

Dartmouth College

## Dartmouth Digital Commons

---

Computer Science Technical Reports

Computer Science

---

12-1-2014

### Information-Theoretic Limits for Density Estimation

James Brofos

*Dartmouth College*

Follow this and additional works at: [https://digitalcommons.dartmouth.edu/cs\\_tr](https://digitalcommons.dartmouth.edu/cs_tr)



Part of the [Computer Sciences Commons](#)

---

#### Dartmouth Digital Commons Citation

Brofos, James, "Information-Theoretic Limits for Density Estimation" (2014). Computer Science Technical Report TR2014-765. [https://digitalcommons.dartmouth.edu/cs\\_tr/365](https://digitalcommons.dartmouth.edu/cs_tr/365)

This Technical Report is brought to you for free and open access by the Computer Science at Dartmouth Digital Commons. It has been accepted for inclusion in Computer Science Technical Reports by an authorized administrator of Dartmouth Digital Commons. For more information, please contact [dartmouthdigitalcommons@groups.dartmouth.edu](mailto:dartmouthdigitalcommons@groups.dartmouth.edu).

# INFORMATION-THEORETIC LIMITS FOR DENSITY ESTIMATION

BY JAMES BROFOS\*

*Dartmouth College\**

This paper is concerned with the information-theoretical limits of density estimation for Gaussian random variables with data drawn independently and with identical distributions. We apply Fano's inequality to the space of densities and an arbitrary estimator. We derive necessary conditions on the sample size for reliable density recovery and for reliable density estimation. These conditions are true simultaneously for both finitely and infinitely dimensional density spaces.

**1. Introduction.** Given a set of  $k$  normal densities  $\{f_1, \dots, f_k\}$  (assumed to be univariate), a single density is selected uniformly at random and  $n$  samples are generated from it. The density that was selected (and in particular its index in  $\{1, \dots, k\}$ ) is assumed to be unknown, though we are provided the samples. Denoting these  $n$  samples collectively as  $\mathbf{X}^n$ , the task is then to estimate, from  $\mathbf{X}^n$ , which of the  $k$  densities was responsible for generating the data. This problem of density recovery from data is one that appears ubiquitously in statistics and related fields. We also treat the problem of density estimation, which, rather than estimating the index of the underlying distribution, estimates the distribution itself.

In this paper, we consider the information-theoretic limits of the density recovery and density estimation problems for Gaussian random variables and samples generated from them in an i.i.d. manner. We treat this problem very generally, and in fact our results hold for any algorithm and associated computational complexity. The key ingredient of this analysis is Fano's inequality, which has been at the source of many recent exciting developments in theoretical statistics and machine learning.

Throughout this paper we will assume that the distributions are identifiable in the sense that if  $f_i$  is parameterized by  $\theta_i = (\mu_i, \sigma_i^2)$ , then  $\theta_i \neq \theta_j$  whenever  $i \neq j$ . It is understood that, throughout this paper, probabilities and expectations are computed with respect to the appropriate probability measure induced by a density.

**2. Main Results.** We establish necessary conditions on the sample size for both reliable density recovery and reliable density estimation associated

with an estimator. We define  $\psi : \mathcal{X}^n \rightarrow \{1, \dots, k\}$  to be an estimator of the density index. We begin by introducing two definitions for  $\psi$  that will be useful in our analysis.

**DEFINITION 2.1** ( $\delta$ -Suspect and  $\delta$ -Inconsistent). We say that an estimator for the index of the true density  $\psi$  is  $\delta$ -suspect if

$$(2.1) \quad \max_{i \in \{1, \dots, k\}} \mathbb{P}[\psi(\mathbf{X}^n) \neq i] \geq \delta.$$

Furthermore, the density selected by  $\psi$  on the basis of  $n$  samples, denoted  $f_\psi$ , is said to be  $\delta$ -inconsistent if

$$(2.2) \quad \max_{i \in \{1, \dots, k\}} \mathbb{E}[\|f_\psi - f_i\|_1] \geq \delta.$$

The notation  $\|\cdot\|_1$  denotes, as usual, the manhattan norm, which is for probability densities twice the total variation distance.

