By Theorem 3.2, we have  $\text{Risk}_{N,*}(\epsilon'') \geq \vartheta^{-1} \tilde{S}^N(\ln(2/\epsilon'')) \geq \vartheta^{-1} \text{Opt}$ . Invoking (51), we conclude that  $R \geq \vartheta^{-1} \text{Opt}$ , which is a desired contradiction (look at the origin of  $R$ ).

### 3.5 Adaptive version of the estimate

Let  $X^1 \subset X^2 \subset \dots \subset X^K$  be a nested collection of nonempty convex compact sets in  $\mathbb{R}^n$ . Consider a modification of Problem I where the set  $X$ , instead of being given in advance, is known to be one of the sets of the collection. Given a linear form  $g^T z$  on  $\mathbb{R}^n$ , let  $\text{Risk}^k(\hat{g}; \epsilon)$  and  $\text{Risk}_*^k(\epsilon)$  be, respectively, the  $\epsilon$ -risk of an estimate  $\hat{g}$  on  $X^k$ , and the minimax optimal  $\epsilon$ -risk of recovering  $g^T x$  on  $X^k$ . Let also  $S_k(\cdot)$  be the function (16) associated with  $X = X^k$ . As it is immediately seen, the functions  $S_k(\cdot)$  grow with  $k$ . Our goal is to modify the estimate  $\hat{g}$  we have built in such a way that the  $\epsilon$ -risk of the modified estimate on  $X^k$  will be “nearly”  $\text{Risk}_*^k(\epsilon)$  for every  $k \leq K$ . This goal can be achieved by a straightforward application of the well-known Lepskii’s adaptation scheme [26] as follows.

Given  $\delta > 0$ , let  $\delta' \in (0, \delta)$ , and let  $\hat{g}^k(\cdot)$  be the simple estimate with the  $(\epsilon/K)$ -risk on  $X^k$  not exceeding  $S_k(\ln(2K/\epsilon)) + \delta'$  given by the construction of Lemma 3.2 applied with  $\epsilon/K$  substituted for  $\epsilon$  and  $X^k$  substituted for  $X$ . Then,

$$\forall(k \leq K) : \sup_{x \in X^k} \text{Prob}_{i^N \sim (A(x))^N} \{ |\hat{g}^k(i^N) - g^T x| > S_k(\ln(2K/\epsilon)) + \delta \} \leq \epsilon'/K < \epsilon/K. \quad (52)$$

Given observation  $i^N$ , let us say that the index  $k \leq K$  is  $i^N$ -good, if

$$\forall(k', k \leq k' \leq K) : |\hat{g}^{k'}(i^N) - \hat{g}^k(i^N)| \leq S_k(\ln(2K/\epsilon)) + S_{k'}(\ln(2K/\epsilon)) + 2\delta. \quad (53)$$

Note that  $i^N$ -good indexes do exist (e.g.,  $k = K$ ). Given  $i^N$ , we can find the smallest  $i^N$ -good index  $k = k(i^N)$ ; our estimate is nothing but  $\hat{g}(i^N) = \hat{g}^{k(i^N)}(i^N)$ .

**Proposition 3.1** *Assume that  $\epsilon \in (0, 1/4)$ , and let*

$$\vartheta = 3 \frac{\ln(2K/\epsilon)}{\ln(2/\epsilon)}. \quad (54)$$

Then

$$\forall(k, 1 \leq k \leq K) : \sup_{x \in X^k} \text{Prob}_{i^N \sim (A(x))^N} \{ |\hat{g}(i^N) - g^T x| > \vartheta S_k(\ln(2/\epsilon)) + 3\delta \} < \epsilon, \quad (55)$$

whence also

$$\forall(k, 1 \leq k \leq K) : \text{Risk}^k(\widehat{g}; \epsilon) \leq \frac{6 \ln\left(\frac{2K}{\epsilon}\right)}{\ln\left(\frac{1}{4\epsilon}\right)} \text{Risk}_*^k(\epsilon) + 3\delta. \quad (56)$$

**Proof.** Setting  $\omega = \ln(2K/\epsilon)$ , let us fix  $\bar{k} \leq K$  and  $x \in X^{\bar{k}}$  and call a realization  $i^N$   $x$ -good, if

$$\forall(k, \bar{k} \leq k \leq K) : |\widehat{g}^k(i^N) - g^T x| \leq S_k(\omega) + \delta. \quad (57)$$

Since  $X^k \supset X^{\bar{k}}$  when  $k \geq \bar{k}$ , (52) implies that

$$\text{Prob}_{i^N \sim \mathcal{A}^N(x)} \{i^N \text{ is good}\} \geq 1 - \epsilon'. \quad (58)$$

Now, when  $x$  is the signal and  $i^N$  is  $x$ -good, relations (57) imply that  $\bar{k}$  is an  $i^N$ -good index, so that  $k(i^N) \leq \bar{k}$ . Since  $k(i^N)$  is an  $i^N$ -good index, we have

$$|\widehat{g}(i^N) - \widehat{g}^{\bar{k}}(i^N)| = |\widehat{g}^{k(i^N)}(i^N) - \widehat{g}^{\bar{k}}(i^N)| \leq S_{k(i^N)}(\omega) + S_{\bar{k}}(\omega) + 2\delta,$$

which combines with (57) to imply that

$$|\widehat{g}(i^N) - g^T x| \leq 2S_{\bar{k}}(\omega) + S_{k(i^N)}(\omega) + 3\delta \leq 3S_{\bar{k}}(\omega) + 3\delta, \quad (59)$$

where the concluding inequality is due to  $k(i^N) \leq \bar{k}$  and to the fact that  $S_k$  grows with  $k$ . The bound (59) holds true whenever  $i^N$  is  $x$ -good, which, as we have seen, happens with probability  $\geq 1 - \epsilon'$ . Since  $\epsilon' < \epsilon$  and  $\bar{x} \in X^{\bar{k}}$  is arbitrary, we conclude that

$$\text{Risk}^{\bar{k}}(\widehat{g}; \epsilon) \leq 3S_{\bar{k}}(\omega) + 3\delta. \quad (60)$$

