

1. Fisher Information

Let $f(x | \theta)$ be a density function with the property that $\log f(x | \theta)$ is differentiable in θ throughout the open p -dimensional parameter set $\Theta \subset \mathbb{R}^p$; then the **score statistic** (or score function) is defined by

$$Z(X) \equiv \nabla_{\theta} \log f(x | \theta) = \frac{\nabla_{\theta} f(x | \theta)}{f(x | \theta)}$$

and the **Fisher** (or **Expected**) **Information** matrix is defined by

$$I(\theta) \equiv \mathbb{E} [Z(X)Z(X)' | \theta];$$

if we may exchange integration with differentiation then we can calculate

$$\begin{aligned} \mathbb{E}[Z_i(X) | \theta] &= \int_{\mathcal{X}} \left[\frac{d}{d\theta_i} \log f(x | \theta) \right] f(x | \theta) dx \\ &= \int_{\mathcal{X}} \frac{\frac{d}{d\theta_i} f(x | \theta)}{f(x | \theta)} f(x | \theta) dx \\ &= \int_{\mathcal{X}} \frac{d}{d\theta_i} f(x | \theta) dx \\ &= \frac{d}{d\theta_i} \int_{\mathcal{X}} f(x | \theta) dx \\ &= 0 \end{aligned}$$

and hence $\mathbb{E}[Z(X) | \theta] = 0$ and $\text{Cov}[Z(X) | \theta] = \mathbb{E}[Z(X)Z(X)' | \theta] = I(\theta)$; taking another derivative with respect to θ_j of the equation $\mathbb{E}[Z_i(X) | \theta] = 0$ gives, by the product rule,

$$\begin{aligned} 0 &= \frac{d}{d\theta_j} \mathbb{E}[Z_i(X) | \theta] \\ &= \frac{d}{d\theta_j} \int_{\mathcal{X}} \left[\frac{d}{d\theta_i} \log f(x | \theta) \right] [f(x | \theta)] dx \\ &= \int_{\mathcal{X}} \left[\frac{d^2}{d\theta_i d\theta_j} \log f(x | \theta) \right] [f(x | \theta)] dx + \int_{\mathcal{X}} \left[\frac{d}{d\theta_i} \log f(x | \theta) \right] \left[\frac{d}{d\theta_j} f(x | \theta) \right] dx \\ &= \mathbb{E} \left[\frac{d^2}{d\theta_i d\theta_j} \log f(x | \theta) \right] + I(\theta), \end{aligned}$$

so we may also compute the Fisher Information as

$$I(\theta) = \mathbb{E} [-\nabla^2 \log f(x | \theta)]$$

as the matrix of expected negative second derivatives of the log likelihood with respect to θ .

The Fisher Information matrix depends on the parametrization chosen. If we rewrite our model in terms of some other parameter η , then by the vector chain rule $\nabla_\theta = J^\top \nabla_\eta$ (J denotes the Jacobian matrix $J = \frac{\partial \eta}{\partial \theta}$ with components $J_{ij} = \partial \eta_i / \partial \theta_j$), the Fisher Information for the two parametrizations are related by

$$I(\theta) = J^\top I(\eta) J. \quad (1)$$

The score statistic and Fisher Information are given in Natural Exponential Families by

$$\begin{aligned} Z(X) &\equiv \nabla_\eta \log f(X|\eta) \\ &= \nabla_\eta [\eta \cdot t(x) - A(\eta) + \log h(x)] \\ &= t(x) - \nabla A(\eta) \\ I(\eta) &\equiv -\nabla_\eta^2 \mathbb{E}[\log f(X|\eta) | \eta] \\ &= -\nabla_\eta^2 (\eta \cdot \mathbb{E}[t(X)]) + \nabla^2 A(\eta) \\ &= \nabla^2 A(\eta) \end{aligned}$$