**REMARK 2.1.** Lower bounds in 2.1 on the error probability and in 2.2 on the total variation distance are useful primarily for establishing necessary conditions for any density index estimator  $\psi$ . In particular, one can establish an upper bound on the right-hand side of either 2.1 or 2.2, say  $\delta < \delta'$ , and deduce a lower bound on the sample size  $n$  that will depend on  $\delta'$ . This lower bound can then be viewed as a necessary (although certainly not sufficient) condition for any estimator  $\psi$  to be reliable in terms of density recovery or density estimation.

We now present our main results as two theorems relating the sample size for density estimation problems for  $k$  normal densities with parameterizations  $\{\theta_i\}_{i=1}^k$  to notions arising from information theory.

**THEOREM 2.1** (Necessary Condition for Density Recovery). *Let  $\Xi$  be distributed uniformly at random from  $\{1, \dots, k\}$ . Then if  $\Xi = i$ , we have that  $\mathbf{X}_j^n \sim f_i$  for all  $j \in \{1, \dots, n\}$ . If  $\psi$  is an estimator of  $\Xi$  from  $\mathbf{X}^n$ , then*

$$(2.3) \quad \max_i \mathbb{P}[\psi(\mathbf{X}^n) \neq i] \geq 1 - \frac{\frac{n}{k^2} \sum_{a=1}^k \sum_{b=1}^k \mathcal{D}_{KL}(f_a, f_b) + \log 2}{\log k}$$

$$(2.4) \quad \geq 1 - \frac{n \cdot \max_{a,b} \{\mathcal{D}_{KL}(f_a, f_b)\} + \log 2}{\log k}.$$

*It follows that a necessary condition for  $\psi$  to not be a  $\delta$ -suspect estimator we must have,*

$$(2.5) \quad n > \frac{\log \frac{k}{2} - \delta \log k}{\max_{a,b} \left\{ \log \frac{\sigma_b}{\sigma_a} + \frac{\sigma_a^2 + (\mu_a - \mu_b)^2}{2\sigma_b^2} - \frac{1}{2} \right\}}.$$

**THEOREM 2.2** (Necessary Condition for Density Estimation). *Under the same assumptions as in Theorem 2.1, we have for any density estimator  $f_\psi$  that the worst case (with respect to  $\Xi = i$ ), is lower bounded as*

$$(2.6) \quad \mathbb{E} [\|f_\psi - f_i\|_1] \geq \frac{\min_{a,b} \left\{ 2 - 2\sqrt{\frac{2\sigma_a\sigma_b}{\sigma_a^2+\sigma_b^2}} e^{\frac{-1}{4} \frac{(\mu_a-\mu_b)^2}{\sigma_a^2+\sigma_b^2}} \right\}}{2} \times \left( 1 - \frac{n \cdot \max_{a,b} \{\mathcal{D}_{KL}(f_a, f_b)\} + \log 2}{\log k} \right).$$

Denote by  $\alpha$  the quantity attaining the minimum in 2.6. Then a necessary condition for  $f_\psi$  to not be a  $\delta$ -inconsistent estimator of the density is

$$(2.7) \quad n > \frac{(\alpha - 2\delta) \cdot \log k - \alpha \cdot \log 2}{\alpha \cdot \left( \max_{a,b} \left\{ \log \frac{\sigma_b}{\sigma_a} + \frac{\sigma_a^2 + (\mu_a - \mu_b)^2}{2\sigma_b^2} - \frac{1}{2} \right\} \right)}.$$

**REMARK 2.2.** It should be noted that this same approach can be applied to subsets of the  $k$  densities in the event that a more satisfying lower bound on  $n$  can be obtained in this fashion. Indeed, if  $\mathbf{F}$  denotes a class of densities with a subclass of densities  $\mathbf{F}'$ , then clearly the number of samples required for reliable density recovery (or estimation) is the maximum of the lower bounds obtained by applying Theorem 2.1 (or Theorem 2.2) to  $\mathbf{F}$  and  $\mathbf{F}'$  individually.