By Lemma 3.3 we have  $S_k(\theta\gamma) \leq \theta S_k(\gamma)$  whenever  $\gamma \geq 0$  and  $\theta \geq 1$ . Recalling the definition of  $\omega$  and  $\vartheta$ , the right hand side in (60) does not therefore exceed  $\vartheta S_{\bar{k}}(\ln(2/\epsilon)) + 3\delta$ . Since  $\bar{k} \leq K$  is arbitrary, we have proved (55). This bound, by Lemma 3.2, implies (56). ■

## 4 From estimating linear forms to signal recovery

We have seen that in the context of Problem I, we know how to recover a linear form of the unknown signal with  $\epsilon$ -risk just by an absolute constant factor larger than the minimax optimal risk. As it was mentioned in Introduction, in the context of Problem II similar fact was established by Donoho [8], who proved that the minimal  $\epsilon$ -risk on  $X$  of recovering  $g^T x$  achievable with affine in  $y$  estimates is within an absolute constant factor of the minimax risk, and this is so for risks associated with all usual loss functions, not only for the  $\epsilon$ -risk (a significantly simplified, as compared to the original one, proof of Donoho's result can be found in the appendix). Our goal is to demonstrate that when  $X$  has a favourable geometry, nearly optimal estimates of linear forms imply "not too bad" estimates of the unknown signal. For the sake of simplicity, we focus on recovery the signal in the standard Euclidean norm.

Our simple approach is equally applicable to the situations of Problems I and II; in order to make the presentation "problem-independent", we from now denote the observations by  $w$  ( $w$  is  $i^N$  in the case of Problem I and  $w$  is  $y$  in the case of Problem II), and denote by  $P_x$  the distribution of observations associated with  $x$ . A candidate estimate is now a Borel function  $\widehat{x}(w)$  taking values

in the space  $\mathbb{R}^n$  where  $x$  lives. Given a tolerance  $\epsilon \in (0, 1)$ , we quantify the quality of such an estimate by the worst-case, over  $x \in X$ , upper  $(1 - \epsilon)$ -quantile of the recovering error as measured in the Euclidean norm:

$$\text{Risk}_2(\hat{x}; \epsilon) = \inf \left\{ \delta : \sup_{x \in X} \text{Prob}_{w \sim P_x} \{ \|\hat{x}(w) - x\| > \delta \} < \epsilon \right\},$$

and denote by  $\text{Risk}_{2,*}(\epsilon)$  the associated optimal risk:

$$\text{Risk}_{2,*}(\epsilon) = \inf_{\hat{x}(\cdot)} \text{Risk}_2(\hat{x}; \epsilon).$$

**The construction.** Let us choose somehow a collection of  $N$  unit vectors  $e_1, \dots, e_N$  in  $\mathbb{R}^n$ . For  $\epsilon \in (0, 0.1)$ , let  $\text{Risk}_*^\ell(\epsilon)$  be the optimal, in the minimax sense,  $\epsilon$ -risk of recovering  $e_\ell^T x$  via observations  $w \sim P_x$ . By the results of Section 3 (case of Problem I) and the just cited Donoho's Theorem (case of Problem II), whenever  $\epsilon \in (0, 0.01)$  is given, we can build efficiently (provided that  $X$  is computationally tractable) estimates  $\tilde{e}_\epsilon^\ell(\cdot)$  of the linear forms  $e_\ell^T x$  and compute efficiently upper bounds  $R^\ell(\epsilon)$  on the  $\epsilon$ -risks of the estimates:

$$R^\ell(\epsilon) > \inf \left\{ \delta : \sup_{x \in X} \text{Prob}_{w \sim P_x} \left\{ |\tilde{e}_\epsilon^\ell(w) - e_\ell^T x| > \delta \right\} < \epsilon \right\}$$

in such a way that  $R^\ell$  are within an absolute constant factor  $C$  of the minimax optimal  $\epsilon$ -risks  $\text{Risk}_*^\ell(\epsilon)$  of recovering  $e_\ell^T x$ ,  $x \in X$ :

$$R^\ell(\epsilon) < C \text{Risk}_*^\ell(\epsilon). \quad (61)$$

Now, given  $\bar{\epsilon} \in (0, 0.1)$ , consider the following estimate  $\hat{x}$  of a signal  $x \in X$  via observations  $w$ . We build the  $N$  estimates  $\tilde{e}^\ell(\cdot) \equiv \tilde{e}_{\bar{\epsilon}/N}^\ell(\cdot)$ ,  $1 \leq \ell \leq N$ . We further take as  $\hat{x}(w)$  (any) vector  $u$  satisfying the relations

$$u \in X \quad \text{and} \quad |e_\ell^T u - \tilde{e}_{\bar{\epsilon}/N}^\ell(w)| \leq R^\ell(\bar{\epsilon}/N), \quad \ell = 1, \dots, N, \quad (62)$$

if such an  $u$  exists, otherwise  $\hat{x}(w)$  is a once for ever fixed point of  $X$ .

**Analysis.** Let  $p_\infty(z) = \max_\ell |e_\ell^T z|$ ,  $z \in \mathbb{R}^n$ , and let  $\text{Risk}_{\infty,*}(\epsilon)$  be the optimal, in the minimax sense,  $\epsilon$ -risk of recovering  $x \in X$  via observations  $w \sim P_x$ , the loss function being  $p_\infty(\cdot)$ :

$$\text{Risk}_{\infty,*}(\epsilon) = \inf_{\tilde{x}(\cdot)} \text{Risk}_\infty(\tilde{x}; \epsilon), \quad \text{Risk}_\infty(\tilde{x}; \epsilon) = \inf \left\{ \delta : \sup_{x \in X} \text{Prob}_{w \sim P_x} \{ p_\infty(\tilde{x}(w) - x) > \delta \} < \epsilon \right\}.$$

Since  $\|\cdot\| \geq p_\infty(\cdot)$  due to  $\|e_\ell\| = 1$ , we have

$$\text{Risk}_{\infty,*}(\epsilon) \leq \text{Risk}_{2,*}(\epsilon). \quad (63)$$