1.1. The Information Inequality

Now let $\Theta \subset \mathbb{R}$ be one-dimensional and let T be any statistic with finite expectation $\psi(\theta) \equiv \mathbb{E}[T(X) | \theta]$, and assume additionally that ψ is differentiable throughout Θ to justify exchanging integration and differentiation as follows:

$$\begin{aligned} \psi'(\theta) &= \frac{d}{d\theta} \int_{\mathcal{X}} T(x) f(x | \theta) dx \\ &= \int_{\mathcal{X}} T(x) \frac{d}{d\theta} f(x | \theta) dx \\ &= \int_{\mathcal{X}} T(x) Z(x) f(x | \theta) dx \\ &= \mathbb{E}[T(X) Z(X) | \theta] = \text{Cov}[T(X) Z(X)], \end{aligned}$$

so the score statistic $Z(X) \equiv \frac{d}{d\theta} \log f(x | \theta)$ has mean zero, variance $I(\theta)$, and covariance $\psi'(\theta) = \text{Cov}[T(X), Z(X)]$ with $T(X)$; by the Covariance

Inequality $|\text{Cov}(T, Z)|^2 \leq \mathbf{V}(T)\mathbf{V}(Z)$ (Minkowski's inequality), we can conclude that $|\psi'(\theta)|^2 \leq I(\theta)\mathbf{V}(T(X))$, or that

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

in particular, any unbiased estimator T of θ must have risk

$$R(\theta, T) \geq \frac{1}{I(\theta)}$$

bounded below by the celebrated Information Inequality. This result was commonly referred to as the Cramèr-Rao inequality, until Frechèt's earlier discovery was widely recognized.

2. Bayesian Central Limit Theorem

The *observed information* for a single observation $X = x$ from the model $X \sim f(x|\theta)$ is

$$i(\theta, x) = -\nabla_{\theta}^2 \log f(x|\theta);$$

evidently the Fisher (expected) information is related to this by $I(\theta) = \mathbf{E}[i(\theta, X)|\theta]$. The likelihood for a sample of size n is just the product of the individual likelihoods, leading to a *sum* for the *log* likelihoods, and observed information

$$i(\theta, x) = \sum_{j=1}^n i(\theta, x_j).$$

If the log likelihood $\log f(x|\theta)$ is differentiable throughout Θ and attains a unique maximum at an interior point $\hat{\theta}_n(x) \in \Theta$, then we can expand $\log f(x|\theta)$ in a second-order Taylor series for $\theta = \hat{\theta}_n(x) + \epsilon/\sqrt{n}$ close to $\hat{\theta}_n(x)$ to find

$$\begin{aligned} \log f(x|\theta) &= \log f(x|\hat{\theta}_n) + \frac{(\epsilon/\sqrt{n})^1}{1!} \nabla_{\theta} \log f(x|\hat{\theta}_n) + \frac{(\epsilon/\sqrt{n})^2}{2!} \nabla_{\theta}^2 \log f(x|\hat{\theta}_n) \\ &\quad + o((\epsilon/\sqrt{n})^2 |\nabla_{\theta}^2 \log f(x|\hat{\theta}_n)|) \\ &= \log f(x|\hat{\theta}_n) + 0 - \frac{\epsilon^2}{2} \frac{1}{n} \sum_{j=1}^n i(\hat{\theta}_n, x_j) + o(1) \\ &\rightarrow \log f(x|\hat{\theta}_n) - \frac{\epsilon^2}{2} \mathbf{E}[i(\hat{\theta}_n, X)] \\ &= \log f(x|\hat{\theta}_n) - \frac{\epsilon^2}{2} I(\theta), \end{aligned}$$

where we have used the consistency of $\hat{\theta}_n$ and have applied the strong law of large numbers for $i(\theta, X)$. Thus we have the likelihood approximation $f(x|\theta) \approx \text{No}(\hat{\theta}_n(x), nI(\hat{\theta}_n))$, normal with mean the MLE $\hat{\theta}_n(x)$ and *precision* $nI(\hat{\theta}_n)$ (or covariance $\frac{1}{n}I(\hat{\theta}_n)^{-1}$).

