

# A computationally motivated definition of parametric estimation and its applications to the Gaussian distribution

Leonard J. Schulman\*

Vijay V. Vazirani†

December 20, 2001

## Abstract

We introduce<sup>1</sup> a treatment of parametric estimation in which optimality of an estimator is measured *in probability* rather than in variance (the measure for which the strongest general results are known in statistics). Our motivation is that the quality of an approximation algorithm is measured by the probability that it fails to approximate the desired quantity within a set tolerance. We concentrate on the Gaussian distribution and show that the sample mean is the unique “best” estimator, in probability, for the mean of a Gaussian distribution. We also extend this method to general penalty functions and to multidimensional spherically symmetric Gaussians.

The algorithmic significance of studying the Gaussian distribution is established by showing that determining the average matching size in a graph is  $\#\mathbf{P}$ -complete, and moreover approximating it reduces to estimating the mean of a random variable that (under some mild conditions) has a distribution closely approximating a Gaussian. This random variable is (essentially) polynomial time samplable, thereby yielding an FPRAS for the problem.

## 1 Introduction

The task of estimating a parameter via sampling lies at the heart of numerous algorithms. In particular, this task is central to approximation algorithms for  $\#\mathbf{P}$ -complete problems. These algorithms rely on an *in probability* statement: establishing that the parameter in question has been estimated within certain bounds with “high” probability. (Inverse polynomial probability suffices for polynomial time computability.) Despite widespread use of this method, so far questions about the *optimality* of the estimator have not been studied.

Such results belong in the area of parametric estimation within the field of statistics. Traditionally, the notion of optimality most studied in this area is minimization of the mean square error of the estimator. Indeed, among the celebrated theorems of statistics is the Cramer-Rao lower bound on the mean square error of an unbiased estimator of a parameter  $\theta$ , in terms of its Fisher entropy. A key application of this theorem is to show that the sample mean is an optimal unbiased estimator of the mean of a Gaussian from a variable-location, fixed-scale (unit variance) family  $\{G_\theta\}$  where  $G_\theta(x) = (2\pi)^{-1/2} \exp(-(x - \theta)^2/2)$  [CT91, Zac81, Fre43, Rao45, Cra46]. Notice, however, that optimality of an estimator in mean square error does not imply optimality in probability. Furthermore, establishing a bound on the mean square error of an estimator is not sufficient for the purposes of deriving an approximation algorithm, for instance for giving a fully polynomial randomized approximation scheme, FPRAS, for a  $\#\mathbf{P}$ -complete problem.

The traditional worst case analysis of algorithms has been particularly successful in unraveling algorithmically relevant combinatorial structure in problems, and in designing powerful algorithmic tools to exploit this structure. We adopt this paradigm in our criteria for parametric estimation.

The current paper is an attempt at initiating a theory of parametric estimation in which optimality is measured in probability in the worst case over the possible values of the parameter. Although our original

---

\*Caltech, Pasadena CA 91125, schulman@cs.caltech.edu

†College of Computing, Georgia Inst. Technology, Atlanta GA 30332-0280, vazirani@cc.gatech.edu

<sup>1</sup>This paper is based on results in [SV99] and [SV01].

motivation is algorithmic, we believe that this theory will find value in other areas as well. We derive results for the Gaussian distribution. The reason for concentrating on this distribution is twofold. First, the Gaussian distribution lends itself to a very precise analysis. Second, this distribution arises naturally in several computational situations. Here is a particularly striking case.

Consider the problem of computing the average matching size in a given graph. In Theorem 5 we show that exact computation of this parameter is #P-complete. Now, consider the random variable that is the size of a random, uniformly chosen, matching in  $G$ . This random variable is (essentially) polynomial time samplable, since there exists an almost uniform generator for matchings in a graph [JS95]. Therefore, it can be used to estimate the average size of a matching in  $G$ ; in fact it even leads to an FPRAS for this problem. The algorithm is straightforward: sample this random variable an appropriate number of times, depending on the error parameter, and output the mean of these samples.

Under some mild conditions, this random variable has essentially a Gaussian distribution. This follows from an exceptionally strong result of Godsil describing the size distribution of the matchings of a graph [LP86, God81]. Hence, the FPRAS stated above is simply estimating the mean of this Gaussian distribution! Moreover, it is using the mean estimator to do so. We are interested in whether this is the best estimator for the mean of a Gaussian distribution, where “best” is, as is customary for algorithmic analysis, defined in terms of the *probability* of failing to estimate the mean within the desired accuracy on a *worst-case* input. In Section 2 we define this question more precisely, and in Theorem 2 we answer it in the affirmative.

## 2 The model and our results

Consider a probability density  $f$  on the real line with first moment 0 and finite second moment. Form the family of densities  $\{f_\theta\}$ , which are the translations of  $f$ , indexed by their means  $\theta$ . (So  $f_0 = f$ .)

Now,  $\theta$  is fixed and unknown to us, and we collect  $n$  samples  $x_1, \dots, x_n$  from the density  $f_\theta$ . We wish to infer an estimate of the parameter  $\theta$ . For each  $\varepsilon > 0$ , we are interested in the probability that our estimator  $S(x_1, \dots, x_n)$  falls within distance  $\varepsilon$  of  $\theta$ . Furthermore, we are interested in the *worst case* (over  $\theta$ ) performance of  $S$ . For this purpose, let us define the  $\varepsilon$ -*quality* of estimator  $S$  to be

$$Q_S^\varepsilon = \inf_\theta [P(|S - \theta| \leq \varepsilon)].$$

**Definition 1** We say that estimator  $T$  majorizes  $S$  if for all  $\varepsilon > 0$ ,  $Q_T^\varepsilon \geq Q_S^\varepsilon$ .

**Theorem 2** For the family  $\{G_\theta\}$ , for any given  $n$ , the mean estimator,  $T(x_1, \dots, x_n) = \frac{1}{n} \sum x_i$ , majorizes every other estimator.

In Theorem 14 we will further establish that  $T$  is the *unique* majorizing estimator.

Let  $X = (x_1, \dots, x_n)$  denote  $n$  independent samples picked from  $G_\theta$ . Let  $\psi : \mathbb{R} \rightarrow \mathbb{R}^+$  (where  $\mathbb{R}^+$  is the nonnegative reals) be a *penalty function* satisfying the following conditions: (a)  $\psi$  is symmetric, (b)  $\psi(x)$  is nondecreasing in  $|x|$ , (c)  $\psi$  is not a constant function, and (d)  $L = \int_{-\infty}^{\infty} \psi(x)G(x)dx < \infty$ .

Without loss of generality we may assume that  $\psi(0) = 0$ . Note also that conditions (a) and (b) imply that  $\psi$  is measurable.

Define the  $\psi^{th}$  central moment of estimator  $S$  at  $\theta$  to be

$$M_\theta^\psi(S) = \int P(X|\theta) \int_{t \in \mathbb{R}} \psi(t - \theta) P(S(X) = t) dt dX.$$

In case  $\psi(x) = x^r$ , this is simply the  $r$ 'th central moment of the estimator  $S$  at  $\theta$ .

be

By an extension of the method of Theorem 2, we show the more general:

**Theorem 3** For the family  $\{G_\theta\}$ , for every  $n$  and every penalty function  $\psi$ , the mean estimator minimizes  $\sup_\theta M_\theta^\psi(S)$  among all estimators  $S$ .