Our goal is to compare  $\text{Risk}_\infty(\hat{x}; \bar{\epsilon})$  and  $\text{Risk}_{\infty,*}(\bar{\epsilon})$ . By the origin of the estimates and  $R^\ell$ , when  $w \sim P_x$  with  $x \in X$ , every one of the  $N$  inequalities  $|e_\ell^T x - \tilde{e}_{\bar{\epsilon}/N}^\ell(w)| \leq R^\ell(\bar{\epsilon}/N)$  takes place with probability  $\geq 1 - \bar{\epsilon}/N + \delta$  with certain independent of  $x$   $\delta > 0$ . If all these inequalities take place (which happens with probability  $\geq 1 - \bar{\epsilon} + \delta$ ), (62) is feasible (since the constraints in (62) are satisfied when  $u = x$ ), and for every feasible solution  $u$  to (62) we have  $|e_\ell^T u - e_\ell^T x| \leq 2R^\ell(\bar{\epsilon}/N)$ . Thus, we have

$$\sup_{x \in X} \text{Prob}_{w \sim P_x} \left\{ |e_\ell^T [\hat{x}(w) - x]| \leq 2R^\ell(\bar{\epsilon}/N), \quad \ell = 1, \dots, n \right\} < \bar{\epsilon},$$

whence

$$\text{Risk}_\infty(\hat{x}; \bar{\epsilon}) \leq 2 \max_{1 \leq \ell \leq N} R^\ell(\bar{\epsilon}/N). \quad (64)$$

Now, in the situation of Problem I we have  $R^\ell(\epsilon) = S(\ln(2/\epsilon))$ , and therefore by Lemma 3.3 one has  $R^\ell(\bar{\epsilon}/N) \leq O(\ln(2N/\bar{\epsilon})/\ln(2/\bar{\epsilon}))R^\ell(\bar{\epsilon})$ . One can easily extract from [8] that the resulting inequality holds true in the situation of Problem II as well. Thus, (64) implies that

$$\text{Risk}_\infty(\hat{x}; \bar{\epsilon}) \leq \vartheta \max_{\ell \leq N} R^\ell(\bar{\epsilon}), \quad \vartheta = O(\ln(N/\bar{\epsilon})/\ln(1/\bar{\epsilon})),$$

which combines with (61) to imply that

$$\begin{aligned} \text{Risk}_\infty(\hat{x}; \bar{\epsilon}) &\leq \widehat{\vartheta} \max_{\ell \leq N} R_*^\ell(\bar{\epsilon}) \leq \widehat{\vartheta} \text{Risk}_{\infty,*}(\bar{\epsilon}) \leq \widehat{\vartheta} \text{Risk}_{2,*}(\bar{\epsilon}), \\ \widehat{\vartheta} &= O(\ln(N/\bar{\epsilon})/\ln(1/\bar{\epsilon})). \end{aligned} \quad (65)$$

Note that  $\widehat{\vartheta}$  is a moderate constant, unless  $N$  is astronomically large. We conclude that *unless  $N$  is astronomically large, the estimate  $\hat{x}$  is nearly optimal in the sense of its  $\epsilon$ -risk on  $X$  associated with the loss function  $p_\infty$ .*

Now assume that the geometry of  $X$  allows to choose a collection  $\{e_\ell\}$  of a “moderate” number  $N$  of unit vectors in such a way that

$$\forall u \in X - X : \|u\| \leq C_X p_\infty^{\chi(X)}(u) \quad (66)$$

where  $C_X > 0$  and  $\chi(X) \in (0, 1]$  are appropriate constants. Since  $\hat{x}$  takes values in  $X$ , (66) combines with (65) to imply that

$$\text{Risk}_2(\hat{x}; \bar{\epsilon}) \leq C_X [\text{Risk}_\infty(\hat{x}, \bar{\epsilon})]^{\chi(X)} \leq C_X \widehat{\vartheta}^{\chi(X)} [\text{Risk}_{2,*}(\bar{\epsilon})]^{\chi(X)}, \quad (67)$$

so that *the  $\bar{\epsilon}$ -risk of the estimate  $\hat{x}$ , the loss function being  $\|\cdot\|$ , can be bounded from above in terms of the corresponding minimax optimal risk.* Ideally, we would like to have  $\chi(X) = 1$ , meaning that our estimate  $\hat{x}$  is “nearly minimax optimal” in terms of  $\|\cdot\|$ -risk (recall that for all practical purposes,  $\widehat{\vartheta}$  is a moderate constant). How “far” we are from this ideal situation (that is, how far is  $\chi(X)$  from 1), it depends solely on the geometry of  $X$  and is completely independent of how good is the affine mapping  $A(x)$  (Problem I) or the matrix  $A$  (Problem II). It should be added that there are important situations where (66) is satisfied with “not too bad” constants  $C_X, \chi(X)$ . Here are two instructive examples:

**Example 1:  $\ell_1$ -ball.** Assume that  $X \subset \Delta_R = \{x \in \mathbb{R}^N : \sum_i |x_i| \leq R\}$  (this is a frequently used model of a “sparse” signal). In this case, choosing as  $e_i$  the standard basic orths, we clearly have  $\sum_i |u_i| \leq 2R$  for every  $u \in X - X$ , whence  $\|u\| = \sqrt{\sum_i u_i^2} \leq \sqrt{p_\infty(u) \sum_i |u_i|} \leq \sqrt{2R} p_\infty^{1/2}(u)$ , that is, (66) holds true with  $C_X = \sqrt{2R}$ ,  $\chi(X) = 1/2$ .

**Example 2: Ellipsoid.** Now assume that  $X$  is a centered at the origin ellipsoid with half-axes  $d_i = Ri^{-\gamma}$ ,  $\gamma > 0$  (this is the standard model of signals from Sobolev balls restricted onto uniform grids). In this case, assuming w.l.o.g. that the directions of the ellipsoid axes are the standard basic orths and choosing these orths as  $e_1, \dots, e_N$ , for  $u \in X - X$  we have

$$\sum_{i=1}^N u_i^2 i^{2\gamma} \leq (2R)^2,$$

whence for every integer  $k \geq 0$  one has  $\sum_{i=k+1}^N u_i^2 \leq (2R)^2(k+1)^{-2\gamma}$ . It follows that for every integer  $k \geq 0$  we have

$$\|u\|^2 \leq kp_\infty^2(u) + \sum_{i=k+1}^n u_i^2 \leq kp_\infty^2(u) + (2R)^2(k+1)^{-2\gamma}.$$

Minimizing the resulting bound in  $k$ , we get  $\|u\| \leq O(1)R^{\frac{1}{2\gamma+1}}p_\infty^{\frac{2\gamma}{2\gamma+1}}(u)$ , that is, in the case in question  $C_X = O(1)R^{\frac{1}{2\gamma+1}}$ ,  $\chi(X) = \frac{2\gamma}{2\gamma+1}$ .

