

Math 6070

A Primer on Statistical Inference

Davar Khoshnevisan
University of Utah

Spring 2014

Contents

1	Statistical Models	1
2	Classical Parametric Inference	2
3	The Information Inequality	4
4	A Glance at Confidence Intervals	7
5	A Glance at Testing Statistical Hypotheses	9

1 Statistical Models

It is convenient to have an abstract framework for discussing statistical theory. The general problem is that there exists an unknown *parameter* θ_0 , which we wish to find out about. To have something concrete in mind, consider for example a population with the $N(\theta_0, 1)$ distribution, where θ_0 is an unknown constant. If we do not have any *a priori* information about θ_0 then it stands to reason that we consider every distribution of the form $N(\theta, 1)$, as θ ranges over \mathbf{R} , and then use data to make inference about the real, unknown θ_0 .

The general framework is this: We have a *parameter space* Θ and the real θ_0 is in Θ , but we do not its value. For every $\theta \in \Theta$, let P_θ denote the underlying probability, which is computed by assuming that $\theta_0 = \theta$. Similarly define E_θ , Var_θ , Cov_θ , etc. Then, the idea is to take a sample—typically an independent sample— $\mathbf{X} = (X_1, \dots, X_n)$ —from P_{θ_0} . If the true (unknown) θ_0 were equal to some (known) $\theta_1 \in \Theta$, then one would expect \mathbf{X} to behave like an independent sample from P_{θ_1} . If so, then we declare that θ_0 might well be θ_1 . Else, we reject the notion that $\theta_0 = \theta_1$. The remainder of these notes make this technique precise in more special settings.

2 Classical Parametric Inference

The typical problem of classical statistics is the following: Given a family of probability densities $\{f_\theta\}_{\theta \in \Theta}$ how can we decide whether or not ours is f_θ ? More precisely, we have an unknown density f_{θ_0} ; we wish to estimate it by choosing one from the family $\{f_\theta\}_{\theta \in \Theta}$ of densities available to us. [Alternatively, you could replace f_θ by a mass function p_θ .] Here, Θ is the “parameter space,” and θ_0 is the unknown “parameter.”

To estimate θ_0 one typically considers an independent sample X_1, \dots, X_n from the true distribution with density f_{θ_0} , and constructs an estimator $\hat{\theta}$.

Example 1 Let $\Theta := \mathbf{R}$, and f_θ the $N(\theta, 1)$ density. The standard approach is to estimate θ_0 with

$$\hat{\theta} := \frac{X_1 + \dots + X_n}{n}. \quad (1)$$

There are many reasons why $\hat{\theta}$ is a good estimate of θ .

1. [Unbiasedness] Evidently,

$$\mathbf{E}_\theta(\hat{\theta}) = \theta, \quad \text{for all } \theta \in \Theta. \quad (2)$$

This is called *unbiasedness*. In general, a random variable T is said to be an *unbiased* estimator of θ if $\mathbf{E}_\theta(T) = \theta$ for all $\theta \in \Theta$.

2. [Consistency] By the law of large numbers, for all $\theta \in \Theta$,

$$\hat{\theta} \xrightarrow{\text{P}} \theta \quad \text{as } n \rightarrow \infty. \quad (3)$$

This is called *consistency*. In general, a random variable T is said to be a *consistent* estimator of θ_0 if $T \xrightarrow{\text{P}} \theta$ for all $\theta \in \Theta$ as the sample size tends to infinity.

3. [MLE] The *maximum likelihood estimate* of θ_0 —in all cases—is an estimator that maximizes $\theta \mapsto f_\theta(X_1, \dots, X_n)$ for an independent sample (X_1, \dots, X_n) , where f_θ here represents the joint density function of n i.i.d. random variables each with density $N(\theta, 1)$. In the present example.

$$f_\theta(X_1, \dots, X_n) = \frac{1}{(2\pi)^{n/2}} \exp \left(-\frac{1}{2} \sum_{j=1}^n (X_j - \theta)^2 \right). \quad (4)$$

To find a MLE, it is easier to maximize the *log likelihood*,

$$L(\theta) := \ln f_\theta(X_1, \dots, X_n), \quad (5)$$

which is the same as minimizing $h(\theta) := \sum_{j=1}^n (X_j - \theta)^2$ over all θ . But $h'(\theta) = -2 \sum_{j=1}^n (X_j - \theta)$ and $h''(\theta) = 2n > 0$. Therefore, the MLE is uniquely $\hat{\theta}$.