We next extend the method to higher dimensions. Let  $X = (x_1, \dots, x_n)$  denote  $n$  independent samples in  $\mathbb{R}^d$  picked from  $G_\theta^d$ , the spherically symmetric Gaussian distribution

$$G_\theta^d(z) = (2\pi)^{-d/2} \exp\left(-\frac{1}{2} \sum_1^d (z(i) - \theta(i))^2\right).$$

Let  $\psi : \mathbb{R}^d \rightarrow \mathbb{R}^+$  (where  $\mathbb{R}^+$  is the nonnegative reals) be a *penalty function* satisfying the following conditions: (a)  $\psi$  is spherically symmetric:  $\psi(x) = \psi(y)$  if  $|x| = |y|$ . (b)  $\psi(x)$  is nondecreasing in the Euclidean norm  $|x|$ . (c)  $\psi$  is not a constant function. (d)  $L = \int_{\mathbb{R}^d} \psi(x) G_0^d(x) dx < \infty$ .

These assumptions imply that  $\psi$  is “unimodal on lines,” meaning that for every set of real parameters  $a_1, \dots, a_d, b_1, \dots, b_d$ , the function  $\psi(a_1 t + b_1, \dots, a_d t + b_d)$  is a unimodal function of the real parameter  $t$ . The central moment of an estimator  $S$ ,  $M_\theta^\psi(S)$ , is defined in  $\mathbb{R}^d$  analogously to  $\mathbb{R}$ . In section 6 we show:

**Theorem 4** *For the family  $\{G_\theta^d\}$ , for every  $n$  and every penalty function  $\psi$ , the mean estimator minimizes  $\sup_\theta M_\theta^\psi(S)$  among all estimators  $S$ .*

## Cramer-Rao

Among the celebrated theorems of statistics is the Cramer-Rao lower bound (due independently to Cramer, Rao and Frechet) on the mean squared error of an unbiased estimator of a parameter  $\theta$ . A key application of that theorem is to show that the sample mean is an optimal unbiased estimator of the mean of a Gaussian from the family  $\{G_\theta\}$  [CT91, Zac81, Fre43, Rao45, Cra46]. Theorem 3 represents an improvement in the sense in which the mean estimator for the Gaussian is shown to be optimal, as it implies optimality of the mean estimator in mean square (variation), as indeed in any central moment, among all (not only unbiased) estimators. However it says this only about the worst-case  $\theta$  while Cramer-Rao enables statements about each  $\theta$ .

## Confidence Intervals

Motivated by the algorithmic applications, we have chosen to measure the quality of an estimator  $T$  of a parameter  $\theta$  by the function  $\inf_\theta [P(|T - \theta| \leq \varepsilon)]$ . A somewhat “dual” notion is studied in the statistical literature. A *confidence interval of level  $p$*  is a pair of estimators  $T_1, T_2$  s.t. for every  $\theta$ , with probability at least  $p$ ,  $T_1 \leq \theta \leq T_2$ . Obviously it is desirable that the intervals  $[T_1, T_2]$  be as short as possible subject to the confidence level  $p$ ; this objective is complementary to our goal of maximizing the estimator’s probability of falling within a fixed width interval,  $\inf_\theta [P(|T - \theta| \leq \varepsilon)]$ .

However, in the case of confidence intervals, there is an additional degree of freedom available in “sliding” both ends of the interval without changing the confidence level. While this flexibility is desirable for some applications (e.g. if the penalties for errors in the two directions are unequal), it reduces the extent to which the quality of estimators can be compared. In particular, there does not exist any family of densities  $\{f_\theta\}$ , and any  $0 < p < 1$ , for which there is an optimal estimator (in the sense that its confidence intervals are contained within those of any other estimator). (And a statement nearly as strong can be made also for families of distributions which do not arise from densities.) To see this, one has only to consider the two optimal estimators subject to the restrictions that the lower or upper endpoints are at  $-\infty$  or  $+\infty$ . Estimators that are optimal subject to these restrictions are termed “uniformly most accurate upper/lower (respectively) confidence limits”; this appears to be the closest definition in the literature to our notion of a majorizing estimator. However, as just implied, no statement resembling theorem 2 can exist for upper and lower confidence limits. Thus one of the contributions of this paper is the introduction of  $Q_S^\varepsilon$  as a measure of the quality of an estimator  $S$ , since this defines a partial order on estimators that is on the one hand, more refined than that defined by commonly used criteria such as mean squared error; and on the other hand, the resulting partial order on estimators is not so weak as to preclude the existence of a greatest element in the partial order.

### 3 Estimating average matching size in a graph

A *counting problem*,  $\Pi$ , consists of:

- A set of *instances*,  $D_\Pi$ .
- The *size* of instance  $I \in D_\Pi$ , denoted by  $|I|$ , is defined as the number of bits needed to write  $I$  under the assumption that all numbers occurring in the instance are written in binary.
- A *solution space*  $S_I$ , typically of size exponential in  $|I|$ , is associated with instance  $I$ . A parameter of  $S_I$ ,  $\theta_I$ , is defined.

For the problem of interest, an instance is an undirected graph,  $G$ , and the solution space is the set of matchings (of all sizes) in  $G$ . The parameter of interest is the average size of a matching in  $G$ . Let us denote it by  $\mu(G)$ .

The interesting case is when  $\Pi \in \mathbf{P}$ , and when computing  $\theta_I$  as a function of  $I$  is complete for the counting class  $\#\mathbf{P}$  introduced by Valiant in [Val79a]. The problem of finding a matching, even a maximum matching, in  $G$  is polynomial time solvable. Below we show that the problem of computing  $\mu(G)$  exactly is  $\#\mathbf{P}$ -complete.

**Theorem 5** *The problem of computing  $\mu(G)$  is  $\#\mathbf{P}$ -complete.*

Note that  $\mu(G)$  is the ratio  $p/q$  of two integers each bounded by  $2^{|E(G)|}$  ( $E$  being the edge set of  $G$ ). Computation of  $\mu(G)$  can be understood to mean either of two things: (a) computation of a pair  $p', q'$  such that  $p'/q' = p/q$ . (b) Computation of any number  $r$  such that  $|r - p/q| < 2^{-2^{|E(G)|-1}}$ . By the theory of continued fractions,  $p/q$  is the unique rational with denominator bounded by  $2^{|E(G)|}$  within this distance of  $r$ . Hence a pair  $p', q'$  as in (a) can be computed from  $r$  (simply by computing enough of its continued fraction expansion).

**Proof:** We reduce from the problem of computing  $\phi(G)$ , the number of matchings in  $G$ , demonstrated  $\#\mathbf{P}$ -complete by Valiant [Val79b].

Consider any edge  $uv$  of  $G$ . Note that

$$\phi(G) = \phi(G - u - v) + \phi(G - uv).$$

Denote by  $s(G)$  the sum of sizes of all matchings in  $G$ ; clearly,  $s(G) = \mu(G)\phi(G)$ . Observe that among matchings of  $G$  which do not use  $uv$ , the average matching size is  $\mu(G - uv) = s(G - uv)/\phi(G - uv)$ . Among matchings of  $G$  which do use  $uv$ , the average matching size is one more than in  $G - u - v$ , hence  $(s(G - u - v) + \phi(G - u - v))/\phi(G - u - v)$ . Considering the matchings of  $G$  according to whether they contain  $uv$ , we see that

$$\mu(G) = \frac{s(G - u - v) + \phi(G - u - v) + s(G - uv)}{\phi(G - u - v) + \phi(G - uv)}. \quad (1)$$

The following is our polynomial time Turing reduction of the computation of  $\phi(G)$  to the computation of  $\mu(G)$ :

*Algorithm which computes  $\phi(G)$  on input  $G$ :*

Pick any edge  $uv$ .

Compute  $\mu(G - uv)$ ,  $\mu(G - u - v)$  and  $\mu(G)$ .

Recursively compute  $\phi(G - uv)$ .