## References

- [1] Ben-Tal, A., and Nemirovski, A., *Lectures on Modern Convex Optimization: Analysis, Algorithms and Engineering Applications*, MPS-SIAM Series on Optimization, SIAM, Philadelphia, 2001.
- [2] Birgé, L., Sur un théorème de minimax et son application aux tests. (French) *Probab. Math. Stat.* **3** (1982), 259-282.
- [3] Birgé, L., Approximation dans les espaces métriques et théorie de l'estimation. *Z. Wahrscheinlichkeitstheorie verw. Geb.*, **65** (1983), 181-237.
- [4] A. W. Bowman, A.W., Azzalini, A., *Applied Smoothing Techniques for Data Analysis: The Kernel Approach with S-Plus Illustrations*. Oxford University Press, Oxford, 1997.
- [5] Donoho D., Johnstone, I., Ideal spatial adaptation via wavelet shrinkage. *Biometrika* **81:3** (1994), 425-455.
- [6] Donoho, D., Liu, R., Geometrizing Rates of Convergence, II. *The Annals of Statistics* **19:2** (1991), 633-667.
- [7] Donoho, D., Liu, R., Geometrizing Rates of Convergence, III. *The Annals of Statistics* **19:2** (1991), 668-701.
- [8] Donoho, D., Statistical estimation and optimal recovery. *The Annals of Statistics* **22:1** (1995), 238-270.
- [9] Donoho, D., Johnstone, I., Adapting to unknown smoothness via wavelet shrinkage. *J. Amer. Statist. Assoc.* **90:432** (1995), 1200-1224.
- [10] Donoho D., Johnstone, I., Kerkyacharian, G., Picard, D., Wavelet shrinkage: Asymptopia? (with discussion and reply by the authors). *J. Royal Statist. Soc. Series B* **57:2** (1995), 301-369.
- [11] Eubank, R., *Spline smoothing and Nonparametric Regression*, Dekker, New York, 1988.
- [12] Fox, J., *Nonparametric Simple Regression: Smoothing Scatterplots*. Sage, Thousand Oaks CA, 2000.
- [13] Fox, J., *Multiple and Generalized Nonparametric Regression*. Sage, Thousand Oaks CA, 2000.

- [14] Goldenshluger, A., Nemirovski, A., On spatially adaptive estimation of nonparametric regression. *Math. Methods of Statistics*, **6:2** (1997), 135 – 170.
- [15] Golubev, Yu., Asymptotic minimax estimation of regression function in additive model. *Problemy peredachi informatsii*, **28:2** (1992), 3-15. (English transl. in *Problems Inform. Transmission* **28**, 1992.)
- [16] Härdle, W., *Applied Nonparametric Regression*. ES Monograph Series 19, Cambridge, U.K., Cambridge University Press, 1990.
- [17] Härdle, W., Kerkyacharian, G., Picard, D., Tsybakov, A.B., *Wavelets, Approximation and Statistical Applications*. Lecture Notes in Statistics **129**, Springer, New York, 1998.
- [18] Hastie, T.J., Tibshirani, R.J., Freedman, J., *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer, New York, 2001.
- [19] Hastie, T.J., Tibshirani, R.J., *Generalized Additive Models*. Chapman and Hall, London, 1990.
- [20] Ibragimov, I.A., Khasminski, R.Z., *Statistical Estimation: Asymptotic Theory*, Springer, 1981.
- [21] Ibragimov, I., Nemirovski, A., Khas'minski, R., Some problems of nonparametric estimation in Gaussian white noise. *Theory Probab. Appl.* **31:3** (1986), 391-406.
- [22] Juditsky, A., Wavelet estimators: Adapting to unknown smoothness. *Math. Methods of Statistics* **6:1** (1997), 1-25.
- [23] Juditski, A., Nemirovski, A., Functional aggregation for nonparametric regression. *Annals of Statistics* **28** (2000), 681-712.
- [24] Klemela, J., Tsybakov, A.B., Sharp adaptive estimation of linear functionals. *Annals of Statistics* **29** (2001), 1567-1600.
- [25] Korostelev, A., Tsybakov, A., *Minimax theory of image reconstruction*. *Lecture Notes in Statistics* **82**, Springer, New York, 1993.
- [26] Lepskii, O., On a problem of adaptive estimation in Gaussian white noise. *Theory of Probability and Its Applications*, **35:3** (1990), 454-466.
- [27] Lepskii, O., Asymptotically minimax adaptive estimation I: Upper bounds. Optimally adaptive estimates. *Theory of Probability and Its Applications*, **36:4** (1991), 682-697.
- [28] Lepskii, O., Mammen, E., Spokoiny, V., Optimal spatial adaptation to inhomogeneous smoothness: an approach based on kernel estimates with variable bandwidth selectors. *Annals of Statistics* **25:3** (1997), 929-947.
- [29] Nemirovski, A., On nonparametric estimation of smooth regression functions. *Sov. J. Comput. Syst. Sci.*, **23:6** (1985), 1-11.
- [30] Nemirovski, A., *Topics in Non-parametric Statistics*, in: M. Emery, A. Nemirovski, D. Voiculescu, *Lectures on Probability Theory and Statistics*, Ecole d'Été de Probabilités de Saint-Flour XXVII – 1998, Ed. P. Bernard. - *Lecture Notes in Mathematics* **1738**, 87–285.



