

# STA 532: Theory of Statistical Inference

Robert L. Wolpert  
Department of Statistical Science  
Duke University, Durham, NC, USA

## 4 Parametric Inference II

In this section we will take  $X = \{X_1, \dots\}$  to be a sequence of independent random variables all with the same probability distribution (so they are “iid”), from some parametric family  $\mathfrak{F} = \{F(x | \theta) : \theta \in \Theta \subset \mathbb{R}^k\}$  of distributions that either have pdfs or pmfs, either of which we will denote by  $f(x | \theta)$ . We will look at the behavior of sequences of estimators  $T_n$  of  $\theta$  based on the first  $n$  observations.

Define the “Kullback-Leibler Divergence” from a distribution with pdf  $f(x)$  to another  $g(x)$  to be

$$\text{KL}(f : g) := \int_{\mathcal{X}} f(x) \log \frac{f(x)}{g(x)} dx. \quad (1)$$

Although it isn’t symmetric and doesn’t obey the triangle inequality, and so can’t be a true “distance” measure, the KL divergence has many of the properties of a metric. In particular,  $\text{KL}(f : f) = 0$  for any  $f$  and  $0 < \text{KL}(f : g) \leq \infty$  for any distinct densities  $f, g$ , because

$$\text{KL}(f : g) = \int_{\mathcal{X}} f(x) \left[ -\log \frac{g(x)}{f(x)} \right] dx \geq \int_{\mathcal{X}} f(x) \left[ 1 - \frac{g(x)}{f(x)} \right] dx = 0$$

since  $-\log y \geq (1 - y)$  for all  $y > 0$  (because  $e^x \geq x + 1$  for all  $x \in \mathbb{R}$ ). For pdfs  $f(x | \theta)$  and  $f(x | \theta')$  from distributions in the parametric family  $\mathfrak{F}$ , we simplify the notation by writing  $\text{KL}(\theta : \theta') := \text{KL}(f(x | \theta) : f(x | \theta'))$ . Some denote this as  $\text{KL}(\theta || \theta')$ .

The family  $\mathfrak{F}$  is called “identifiable” if  $\text{KL}(\theta : \theta') > 0$  whenever  $\theta \neq \theta'$ , *i.e.*, when different parameters lead to different distributions. This would *not* hold for the Normal  $\text{No}(\mu, \sigma^2)$  distribution if we took  $\sigma \in \mathbb{R}$ , for example, even though the normal pdf is well-defined for all  $(\mu, \sigma) \in \mathbb{R}^2$ . From now on we will assume all our parametric families are identifiable.

### Asymptotic Properties of MLEs

Suppose the observations  $\{X_i\} \stackrel{\text{iid}}{\sim} F(x | \theta_*)$  for some  $\theta_* \in \Theta$ , and let  $\hat{\theta}_n$  be the MLE based on the first  $n$  observations, *i.e.*, the value of  $\theta \in \Theta$  that maximizes  $\mathcal{L}_n(\theta) := \prod_{j \leq n} f(x_j | \theta)$  or, equivalently, that maximizes  $\ell_n(\theta) := \sum_{j \leq n} \log f(x_j | \theta)$ .

### Consistency

Since  $\theta_*$  and  $n$  are just constants, the MLE  $\hat{\theta}_n$  also *minimizes*

$$n^{-1}[\ell_n(\theta_*) - \ell_n(\theta)] = \frac{1}{n} \sum \log \frac{f(x_i | \theta_*)}{f(x_i | \theta)}$$

which, by the law of large numbers, converges as  $n \rightarrow \infty$  to

$$\int_{\mathcal{X}} \log \frac{f(x | \theta_*)}{f(x | \theta)} f(x | \theta_*) dx = \text{KL}(\theta_* : \theta).$$

But this limit is minimized at  $\theta = \theta_*$ , where  $\text{KL}(\theta_* : \theta)$  attains its minimum value of zero. The details of proving  $\hat{\theta}_n \rightarrow \theta_*$  for every model under very weak conditions are a bit more involved, ask me if you'd like more details.

### Asymptotic Normality

#### Fisher Information

The  $k$ -dimensional random vector  $Z(X; \theta)$  with components  $Z_i(X; \theta) = (\partial/\partial\theta_i) \log f(X | \theta)$  for  $1 \leq i \leq k$  is called the “score function”. It isn't a *statistic*, because it depends explicitly on the parameter  $\theta$ . The mean is  $\mathbb{E}[Z(X; \theta)] = 0$  because

$$\begin{aligned} \mathbb{E}Z_i(X; \theta) &= \int_{\mathcal{X}} \{(\partial/\partial\theta_i) \log f(x | \theta)\} f(x | \theta) dx \\ &= \int_{\mathcal{X}} \frac{(\partial/\partial\theta_i)f(x | \theta)}{f(x | \theta)} f(x | \theta) dx \\ &= (\partial/\partial\theta_i) \int_{\mathcal{X}} f(x | \theta) dx = (\partial/\partial\theta_i)1 = 0, \end{aligned} \tag{2}$$

since any pdf integrates identically to one. The  $(k \times k)$  *covariance* matrix