Set  $s(G - uv) := \mu(G - uv)\phi(G - uv)$ .

Set  $\phi(G - u - v) := \frac{\phi(G - uv)(\mu(G) - \mu(G - uv))}{1 + \mu(G - u - v) - \mu(G)}$  which is a consequence of equation 1.

Output  $\phi(G) := \phi(G - uv) + \phi(G - u - v)$ . □



We will say that an algorithm  $A$  is a fully polynomial randomized approximation scheme (FPRAS) for computing  $\theta_I$  if for each instance  $I$  and error parameter  $\varepsilon > 0$ ,

$$P(|A(I) - \theta_I| \leq \varepsilon \theta_I) \geq \frac{3}{4},$$

and the running time of  $A$  is polynomially bounded in  $|I|$  and  $\frac{1}{\varepsilon}$ . Once this is achieved, the probability of success can be amplified using the “median trick”: Run algorithm  $A$  a number of times and output the median answer. It is easy to show that to achieve a probability of success of  $1 - \delta$  it suffices to run  $A$   $O(\log(1/\delta))$  times.

Typically, an FPRAS for computing  $\theta_I$  is constructed as follows: A polynomial time samplable probability distribution is defined on  $S_I$ , together with a random variable  $X_I$ , which is shown to be an unbiased estimator of  $\theta_I$ , or nearly so, the error being  $< \varepsilon$ . A specified number of sample points are picked from the probability distribution, the random variable is computed at these points, and the mean of these values is output. It is a consequence of work of Canetti, Even and Goldreich [CEG95], and also implied by the present paper, that generally there is not much more that one can do: Specifically, if all we know about the random variable  $X_I$  is its standard deviation  $\sigma(X_I)$ , then a necessary and sufficient condition for the existence of a FPRAS for  $E(X_I)$  is that  $\sigma(X_I)/E(X_I) \leq p(|I|)$  for some polynomial  $p$ .

Let  $X$  be the random variable that on random, uniformly chosen, matching in  $G$  is the size of the matching. Clearly,  $E(X) = \mu(G)$ , and therefore, estimating  $\mu(G)$  amounts to estimating the mean of  $X$ . Jerrum and Sinclair [JS95] use the Markov chain Monte Carlo method to give an almost uniform generator for matchings in a graph, thereby showing that a random variable having probability distribution arbitrarily close to that of  $X$  is polynomial time samplable. As shown in Theorem 7 this yields an FPRAS for estimating  $\mu(G)$ .

Under some mild conditions, random variable  $X$  has essentially a Gaussian distribution. This follows from an exceptionally strong result of Godsil describing the size distribution of the matchings of a graph [LP86, God81]. Let  $G_1, \dots$  be a family of graphs. Let  $\phi_k(G_n)$  be the number of matchings with  $k$  edges in  $G_n$  and let  $\phi(G_n) = \sum_k \phi_k(G_n)$  be the total number of matchings of  $G_n$ . Let  $\mu(G_n)$  be the average size of a matching of  $G_n$ ,  $\mu(G_n) = (\sum_k k \phi_k(G_n)) / \phi(G_n)$ . Let  $\sigma^2(G_n)$  be the variance of the distribution of sizes of matchings of  $G_n$ ,  $\sigma^2(G_n) = (\sum_k (k - \mu(G_n))^2 \phi_k(G_n)) / \phi(G_n)$ . Suppose that  $\sigma(G_n) \rightarrow \infty$ . Then:

**Theorem 6 (Godsil)** *The distribution of matching sizes is asymptotically locally normal, meaning that if we fix any real  $x$  and let  $n \rightarrow \infty$ , then*

$$\frac{\phi_k(G_n)}{\phi(G_n)} \sigma(G_n) \rightarrow (2\pi)^{-1/2} e^{-x^2/2}$$

for  $k$  such that  $k \sim \mu(G_n) + x\sigma(G_n)$ .

**Theorem 7** *There exists an FPRAS for estimating  $\mu(G)$ .*

**Proof:** Let  $n$  denote the number of vertices in the given graph  $G$ . First consider the random variable  $X$  defined above.  $X$  has polynomially bounded standard deviation, since it takes values only in the polynomial range:  $\{1, 2, \dots, n/2\}$ . Therefore, it suffices to sample it polynomially many times, in  $n$  and  $1/\varepsilon$ , and output the mean value. We formalize this first, and then deal with the fact that whereas we do not know how to efficiently sample  $X$  itself, we do know how to sample a close, in variation distance, random variable.

Let  $X_k$  denote the mean of  $k$  samples of  $X$ . Using Chebyshev’s inequality it is easy to see that for  $k = n^2/\varepsilon^2$ ,

$$P[|X_k - \mu(G)| \geq \varepsilon \mu(G)] \leq \frac{1}{4}.$$

As before, let  $\phi(G)$  denote the number of matchings in  $G$ . An *almost uniform generator* for matchings is a randomized polynomial time algorithm  $\mathcal{A}$  that for any  $\delta > 0$  and graph  $G$  outputs a matching satisfying: for each matching  $M$  in  $G$ ,

$$P[\mathcal{A} \text{ outputs } M] \in [(1 - \delta) \frac{1}{\phi(G)}, (1 + \delta) \frac{1}{\phi(G)}].$$

Furthermore, the running time of  $\mathcal{A}$  is polynomial in  $n$  and  $\log(1/\delta)$ .

Let  $Y$  be the random variable that is the size of the matching generated by the almost uniform generator of Jerrum and Sinclair [JS95]. Observe that the error parameter  $\delta$  can be made inverse exponential in polynomial time, therefore giving

$$|E(Y) - \mu(G)| \leq \varepsilon \mu(G).$$

As before, an FPRAS follows by sampling  $Y$  polynomially many times and outputting the mean.  $\square$

**Remark 8** Typically, FPRAS's for  $\#\mathbf{P}$ -complete problems, obtained using the Markov chain Monte Carlo method, reduce the approximate counting problem to random generation using self-reducibility of the problem [JV86]. An interesting feature of the FPRAS derived above is that we did not need to resort to self-reducibility.

## 4 The mean is a majorizing estimator for the family $\{G_\theta\}$

Let  $T$  be the mean estimator. Clearly, an arbitrary estimator  $S$  may be able to do better than  $T$  on certain specific values of  $\theta$ . We wish to show that even so, in the worst case,  $T$  must be doing at least as well as  $S$ . An important observation is that  $T$  commutes with translation, i.e.,  $T[X + a] = T[X] + a$ , where  $X + a$  denotes the  $n$  samples  $(x_1 + a, x_2 + a, \dots, x_n + a)$ . Therefore, its probability of falling within an  $\varepsilon$  distance of  $\theta$ ,  $P(|T - \theta| \leq \varepsilon)$ , is independent of  $\theta$ .

Thus, the worst case performance of  $T$  is the same as its performance at any  $\theta$ . The worst case performance of a general estimator  $S$ , however, is difficult to characterize. Instead, we will show that in the limit, the average performance of  $T$  over a large range of  $\theta$ 's must be at least as good as that of  $S$ . This will lead to the majorization result.

### Proof of theorem 2:

We begin with a fact that substantially simplifies the matter.