The statistics $\hat{\theta}$ has other optimality features too. See for instance Example 8 (page 7) below.

Example 2 Suppose $\Theta := \mathbf{R} \times (0, \infty)$. Then, we can write $\theta \in \Theta$ as $\theta = (\mu, \sigma^2)$ where $\mu \in \mathbf{R}$ and $\sigma > 0$. Suppose f_θ is the $N(\mu, \sigma^2)$ density. Then the usual estimator for the true parameter $\theta_0 = (\mu_0, \sigma_0^2)$ is $\hat{\theta} := (\hat{\mu}, \hat{\sigma}^2)$, where

$$\begin{aligned}\hat{\mu} &:= \frac{1}{n} \sum_{j=1}^n X_j, \\ \hat{\sigma}^2 &:= \frac{1}{n} \sum_{j=1}^n (X_j - \hat{\mu})^2.\end{aligned}\tag{6}$$

[As before, X_1, \dots, X_n is an independent sample.] As in the previous example, $\hat{\theta}$ is the unique MLE, and is consistent. However, it is not unbiased. Indeed,

$$E_\theta(\hat{\theta}) = \left(\begin{array}{c} \mu \\ [1 - \frac{1}{n}]^2 \sigma^2 \end{array} \right), \quad \text{for all } \theta = (\mu, \sigma^2) \in \Theta.\tag{7}$$

So $\hat{\theta}$ is “biased,” although it is *asymptotically unbiased*; i.e., $E_\theta(\hat{\theta}) \rightarrow \theta$ as $n \rightarrow \infty$.

Example 3 Suppose $\Theta = (0, \infty)$, and f_θ is the uniform- $(0, \theta)$ density for all $\theta \in \Theta$. Given an independent sample X_1, \dots, X_n , we consider

$$\hat{\theta} := \max_{1 \leq j \leq n} X_j.\tag{8}$$

The distribution of $\hat{\theta}$ is easily computed, viz.,

$$P_\theta \left\{ \hat{\theta} \leq a \right\} = [P_\theta \{X_1 \leq a\}]^n = (a/\theta_0)^n, \quad 0 \leq a \leq \theta_0.\tag{9}$$

This gives the density $f_{\hat{\theta}}(a) = n\theta_0^{-n}a^{n-1}$ for $0 \leq a \leq \theta_0$. Consequently,

$$E_\theta(\hat{\theta}) = \theta_0^{-n} \int_0^{\theta_0} na^n da = \frac{n\theta_0}{n+1}.\tag{10}$$

Therefore: (i) $\hat{\theta}$ is biased; but (ii) it is asymptotically unbiased. Next we show that $\hat{\theta}$ is consistent. Note that $\hat{\theta} \leq \theta_0$, by force. So it is enough to show that with high probability $\hat{\theta}$ is not too much smaller than θ_0 . Fix $\epsilon > 0$, and note that

$$P_\theta \left\{ \hat{\theta} \leq (1 - \epsilon)\theta_0 \right\} = \int_0^{(1-\epsilon)\theta_0} n\theta_0^{-n}a^{n-1} da = (1 - \epsilon)^n.\tag{11}$$

Thus,

$$P_\theta \left\{ \left| \frac{\hat{\theta}}{\theta_0} - 1 \right| > \epsilon \right\} \leq 1 - (1 - \epsilon)^n \rightarrow 0.\tag{12}$$

That is, $\hat{\theta}$ is consistent, as asserted earlier. To complete the example let us compute the MLE for θ_0 . Evidently,

$$f_{\theta}(X_1, \dots, X_n) = \frac{1}{\theta^n} \mathbf{I}\{\theta > \hat{\theta}\}, \quad (13)$$

where $\mathbf{I}\{A\}$ is the indicator of A . So to find the MLE we observe that $\mathbf{I}\{A\} \leq 1$, so that $f_{\theta}(X_1, \dots, X_n) \leq 1/\hat{\theta}^n$. The MLE is $\hat{\theta}$ uniquely.