$$I(\theta) = \mathbb{E}_{\theta}[Z(X; \theta)Z(X; \theta)'] \tag{3a}$$

(here  $Z(X; \theta)'$  denotes the  $(1 \times k)$  transpose of the  $(k \times 1)$  vector  $Z(X; \theta)$ ) is called the “Fisher Information”. Differentiating (2) w.r.t  $\theta_j$  gives

$$\begin{aligned} 0 &= (\partial/\partial\theta_j) \int_{\mathcal{X}} \{(\partial/\partial\theta_i) \log f(x | \theta)\} f(x | \theta) dx \\ &= \int_{\mathcal{X}} \{(\partial^2/\partial\theta_i\partial\theta_j) \log f(x | \theta)\} f(x | \theta) dx + \int_{\mathcal{X}} \{(\partial/\partial\theta_i) \log f(x | \theta)\} (\partial/\partial\theta_j)f(x | \theta) dx \\ &= \int_{\mathcal{X}} \{(\partial^2/\partial\theta_i\partial\theta_j) \log f(x | \theta)\} f(x | \theta) dx + \mathbb{E}_{\theta}[Z_i(X; \theta)Z_j(X; \theta)], \end{aligned}$$

showing that the  $(k \times k)$  Fisher Information can also be written as

$$I_{ij}(\theta) = \mathbb{E}_{\theta} \left[ \frac{-\partial^2}{\partial\theta_i\partial\theta_j} \log f(X | \theta) \right]; \tag{3b}$$

sometimes one is easier to compute, sometimes the other. For a multivariate normal  $\text{No}(\mu, \Sigma)$  distribution with uncertain mean  $\mu \in \mathbb{R}^k$  and known  $(k \times k)$  covariance matrix  $\Sigma$ , for example, the log likelihood is

$$-\frac{1}{2} \log \det(2\pi\Sigma) - \frac{1}{2}(X - \mu)' \Sigma^{-1}(X - \mu)$$

whose negative second partial derivative wrt  $\mu_i$  and  $\mu_j$  is  $\Sigma_{ij}^{-1}$ , so  $I(\mu) := \Sigma^{-1}$  for the MVN distribution. This is worth remembering: in some sense  $I(\theta)$  acts like a *precision matrix* (inverse of a covariance) for  $\hat{\theta}_n$  in almost all problems, as we'll see below.

The score function for the one-parameter Poisson  $\text{Po}(\lambda)$  distribution is

$$Z(X; \lambda) = (\partial/\partial\lambda)[X \log \lambda - \log X! - \lambda] = X/\lambda - 1,$$

so  $E_\lambda[Z(X; \lambda)] = 0$  (as always) and the Fisher information is  $I(\lambda) = E_\lambda Z(X; \lambda)^2 = 1/\lambda$ , again one over the variance of  $X$ .

The Fisher Information for a sample of size  $n$  is just  $n$  times that for a sample of size one,  $I_n(\theta) = nI(\theta)$ . Fisher Information depends on the particular choice of parameterization: if  $\eta = H(\theta)$  is another parameterization, then the Fisher informations for the two parameterization are related by

$$I^\theta(\theta) = J(\theta)^2 I^\eta(H(\theta))$$

where  $J(\theta) = \partial H(\theta)/\partial\theta$  is the Jacobian of the transformation  $\theta \mapsto \eta = H(\theta)$ . In  $k > 1$  a matrix version of the same formula applies:  $I^\theta(\theta) = J(\theta)' I^\eta(H(\theta)) J(\theta)$ .

### MLE and Information

Suppose that the likelihood function  $\mathcal{L}$  or equivalently the joint pdf  $f(X | \theta)$  has at least two continuous derivatives in  $\theta$ , and that for increasing samples of size  $n \in \mathbb{N}$  from  $X \sim f(x | \theta_*)$  with  $\theta_*$  in the interior of  $\Theta$  it takes its maximum  $\hat{\theta}_n(X)$  for each  $n$  at a sequence of points  $\hat{\theta}_n \rightarrow \theta_*$  (mild regularity conditions on the statistical model will ensure this). By Taylor's theorem the log likelihood function in a neighborhood of  $\theta_*$  will be approximately

$$0 = \ell'_n(\hat{\theta}) \approx \ell'_n(\theta_*) + [\hat{\theta}_n - \theta_*] \ell''_n(\theta_*), \quad \text{so}$$

$$\sqrt{n}[\hat{\theta}_n - \theta_*] \approx \frac{n^{-1/2} \ell'_n(\hat{\theta})}{-n^{-1} \ell''_n(\theta_*)} \approx \frac{\sqrt{n} \overline{Z(X; \theta)}}{I(\theta)}$$

by (3b), since the gradient of  $\ell_n$  will vanish at  $\theta_*$ . By the CLT,  $\overline{Z(X; \theta)}$  is approximately normally distributed with mean zero and variance  $I(\theta)/n$ , so

$$\sqrt{n}[\hat{\theta}_n - \theta_*] \approx \text{No}(0, 1/I(\theta_*)). \quad (4)$$

A similar result holds for multivariate  $\theta \in \Theta \subset \mathbb{R}^k$ , where  $\sqrt{n} [\hat{\theta}_n - \theta_*]$  has approximately a  $k$ -variate  $\text{No}(0, I(\theta_*)^{-1})$  distribution.