**Lemma 9** It suffices to consider the single-sample case.

**Proof:** This is because: (a) The mean of several iid Gaussian random variables is also a Gaussian random variable. (b) The mean is a *sufficient statistic* for samples drawn from the family  $\{G_\theta\}$ . This means that for every  $\theta$ , the samples  $x_1, \dots, x_n$  drawn from  $G_\theta$  are independent of  $\theta$  given  $\bar{x} = \frac{1}{n} \sum x_i$ , or in other words that there is a distribution  $P((x_1, \dots, x_n) | \bar{x})$  such that

$$P((x_1, \dots, x_n) | \theta) = P((x_1, \dots, x_n) | \bar{x}) P(\bar{x} | \theta).$$

Consequently the performance of any estimator will be unchanged if, given  $x_1, \dots, x_n$ , we first compute the mean  $\bar{x} = \frac{1}{n} \sum x_i$ , then choose a list of differences  $(x'_i - \bar{x})_1^n$  from the same distribution as for the Gaussian (note in particular that the distribution is supported only on lists whose sum is 0), then supply the estimator with the list  $x'_1, \dots, x'_n$ . The distribution of the lists produced this way is the same as that of the lists  $x_1, \dots, x_n$ , whence the conclusion that the performance is unaffected. Now, the process of substitution followed by application of the estimator, may be viewed jointly as a (randomized) estimator that takes as its input only the mean  $\bar{x}$ .  $\square$

Now consider the following process: Let  $\varepsilon > 0$  be fixed. For fixed  $\alpha > 0$ ,  $\theta$  is picked uniformly at random from the interval  $I_\alpha = [-\alpha, \alpha]$ , and then a sample  $x$  is picked from the distribution  $G_\theta$ . (We will call this the finite- $\alpha$  experiment.)

Let  $S : \mathbb{R} \rightarrow \mathbb{R}$  be an estimator of  $\theta$ . In general,  $S$  may be randomized;  $P(S(x) = y)$  denotes the probability (density) with which the estimator  $S$  outputs  $y$  on input  $x$ . Let  $\varepsilon > 0$  be fixed. We will say that  $S$  succeeds if  $\theta \in [S(x) - \varepsilon, S(x) + \varepsilon]$ . The probability of success of  $S$  over the entire finite- $\alpha$  experiment is given by

$$\int \int_{S(x) - \varepsilon}^{S(x) + \varepsilon} P(\theta) P(x | \theta) d\theta dx$$

if  $S$  is deterministic, and by

$$\int \int_{-\infty}^{\infty} P(S(x) = y) \int_{y-\varepsilon}^{y+\varepsilon} P(\theta)P(x|\theta)d\theta dy dx$$

if  $S$  is randomized.

In the single-sample case, the mean estimator is simply the identity estimator  $T(x) = x$ . For  $\alpha > \varepsilon$  let  $I'_\alpha$  denote the interval  $[-(\alpha - \varepsilon), (\alpha - \varepsilon)]$ .

**Lemma 10** For  $x \in I'_\alpha$ ,

$$\int_{-\infty}^{\infty} P(T(x) = y) \int_{y-\varepsilon}^{y+\varepsilon} P(\theta)G_\theta(x)d\theta dy$$

is uniquely maximized for the identity estimator.

**Proof:**

$$\int_{y-\varepsilon}^{y+\varepsilon} P(\theta)G_\theta(x)d\theta$$

is uniquely maximized at  $y = x$ . The lemma follows.  $\square$

For an estimator  $S$ , let  $P_S^{\alpha, \varepsilon}$  denote the probability of success of  $S$  in the finite- $\alpha$  experiment. Since the identity estimator commutes with translation, we find:

**Observation 11**  $Q_T^\varepsilon = P_T^{\alpha, \varepsilon}$ .

Let  $M(\alpha, \varepsilon)$  denote the supremum over all estimators  $S$  of  $P_S^{\alpha, \varepsilon}$ . Let  $B(\alpha, \varepsilon)$  be the event that, after picking  $\theta$  at random from  $I_\alpha$  and  $x$  at random using the distribution  $G_\theta$ ,  $x \notin I'_\alpha$ . By Lemma 10, we get:

**Corollary 12**

$$P_T^{\alpha, \varepsilon} \geq M(\alpha, \varepsilon) - P(B(\alpha, \varepsilon)).$$

$\square$

Finally, let  $Q(\varepsilon) = \sup_S Q_S^\varepsilon$ . We wish to show that  $Q_T^\varepsilon = Q(\varepsilon)$ , and thus prove the theorem. By Observation 11 and Corollary 12,

$$Q_T^\varepsilon = P_T^{\alpha, \varepsilon} \geq \liminf_{\alpha} M(\alpha, \varepsilon) - \limsup_{\alpha} P(B(\alpha, \varepsilon)).$$

Since any estimator can be employed without modification in the finite- $\alpha$  experiment,  $M(\alpha, \varepsilon) \geq Q(\varepsilon)$ . Therefore,

$$Q_T^\varepsilon \geq Q(\varepsilon) - \limsup_{\alpha} P(B(\alpha, \varepsilon)).$$

Now,

$$\begin{aligned} & \limsup_{\alpha} P(B(\alpha, \varepsilon)) \\ & \leq \limsup_{\alpha} [P(|\theta| > \alpha - \alpha^{1/2}) + P(x \notin I'_\alpha \mid |\theta| \leq \alpha - \alpha^{1/2})] \\ & \leq 0 + \limsup_{\alpha} P(|x - \theta| > \alpha^{1/2} - \varepsilon). \end{aligned}$$

Since  $x$  is normally distributed with variance  $1/n$ , this is bounded above by

$$\limsup_{\alpha} \exp(-n(\alpha^{1/2} - \varepsilon)^2/2) = 0.$$

Hence  $Q_T^\varepsilon \geq Q(\varepsilon)$ .  $\square$

## 5 Optimality of the mean estimator for $\{G_\theta\}$ with respect to general penalty functions

### Proof of theorem 3:

The proof of theorem 3 is similar to that of theorem 2. Instead of providing an upper bound on a “benefit function” (e.g. 1 if  $|T(x) - \theta| < \varepsilon$  and 0 otherwise), we provide a lower bound on a penalty function (e.g.  $|T(x) - \theta|^r$ ). Again as in theorem 2, because the mean  $\bar{x} = \frac{1}{n} \sum x_i$  is a sufficient statistic for the parameter, we may assume that an estimator is a function only of the mean; and because the mean has a Gaussian distribution, it suffices to show that the identity estimator performs at least as well (in supremum over all  $\theta$  of expected penalty) as any other estimator.

Recall that we have defined  $L = \int_{-\infty}^{\infty} \psi(x)G(x)dx$  and that  $L < \infty$ . Note that for any  $\theta$ , the expected penalty of the identity estimator is  $L$ .

Now fix any  $\varepsilon > 0$ . Select  $\alpha$  large enough so that the following condition is satisfied:

$$\int_{\varepsilon\alpha}^{\infty} G(\theta)\psi(\theta)d\theta < L\varepsilon/2.$$

Let  $S$  be an arbitrary estimator. Pick  $\theta$  uniformly in the interval  $[-\alpha, \alpha]$ . The expected penalty of  $S$  is

$$\frac{1}{2\alpha} \int_{-\alpha}^{\alpha} \int_{-\infty}^{\infty} G(x - \theta)\psi(S(x) - \theta)dx d\theta = \frac{1}{2\alpha} \int_{-\infty}^{\infty} \int_{-\alpha}^{\alpha} G(x - \theta)\psi(S(x) - \theta)d\theta dx$$

Now we ignore  $x$ 's outside the range  $[-\alpha, \alpha]$ , and treat negative and positive  $x$ 's separately. For each  $x$  we consider the penalty due only to a limited range of  $\theta$ 's.

$$\dots \geq \frac{1}{2\alpha} \int_{-\alpha}^0 \int_{-\alpha}^{2x+\alpha} G(x - \theta)\psi(S(x) - \theta)d\theta dx + \frac{1}{2\alpha} \int_0^{\alpha} \int_{2x-\alpha}^{\alpha} G(x - \theta)\psi(S(x) - \theta)d\theta dx$$