**EXAMPLE 2.1.** Suppose we have a set of three normal densities with different means and variances  $\{\mathcal{N}(0.89, 0.59), \mathcal{N}(0.90, 0.60), \mathcal{N}(0.91, 0.61)\}$ . The bounds in Theorems 2.1 and 2.2 are most useful in “difficult” estimation problems, where the Kullback-Leibler divergence is small. In this example,

$$(2.8) \quad \max_{a,b} \left\{ \log \frac{\sigma_b}{\sigma_a} + \frac{\sigma_a^2 + (\mu_a - \mu_b)^2}{2\sigma_b^2} - \frac{1}{2} \right\} = 0.017$$

$$(2.9) \quad \min_{a,b} \left\{ 2 - 2\sqrt{\frac{2\sigma_a\sigma_b}{\sigma_a^2+\sigma_b^2}} e^{\frac{-1}{4} \frac{(\mu_a-\mu_b)^2}{\sigma_a^2+\sigma_b^2}} \right\} = 7.55 \times 10^{-5}.$$

Therefore, the necessary condition from Theorem 2.1 is  $n \geq 18$  for  $\delta = 0.10$ . Additionally, the necessary condition from Theorem 2.2 is  $n \geq 22$  for  $\delta = 1 \times 10^{-6}$ .

**3. Proof of the Main Results.** Theorems 2.1 and 2.2 characterize the number of samples required for reliable density recovery and estimation using bounds on information-theoretic quantities. The remainder of the paper is devoted to proving these two theorems. We begin with a quantity well-understood in statistics and information theory that quantifies the similarity between two probability distributions.

**DEFINITION 3.1** (Squared Hellinger Distance). The squared Hellinger distance between two probability measures  $\pi_1$  and  $\pi_2$  is given as

$$(3.1) \quad H^2(\pi_1, \pi_2) = \frac{1}{2} \int \left( \sqrt{\frac{d\alpha_1}{d\gamma}} - \sqrt{\frac{d\alpha_2}{d\gamma}} \right)^2 d\gamma.$$

The variable  $\gamma$  is a probability measure. In particular, if  $\pi_1 \sim \mathcal{N}(\mu_1, \sigma_1^2)$  and  $\pi_2 \sim \mathcal{N}(\mu_2, \sigma_2^2)$  then the squared Hellinger distance becomes

$$(3.2) \quad H^2(\pi_1, \pi_2) = 1 - \sqrt{\frac{2\sigma_1\sigma_2}{\sigma_1^2 + \sigma_2^2}} \exp \left\{ \frac{-1}{4} \cdot \frac{(\mu_1 - \mu_2)^2}{\sigma_1^2 + \sigma_2^2} \right\}.$$

**REMARK 3.1.** The squared Hellinger distance will prove to be an indispensable quantity for our application. It is a well-known result that the squared Hellinger distance is always less than or equal to the total variation distance for two probability measures. Since the total variation distance for normally distributed random variables does not have closed-form, the tractability of the squared Hellinger distance makes it a natural choice in practical applications as well.

It is easy to see that since the total variation is exactly one-half the Manhattan distance, we can introduce a factor of two into the squared Hellinger distance and still obtain a valid lower bound on the Manhattan distance.

We have assumed that each of the  $k$  distributions are parameterized by vectors  $\{\theta_i\}_{i=1}^k$ . For our purposes it will suffice to lower bound the (non-trivial) total variation for all distributions using the squared Hellinger distance. We see that a global lower bound, and the corresponding indices in  $\{1, \dots, k\}$ , is given as

$$(3.3) \quad (i^*, j^*) = \arg \min_{i, j \in \{1, \dots, k\}: i \neq j} 1 - \sqrt{\frac{2\sigma_i\sigma_j}{\sigma_i^2 + \sigma_j^2}} \exp \left\{ \frac{-1}{4} \frac{(\mu_i - \mu_j)^2}{\sigma_i^2 + \sigma_j^2} \right\}.$$

It makes no difference as to how ties are arbitrated and we do not discuss the possibility further.