This leads to an asymptotic  $100\gamma\%$  confidence intervals of the form

$$\left[ \hat{\theta}_n - \frac{z_{\alpha/2}}{\sqrt{nI(\hat{\theta}_n)}}, \quad \hat{\theta}_n + \frac{z_{\alpha/2}}{\sqrt{nI(\hat{\theta}_n)}} \right]$$

for  $\theta$ , where  $\alpha = 1 - \gamma$  and  $\Phi(-z_{\alpha/2}) = \alpha/2$ .

### Optimality

One form of the “Cramér-Rao” or “Information” Inequality asserts that *any* statistic  $T : \mathcal{X} \rightarrow \mathbb{R}$  with mean  $\psi(\theta) = \mathbf{E}_\theta T(X)$  satisfies

$$\mathbf{V}_\theta T \geq \frac{|\psi'(\theta)|^2}{I(\theta)}$$

for one-dimensional parametric families. Any unbiased estimator  $T(X)$  has mean  $\psi(\theta) = \theta$  with  $|\psi'(\theta)| \equiv 1$ , so every unbiased estimator based on a sample of size  $n$  has  $\text{mse}[T_n] \geq 1/nI(\theta)$ . By (4) the MLE achieves this lower bound asymptotically, so it is optimal in the sense that no other asymptotically unbiased estimator has strictly lower mse. Also we see that, under mild regularity,  $\text{mse}(T_n)$  decreases no faster than rate  $\propto 1/n$  for *any* sequence  $T_n$  of estimators.

### Cautionary Note

Sometimes “mild regularity” doesn’t hold. For example, consider estimating the parameter  $\theta$  for  $\{X_i\} \stackrel{\text{iid}}{\sim} \text{Un}(0, \theta)$  on the basis of the first  $n$  observations  $X = (X_1, \dots, X_n)$ . The likelihood function

$$\mathcal{L}_n(\theta | X) = \prod_{i=1}^n \{\theta^{-1} \mathbf{1}_{[0, \theta]}(X_i)\} = \theta^{-n} \mathbf{1}_{[0, \theta]}(X_n^*)$$

vanishes for  $\theta < X_n^*$  and decreases monotonically for  $\theta > X_n^*$ , the sufficient statistic  $X_n^* := \max(X_i : 1 \leq i \leq n)$ , so the MLE is in fact  $\hat{\theta}_n = X_n^*$ . The CDF for  $X_n^*$  is

$$F_n(x) = \mathbf{P}_\theta[X_n^* \leq x] = \{\mathbf{P}_\theta[X_1 \leq x]\}^n = (x/\theta)^n$$

for  $0 \leq x \leq \theta$ , and zero for  $x < 0$ , one for  $x > \theta$ . Thus the MLE has mean square error

$$\begin{aligned} \text{mse}(\theta) &= \mathbf{E}_\theta |X_n^* - \theta|^2 \\ &= \int_0^\theta (x - \theta)^2 n x^{n-1} \theta^{-n} dx \\ &= 2\theta^2 / (n+1)(n+2), \end{aligned}$$

shrinking at a rate faster than  $\propto 1/n$ . The error  $[\hat{\theta}_n - \theta]$  does not have a Normal asymptotic distribution. Rather, it is always negative (since  $\mathbf{P}_\theta[\hat{\theta}_n < \theta] = 1$ ) and, for  $z > 0$ , satisfies

$$\begin{aligned} \mathbf{P}_\theta \left\{ n[\hat{\theta}_n - \theta] \leq -z \right\} &= \mathbf{P}_\theta \left\{ \hat{\theta}_n \leq \theta - z/n \right\} \\ &= (1 - z/n\theta)^n \rightarrow \exp(-z/\theta). \end{aligned}$$

Thus  $n[\theta - \hat{\theta}_n]$  has approximately the Weibull  $\text{We}(1, 1/\theta)$  distribution for large  $n$ , with mean  $\theta$  and variance  $\theta^2$ . For more on the Information Inequality, see (or click on)

<https://stat.duke.edu/courses/Spring16/sta532/lec/fish.pdf>

### The Delta Method

Let  $\tau = g(\theta)$  for some smooth function  $g : \Theta \rightarrow \mathbb{R}$ , and let  $\{T_n\}$  be a consistent sequence of estimators of  $\theta$  (like the MLEs  $\hat{\theta}_n$  based on the first  $n$  observations) that are asymptotically normal in the sense that

$$\sqrt{n}[T_n - \theta] \approx \text{No}(0, \sigma^2)$$

for large  $n$ , for some constant  $\sigma^2 > 0$ . By Taylor's theorem with remainder, there exist numbers  $\tilde{\theta}_n$  between  $\theta$  and  $T_n$  such that

$$\begin{aligned} g(T_n) &= g(\theta) + g'(\tilde{\theta}_n)[T_n - \theta], \quad \text{so} \\ \sqrt{n}[g(T_n) - g(\theta)] &= g'(\tilde{\theta}_n)\sqrt{n}[T_n - \theta] \\ &\approx \text{No}(0, |g'(\theta)|^2\sigma^2), \end{aligned}$$

*i.e.*,  $g(T_n)$  is approximately normally-distributed with mean  $g(\theta)$  and variance  $|g'(\theta)|^2\sigma^2/n$ .