- [31] Nemirovski, A., Shapiro, A., Convex Approximations of Chance Constrained Programs. *SIAM Journal on Optimization* **17:4** (2006), 969-996.
- [32] Pinsker M., Optimal filtration of square-integrable signals in Gaussian noise. *Problemy peredachi informatsii*, **16:2** (1980), 120-133. (English transl. in *Problems Inform. Transmission* **16**, 1980.)
- [33] Pinsker M., S. Efroimovich. Learning algorithm for nonparametric filtering. *Automation and Remote Control*, **45:11** (1984), 1434-1440.
- [34] Prakasa Rao, B.L.S., *Nonparametric functional estimation*. Academic Press, Orlando, 1983.
- [35] Rosenblatt, M., *Stochastic curve estimation*. Institute of Mathematical Statistics, Hayward, California, 1991.
- [36] Simonoff, J.S., *Smoothing Methods in Statistics*. Springer, New York, 1996.
- [37] Takezawa, K., *Introduction to Nonparametric Regression*. Wiley Series in Probability and Statistics, 2005.
- [38] Tsybakov, A.B. *Introduction a l'estimation non-paramétrique*. Springer, 2004.
- [39] Wahba, G., *Spline models for observational data*. SIAM, Philadelphia, 1990.

## Appendix: Donoho's Theorem revisited

We are about to present an alternative proof of the (finite-dimensional version of the) main result of Donoho [8]. While the underlying ideas are exactly the same as in [8], we believe that our proof is much simpler and more transparent than the original one.

The fact we intend to demonstrate is as follows

**Proposition 4.1** [D. Donoho, [8]] *Let  $\epsilon \in (0, 0.1]$ , and consider the problem of estimating a linear form  $g^T x$  via observations (1). The minimal, over all affine in  $y$  estimates, worst-case  $\epsilon$ -risk  $\text{RiskAff}_*(\epsilon)$  achievable with affine in  $y$  estimates is within the factor 2.05 of the minimax optimal risk  $\text{Risk}_*(\epsilon)$ .*

**Proof.** We lose nothing by assuming that  $X$  possesses a nonempty interior (otherwise we could replace  $\mathbb{R}^n$  with the affine hull of  $X$ ). By scaling the observations, we can assume further that  $\sigma = 1$ . Finally, eliminating the trivial case when  $g$  is constant on  $X$  and scaling  $g$ , we can normalize the situation by the requirement that  $\text{Risk}_*(\epsilon) = 0.49$ ; all we need to prove is that then  $\text{RiskAff}_*(\epsilon) \leq 1$ . To this end, let us set

$$\alpha = 0.49, \beta = \frac{\alpha}{\text{ErfInv}(\epsilon)}$$

where  $\text{ErfInv}(\epsilon)$ ,  $0 < \epsilon < 1$ , is the inverse of the error function

$$\text{Erf}(s) = \frac{1}{\sqrt{2\pi}} \int_s^\infty \exp\{-s^2/2\} ds.$$

Let  $X_* = X - X$ , so that  $X_*$  is a convex compact symmetric w.r.t. origin set with  $0 \in \text{int } X_*$ , that is,  $X_*$  is the unit ball of certain norm  $p(\cdot)$  on  $\mathbb{R}^n$ . Let  $p_*(\cdot)$  be the conjugate norm:  $p_*(x) = \max_u \{x^T u : p(u) \leq 1\}$ . We claim that there exists  $h \in \mathbb{R}^m$  such that

$$\|h\| \leq \beta \quad \text{and} \quad p_*(g - A^T h) \leq 2\alpha. \tag{68}$$

Our claim immediately implies the desired relation  $\text{RiskAff}_*(\epsilon) \leq 1$ . Indeed, let  $h$  solve (68). The relation  $p_*(g - A^T h) \leq 2\alpha$  reads

$$2\alpha \geq \max_{x \in X, y \in X} (g - A^T h)^T(x - y) = \max_{x \in X} (g - A^T h)^T x - \min_{x \in X} (g - A^T h)^T x.$$

Setting  $c = \frac{1}{2} [\max_{x \in X} (g - A^T h)^T x - \min_{x \in X} (g - A^T h)^T x]$  and  $H(x) = h^T A x + c$ , we see that for the affine form  $a(x) = g^T x - H(x)$  one has

$$\max_{x \in X} |a(x)| \leq \alpha.$$

Now consider the affine in  $y$  estimate  $\widehat{g}(y) = h^T y + c$ . For this estimate and  $x \in X$ , we have

$$|g^T x - \widehat{g}(Ax + \xi)| = |g^T x - h^T(Ax + \xi) - c| = |a(x) - h^T \xi|;$$

denoting  $\eta = h^T \xi$ , note that  $\eta \sim \mathbb{N}(0, \bar{\beta}^2)$  with  $\bar{\beta} = \|h\| \leq \beta$ . It follows that

$$\text{Prob}\{|g^T x - \widehat{g}(Ax + \xi)| > 1\} \leq \text{Erf}((1 - a(x))/\bar{\beta}) + \text{Erf}((1 + a(x))/\bar{\beta}) \leq \text{Erf}((1 - \alpha)/\beta) + \text{Erf}((1 + \alpha)/\beta);$$

it is easily seen that the concluding quantity with our  $\alpha$  and  $\beta$  is  $< \epsilon$ , provided  $\epsilon \leq 0.1$ .

It remains to justify the claim. Assume that no  $h$  satisfying (68) exists. Then the convex compact sets  $U = \{u \in \mathbb{R}^n : p_*(g - u) \leq 2\alpha\}$  and  $V = \{v = A^T h : \|h\| \leq \beta\}$  do not intersect and thus can be strongly separated: there exists  $z \in \mathbb{R}^n$  such that  $\min_{u \in U} z^T u > \max_{v \in V} z^T v$ . W.l.o.g. we can assume that  $p(z) = 1$  (i.e.,  $z \in X_*$ ); with this normalization, we have  $\min_{u \in U} z^T u = g^T z - \max\{z^T u : p_*(u) \leq 2\alpha\} = g^T z - 2\alpha p(z) = g^T z - 2\alpha$  and  $\max_{v \in V} z^T v = \max_{w: \|w\| \leq \beta} z^T A^T w = \beta \|Az\|$ . Since  $z \in X_*$ , we have  $z_* = r - s$  with  $r, s \in X$ . Thus, there exist  $r, s \in X$  such that