Now we apply unimodality of  $G$  and  $\psi$  to conclude that

$$\geq \frac{1}{2\alpha} \int_{-\alpha}^0 \int_{-\alpha}^{2x+\alpha} G(x - \theta)\psi(x - \theta)d\theta dx + \frac{1}{2\alpha} \int_0^{\alpha} \int_{2x-\alpha}^{\alpha} G(x - \theta)\psi(x - \theta)d\theta dx$$

Since these two expressions are identical,

$$= \frac{1}{\alpha} \int_0^{\alpha} \int_{2x-\alpha}^{\alpha} G(x - \theta)\psi(x - \theta)d\theta dx$$

Now make the change of variables  $z = \alpha - x$ .

$$= \frac{1}{\alpha} \int_0^{\alpha} \int_{-z}^z G(\theta)\psi(\theta)d\theta dz \geq \frac{1}{\alpha} \int_{\varepsilon\alpha}^{\alpha} \int_{-\varepsilon\alpha}^{\varepsilon\alpha} G(\theta)\psi(\theta)d\theta dz = (1 - \varepsilon) \int_{-\varepsilon\alpha}^{\varepsilon\alpha} G(\theta)\psi(\theta)d\theta$$

Now apply the assumption on  $\alpha$

$$\geq (1 - \varepsilon)^2 L.$$

The supremum penalty of  $S$  over  $\theta$ 's in the interval  $[-\alpha, \alpha]$  is therefore at least  $L(1 - \varepsilon)^2$ ; since  $\varepsilon$  was arbitrary, this means that the supremum of  $S$  over  $\theta \in \mathbb{R}$  is at least  $L$ , and therefore at least as great as the supremum penalty of the identity estimator.  $\square$

## 6 Optimality of the mean estimator for $\{G_\theta^d\}$ with respect to general penalty functions

### Proof of theorem 4:

The mean  $\frac{1}{n} \sum x_i$  is, just as in one dimension, a sufficient statistic for the parameter  $\theta$ , so we may assume that an estimator is a function only of the mean; and because the mean has a spherically symmetric Gaussian distribution, it suffices to show that the identity estimator performs at least as well (in supremum over all  $\theta$  of expected penalty) as any other estimator.

For  $x \in \mathbb{R}^d$  and  $r \in \mathbb{R}$  let  $b(x, r)$  be the open ball of radius  $r$  about  $x$ ; if  $r < 0$   $b(x, r)$  is empty.

Recall that we have defined  $L = \int_{\mathbb{R}^d} \psi(x)G(x)dx$  and that  $L < \infty$ . Note that for any  $\theta$ , the expected penalty of the identity estimator is  $L$ .

Now fix any  $\varepsilon > 0$ . Select  $\alpha$  large enough so that the following condition is satisfied:

$$\int_{\mathbb{R}^d - b(0, \varepsilon \alpha)} G(\theta)\psi(\theta)d\theta < L\varepsilon.$$

Let  $S$  be an arbitrary estimator. Pick  $\theta$  uniformly in the ball  $b(0, \alpha)$ . The expected penalty of  $S$  is (with  $c_d r^d = \int_{b(0, r)} 1d\theta$ )

$$\frac{1}{c_d \alpha^d} \int_{b(0, \alpha)} \int_{\mathbb{R}^d} G(x - \theta)\psi(S(x) - \theta)dx d\theta = \frac{1}{c_d \alpha^d} \int_{\mathbb{R}^d} \int_{b(0, \alpha)} G(x - \theta)\psi(S(x) - \theta)d\theta dx$$

Now we ignore  $x$ 's outside the ball  $b(0, \alpha)$ . For each  $x$  we consider the penalty due only to a limited range of  $\theta$ 's.

$$\dots \geq \frac{1}{c_d \alpha^d} \int_{b(0, \alpha)} \int_{b(x, \alpha - |x|)} G(x - \theta)\psi(S(x) - \theta)d\theta dx$$

Let  $u = \theta - x$ .

$$\dots = \frac{1}{c_d \alpha^d} \int_{b(0, \alpha)} \int_{b(0, \alpha - |x|)} G(u)\psi(S(x) - x - u)dudx$$

We show below that unimodality (and symmetry) of  $G$  and  $\psi$ , applied to the inner integral, imply  $\int_{b(0, \alpha - |x|)} G(u)\psi(S(x) - x - u)du \geq \int_{b(0, \alpha - |x|)} G(u)\psi(u)du$ . Thus

$$\dots \geq \frac{1}{c_d \alpha^d} \int_{b(0, \alpha)} \int_{b(0, \alpha - |x|)} G(u)\psi(u)dudx \geq \frac{1}{c_d \alpha^d} \int_{b(0, (1-\varepsilon)\alpha)} \int_{b(0, \varepsilon\alpha)} G(u)\psi(u)dudx$$

And, applying the assumption on  $\alpha$ :

$$\dots > \frac{1}{c_d \alpha^d} \int_{b(0, (1-\varepsilon)\alpha)} (1-\varepsilon)L dx = \frac{c_d((1-\varepsilon)\alpha)^d(1-\varepsilon)L}{c_d \alpha^d} = (1-\varepsilon)^{d+1}L.$$

The supremum penalty of  $S$  over  $\theta$ 's in the ball  $b(0, \alpha)$  is therefore at least  $L(1-\varepsilon)^{d+1}$ ; since  $\varepsilon$  was arbitrary, this means that the supremum of  $S$  over  $\theta \in \mathbb{R}$  is at least  $L$ , and therefore at least as great as the supremum penalty of the identity estimator.  $\square$

It remains to show that:

**Lemma 13**  $\int_{b(0, \alpha - |x|)} G(u)\psi(S(x) - x - u)du \geq \int_{b(0, \alpha - |x|)} G(u)\psi(u)du$ .

**Proof:** Applying symmetry of  $\psi$ , letting  $\chi = S(x) - x$ , and introducing a real parameter  $s$ , it suffices to show that

$$F(s) = \int_{b(0, \alpha - |x|)} G(u)\psi(u - s\chi)du$$

is minimized at  $s = 0$ . For  $u \in \mathbb{R}^d$  and  $t > 0$ , let  $\delta_{u,t} = 1$  if  $G(u) \geq t$ , and  $\delta_{u,t} = 0$  otherwise. Let  $\eta(t) = \sup\{|u| : \delta_{u,t} = 1\}$ ; if the set is empty let  $\eta(t) = 0$ . Now

$$\begin{aligned} F(s) &= \int_{b(0, \alpha - |x|)} G(u)\psi(u - s\chi)du = \int_{b(0, \alpha - |x|)} \int_0^\infty \delta_{u,t}\psi(u - s\chi)dt du \\ &= \int_0^\infty \int_{b(0, \alpha - |x|)} \delta_{u,t}\psi(u - s\chi)dudt = \int_0^\infty \int_{b(0, \min\{\alpha - |x|, \eta(t)\})} \psi(u - s\chi)dudt \end{aligned}$$

Let  $\chi^\perp = \{v \in \mathbb{R}^d : v \cdot \chi = 0\}$ . For  $v \in \chi^\perp$  let  $k(x, t, v) = \{a \in \mathbb{R} : v + a\chi \in b(0, \min\{\alpha - |x|, \eta(t)\})\}$ . Observe that  $k(x, t, v)$  is a (possibly empty) interval that is symmetric about the origin. Now

$$F(s) = |\chi| \int_0^\infty \int_{\chi^\perp} \int_{k(x,t,v)} \psi(v + (a-s)\chi) da dv dt$$