Example: Let  $S_n \sim \text{Po}(n\lambda)$  be the sum of  $n$  iid RVs  $X_i \sim \text{Po}(\lambda)$  and let  $\bar{X}_n := S_n/n$  be their sample mean. Then  $\bar{X}_n \rightarrow \lambda$  by the LLN and  $\sqrt{n}(\bar{X}_n - \lambda) \approx \text{No}(0, \lambda)$  by the CLT. Set  $g(x) := \log(1 + x)$ . By the Delta Method, the log of one plus the sample mean  $g(X_n/n) = \log(1 + \bar{X}_n)$  is approximately normally-distributed with mean and variance

$$\mathbb{E}[g(\bar{X}_n)] \approx \log(1 + \lambda) \quad \mathbb{V}[g(\bar{X}_n)] \approx \frac{\lambda}{n(1 + \lambda)^2}.$$

A similar result holds for multivariate distributions with  $\Theta \subset \mathbb{R}^k$ .

**Warning:** What goes wrong in a similar argument with the functions  $g(x) = \log(x)$  or  $g(x) = 1/x$ ? Note that success of the Delta Method hinges on the assumption that, for large  $n$ , all the “action” is close to  $T_n \approx \theta$ , and that nothing too bad happens away from there.

### The Parametric Bootstrap

Earlier we computed approximate features of the sampling distributions of statistics  $T : \mathcal{X} \rightarrow \mathbb{R}$  by replicating draws  $X_j^{(b)}$  from the empirical distribution  $\hat{F}_n$  (*i.e.*, by draws with replacement from the random sample). In parametric models an alternative is to estimate  $\theta \in \Theta$  by some estimator  $\tilde{\theta}_n(X)$  (for example, the MLE or a Bayes estimator), then use simulation to approximate the sampling distribution of  $T(X)$  by that of the “parametric bootstrap sample”  $\{T(X_i^{(b)}) : X_i^{(b)} \stackrel{\text{iid}}{\sim} \hat{F}_n(x | \tilde{\theta}_n)\}$ . This approach is typically more efficient than the usual nonparametric bootstrap when the parametric model is correct, but can be grossly misleading if the model is misspecified. It is prudent to check model assumptions, at least informally.

### Checking Model Assumptions

Famously George Box quipped “All models are wrong, but some are useful” and, less famously,

Remember that all models are wrong; the practical question is how wrong do they have to be to not be useful.

One way to explore “how wrong” a model is is to make an overlay plot of a histogram of the data along with the pdf  $f(x | \hat{\theta}_n)$  at the estimated value of  $\theta \in \Theta$ . If these differ grossly, the model may not be useful. Another is to make a “Q-Q plot” of data quantiles (on the vertical axis) *vs.* theoretical quantiles (on the horizontal axis) for the best-fitting distribution from the model family. For example, Figure (1) shows a Q-Q plot for a sample of  $N = 100$  replicates for a normal  $\text{No}(\mu, \sigma^2)$  model, along with a histogram (same sample) with best  $\text{No}(\mu, \sigma^2)$  pdf fit overlaid. In fact these

data were generated from the  $t_2$  distribution, with substantially heavier tails; the sigmoid shape of the curve is a tip-off about this problem, as are the outliers above the  $x = y$  line on the right and below it on the left.

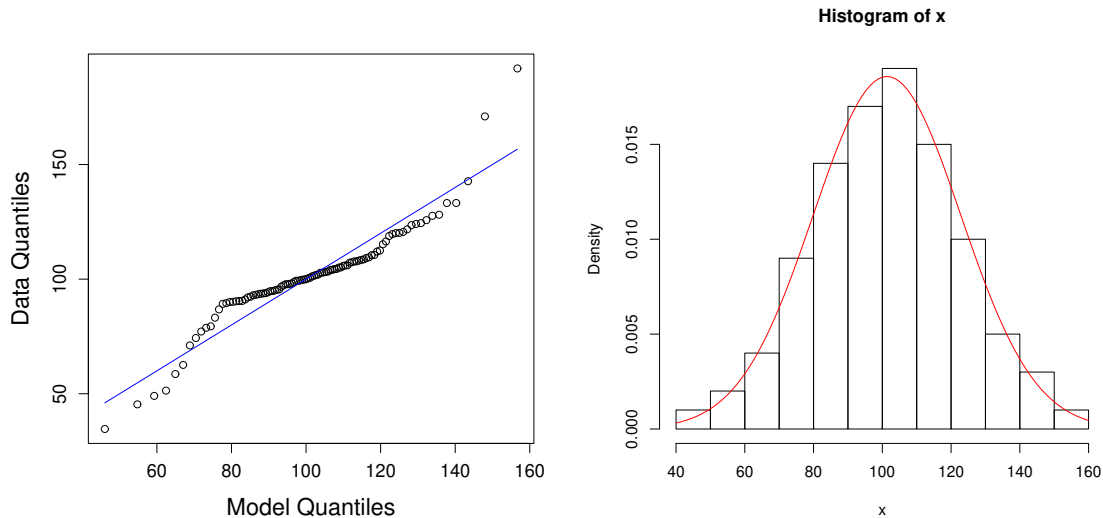


Figure 1: Q-Q plot and histogram to explore how well a  $\text{No}(\mu, \sigma^2)$  model fits data  $\{X_i\}$ . In fact, data were generated with  $t_2(100, 10)$  distribution with heavier tails than the normal, but it's hard to see that from the histogram.