One can consider a variant of $\hat{\theta}$, here, that is unbiased and consistent, but only “approximately” MLE for large n . Namely, we can consider the statistic $\tilde{\theta} := (n+1) \max_{1 \leq j \leq n} X_j / n = (1 + \frac{1}{n}) \max_{1 \leq j \leq n} X_j$.

3 The Information Inequality

Let us concentrate on the case where every $\theta \in \Theta$ is one-dimensional, and hence so is θ_0 .

Let $\mathbf{X} := (X_1, \dots, X_n)$ be a random vector with joint density $f_{\theta}(\mathbf{x})$. The *Fisher information* of the family $\{f_{\theta}\}_{\theta \in \Theta}$ is defined as the function $I(\theta)$, where

$$I(\theta) := \mathbb{E}_{\theta} \left[\left(\frac{\partial}{\partial \theta} \ln f_{\theta}(\mathbf{X}) \right)^2 \right], \quad (\theta \in \Theta), \quad (14)$$

provided that the expectation exists and is finite. If \mathbf{X} is discrete we define $I(\theta)$ in the same way, but replace f_{θ} by the joint mass function p_{θ} .

In the continuous case, for example, the Fisher information is computed as follows:

$$\begin{aligned} I(\theta) &= \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} \left(\frac{\partial}{\partial \theta} \ln f_{\theta}(\mathbf{x}) \right)^2 f_{\theta}(\mathbf{x}) d\mathbf{x} \\ &= \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} \frac{1}{f_{\theta}(\mathbf{x})} \left(\frac{\partial}{\partial \theta} f_{\theta}(\mathbf{x}) \right)^2 d\mathbf{x}. \end{aligned} \quad (15)$$

So in fact $I(\theta)$ is always defined, but could be any number in $[0, \infty]$.

Example 4 In the case of independent $N(\theta, 1)$ ’s,

$$\ln f_{\theta}(\mathbf{x}) = -\frac{n}{2} \ln(2\pi) - \frac{1}{2} \sum_{j=1}^n (x_j - \theta)^2. \quad (16)$$

The θ -derivative is $\sum_{j=1}^n (x_j - \theta)$. Therefore,

$$I(\theta) = \mathbb{E}_{\theta} \left[\left(\sum_{j=1}^n X_j - n\theta \right)^2 \right] = \text{Var}_{\theta} \left(\sum_{j=1}^n X_j \right) = n. \quad (17)$$

[Here it does not depend on θ .]

Example 5 Suppose $X_1, \dots, X_n \sim \text{Poisson}(\theta)$ are independent, where $\theta \in \Theta := (0, \infty)$. [Remember that “ $Y \sim D$ ” means that “ Y is distributed as D .”] Now we have the joint mass function $p_\theta(\mathbf{x})$ instead of densities. Then,

$$\ln p_\theta(\mathbf{x}) = -n\theta + \ln \theta \sum_{j=1}^n x_j - \sum_{j=1}^n \ln(x_j!). \quad (18)$$

Differentiate with respect to θ in order to obtain

$$\frac{\partial}{\partial \theta} \ln p_\theta(\mathbf{x}) = -n + \frac{1}{\theta} \sum_{j=1}^n x_j. \quad (19)$$

Therefore,

$$I(\theta) = \frac{1}{\theta^2} \mathbb{E} \left[\left(\sum_{j=1}^n X_j - n\theta \right)^2 \right] = \frac{\text{Var}(\sum_{j=1}^n X_j)}{\theta^2} = \frac{n}{\theta}. \quad (20)$$

The following is due to Fréchet originally, and was rediscovered independently, and later on, by Crámer and Rao.

Theorem 6 (The Information Inequality) *Suppose T is a non-random function of n variables. Then, under “mild regularity conditions,”*

$$\text{Var}_\theta(T(\mathbf{X})) \geq \frac{[h'(\theta)]^2}{I(\theta)}, \quad (21)$$

for all θ , where $h(\theta) := \mathbb{E}_\theta[T(\mathbf{X})]$.

The regularity conditions are indeed mild; they guarantee that certain integrals and derivatives commute. See (24) and (27) below.