$$g^T(r - s) > 2\alpha + \beta \|A(r - s)\|. \quad (69)$$

It may happen that  $\|A(r - s)\| \leq 2 \text{ErfInv}(\epsilon)$ . In this case, for every decision rule for distinguishing between the hypotheses  $\Pi_r, \Pi_s$  stating that the distribution of observation (1) is  $\mathcal{N}(Ar, I)$  and  $\mathcal{N}(As, I)$ , respectively, the sum of error probabilities is at least  $2\epsilon$ ; since  $g^T r$  and  $g^T s$  differ by more than  $2\alpha$ , it follows immediately that the worst-case, over  $x \in \{r, s\}$ ,  $\epsilon$ -risk of an arbitrary estimate  $\widehat{g}(\cdot)$  of  $g^T x$  is  $> \alpha$  – otherwise we could use the estimate to distinguish between our two hypotheses by accepting the first of them when  $\widehat{g}(y) \geq \frac{1}{2} g^T(r + s)$  and otherwise accepting the second; if the  $\epsilon$ -risk of our estimate were  $\leq \alpha$ , the probability to reject the true hypothesis would be  $< \epsilon$ . Thus, in the situation in question  $\text{Risk}_*(\epsilon) > 0.49$ . Now consider the case when  $\|A(r - s)\| > 2 \text{ErfInv}(\epsilon)$ . Setting  $\lambda = 2 \text{ErfInv}(\epsilon) / \|A(r - s)\|$  and  $r' = s + \lambda(r - s)$ , we get  $r', s \in X$  and  $\|Ar' - As\| = 2 \text{ErfInv}(\epsilon)$ , while  $g^T(r' - s) = \lambda g^T(r - s) > \lambda[2\alpha + \beta \|A(r - s)\|] \geq 2\beta \text{ErfInv}(\epsilon) \equiv 2\alpha$ . The same reasoning as above, with  $r', s$  in the role of  $r, s$  results in  $\text{Risk}_*(\epsilon) > \alpha = 0.49$ . Thus, we end up with  $\text{Risk}_*(\epsilon) > 0.49$ , which is a contradiction.

#### 4.1 The case of the expected squared error risk

The outlined reasoning can be easily modified to cover the ESE risks, that is, the risks associated with the expected squared error of estimation. The result is, that the worst-case, over  $x \in X$ , ESE-risk of estimating linear form achievable with affine estimates is within the factor 1.25 of the corresponding minimax optimal risk (this is exactly the constant appearing in [8]). To get the announced result, we act as follows.

**Preliminaries.** For  $D > 0$ , let  $R_*(D)$  be the minimax ESE-risk of recovering a parameter  $\theta$  known to belong to  $[-D, D]$  from observation  $\theta + \zeta$ ,  $\zeta \sim \mathcal{N}(0, 1)$ :

$$R_*(D) = \inf_{\phi(\cdot)} \max_{|\theta| \leq D} \mathbf{E}\{[\phi(\theta + \zeta) - \theta]^2\},$$

and let

$$\Omega = \inf_{D > 0} \frac{(1 + D^2)R_*(D)}{D^2}.$$

We claim that the minimal worst-case, over  $x \in X$ , ESE-risk of an affine in  $y = Ax + \xi$ ,  $\xi \sim \mathcal{N}(0, \sigma^2 I)$ , estimate of  $g^T x$  is within the factor  $1/\Omega$  from the minimax optimal ESE-risk  $R_*$  of recovering  $g^T x$ . We first justify our claim, and then explain how to bound  $\Omega$  from below to get the desired numbers.

It clearly suffices to prove our claim for  $\sigma = 1$ . By scaling  $g$ , it suffices to show that when  $\sigma = 1$  and  $R_* < \Omega$ , then there exists an affine in  $y$  estimate with the worst-case, over  $x \in X$ , ESE-risk not exceeding 1. Thus, let  $R_* \leq \Omega$ , and let us build an affine estimate with risk 1.

**Temporary assumption.** We assume for the time being that  $A$  has a trivial kernel; eventually, this assumption will be eliminated.

**The strategy.** There are two possibilities:

(a) There exist  $\alpha, \beta \geq 0$  with  $\alpha^2 + \beta^2 = 1$  such that for certain vector  $h$  and real  $c$  we have

$$\|h\| \leq \beta \quad \text{and} \quad |g^T x - h^T Ax - c| \leq \alpha \quad \forall x \in X. \quad (*)$$

In this case we are done: the squared bias of the affine estimate  $\widehat{g}(Ax + \xi) = h^T(Ax + \xi) + c$  does not exceed  $\alpha^2$  for every  $x \in X$ , and the expectation of the squared stochastic term does not exceed  $\beta^2$ , whence with our estimate, the expected square error is  $\leq 1$  for every  $x \in X$ .

(b) (a) does not take place. All we need is to prove that in fact (b) is impossible – it contradicts the assumption that  $R_* < \Omega$ . Thus, our goal is to prove that if (a) does not take place, then

$$R_* \geq \Omega = \inf_{D>0} \frac{1+D^2}{D^2} R_*(D). \quad (70)$$

The proof goes as follows. As we know, when (a) does not take place, for every  $\alpha, \beta \geq 0$  with  $\alpha^2 + \beta^2 = 1$ , there is no  $h$  satisfying (68), whence, as we have seen, there exists  $z \in X_* = X - X$  such that

$$g^T z > 2\alpha + \beta \|Az\|.$$

Note that this inequality implies that  $z \neq 0$ , and therefore  $D \equiv \|Az\|/2 > 0$ , since  $A$  has a trivial kernel. Representing  $z = r - s$  with  $r, s \in X$  and setting  $\Delta = [s, r]$ ,  $e = \frac{r-s}{2}$ ,

$$\gamma(\theta) = g^T \left[ \frac{r+s}{2} + \frac{\theta}{D} e \right] = \text{const} + \gamma\theta, \quad -D \leq \theta \leq D,$$

observe that an estimate of  $g^T x$  via observations (1) with the ESE-risk  $\leq R$  on  $X$  induces an estimate of the value of  $\gamma(\theta)$  via a single observation of a random variable  $\zeta \sim \mathcal{N}(\theta, 1)$ , the ESE of the estimate being  $\leq R$  when  $-D \leq \theta \leq D$ . The slope  $\gamma$  of  $g(\theta)$  is  $> (\alpha + \beta D)/D$ , so that the latter estimate induces an estimate of  $\theta \in [-D, D]$  from a single observation of  $\zeta$  with the expected squared error less than  $\gamma^{-2} R$ , that is,  $R \geq \gamma^2 R_*(D) = \frac{(\alpha + \beta D)^2}{D^2} R_*(D)$ . Since this bound holds true for every estimate of  $g^T x$ , we conclude that

$$R_* \geq \frac{(\alpha + \beta D)^2}{D^2} R_*(D),$$

whence

$$R_* \geq \inf_{D>0} \frac{(\alpha + \beta D)^2}{D^2} R_*(D)$$

and thus

$$R_* \geq \sup_{\substack{\alpha, \beta > 0: \\ \alpha^2 + \beta^2 = 1}} \inf_{D>0} \frac{(\alpha + \beta D)^2}{D^2} R_*(D) \quad (71)$$