As it will turn out, we will also require an upper bound on the Kullback-Leibler divergence between all the  $k$  densities. In the case of Gaussian densities, the Kullback-Leibler has a simple closed-form. In particular, we have for Gaussian probability measures  $\pi_1$  and  $\pi_2$ ,

$$(3.4) \quad \mathcal{D}_{\text{KL}}(\pi_1, \pi_2) = \log \frac{\sigma_2}{\sigma_1} + \frac{\sigma_1^2 + (\mu_1 - \mu_2)^2}{2\sigma_2^2} - \frac{1}{2}.$$

This quantity is easy to maximize as a function of the densities  $f_i$  and  $f_j$  (noting, of course, that due to the asymmetry of the Kullback-Leibler divergence,  $\mathcal{D}_{\text{KL}}(f_i, f_j) \neq \mathcal{D}_{\text{KL}}(f_j, f_i)$ ). We denote the indices achieving this maximum as the *ordered* tuple  $(i^\circ, j^\circ)$ .

**THEOREM 3.1** (Fano's Inequality for Statistics). *Let  $\Xi$  be distributed uniformly at random in the set  $\{1, \dots, k\}$ . Let a realization of  $\Xi$  be the index  $i$ . Then the density of  $\mathbf{X}_j^n$  is  $f_i$  which has parameterization  $\theta_i$  for all  $j \in \{1, \dots, n\}$ . Then if  $\psi(\mathbf{X}^n)$  is an estimator of  $\Xi$  based on the available data, we have*

$$(3.5) \quad \max_i \mathbb{P}[\psi(\mathbf{X}^n) \neq i] \geq 1 - \frac{\mathcal{I}(\Xi : \mathbf{X}^n) + \log 2}{\log k}$$

$$(3.6) \quad \geq 1 - \frac{n\beta + \log 2}{\log k}.$$

Here the parameter  $\beta$  is chosen such that  $\forall i, j$  the Kullback-Leibler divergence for densities  $f_i$  and  $f_j$  is upper bounded like  $\mathcal{D}_{\text{KL}}(f_i, f_j) \leq \beta$ . It should be noted that the bound in 3.6 follows from standard entropy bounds for collections of random variables.

**COROLLARY 3.1.** *A more general, but markedly less popular, version of Fano's inequality emerges when one considers the total variation distance. The only additional requirement compared to the “vanilla” version in Theorem 3.1 is a lower bound on the (non-trivial) total variation distance across all density pairs. If we denote this lower bound by  $\alpha$  the general form of Fano's inequality becomes,*

$$(3.7) \quad \max_i \mathbb{E}[\|f_\psi - f_i\|_1] \geq \frac{\alpha}{2} \left( 1 - \frac{n\beta + \log 2}{\log k} \right).$$

We use the notation  $f_\psi$  to indicate a density estimated by the estimator  $\phi$  on the basis of  $n$  samples drawn in an i.i.d. fashion from  $f_i$ .

At this point the proof of Theorems 2.1 and 2.2 requires only a little algebra to derive.

PROOF. Notice from first principles that both 3.5 and 3.7 are decreasing functions of  $n$ . It should be apparent that both can be crushed to zero in the limit as  $n \rightarrow \infty$ . Plainly, if we set either of the right-hand side quantities in 3.5 or 3.7 to be less than any  $\delta > 0$ , then the bound can be satisfied by a particular choice of  $n$ . By rearranging we have for the density recovery case,

$$(3.8) \quad n > \frac{\log \frac{k}{2} - \delta \log k}{\beta}.$$

Similarly for the density estimation case,

$$(3.9) \quad n > \frac{(\alpha - 2\delta) \cdot \log k - \alpha \log 2}{\alpha \beta}.$$

The desired inequalities can be obtained by the substitutions,

$$(3.10) \quad \alpha = 2 - 2\sqrt{\frac{2\sigma_{i^*}\sigma_{j^*}}{\sigma_{i^*}^2 + \sigma_{j^*}^2}} \exp \left\{ \frac{-1}{4} \cdot \frac{(\mu_{i^*} - \mu_{j^*})^2}{\sigma_{i^*}^2 + \sigma_{j^*}^2} \right\}$$

$$(3.11) \quad \beta = \log \frac{\sigma_{j^\circ}}{\sigma_{i^\circ}} + \frac{\sigma_{i^\circ}^2 + (\mu_{i^\circ} - \mu_{j^\circ})^2}{2\sigma_{j^\circ}^2} - \frac{1}{2},$$

This gives the claimed necessary conditions on the sample size  $n$ . □