The proof requires the following form of the Cauchy–Schwarz inequality:

Lemma 7 (Cauchy–Schwarz Inequality) *For all rv’s X and Y ,*

$$|\text{Cov}(X, Y)|^2 \leq \text{Var}(X) \cdot \text{Var}(Y), \quad (22)$$

provided that all the terms inside the expectations are integrable.

Proof. Let $X' := (X - \mathbb{E}X)/\sqrt{\text{Var}(X)}$ and $Y' := (Y - \mathbb{E}Y)/\sqrt{\text{Var}(Y)}$. Then,

$$\begin{aligned} 0 \leq \mathbb{E} \left[(X' - Y')^2 \right] &= \mathbb{E}[(X')^2] + \mathbb{E}[(Y')^2] - 2\mathbb{E}[X'Y'] \\ &= 2 \left[1 - \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X) \cdot \text{Var}(Y)}} \right]. \end{aligned} \quad (23)$$

This proves the result when $\text{Cov}(X, Y) \geq 0$. When $\text{Cov}(X, Y) < 0$, we consider instead $E[(X' + Y')^2]$. \square

Proof of the Information Inequality in the Continuous Case. Note that if f_θ is nice then

$$\int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} \frac{\partial}{\partial \theta} f_\theta(\mathbf{x}) d\mathbf{x} = \frac{\partial}{\partial \theta} \left[\int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} f_\theta(\mathbf{x}) d\mathbf{x} \right] = 0. \quad (24)$$

This is so simply because $[\cdots] = 1$. Therefore,

$$E_\theta \left[\frac{\partial}{\partial \theta} \ln f_\theta(\mathbf{X}) \right] = \int_{-\infty}^{\infty} f_\theta(\mathbf{x}) \frac{\partial}{\partial \theta} \ln f_\theta(\mathbf{x}) d\mathbf{x} = 0. \quad (25)$$

This proves that

$$I(\theta) = \text{Var}_\theta \left(\frac{\partial}{\partial \theta} \ln f_\theta(\mathbf{X}) \right). \quad (26)$$

Similarly, if things are nice then

$$\begin{aligned} E_\theta \left[T(\mathbf{X}) \frac{\partial}{\partial \theta} \ln f_\theta(\mathbf{X}) \right] &= \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} T(\mathbf{x}) \frac{\partial}{\partial \theta} f_\theta(\mathbf{x}) d\mathbf{x} \\ &= \frac{\partial}{\partial \theta} \left[\int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} T(\mathbf{x}) f_\theta(\mathbf{x}) d\mathbf{x} \right] \\ &= \frac{\partial}{\partial \theta} E_\theta[T(\mathbf{X})] = h'(\theta). \end{aligned} \quad (27)$$

Combine (24) and (27) to find that

$$\text{Cov}_\theta \left(T(\mathbf{X}), \frac{\partial}{\partial \theta} \ln f_\theta(\mathbf{X}) \right) = h'(\theta). \quad (28)$$

Thanks to Lemma 7,

$$|h'(\theta)|^2 \leq \text{Var}_\theta(T(\mathbf{X})) \cdot \text{Var}_\theta \left(\frac{\partial}{\partial \theta} \ln f_\theta(\mathbf{X}) \right) = \text{Var}_\theta(T(\mathbf{X})) \cdot I(\theta). \quad (29)$$

See (26). This proves the information inequality. \square

A useful consequence of the information inequality is that, under mild conditions, any **unbiased** estimator $T(\mathbf{X})$ has the property that

$$\text{Var}_\theta(T(\mathbf{X})) \geq \frac{1}{I(\theta)}. \quad (30)$$

This leads to the notion of MVU estimators: These are unbiased estimators that have minimum variance. Thanks to (30), if we can find a function T such that $\text{Var}(T(\mathbf{X})) = 1/I(\theta_0)$, then we have found an MVU estimator of θ .