Already this bound allows to prove that the minimum of the worst-case, over  $x \in X$ , ESE of an affine estimate is within  $O(1)$ -factor of  $R_*$ . Indeed, it is immediately seen that  $R_*(D) \geq O(1) \min[1, D^2]$ , and therefore the announced claim is given by (71) where we replace the maximum over  $\alpha, \beta$  by the value of the maximized function at, say,  $\alpha = \beta = 1/\sqrt{2}$ .

In fact, our target relation (70) is nothing but a strengthened version of (71) obtained by swapping inf and sup in the latter relation. We are about to demonstrate that this swap is indeed legitimate.

**The proof.** Assume that (a) does not take place, and let us prove (70). Since (a) does not take place, for every  $\alpha \in (0, 1]$  the conic quadratic inequality

$$g^T z \geq 2\alpha + \sqrt{1 - \alpha^2} \|Az\| \quad (!)$$

has a solution in  $X_*$ . Since  $X_*$  is bounded, it follows that for every  $\epsilon > 0$  there exists  $\delta(\epsilon) > 0$  such that the convex inequality

$$g^T z \geq (2 - \epsilon)\alpha + \delta(\epsilon)z^T z + \sqrt{1 - \alpha^2} \|Az\| \quad (+_\epsilon)$$

has a solution in  $X_*$ , provided that  $\alpha \in [\epsilon, 1]$ . Let us fix  $\epsilon \in (0, 1)$ . The set  $X_\epsilon(\alpha)$  of all belonging to  $X_*$  solutions to  $(+_\epsilon)$  is a nonempty closed convex subset of  $X_*$ , and its image  $Y_\epsilon(\alpha)$  under the mapping  $x \mapsto Ax$  is a nonempty closed convex subset of  $\mathbb{R}^m$ . Let  $y_\epsilon(\alpha)$  be the (uniquely defined) minimal norm point from the set  $Y_\epsilon(\alpha)$ , let  $z_\epsilon(\alpha) \in X_\epsilon(\alpha)$  be such that  $y_\epsilon(\alpha) = Az_\epsilon(\alpha)$ , and let  $D_\epsilon(\alpha) = \frac{1}{2} \|y_\epsilon(\alpha)\|$ . Note that since  $A$  has a trivial kernel, we have  $D_\epsilon(\alpha) > 0$  for all  $\alpha \in [\epsilon, 1]$ .

We claim that  $z_\epsilon(\alpha)$  is continuous in  $\alpha \in [\epsilon, 1]$ . Indeed, let  $[\epsilon, 1] \ni \alpha_t \rightarrow \bar{\alpha}$  as  $t \rightarrow \infty$ , and let  $z_t = z_\epsilon(\alpha_t)$ ,  $\bar{z} = z_\epsilon(\bar{\alpha})$ ; we should prove that  $z_t \rightarrow \bar{z}$  when  $t \rightarrow \infty$ . Assume that this is not the case. Passing to a subsequence, we may assume that  $z_t \rightarrow \tilde{z} \neq \bar{z}$  as  $t \rightarrow \infty$ ; clearly,  $\tilde{z}$  belongs to  $X_*$  and is a solution to the inequality  $(+_\epsilon)$  associated with  $\alpha = \bar{\alpha}$ . Since we are in the situation  $\tilde{z} \neq \bar{z}$ , we have  $\|A\tilde{z}\| > \|A\bar{z}\|$ . Let  $\Delta$  be the segment  $[\bar{z}, \tilde{z}]$ , and let  $f(z) = (2 - \epsilon)\bar{\alpha} + \delta(\epsilon)z^T z + \sqrt{1 - \bar{\alpha}^2} \|Az\| - g^T z$ .  $f$  is strongly convex on  $\Delta$  and is nonpositive at the endpoints of  $\Delta$ , whence it is negative at the midpoint  $\hat{z}$  of  $\Delta$  (since  $\Delta$  is not a singleton). But then  $\hat{z}$  is a solution, belonging to  $X_*$ , of the inequalities  $(+_\epsilon)$  associated with  $\alpha = \alpha_t$  for all large enough values of  $t$ , which is a contradiction (note that  $\|A\hat{z}\| < \|A\tilde{z}\|$  and thus  $\|A\hat{z}\| < \|Az_t\|$  for all large enough values of  $t$ ).

We now claim that there exist  $\alpha_{\min}, \alpha_{\max}$ ,  $0 < \alpha_{\min} \leq \alpha_{\max} < 1$ , such that for all small enough values of  $\epsilon$ , there exists  $\alpha = \alpha_\epsilon \in [\alpha_{\min}, \alpha_{\max}]$  such that

$$\sqrt{1 - \alpha^2}/\alpha = D_\epsilon(\alpha). \quad (:)$$

Indeed, since  $X_*$  is bounded,  $D_\epsilon(\alpha) \leq C < \infty$  for all  $\epsilon \in (0, 1)$  and all  $\alpha \in [\epsilon, 1]$ ; setting  $\alpha_{\min} = 1/\sqrt{1 + C^2}$ , we conclude that when  $\alpha = \alpha_{\min}$ , the left hand side in  $(:)$  is  $\geq$  the right hand side one. On the other hand, since  $A$  is with trivial kernel, there clearly exists  $c > 0$  (which is independent of  $\epsilon$ ) such that whenever  $\epsilon \in (0, 1/2)$  and  $\alpha \in [1/2, 1]$ , we have  $D_\epsilon(\alpha) \geq c$ . It follows that setting  $\alpha_{\max} = \max[1/\sqrt{1 + c^2}, 1/2]$ , we ensure that when  $\epsilon \in (0, 1/2)$  and  $\alpha = \alpha_{\max}$ , the left hand side in  $(:)$  is  $\leq$  the right hand side one. Recalling that  $D_\epsilon(\alpha)$  is continuous in  $\alpha \in [\epsilon, 1]$  along with  $z_\epsilon(\alpha)$ , the bottom line is, that whenever  $\epsilon > 0$  is small enough, equation  $(:)$  has a solution  $\alpha_\epsilon$  in the segment  $[\alpha_{\min}, \alpha_{\max}]$ , as claimed.