A more formal approach is to construct a “goodness of fit” test to check the hypothesis that the data come from the parametric family; we’ll discuss that more when we study hypothesis testing in a few weeks.

## Sufficiency

In some statistical models one can find a low-dimensional statistic  $T(X)$  that captures all the information that a sample  $X = \{X_1, \dots, X_n\}$  has about the uncertain parameter  $\theta \in \Theta$ . For example, if  $\{X_i\} \stackrel{\text{iid}}{\sim} \text{Po}(\lambda)$ , then the likelihood function

$$\begin{aligned} \mathcal{L}(\lambda) &\propto \prod_{i=1}^n \left[ \lambda^{X_i} e^{-\lambda} \right] \\ &= \lambda^{n\bar{X}_n} e^{-n\lambda} \end{aligned}$$

depends on the data *only* through  $\bar{X}_n$  (or, equivalently, through  $\sum X_i$ ). A statistic  $T(X)$  is said to be *sufficient* (for  $\theta \in \Theta$ ) if the conditional distribution of the data vector  $X$ , given the value  $t = T(X)$  of the statistic, does not depend on  $\theta$ . In that case there is no more to be learned about  $\theta$  from the entire vector than there is from  $T(X)$  alone.

**More examples:**

If  $X$  is a vector of  $n$  iid Bernoulli  $\text{Bi}(1, p)$  random variables, then its joint pmf is

$$f(x | p) = \prod_{i=1}^n p^{x_i} (1-p)^{1-x_i} = p^{S_n} (1-p)^{n-S_n},$$

a function that depends on  $X$  only through  $S_n = \sum_{i=1}^n X_i$  or (equivalently) of  $\bar{X}_n = S_n/n$ .

If  $X$  is a vector of  $n$  iid Beta  $\text{Be}(\alpha, \beta)$  random variables, then you can show that the two-dimensional statistic  $T(X) := [\sum \log X_i, \sum \log(1 - X_i)]$  is sufficient (you should do that).

If  $X$  is a vector of  $n$  iid Gamma  $\text{Ga}(\alpha, \lambda)$  random variables, then you can show that the two-dimensional statistic  $T(X) := [\sum \log X_i, \sum X_i]$  is sufficient (do that too). If  $\alpha$  is known then  $\bar{X}_n$  is sufficient for  $\lambda$ ; if  $\lambda$  is known then  $\log \bar{X}_n$  is sufficient for  $\alpha$ .

For iid Normal random variables  $\{X_i\} \stackrel{\text{iid}}{\sim} \text{No}(\mu, \sigma^2)$ , the sample mean  $\bar{X}_n$  and sample variance  $S_n^2$  are easily shown to be sufficient.

**Two caveats:**

If  $X$  is a vector of  $n$  iid Student  $t_\nu(m, s)$  random variables with  $\nu$  degrees of freedom, with location  $m$  and scale  $s$ , then there is no sufficient statistic for  $m$  (whether or not  $\nu$  and  $s$  are known) with dimension less than  $n$ . The ordered sample  $[X_{(1)} < X_{(2)} < \dots < X_{(n)}]$  is sufficient, for example, but not any low-dimensional collection of moments as there were for the Poisson, Bernoulli, Beta, Gamma, and Normal examples. Thus, (1) sufficient statistics don't always exist, and (2) they are highly model-dependent. If one investigator models  $\{X_i\}$  as Normal and reports only the (sufficient) sample mean and variance, another investigator concerned that the data may have heavier tails will be unable to estimate the parameters of a  $t_\nu$  distribution from these statistics.

Finally— sometimes reporting only sufficient statistics will lose the opportunity to check the model. In the Bernoulli example, the likelihood function for a result of ten successes and four failures will be  $\mathcal{L}(p) \propto p^{10}(1-p)^4$ , for example, but if the original data were  $[0, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1]$  then there is reason to doubt the iid model, and ask questions like “did something happen between the fourth and fifth observations?”

**Is  $T$  Sufficient?**

A statistic  $T$  is sufficient for  $\theta$  in a family  $\{f(x | \theta) : \theta \in \Theta\}$  if and only if the pdf or pmf factors in the form of a product

$$f(x | \theta) = h(x) g(T(x), \theta)$$

of a function of  $x$  (that doesn't depend on  $\theta$ ) and a function of  $T(x)$  and  $\theta$  (that doesn't depend on  $x$  except through  $T(x)$ ). This happens if and only if the *likelihood* function can be written as a function of only  $T(X)$  and  $\theta$ . Here's the most common place this happens:

### Examples: Exponential Families

The pdf  $f(x | \theta)$  for many parametric families of distributions can be written in the form

$$f(x | \theta) = h(x) \exp(\eta(\theta)T(x) - B(\theta)) \quad (5)$$

for some functions  $h(x)$ ,  $T(x)$  of  $x \in \mathcal{X}$  and  $\eta(\theta)$ ,  $B(\theta)$  of  $\theta \in \Theta$ . For example:

Po( $\lambda$ )	$h(x) = \frac{1}{x!}$	$T(x) = x$	$\eta(\lambda) = \log \lambda$	$B(\lambda) = \lambda$
Ex( $\lambda$ )	$h(x) = 1$	$T(x) = x$	$\eta(\lambda) = -\lambda$	$B(\lambda) = -\log \lambda$
No( $\theta, 1^2$ )	$h(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$	$T(x) = x$	$\eta(\theta) = \theta$	$B(\theta) = \theta^2/2$
No( $0, \theta^2$ )	$h(x) = \frac{1}{\sqrt{2\pi}}$	$T(x) = x^2$	$\eta(\theta) = -1/2\theta^2$	$B(\theta) = \log \theta$
Bi( $m, p$ )	$h(x) = \binom{m}{x}$	$T(x) = x$	$\eta(p) = \log \frac{p}{1-p}$	$B(p) = -m \log(1-p)$

Note that if iid  $\{X_i\}$  come from an exponential family with functions  $h$ ,  $T$ ,  $\eta$ ,  $B$  then the vector  $X = (X_1, \dots, X_n)$  has an exponential family distribution too, with

$$h_n(x) = \prod h(x_i) \quad T_n(x) = \sum T(x_i) \quad \eta_n(\theta) = \eta(\theta) \quad B_n(\theta) = nB(\theta)$$

Exponential families are especially pleasant to work with, for several reasons:

- They have a natural sufficient statistic,  $T$  (or  $\sum T(x_i)$  for vectors), whose mean  $B'(\theta)/\eta'(\theta)$  is easy to compute from the score  $Z_\theta = \eta'(\theta)T(X) - B'(\theta)$ ;
- Their Fisher information  $I(\theta) = B''(\theta) - \eta''(\theta)B'(\theta)/\eta'(\theta)$  is easy to compute, especially for “natural” exponential families where  $\eta(\theta) = \theta$  and so  $\eta'' \equiv 0$  and  $I(\theta) = B''(\theta)$ ;
- The MLE  $\hat{\theta}_n$  can be found by solving  $\bar{T}_n = E_{\hat{\theta}_n} T = B'(\hat{\theta}_n)/\eta'(\hat{\theta}_n)$ ;
- They have conjugate families of prior distributions (see below).

Note the name “exponential family” has nothing to do with the “exponential distribution” (except that Ex( $\lambda$ ) is one). Exponential families with multivariate parameter  $\theta \in \Theta \subset \mathbb{R}^k$  arise too (usually with  $k = 2$  or so); examples include No( $\mu, \sigma^2$ ), Ga( $\alpha, \lambda$ ), and Be( $\alpha, \beta$ ). The only change in the definition (5) is that now  $T(x)$  and  $\eta(\theta)$  are  $k$ -dimensional vector functions, and that the scalar product “ $\eta(\theta)T(x)$ ” must be replaced with a dot product “ $\eta(\theta) \cdot T(x)$ ”. The sufficient statistic  $T$  has the same dimension ( $k$ ) as the parameter space  $\Theta$ .

Many distributions commonly used in modeling do have an exponential family representation, for fixed values of all but one (or sometimes two) of their parameters. Those include the beta, binomial, exponential, gamma, geometric, negative binomial, normal, Pareto, Poisson, and Weibull. Sufficient statistics and conjugate prior distributions are available for each of these. Many other useful distributions do *not* have exponential family form, and typically these will not have sufficient statistics or conjugate priors—including the logistic, Student  $t$ ,  $\alpha$ -stable, and others. One way to check is to write out the log likelihood function and see if its dependence on  $\theta$  and  $x$  includes only terms that are functions of either  $\theta$  or  $x$  but not both, plus a product (or dot-product) of a function of  $\theta$  (that will be  $\eta$ ) and a function of  $x$  (that will be  $T(x)$ ).



### Efficiency and Sufficiency

Let  $T(\mathbf{x})$  be an estimator of a parameter  $\theta$ , with MSE

$$m(\theta) = \mathbb{E}_\theta [ |T(\mathbf{x}) - \theta|^2 ].$$

The celebrated theorem of Raô and Blackwell asserts that if  $T$  does *not* depend on the data  $\mathbf{x}$  through a sufficient statistic, then it can be improved by replacing it with one that does, in the following sense.

**Theorem 1** *Let  $\mathcal{F} = \{f_\theta(\mathbf{x}) : \theta \in \Theta \subset \mathbb{R}, \mathbf{x} \in \mathcal{X}\}$  be a parametric statistical model. Let  $T : \mathcal{X} \rightarrow \mathbb{R}$  be any statistic and let  $S : \mathcal{X} \rightarrow \mathcal{S}$  be a sufficient statistic. Set*

$$T^*(\mathbf{x}) := \mathbb{E}_\theta [ T(\mathbf{x}) \mid S(\mathbf{x}) ],$$

*the conditional expectation of  $T(\mathbf{x})$  given the value  $s = S(\mathbf{x})$  of the sufficient statistic. Then  $T^*$  is a statistic (i.e., does not depend on  $\theta$ ) with mse*

$$\mathbb{E}_\theta [ |T^*(\mathbf{x}) - \theta|^2 ] \leq \mathbb{E}_\theta [ |T(\mathbf{x}) - \theta|^2 ],$$

*with strict inequality unless  $T(\mathbf{x})$  already depends only on  $S(\mathbf{x})$ , almost-surely.*