Example 8 Suppose X_1, \dots, X_n are i.i.d. $N(\theta, 1)$'s. Let T be such that $T(\mathbf{X})$ is an unbiased estimator of θ . According to Example 4, $I(\theta) = n$, so that $\text{Var}_\theta(T(\mathbf{X})) \geq 1/n = \text{Var}_\theta(\bar{X}_n)$. That is, $\hat{\theta} := (X_1 + \dots + X_n)/n$ has the smallest variance among all unbiased estimators of θ . This is the “MVU” property. More precisely, any estimator $\hat{\theta}$ is said to be *MVUE* when it is a (often, “the”) *minimum variance unbiased estimator* of θ_0 .

Example 9 Suppose X_1, \dots, X_n are $\text{Poisson}(\theta)$, where $\theta > 0$ is an unknown parameter. [The true parameter is some unknown θ_0 , so we model it this way.] Because $E_\theta(X_1) = \theta$, the law of large numbers implies that

$$\bar{X}_n := \frac{X_1 + \dots + X_n}{n} \xrightarrow{P_\theta} \theta. \quad (31)$$

So, \bar{X}_n is a consistent estimator of θ_0 . Recall also that $\text{Var}_\theta(X_1) = \theta$, so that $\text{Var}_\theta(\bar{X}_n) = \theta/n$. We claim that \bar{X}_n is a minimum variance unbiased estimator. In order to prove it it suffices to show that $I(\theta) = n/\theta$. But this was shown to be the case already; see Example 5 on page 5.

4 A Glance at Confidence Intervals

Choose and fix $\alpha \in (0, 1)$. A *confidence set* C with *level* $(1 - \alpha)$ is a random set that depends on the sample \mathbf{X} , and has the property that $P_\theta\{\theta \in C\} \geq 1 - \alpha$ for all $\theta \in \Theta$. If C varies with n , and $\lim_{n \rightarrow \infty} P_\theta\{\theta \in C\} \geq 1 - \alpha$ for all $\theta \in \Theta$, then we say that C is a confidence interval for θ_0 with *asymptotic level* $(1 - \alpha)$.

Example 10 Consider the model $N(\theta, 1)$ where $\theta \in \Theta := \mathbf{R}$. Then, it easy to see that

$$\frac{\bar{X}_n - \theta}{1/\sqrt{n}} \sim N(0, 1) \quad \text{under } P_\theta. \quad (32)$$

Here, “Under P_θ ” is short-hand for “If θ were the true parameter, for all $\theta \in \Theta$.” Consider the random set

$$C(z) := \left[\bar{X}_n - \frac{z}{\sqrt{n}}, \bar{X}_n + \frac{z}{\sqrt{n}} \right], \quad (33)$$

where $z \geq 0$ is fixed. Then,

$$\begin{aligned} P_\theta\{\theta \in C(z)\} &= P_\theta\left\{|\bar{X}_n - \theta| \leq \frac{z}{\sqrt{n}}\right\} \\ &= P_\theta\left\{\frac{|\bar{X}_n - \theta|}{1/\sqrt{n}} \leq z\right\} \\ &= P\{|N(0, 1)| \leq z\} = 2\Phi(z) - 1. \end{aligned} \quad (34)$$

See (32) for the last identity. Choose $z = z_{\alpha/2}$ such that $2\Phi(z_{\alpha/2}) - 1 = 1 - \alpha$ to see that $P_\theta\{\theta \in C(z_{\alpha/2})\} = 1 - \alpha$. That is, $C(z_{\alpha/2})$ is a confidence interval

for θ_0 with level $1 - \alpha$. Note that $z_{\alpha/2}$ is defined by $\Phi(z_{\alpha/2}) = 1 - (\alpha/2)$. The numbers $z_{\alpha/2}$ are called “normal quantiles,” because $P\{N(0, 1) \leq z_{\alpha/2}\} = \Phi(z_{\alpha/2}) = 1 - (\alpha/2)$.

Example 11 Consider the model $\text{Binomial}(n, p)$, where n is a known integer, but $p \in [0, 1]$ is an unknown constant. Here, $\Theta = [0, 1]$, and every $p \in \Theta$ is a parameter. We consider the estimate

$$\hat{p} := \frac{S_n}{n}, \quad (35)$$

where S_n denotes the total number of successes in n independent samples. Evidently, $S_n \sim \text{Binomial}(n, p)$ under P_p . Therefore, $E_p(\hat{p}) = p$ and $\text{Var}_p(\hat{p}) = p(1 - p)/n$.