Now let us choose a sequence  $\epsilon_t \rightarrow +0$  and set  $\alpha_t = \alpha_{\epsilon_t}$ ,  $z_t = z_{\epsilon_t}(\alpha_t)$ . Passing to a subsequence, we may assume that  $\alpha_t \rightarrow \bar{\alpha} \in [\alpha_{\min}, \alpha_{\max}]$  and  $z_t \rightarrow \bar{z} = r - s$ ,  $r, s \in X_*$ , as  $t \rightarrow \infty$ . Recalling that  $D_\epsilon(\alpha) = \frac{1}{2} \|z_\epsilon(\alpha)\|$ , we clearly have

$$\begin{aligned} (a) \quad & g^T(r - s) \geq 2\bar{\alpha} + \bar{\beta} \|A(r - s)\|, \quad \bar{\beta} = \sqrt{1 - \bar{\alpha}^2} \in (0, 1) \\ (b) \quad & \bar{D} \equiv \frac{1}{2} \|A(r - s)\| = \bar{\beta}/\bar{\alpha}. \end{aligned} \quad (72)$$

Same as in our preliminary reasoning, this relation implies that

$$R_* \geq \frac{(\bar{\alpha} + \bar{\beta}\bar{D})^2}{\bar{D}^2} R_*(\bar{D}) = \frac{1 + \bar{D}^2}{\bar{D}^2} R_*(\bar{D}),$$

where the concluding equality is given by (72.b). Since  $\bar{D} > 0$ , (70) follows.

**Eliminating the assumption that  $A$  has trivial kernel.** Let  $A$  be arbitrary. Consider the artificial problem where we want to recover  $g^T x$  from observations

$$[Ax + \xi; \epsilon x + \eta],$$

where  $\epsilon > 0$  is a parameter, and  $\eta \sim \mathcal{N}(0, I_n)$  is independent of  $\xi$ . The minimax risk  $R_*(\epsilon)$  of recovering  $g^T x$  via these “augmented” observations clearly is  $\leq$  the original minimax risk  $R_*$ . As about the optimal worst-case ESE-risk  $R_\epsilon$  achievable with affine estimates based on the augmented observations, we clearly have that  $\liminf_{\epsilon \rightarrow +0} R_\epsilon$  is the optimal ESE-risk  $R_0$  achievable with estimates affine in the original observations  $y$ . Indeed, if  $\hat{g}_\epsilon(y) = h_\epsilon^T y + c_\epsilon + d_\epsilon^T(\epsilon x + \eta)$  is the estimate underlying  $R_\epsilon$ , then  $\|d_\epsilon\|^2$  remains bounded (by a risk of a whatever once for ever fixed affine estimate which does not use the artificial observations). It follows that when replacing  $\hat{g}_\epsilon(y)$  with the estimate  $h_\epsilon^T y + c_\epsilon$ , the resulting increase in the risk goes to 0 as  $\epsilon \rightarrow 0$ , meaning that  $\liminf_{\epsilon \rightarrow +0} R_\epsilon = R_0$ , as claimed. Now, since the matrix  $[A; \epsilon I_n]$  underlying the augmented observations has a trivial kernel, we have  $R_\epsilon/R_*(\epsilon) \leq 1/\Omega$  by our previous results, whence  $R_\epsilon/R_* \leq 1/\Omega$  due to  $R_* \geq R_*(\epsilon)$ . Thus,  $R_0/R_* = \liminf_{\epsilon \rightarrow 0} R_\epsilon/R_* \leq 1/\Omega$ , as claimed.

**Bounding  $R_*(D)$  from below.** For every probability measure  $\mu$  on  $[-D, D]$ , we clearly have

$$R_*(D) \geq \inf_{g(\cdot)} \int \left\{ \int (g(s) - \theta)^2 p(s - \theta) ds \right\} d\mu(\theta),$$

where  $p(\cdot)$  is the standard Gaussian density on the axis. It follows that if  $\Gamma$  is a finite subset of  $[-D, D]$ , then

$$R_*(D) \geq \max_{\substack{\{\mu_\gamma \geq 0\}_{\gamma \in \Gamma} \\ \sum_\gamma \mu_\gamma = 1}} \inf_{g(\cdot)} \sum_{\gamma \in \Gamma} \mu_\gamma \int (g(s) - \gamma)^2 p(s - \gamma) ds.$$

Note that the quantity

$$\inf_{g(\cdot)} \sum_{\gamma \in \Gamma} \mu_\gamma \int (g(s) - \gamma)^2 p(s - \gamma) ds = \sum_{\gamma \in \Gamma} \int \left( \frac{\sum_{\gamma \in \Gamma} \gamma \mu_\gamma p(s - \gamma)}{\sum_{\gamma \in \Gamma} \mu_\gamma p(s - \gamma)} - \gamma \right)^2 \mu_\gamma p(s - \gamma) ds$$

can be easily computed to whatever accuracy, and that this quantity, due to its origin, is a concave function of  $\{\mu_\gamma\}_{\gamma \in \Gamma}$  and thus can be efficiently maximized numerically. We have carried out this maximization with the equidistant 31-element grid in  $[-D, D]$  for a “fine resolution” finite set  $\mathcal{D}$  of values of  $D$ , thus obtaining a valid lower bound on  $R_*(D)$  along the set and used this calculations, combined with the evident formula  $R_*(D) \geq (1 + o(1))D^2$  as  $D \rightarrow 0$  (with easily quantifiable  $o(1)$ ) to conclude that  $\Omega \geq 0.802$ , whence  $1/\Omega \leq 1.25$ .