**Proof.** Because  $S$  is sufficient, the conditional distribution of  $\mathbf{x}$  given  $S(\mathbf{x})$  doesn't depend on  $\theta$ , so for any  $t \in \mathbb{R}$

$$\mathbb{P}_\theta [ T^*(\mathbf{x}) \leq t \mid S(\mathbf{x}) ]$$

is a function  $g(s)$  of the value  $s = S(\mathbf{x})$  that doesn't depend on  $\theta$  either— *i.e.*,  $T^*$  is a statistic. The MSE of  $T$  is

$$\begin{aligned} \mathbb{E}_\theta [ |T(\mathbf{x}) - \theta|^2 ] &= \mathbb{E}_\theta [ |T(\mathbf{x}) - T^*(\mathbf{x}) + T^*(\mathbf{x}) - \theta|^2 ] \\ &= \mathbb{E}_\theta [ |T(\mathbf{x}) - T^*(\mathbf{x})|^2 ] + \mathbb{E}_\theta [ |T^*(\mathbf{x}) - \theta|^2 ] \end{aligned}$$

because the cross-product  $\mathbb{E}_\theta (T(\mathbf{x}) - T^*(\mathbf{x})) (T^*(\mathbf{x}) - \theta)$  vanishes

$$\geq \mathbb{E}_\theta [ |T^*(\mathbf{x}) - \theta|^2 ],$$

with equality only if  $\mathbb{E}_\theta [ |T(\mathbf{x}) - T^*(\mathbf{x})|^2 ] = 0$ — *i.e.*, only if  $\mathbb{P}_\theta [ T = T^* ] = 1$ , a function of  $S(\mathbf{x})$ . □

Thus it's always best to use estimators that depend on the data only through sufficient statistics, when possible.

#### Examples

Let  $\mathbf{x} = \{X_i\}_{1 \leq i \leq 9} \stackrel{\text{iid}}{\sim} \text{No}(\theta, 1)$  be a sample of size  $n = 9$  from the normal distribution with unit variance. Then  $\bar{X}_9$  and  $X_1$  are both unbiased estimators of  $\theta$ , but  $\bar{X}_9$  is a function of the sufficient statistic  $S_9 = \sum X_i$  and  $X_1$  is not. This is reflected in their MSEs,

$$\mathbb{E} | \bar{X}_9 - \theta |^2 = 1/9 < 1 = \mathbb{E} | X_1 - \theta |^2.$$

The inferior estimator  $X_1$  can be improved, however, by the Raô-Blackwell procedure, to

$$X_1^* := \mathbb{E}_\theta[X_1 | S_9] = S_9/9 = \bar{X}_9$$

where we evaluated  $\mathbb{E}_\theta[X_1 | S_9] = S_9/9$  by invoking symmetry. Similarly, the sample median  $\hat{X}_9$  has conditional expectation  $\hat{X}_9^* := \mathbb{E}_\theta[\hat{X}_9 | S_9] = \bar{X}_9$  with lower MSE. The sample mean  $\bar{X}_n$  cannot be improved, however; a naïve attempt at doing so yields

$$\bar{X}_9^* := \mathbb{E}_\theta[\bar{X}_9 | S_9] = S_9/9 = \bar{X}_9,$$

unchanged.

### Conjugate Prior Distributions

A family  $\{\pi_\xi(\theta) : \xi \in \Xi\}$  of prior density (or pmf) functions on  $\Theta$  is called “conjugate” for a statistical model  $\{f(x | \theta) : \theta \in \Theta\}$  if the posterior distribution upon observing  $X \in \mathcal{X}$  is again in the family  $\{\pi_\xi(\theta) : \xi \in \Xi\}$ , *i.e.*, if for each  $\xi \in \Xi$  and  $x \in \mathcal{X}$  there is some  $\xi^{(x)} \in \Xi$  such that

$$\pi_\xi(\theta)f(x | \theta) \propto \pi_{\xi^{(x)}}(\theta).$$

For example, the family  $\text{Be}(\alpha, \beta)$  of prior distributions for  $\theta \in \Theta = [0, 1]$  is conjugate for the Binomial statistical model  $X \sim \text{Bi}(n, \theta)$  for any fixed  $n \in \mathbb{N}$ : if  $\theta \sim \text{Be}(\alpha, \beta)$  and  $X \sim \text{Bi}(n, \theta)$ , then the posterior distribution is  $\theta | X \sim \text{Be}(\alpha + X, \beta + n - X)$ , *i.e.*,

$$\begin{aligned} \left\{ \pi_\xi(\theta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1 - \theta)^{\beta-1} \right\} \left\{ f(x | \theta) = \binom{n}{x} \theta^x (1 - \theta)^{n-x} \right\} \\ \propto \theta^{\alpha+x-1} (1 - \theta)^{\beta+n-x-1} \propto \left\{ \pi_{\xi^{(x)}}(\theta) = \frac{\Gamma(\alpha + \beta + n)}{\Gamma(\alpha + x)\Gamma(\beta + n - x)} \theta^{\alpha+x-1} (1 - \theta)^{\beta+n-x-1} \right\} \end{aligned}$$