Define  $\psi_v : \mathbb{R} \rightarrow \mathbb{R}^+$  by  $\psi_v(y) = \psi(v + y\chi)$ . By the assumptions on  $\psi$ ,  $\psi_v$  is unimodal and symmetric. We have

$$F(s) = |\chi| \int_0^\infty \int_{\chi^\perp} \int_{k(x,t,v)} \psi_v(a-s) da dv dt$$

Let  $K(x, t) \subseteq \chi^\perp$  be the ball for which  $k(x, t, v)$  is nonempty, and for  $v \in K(x, t)$  let  $\bar{k}(x, t, v) = \sup\{a : a \in k(x, t, v)\}$ . Then

$$F(s) = |\chi| \int_0^\infty \int_{K(x,t)} \int_{-\bar{k}(x,t,v)}^{\bar{k}(x,t,v)} \psi_v(a-s) da dv dt$$

We differentiate with respect to  $s$ .

$$\begin{aligned} \frac{d}{ds} F(s) &= |\chi| \frac{d}{ds} \int_0^\infty \int_{K(x,t)} \int_{-\bar{k}(x,t,v)}^{\bar{k}(x,t,v)} \psi_v(a-s) da dv dt \\ &= |\chi| \int_0^\infty \int_{K(x,t)} (\psi_v(-\bar{k}(x,t,v) - s) - \psi_v(\bar{k}(x,t,v) - s)) dv dt \end{aligned}$$

If  $s$  is negative, then the symmetry and unimodality (increasing away from the origin) of  $\psi$ , along with the nonnegativity of  $\bar{k}(x, t, v)$ , imply that  $\psi_v(\bar{k}(x, t, v) - s) \geq \psi_v(-\bar{k}(x, t, v) - s)$  and therefore that  $\frac{d}{ds} F(s)$  is nonpositive. Similarly if  $s$  is positive,  $\frac{d}{ds} F(s)$  is nonnegative. Hence  $F(s)$  achieves its minimum at  $s = 0$ .  $\square$

## 7 Uniqueness of the mean as a majorizing estimator for $\{G_\theta\}$

We now strengthen Theorem 2 by showing that the mean is the *unique* majorizing estimator for the family  $\{G_\theta\}$ . This requires a more delicate argument than the earlier theorem. In the earlier case we did not have to rule out an estimator which improved its odds of success at some values of  $\theta$ , so long as we could rule out its doing better, by an amount bounded away from 0, everywhere; for this purpose it was sufficient to look at a long enough segment of  $\theta$ 's, show that not much benefit could be contributed to this interval by samples from outside of it, and then average uniformly the probability of success within the interval, showing that this average could improve over the mean estimator only by a quantity tending to zero in the length of the interval. There was nothing to prevent the estimator differing from the mean estimator, and indeed improving on the mean estimator locally, so long as it compensated for that change by “importing” estimates toward the values of  $\theta$  that were “neglected”. Now, however, we have to show that if the estimator differs from the mean estimator anywhere, then such a compensation mechanism, while easy to construct in the neighborhood of a small difference, must ultimately fail. The reason for this failure is that the needed compensations in the estimator themselves require compounding compensations, and that this process “diverges”.

We begin with some notation:  $\mathcal{L}$  is the set of Lebesgue measurable sets in  $\mathbb{R}^j$  and  $\mu$  is the usual Lebesgue measure on  $\mathbb{R}^j$  (we write  $\mathcal{L}$  and  $\mu$  regardless of  $j$ ). For an interval  $B \subseteq \mathbb{R}$  we also write  $|B| = \mu(B)$ . Let  $G(y) = (2\pi)^{-1/2} \exp(-y^2/2)$ , and let  $\mathcal{N}(y) = \int_{-\infty}^y G(z) dz$ . The uniqueness theorem is proven in the following generality: an estimator  $S$  is a measure on  $(\mathbb{R}^{n+1}, \mathcal{L})$  (arguments 2, ...,  $n+1$  are the samples  $x_1, \dots, x_n$ , the first argument is the estimate for  $\theta$ ), that satisfies the following condition: for all measurable sets  $U \subseteq \mathbb{R}^n$ ,  $S(\mathbb{R} \times U) = \mu(U)$ .

**Theorem 14** *If there is a measurable set  $A$  such that  $S(A) \neq T(A)$  then for every  $\varepsilon$ ,  $Q_S^\varepsilon < Q_T^\varepsilon$ .*

**Proof:** Again as in theorem 2, it suffices to consider the single-sample case, with  $T$  the identity estimator.

More precisely  $T$  is the diagonal measure: if  $J = \{(t, x) : t = x\}$  and  $\pi_2$  is the projection of  $\mathbb{R}^2$  on its second coordinate then  $T(A) = \mu(\pi_2(A \cap J))$ .

Define the  $\varepsilon$ -quality of estimator  $S$  at  $\theta$  to be

$$Q_S^\varepsilon(\theta) = \int_{x \in \mathbb{R}} G(x - \theta) \int_{t=\theta-\varepsilon}^{\theta+\varepsilon} dS(t, x).$$

For any  $\theta \in \mathbb{R}$ , the  $\varepsilon$ -quality of  $T$  at  $\theta$  is  $\mathcal{N}(\varepsilon) - \mathcal{N}(-\varepsilon) = 2\mathcal{N}(\varepsilon) - 1$ . For the rest of the discussion, assume that  $\varepsilon > 0$  is fixed.

The quantity of interest for us is  $Q_S^\varepsilon = \inf_\theta Q_S^\varepsilon(\theta)$ . As in the proof of Theorem 2, we will need to consider the *average* performance of  $S$  in order to characterize its worst case performance. Thus, for a measurable set  $B$ , we will be interested in

$$F_S(B; \mathbb{R}) = \int_{\theta \in B} \int_{x \in \mathbb{R}} G(x - \theta) \int_{t=\theta-\varepsilon}^{\theta+\varepsilon} dS(t, x) d\theta.$$

Let us define this to be the *estimation total for  $\theta$  in  $B$* . For convenience, let us first express this as a double integral: Let  $u_B$  be the characteristic function for  $B$ . For  $x, t \in \mathbb{R}$ , define

$$\alpha(x, t, B) = \int_{t-\varepsilon}^{t+\varepsilon} G(s - x) u_B(s) ds,$$

For instance, if  $B = \mathbb{R}$ , then this is simply  $\mathcal{N}(-x + t + \varepsilon) - \mathcal{N}(-x + t - \varepsilon)$ . The reader can now verify that

$$F_S(B; \mathbb{R}) = \int_{x \in \mathbb{R}} \int_{t \in \mathbb{R}} \alpha(x, t, B) dS(t, x).$$

More generally, for two measurable sets  $B$  and  $D$ , let us define the *estimation total for  $\theta$  in  $B$  due to  $x$  in  $D$*  to be

$$F_S(B; D) = \int_{x \in D} \int_{t \in \mathbb{R}} \alpha(x, t, B) dS(t, x).$$

A quantity of special interest is the total amount accrued due to  $x$  in  $D$ ,  $F_S(\mathbb{R}; D)$ . Notice that this is maximized by the mean estimator; in particular,

$$F_T(\mathbb{R}; D) = \mu(D)(2\mathcal{N}(\varepsilon) - 1).$$

Finally, define the *deficit of estimator  $S$  on set  $B$* ,

$$\Delta_S(B) = \int_{x \in B} \int_{t \in \mathbb{R}} \alpha(x, t, \mathbb{R})(dT(t, x) - dS(t, x)).$$

For the special case of finite measure  $B$  this is the same as

$$\Delta_S(B) = F_T(\mathbb{R}; B) - F_S(\mathbb{R}; B).$$

**Lemma 15** *If  $S(A) \neq T(A)$  for a measurable set  $A$ , then there is a finite interval  $B$  for which  $\Delta_S(B) > 0$ .*

**Proof:** By countable additivity, we may assume that there is a finite interval  $B$  such that  $A \subseteq \mathbb{R} \times B$ . We first claim that  $S((\mathbb{R} \times B) - J) > T((\mathbb{R} \times B) - J)$ . Clearly,  $S(A \cap J) \leq T(A \cap J)$ . If  $S(A \cap J) < T(A \cap J)$ , the claim follows since  $S(\mathbb{R} \times B) = \mu(B)$ . On the other hand, if  $S(A \cap J) = T(A \cap J)$ , then  $S(A - J) > T(A - J) = 0$ . Since  $T((\mathbb{R} \times B) - J) = 0$ , the claim follows again.