By the central limit theorem, as n tends to infinity,

$$\frac{S_n - np}{\sqrt{np(1 - p)}} \xrightarrow{d} N(0, 1), \quad (36)$$

under P_p . (Why?) Equivalently,

$$\frac{\hat{p} - p}{\sqrt{p(1 - p)/n}} \xrightarrow{d} N(0, 1), \quad (37)$$

under P_p . Also, by the law of large numbers, $\hat{p} \xrightarrow{P_g} p$. (Why?) Apply the latter two results, via Slutsky’s theorem, to find that under P_p ,

$$\frac{\hat{p} - p}{\sqrt{\hat{p}(1 - \hat{p})/n}} \xrightarrow{d} N(0, 1). \quad (38)$$

Now consider

$$C_n(z) := \left[\hat{p} - z \sqrt{\frac{\hat{p}(1 - \hat{p})}{n}}, \hat{p} + z \sqrt{\frac{\hat{p}(1 - \hat{p})}{n}} \right]. \quad (39)$$

Then, we have shown that

$$\lim_{n \rightarrow \infty} P_p \{p \in C_n(z)\} = P\{|N(0, 1)| \leq z\} = 2\Phi(z) - 1. \quad (40)$$

Therefore, $C_n(z_{\alpha/2})$ is asymptotically a level- $(1 - \alpha)$ confidence interval for p .

There are many variants of confidence intervals that are also useful. For instance, a *one-sided confidence interval* is a half-infinite random interval that should contain the parameter of interest with a pre-described probability. Similarly, there are one-sided confidence intervals that have a given asymptotic level. Finally, there are higher-dimensional generalizations. For example, there are confidence ellipsoids, confidence bands, etc. All of them are random sets—often with a pre-described geometry—that have exact or asymptotic level $(1 - \alpha)$ for a pre-described level $\alpha \in (0, 1)$.

5 A Glance at Testing Statistical Hypotheses

Someone proposes the theory that a certain coin is fair. To test this hypothesis, a statistician can flip the said coin n times, independently. Record the number of heads S_n . In any event, we know that $S_n \sim \text{binomial}(n, p)$ for some p . Thus, we write the proposed hypothesis as the *null hypothesis*, $H_0 : p = \frac{1}{2}$, versus the *alternative*, $H_1 : p \neq \frac{1}{2}$. If the null hypothesis is correct, then $\hat{p} := S_n/n$ is close to $p = 1/2$ with high probability. Fix $\alpha \in (0, 1)$, and consider the confidence interval $C_n(z_{\alpha/2})$ from Example 11 on page 8. It is more convenient to write P_{H_0} here instead of P_p . With this in mind, we know then that for large n ,

$$P_{H_0} \{p \notin C_n(z_{\alpha/2})\} \approx \alpha. \quad (41)$$

Here is how we make an inference about H_0 : If $p \notin C_n(z_{\alpha/2})$, then we reject the null hypothesis H_0 . Else, we accept H_0 , but only in the sense that we do not reject it. There are two sources of error in testing statistical hypotheses:

1. Type-I Error: This is the probability of incorrect rejection of H_0 . In our example, (41) shows that the type-I error is asymptotically α .
2. Type-II Error: This is the probability of incorrect acceptance of H_1 . In our example, type-II error is

$$\beta = P_{H_1} \{p \in C_n(z_{\alpha/2})\}, \quad (42)$$

which goes to zero as $n \rightarrow \infty$.

A slightly more general parametric testing problem is to decide between $H_0 : \theta \in \Theta_0$ versus $H_1 : \theta \in \Theta_1$, where Θ_0 and Θ_1 are subsets of Θ . It need not be the case that $\Theta_0 \cup \Theta_1 = \Theta$, but it must be that $\Theta_0 \cap \Theta_1 = \emptyset$. Our answer is typically found by finding a confidence interval (or set, or ...) C of a prescribed asymptotic level $(1 - \alpha)$ such that $P_{H_0}\{\theta \in C\} \approx 1 - \alpha$, and hopefully $P_{H_1}\{\theta \in C\}$ is small. If $C \cap \Theta_0 = \emptyset$ then reject H_0 , else accept H_1 .