3. Exponential Families

Consider a sample $X = (X_1, \dots, X_n)$ of some number n of independent replicates, all from the same probability distribution with pdf or pmf $f(x | \theta)$ of *exponential family* form

$$f(x | \theta) = \exp \left[\sum_{i=1}^q \eta_i(\theta) t_i(x) - B(\theta) \right] h(x);$$

since $1 \equiv \int_{\mathcal{X}} f(x | \theta) dx = e^{-B(\theta)} \int_{\mathcal{X}} e^{\eta(\theta) \cdot T(x)} h(x) dx$, $B(\theta)$ must be given by

$$B(\theta) \equiv \log \left(\int_{\mathcal{X}} e^{\eta(\theta) \cdot T(x)} h(x) dx \right).$$

Many of the commonly-considered distributions can be written as exponential families with $q = 1$ or 2 , for suitable h , B , and $\{\eta_i, t_i\}_{i \leq q}$. The likelihood function for a random sample of size n from the exponential family is

$$L(\theta) = \exp \left[\sum_{i=1}^q \eta_i(\theta) \sum_{j=1}^n t_i(x_j) - nB(\theta) \right],$$

which depends on the data only through the q -dimensional statistic T with components $T_i = \sum_{j \leq n} t_i(x_j)$. This *natural sufficient statistic* (see below) summarizes the data completely for any inference about θ .

It is often convenient to reparametrize exponential families to the *natural parameter* $\eta = \eta(\theta) \in \mathbb{R}^q$, leading (after rewriting the normalizing constant $B(\theta)$ as $A(\eta)$) to

$$f(x | \eta) = e^{\eta \cdot T(x) - nA(\eta)} h(x)$$

for a sample of size n , where again $A(\eta) \equiv \log \left(\int_{\mathcal{X}} e^{\eta \cdot T(x)} h(x) dx \right)$. We can calculate the moment generating function (MGF) for $T(X)$ as

$$\begin{aligned}
M_T(s) &= \mathbf{E} \left[e^{s \cdot T(X)} \right] \\
&= \int_{\mathcal{X}} e^{s \cdot T(x) - n A(\eta)} e^{\eta \cdot T(x)} h(x) dx \\
&= e^{-n A(\eta)} \int_{\mathcal{X}} e^{(\eta+s) \cdot T(x)} h(x) dx \\
&= e^{n [A(\eta+s) - A(\eta)]},
\end{aligned}$$

so we can find its mean and (co)variance as

$$\begin{aligned}
\mathbf{E}[T] &= \nabla_{\eta} \log M_T(0) = n \nabla_{\eta} A(\eta) \\
\mathbf{V}[T] &= \nabla_{\eta}^2 \log M_T(0) = n \nabla_{\eta}^2 A(\eta).
\end{aligned}$$

3.1. The Information Inequality

The Fisher Information about the natural parameter η from a single observation ($n = 1$) from an exponential family $f(x | \eta) = \exp(\eta \cdot T(x) - A(\eta)) h(x)$ is given by

$$I(\eta) = -\nabla_{\eta}^2 \log f(x | \eta) = \nabla_{\eta}^2 A(\eta),$$

so by Equation (1) the information in any parametrization is given by

$$I(\theta) = J^T \nabla_{\eta}^2 A(\eta) J$$

for $J = \partial \eta / \partial \theta$. The natural sufficient statistic $T(x)$ has mean $\psi(\theta) \equiv \mathbf{E}[T(X) | \theta] = \nabla_{\eta} A(\eta(\theta))$.