Partition  $(\mathbb{R} \times B) - J$  into regions  $K_j = \{(t, x) : 2^j \leq |t - x| < 2^{j+1}\} \cap (\mathbb{R} \times B)$ , for each integer  $j$ . Again, using countable additivity, there is a  $j$  such that  $S(K_j) > 0$ .

Then  $\Delta_S(B) \geq S(K_j)[(\mathcal{N}(\varepsilon) - \mathcal{N}(-\varepsilon)) - (\mathcal{N}(\varepsilon + 2^j) - \mathcal{N}(-\varepsilon + 2^j))] > 0$ .  $\square$

Let  $B'$  denote the interval obtained by extending interval  $B$  by  $\varepsilon$  on each side. The next lemma shows that deficit must lead to a smaller estimation total for  $S$  (as compared to  $T$ ).

**Lemma 16** For a finite interval  $B$ ,  $F_T(B'; B) - F_S(B'; B) \geq \Delta_S(B)$ .

**Proof:** Observe that  $F_T(\mathbb{R}; B) = F_T(B'; B)$ . Furthermore, since  $\alpha(x, t, \mathbb{R}) \geq \alpha(x, t, B')$  for any  $x, t \in \mathbb{R}$ ,  $F_S(\mathbb{R}; B) \geq F_S(B'; B)$ . Therefore,

$$\Delta_S(B) = F_T(\mathbb{R}; B) - F_S(\mathbb{R}; B) \leq F_T(B'; B) - F_S(B'; B).$$

□

**Lemma 17** If  $\Delta_S(\mathbb{R}) > 0$  then there is a set  $D$  of finite measure such that  $F_S(D; \mathbb{R}) < F_T(D; \mathbb{R})$ .

**Proof:** There are two cases:

**Case (i):**  $\Delta_S(\mathbb{R})$  is infinite.

Let  $B$  be a finite interval such that  $\Delta_S(B) > \varepsilon$ . By Lemma 16,  $F_S(B'; B) \leq F_T(B'; B) - \Delta_S(B) < F_T(B'; B) - \varepsilon$ . Clearly,  $F_S(B'; \mathbb{R} - B)$  is maximized by the estimator that, for each  $x \in \mathbb{R} - B$ , guesses the closest endpoint of  $B$ . It is easy to verify that for such an estimator,  $F_S(B'; \mathbb{R} - B) \leq \varepsilon$ . Therefore,

$$F_S(B'; \mathbb{R}) \leq F_S(B'; B) + \varepsilon < F_T(B'; B) < F_T(B'; \mathbb{R}).$$

**Case (ii):**  $\Delta_S(\mathbb{R})$  is finite.

Let  $B$  be a finite interval such that  $g(\Delta_S(\mathbb{R} - B)) \leq \Delta_S(B)/2$ , where  $g$ , to be defined below, is a monotone increasing continuous function on the nonnegative reals, with  $g(0) = 0$ . Define  $B'$  as above.

By Lemma 16,

$$F_S(B'; B) \leq F_T(B'; B) - \Delta_S(B) \leq |B|(2\mathcal{N}(\varepsilon) - 1) - \Delta_S(B),$$

we get

$$\begin{aligned} F_S(B'; \mathbb{R}) &= F_S(B'; B) + F_S(B'; \mathbb{R} - B) \\ &\leq |B|(2\mathcal{N}(\varepsilon) - 1) - \Delta_S(B) + F_S(B'; \mathbb{R} - B). \end{aligned}$$

In the simplest case, that  $S$  is identical to  $T$  on  $\mathbb{R} - B$ , the last term equals  $2\varepsilon(2\mathcal{N}(\varepsilon) - 1)$  and so  $F_S(B'; \mathbb{R}) \leq |B'|(2\mathcal{N}(\varepsilon) - 1) - \Delta_S(B) < |B'|(2\mathcal{N}(\varepsilon) - 1) = F_T(B'; \mathbb{R})$ . However,  $\Delta_S(\mathbb{R} - B)$  may be nonzero. This allows estimates to be shifted so as to increase  $F_S(B'; \mathbb{R})$ . The remainder of the argument is devoted to showing that this increase, which we call  $DF_S(B'; \mathbb{R})$ , is less than  $\Delta_S(B)$ , provided  $\Delta_S(\mathbb{R} - B)$  is sufficiently small as specified above.

If, at distance  $y$  from  $B$ , the estimator is shifted by distance  $r$  toward  $B$ , then the contribution toward  $\Delta_S(\mathbb{R} - B)$  is proportional to  $\int_0^r (G(-\varepsilon + s) - G(\varepsilon + s)) ds = -\mathcal{N}(\varepsilon + r) + \mathcal{N}(\varepsilon) + \mathcal{N}(-\varepsilon + r) - \mathcal{N}(-\varepsilon)$ . Meanwhile,  $DF_S(B'; \mathbb{R})$  is  $\mathcal{N}(\varepsilon + r) - \mathcal{N}(\varepsilon)$  for  $y \leq 2\varepsilon$ , provided  $0 \leq r \leq y$  (greater values of  $r$  contribute less to  $F_S(B'; \mathbb{R})$ ); while for  $y \geq 2\varepsilon$   $DF_S(B'; \mathbb{R})$  is 0 for  $0 \leq r \leq y - 2\varepsilon$ , and  $\mathcal{N}(\varepsilon + r) - \mathcal{N}(y - \varepsilon)$  for  $y - 2\varepsilon \leq r \leq y$  (again, greater values of  $r$  contribute less to  $F_S(B'; \mathbb{R})$ ).

First, we claim that the best gain in  $F_S(B'; \mathbb{R})$  (greatest value of  $DF_S(B'; \mathbb{R})$ ) given the limit on  $\Delta_S(\mathbb{R} - B)$  is achieved by a “deterministic” estimator, i.e. one which for any  $y$ , places the entire measure on a particular value of  $r$ . This is for the following reason. Let the equation

$$-\mathcal{N}(\varepsilon + r) + \mathcal{N}(\varepsilon) + \mathcal{N}(-\varepsilon + r) - \mathcal{N}(-\varepsilon) = z$$

implicitly define  $r$  as a function of  $z$ , and let  $h$  denote the function such that  $h(z)$  equals  $\mathcal{N}(\varepsilon + r) - \mathcal{N}(\varepsilon)$  for the  $r$  corresponding to  $z$ . Then calculation shows that for  $y \leq 2\varepsilon$ ,  $h$  is a convex cap, increasing function, hence a convex combination  $\sum p_i h(z_i)$  is maximized, given an upper bound on  $\sum p_i z_i$  (the local deficit), by choosing a singular distribution, i.e. a deterministic estimator. A similar argument yields the same conclusion for  $y \geq 2\varepsilon$ .