Similarly if  $\lambda \sim \text{Ga}(\alpha, \beta)$  and  $X \sim \text{Po}(\lambda)$ , then  $\lambda | X \sim \Gamma(\alpha + X, \beta + 1)$ . Note that both binomial and Poisson are exponential family distributions— *every* exponential family (see (5)) has a conjugate family of prior distributions of the form

$$\pi_\xi(\theta) = c_\xi^{-1} \pi_\star(\theta) \exp(\eta(\theta)\xi_1 - \xi_2 B(\theta))$$

for any function  $\pi_\star(\theta) \geq 0$  for which the set  $\Xi$  of  $\xi = (\xi_1, \xi_2)$  for which

$$0 < c_\xi := \int_\Theta \pi_\star(\theta) \exp(\eta(\theta)\xi_1 - \xi_2 B(\theta)) d\theta < \infty$$

is non-empty and closed under the addition of  $(t, 1)$  for any possible value  $t$  of  $T(X)$ . Evidently  $\xi^{(x)} = (\xi_1 + T(X), \xi_2 + 1)$  or, for a sample of size  $n$ ,  $\xi^{(\mathbf{x})} = (\xi_1 + \sum T(X_i), \xi_2 + n)$ .

Conjugate prior distributions simplify analyses by making it unnecessary to do any numerical integration to perform Bayesian statistical inference, *i.e.*, to evaluate features of posterior probability distributions (means, credible sets, *etc.*) For the first two centuries after the emergence of Bayesian methods, this simplification was critical— and led investigators to limit themselves to conjugate prior distributions or low-dimensional models.

With the emergence in the 1990s of hardware and algorithms capable of evaluating essentially any posterior distribution it’s no longer necessary to limit attention to these families. This permits modelers to use more realistic models with both more realistic sampling distributions and more realistic prior distributions.

## 5 Summary: Desirable Properties for Estimators

We have seen a number of properties that an estimator  $T(X)$ , or a sequence of estimators  $T_n(\mathbf{x})$ , might have in a parametric statistical model  $\mathfrak{F} = \{f(x | \theta) : \theta \in \Theta\}$ . Some of these are:

**Consistent:** A sequence  $T_n(\mathbf{x})$  of estimators of a feature  $\tau = g(\theta)$  of  $\theta \in \Theta$  is *consistent* if

$$\lim_{n \rightarrow \infty} \mathbb{P}_\theta[ |T_n(\mathbf{x}) - \tau| > \epsilon ] = 0$$

for every  $\epsilon > 0$ .

**MS Consistent:** A sequence  $T_n(\mathbf{x})$  of estimators of a feature  $\tau = g(\theta)$  of  $\theta \in \Theta$  is *mean square consistent* (or MS consistent) if

$$\lim_{n \rightarrow \infty} \mathbb{E}_\theta |T_n(\mathbf{x}) - \tau|^2 = 0.$$

This is a stronger condition than consistency, *i.e.*, if  $T_n(\mathbf{x})$  is MS consistent then it is also consistent, since  $\mathbb{P}_\theta[ |T_n(\mathbf{x}) - \tau| > \epsilon ] \leq \mathbb{E}_\theta |T_n(\mathbf{x}) - \tau|^2 / \epsilon^2$  for any  $\epsilon > 0$ .

**Efficient:** An estimator  $T(X)$  with mean  $m(\theta) = \mathbb{E}_\theta T(X)$  is *efficient* if it satisfies the Information Inequality lower bound,

$$\mathbb{E}_\theta |T(X) - \tau|^2 = \frac{[m'(\theta)]^2}{I(\theta)} + [m(\theta) - \tau]^2.$$

In particular, an unbiased estimator  $T(X)$  of  $\theta$  is efficient if  $\mathbb{E}_\theta |T(X) - \theta|^2 = 1/I(\theta)$ . Some authors limit the definition of “efficient” to this case.

**Asymptotically Efficient:** A sequence  $T_n(\mathbf{x})$  of estimators of  $\theta$  is *asymptotically efficient* if it satisfies

$$\lim_{n \rightarrow \infty} n \mathbb{E}_\theta |T_n(\mathbf{x}) - \theta|^2 = \frac{1}{I(\theta)}$$

**Unbiased:** An estimator  $T(X)$  of a feature  $\tau = g(\theta)$  is *unbiased* if  $\mathbb{E}_\theta T(X) = \tau$  for every  $\theta \in \Theta$ .

**Asymptotically Unbiased:** A sequence  $T_n(\mathbf{x})$  of estimators of  $\tau = g(\theta)$  is *asymptotically unbiased* if  $\lim_{n \rightarrow \infty} \mathbb{E}_\theta T_n(\mathbf{x}) = \tau$  for every  $\theta \in \Theta$ .

**Sufficient:** A statistic  $S$  is *sufficient* if the conditional distribution of  $X$ , conditional on  $S(X) = s$ , does not depend on  $\theta$ . Equivalently, the pdf can be written as a product

$$f(x | \theta) = h(x) g(S(x), \theta)$$

of a function  $h(x)$  of the data only, and a function  $g(S(x), \theta)$  that depends on  $\theta$  and also on the data, but only through the value of the statistic  $S(X)$ .

**Admissible:** An estimator  $T(X)$  of a feature  $\tau = g(\theta)$  is *admissible* if there does *not* exist an estimator  $S(X)$  that is “better” than  $T(X)$  in the sense that

$$(\forall \theta \in \Theta) \mathbb{E}_\theta |S(X) - \tau|^2 \leq \mathbb{E}_\theta |T(X) - \tau|^2,$$

with strict inequality for at least one  $\theta \in \Theta$ .