In particular, for scalar ($p = 1$) exponential families, the Information Inequality takes the form

$$\begin{aligned}
\mathbf{V}[T] &\geq \frac{|\psi'(\theta)|^2}{I(\theta)} \\
&= \frac{|c''(\eta(\theta)) \eta'(\theta)|^2}{\eta'(\theta) c''(\eta) \eta'(\theta)} \\
&= c''(\eta(\theta)) \\
&= \mathbf{V}[T],
\end{aligned}$$

so the lower bound is attained. This is not terribly surprising, since the inequality was based on the covariance inequality for the random variables T and $Z \equiv \nabla_{\theta} \log f(x | \theta) = \nabla_{\theta} \eta(\theta) \cdot T(X) - \nabla_{\theta} B(\theta)$, which are related by an affine transformation for scalar exponential families and hence are perfectly correlated.

4. Objective Bayesian Analysis

Laplace in the 1700's used the *uniform* prior distribution $\pi(\theta) \equiv 1$ in his Bayesian statistical analysis, intending it to represent a complete absence of knowledge about θ before observing a data vector $x \in \mathcal{X}$, and leading to a posterior density function

$$\pi(\theta | x) \propto f(x | \theta)$$

proportional to the likelihood. As appealing as this is for a non-subjective analysis, it is not invariant under reparametrization; for example, if we use the uniform distribution $\pi_1(\theta) \equiv 1$ for a binomial success probability $\theta \in \Theta = (0, 1)$ then upon observing y successes in n tries this leads to

$$\pi_1(\theta | x) \propto \theta^y (1 - \theta)^{n-y}$$

the $\text{Be}(1 + y, 1 + n - y)$ posterior distribution, with mean $T_1(Y) = \mathbb{E}^{\pi_1}[\theta | y] = \frac{1+y}{2+n}$, while a similar analysis using a uniform prior density for the (natural) logistic parameter $\eta = \log \frac{\theta}{1-\theta}$ leads to the beta $\text{Be}(y, n - y)$ posterior distribution, with mean $T_2(Y) = \mathbb{E}^{\pi_2}[\theta | y] = \frac{y}{n}$. When η is accorded a uniform prior the *implicit* prior for θ is $\pi_2(\theta) = 1(\eta)\eta'(\theta) = \frac{1}{\theta(1-\theta)}$, the beta $\text{Be}(0, 0)$, while the uniform density $\pi_1(\theta) \equiv 1$ is also the beta $\text{Be}(1, 1)$. For large numbers of success y and failure $n - y$ these two *reference* posterior distributions are very close, but for small y or $n - y$ they are not; which should we use, and why?

Evidently the problem is that in transforming from η to θ any prior density $\pi^\eta(\eta)$ will be transformed to $\pi^\theta(\theta) = |\det J| \pi^\eta(\eta(\theta))$ where $J = \partial\eta/\partial\theta$; the functional form of the priors in these two parametrizations cannot be the same, because of the Jacobian. Harold Jeffreys noticed that, since the Fisher Information $I(\theta) = J^\top \nabla_\eta^2 A(\eta) J$ transforms bilinearly in J , the recipe

$$\pi_J(\theta) \propto \sqrt{\det I(\theta)}$$

will give a prior density that transforms consistently to any parametrization:

$$\begin{aligned} \pi_J(\theta) &\propto \sqrt{\det I(\theta)} \\ &= \sqrt{\det J^\top I(\eta) J} \\ &= |\det J| \sqrt{\det I(\eta)} \\ &= |\det J| \pi_J(\eta), \end{aligned}$$

as required. In the case of Binomial data, for example, $I(\theta) = \frac{1}{\theta(1-\theta)}$ so $\pi_J(\theta) \propto \theta^{\frac{1}{2}-1}(1-\theta)^{\frac{1}{2}-1}$ is the $\text{Be}(\frac{1}{2}, \frac{1}{2})$ distribution (also called the *arc-sin* distribution, since it has CDF $\text{P}[\theta < t] = \frac{2}{\pi} \sin^{-1}(\sqrt{t})$), leading to a $\text{Be}(\frac{1}{2} + y, \frac{1}{2} + n - y)$ posterior distribution for the success probability θ upon observing y successes in n tries, with posterior mean $\text{E}[\theta \mid y] = (y + \frac{1}{2}) / (n + 1)$.