Moreover, the ratio of “gain” to “cost”

$$\frac{\mathcal{N}(\varepsilon + r) - \mathcal{N}(\varepsilon)}{-\mathcal{N}(\varepsilon + r) + \mathcal{N}(\varepsilon) + \mathcal{N}(-\varepsilon + r) - \mathcal{N}(-\varepsilon)} \tag{2}$$



does not depend on  $y$ , for  $y \leq 2\varepsilon$ ; hence it is optimal to use the same shift  $r$  for all  $y \leq 2\varepsilon$ . Moreover since the ratio is only worse for  $y \geq 2\varepsilon$ , where it is given by the equation

$$\frac{\mathcal{N}(\varepsilon + r) - \mathcal{N}(y - \varepsilon)}{-\mathcal{N}(\varepsilon + r) + \mathcal{N}(\varepsilon) + \mathcal{N}(-\varepsilon + r) - \mathcal{N}(-\varepsilon)} \quad (3)$$

it follows that in an optimal estimator the shift used at that range can be no greater. We therefore obtain an upper bound on  $DF_S(B'; \mathbb{R})$  in the following way: considering only  $y \leq 2\varepsilon$ , find the shift  $r_0$  such that  $DF_S(B'; \mathbb{R})$  is maximized without the deficit exceeding  $\Delta_S(\mathbb{R} - B)$ . Observe that  $r_0$  is at least as great as the shift used by the optimal estimator for  $y \leq 2\varepsilon$  (the optimal estimator may not use all of the deficit on these values of  $y$ , and so may not be able to “afford” as great a shift.) Now since  $r_0$  can be at most  $2\varepsilon$ , and since the optimal estimator uses a shift of at most  $r_0$  for  $y \geq 2\varepsilon$ , it follows that the optimal estimator does not introduce any shift at all for any  $y > 4\varepsilon$ . So we can upper bound  $DF_S(B'; \mathbb{R})$  by  $8\varepsilon(\mathcal{N}(\varepsilon + r_0) - \mathcal{N}(\varepsilon))$ . (A factor of 2 has been introduced to account for both sides of  $B$ .)

The equation defining  $r_0$  is  $\Delta_S(\mathbb{R} - B) = 4\varepsilon[-\mathcal{N}(\varepsilon + r_0) + \mathcal{N}(\varepsilon) + \mathcal{N}(-\varepsilon + r_0) - \mathcal{N}(-\varepsilon)]$ . Let  $g_1$  denote the implicitly defined function on  $\mathbb{R}_{\geq 0}$  giving  $r_0$  as a function of  $\Delta_S(\mathbb{R} - B)$ ; note that  $g_1$  is monotone increasing, continuous and that  $\lim_{x \rightarrow 0} g_1(x) = 0$ . Next, let  $g_2(x) = 8\varepsilon(\mathcal{N}(\varepsilon + x) - \mathcal{N}(\varepsilon))$ ; note that  $g_2$  is monotone increasing, continuous and that  $\lim_{x \rightarrow 0} g_2(x) = 0$ . The composite function  $g(x) = g_2(g_1(x))$  is an upper bound on  $DF_S(B'; \mathbb{R})$  as a function of  $\Delta_S(\mathbb{R} - B)$ ; note that  $g$  is monotone increasing, continuous and that  $\lim_{x \rightarrow 0} g(x) = 0$ . This is the function  $g$  required at the outset of the proof in the selection of  $B$ ; and now, using the assumption that  $g(\Delta_S(\mathbb{R} - B)) \leq \Delta_S(B)/2$ , we find that  $DF_S(B'; \mathbb{R}) \leq \Delta_S(B)/2$  and therefore (by comparing with the estimator which is equal to the mean outside  $B$ ), we find that  $F_S(B'; \mathbb{R}) \leq |B'|(2\mathcal{N}(\varepsilon) - 1) - \Delta_S(B) + DF_S(B'; \mathbb{R}) \leq |B'|(2\mathcal{N}(\varepsilon) - 1) - \Delta_S(B)/2 < |B'|(2\mathcal{N}(\varepsilon) - 1) = F_T(B'; \mathbb{R})$ .  $\square$

## 8 Discussion

The main open issue suggested by our work is whether the concept of a majorizing estimator, as well as the techniques we use for the Gaussian family, can be useful in establishing optimality of estimators for other families of distributions.

Regarding the Gaussian distribution we conjecture that the mean estimator  $T$  is the unique penalty-minimizing estimator for any nonzero penalty function  $\psi$ , i.e., for all nonzero penalty functions  $\psi$ , if  $S$  is an estimator and there is a measurable set  $A$  such that  $S(A) \neq T(A)$ , then  $M^\psi(T) < M^\psi(S)$ .

Another interesting question has to do with the fact that the Cramer-Rao lower bound (on the variance of any unbiased estimator) varies at the parameter values  $\theta$ , depending on the sensitivity of the parametric family to change about  $\theta$ . In the same spirit one may ask (for a general parametric family) for a lower bound  $p(\theta, \varepsilon)$  on the probability that an estimator of  $\theta$  falls outside of the interval  $(\theta - \varepsilon, \theta + \varepsilon)$ . The bound should have the property that  $\lim_{\varepsilon \rightarrow 0} p(\theta, \varepsilon) = 1$ . Some assumption must be made to keep the estimator “honest” (to rule out a constant function for example), such as unbiasedness, or an assumption about the estimator achieving some minimal in-probability performance for some interval length at all  $\theta$ .

## 9 Acknowledgments

We wish to thank Prof. D. Blackwell and Prof. C. R. Rao for helping us confirm the status of Theorem 2.

## References

- [CEG95] R. Canetti, G. Even, and O. Goldreich. Lower bounds for sampling algorithms for estimating the average. *Information Processing Letters*, 53:17–25, 1995.
- [Cra46] H. Cramer. A contribution to the theory of statistical estimation. *Skandinavisk Aktuarietidskrift*, 29:85–94, 1946.

- [CT91] T. M. Cover and J. A. Thomas. *Elements of Information Theory*. Wiley, 1991.
- [Fre43] M. Frechet. sur l'extension de certain evaluations statistique au cas des petit echantillons. *Rev. Inst. Stat.*, 11:182–205, 1943.
- [God81] C. D. Godsil. Matching behaviour is asymptotically normal. *Combinatorica*, 1:369–376, 1981.
- [JS95] M. Jerrum and A. Sinclair. The markov chain monte carlo method: an approach to approximate counting and integration. In D. Hochbaum, editor, *Approximation Algorithms for NP-hard problems*. PWS Publishing Co., 1995.
- [JVV86] M.R. Jerrum, L.G. Valiant, and V.V. Vazirani. Random generation of combinatorial structures from a uniform distribution. *Theoretical Computer Science*, 43:169–188, 1986.
- [LP86] L. Lovász and M. D. Plummer. *Matching Theory*. Akadémiai Kiadó, 1986.
- [Rao45] C. R. Rao. Information and accuracy attainable in estimation of statistical parameters. *Bull. Cal. Math. Soc.*, 37:81–91, 1945.
- [SV99] L. J. Schulman and V. V. Vazirani. Majorizing estimators and the approximation of  $\#\mathbf{P}$ -complete problems. In *31st Annual ACM Symposium on Theory of Computing*, pages 288–294, 1999.
- [SV01] L. J. Schulman and V. V. Vazirani. Majorizing estimators II, with applications to estimating average matching size. Manuscript, 2001.
- [Val79a] L. G. Valiant. The complexity of computing the permanent. *Theoretical Computer Science*, 8:189–201, 1979.
- [Val79b] L. G. Valiant. The complexity of enumeration and reliability problems. *SIAM Journal of Computing*, 8(3):410–421, 1979.
- [Zac81] S. Zacks. *Parametric Statistical Inference*. Pergamon Press, 1981.