

Sufficient Statistics and the Exponential Family

1

Statistic: (informal) A summary of the data.

Statistic: (formal) Let X_1, \dots, X_n be a random sample of size n from a population and let $T(X_1, \dots, X_n)$ be a real-valued or vector-valued function whose domain includes the sample space of (X_1, \dots, X_n) . Then the random variable or random vector $Y = T(X_1, \dots, X_n)$ is called a statistic.

Sufficient Statistic: (informal). A summary of the data that contains all of the information about ~~any~~ any unknown parameters Θ in the data

Sufficient Statistic: (formal, not that useful): A statistic $Y = T(\underline{X})$ is a sufficient statistic for Θ if the conditional distribution of the sample data \underline{X} given the value of $T(\underline{X})$ does not depend on Θ .

Intuition: $T(\underline{X})$ follows a distribution that depends on Θ .

\Rightarrow we could use $T(\underline{X})$ to estimate Θ

Once we know $T(\underline{X})$, the distribution of \underline{X} does not depend on Θ

\Rightarrow knowing \underline{X} doesn't give you any more information about Θ than what we had in $T(\underline{X})$.

Factorization Theorem: Let $f(\underline{x} | \Theta)$ denote the joint pd or pmf of a sample \underline{X} . A statistic $T(\underline{X})$ is a sufficient statistic for Θ if and only if there exist functions $g(t | \Theta)$ and $h(\underline{x})$ such that for all sample points \underline{x} and all parameter values Θ ,

$$f(\underline{x} | \Theta) = g(T(\underline{x}) | \Theta) \cdot h(\underline{x})$$

Pf: See book

Note 1: Suppose we want to obtain a maximum likelihood estimate, and $T(\underline{x})$ is a sufficient statistic. Then

$$\begin{aligned} L(\theta | \underline{x}) &= f(\underline{x} | \theta) \\ &= g(T(\underline{x}) | \theta) \cdot h(\underline{x}) \end{aligned}$$

\Rightarrow log-likelihood is

$$\begin{aligned} l(\theta | \underline{x}) &= \log \{L(\theta | \underline{x})\} \\ &= \log \{g(T(\underline{x}) | \theta)\} + \log \{h(\underline{x})\} \end{aligned}$$

$$\Rightarrow \frac{d}{d\theta} l(\theta | \underline{x}) = \frac{d}{d\theta} \log \{g(T(\underline{x}) | \theta)\}$$

$\hat{\theta}_{MLE}$ depends only on $T(\underline{x})$, not the full data vector \underline{x}

Note 2: Suppose we have a prior distribution $\theta \sim f_{\theta}(\theta)$,

The posterior for θ has pdf/pmf:

$$\begin{aligned} f_{\theta|x}(\theta | \underline{x}) &= C \cdot f_{\theta}(\theta) \cdot f_{x|\theta}(\underline{x} | \theta) \\ &= C \cdot f_{\theta}(\theta) \cdot g(T(\underline{x}) | \theta) \cdot h(\underline{x}) \\ &= C_2 \cdot f_{\theta}(\theta) \cdot g(T(\underline{x}) | \theta) \quad (\text{where } C_2 = C \cdot h(\underline{x})) \end{aligned}$$

The posterior distribution of θ

$$\begin{aligned} f_{\theta|x}(\theta | \underline{x}) &= \frac{f_{\theta,x}(\theta, \underline{x})}{f_{\underline{x}}(\underline{x})} = \frac{f_{\theta,x}(\theta, \underline{x})}{\int f_{\theta,x}(\theta, \underline{x}) d\theta} = \frac{f_{\theta}(\theta) f_{x|\theta}(\underline{x} | \theta)}{\int f_{\theta}(\theta) f_{x|\theta}(\underline{x} | \theta) d\theta} \\ &= \frac{f_{\theta}(\theta) \cdot g(T(\underline{x}) | \theta) h(\underline{x})}{\int f_{\theta}(\theta) \cdot g(T(\underline{x}) | \theta) h(\underline{x}) d\theta} = \frac{f_{\theta}(\theta) \cdot g(T(\underline{x}) | \theta)}{\int f_{\theta}(\theta) \cdot g(T(\underline{x}) | \theta) d\theta} \quad \begin{matrix} \text{posterior depends} \\ \text{only on } T(\underline{x}), \\ \text{not full} \\ \text{data vector } \underline{x} \end{matrix} \end{aligned}$$

Example: Suppose $X_1, \dots, X_n \sim \text{iid Normal}(\mu, \sigma^2)$,
both μ and σ^2 unknown. $\Theta = (\mu, \sigma^2)$.

The joint pdf of \underline{X} is

$$\begin{aligned} f_{\underline{X}}(\underline{x} | \mu, \sigma^2) &= \prod_{i=1}^n (2\pi\sigma^2)^{-1/2} \exp \left\{ -\frac{(x_i - \mu)^2}{2\sigma^2} \right\} \\ &= (2\pi\sigma^2)^{-n/2} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n \frac{(x_i - \mu)^2}{\cancel{\sigma^2}} \right\} \\ &= (2\pi\sigma^2)^{-n/2} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \bar{x} + \bar{x} - \mu)^2 \right\} \\ &= (2\pi\sigma^2)^{-n/2} \exp \left[-\frac{1}{2\sigma^2} \left\{ \sum_{i=1}^n (x_i - \bar{x})^2 + n(\bar{x} - \mu)^2 \right\} \right] \end{aligned}$$

Suppose μ is unknown, σ^2 known.

$$\begin{aligned} f_{\underline{X}}(\underline{x} | \mu, \sigma^2) &= \underbrace{\exp \left\{ -\frac{1}{2} \cdot \frac{n(\bar{x} - \mu)^2}{\sigma^2} \right\}}_{\text{involves } \mu} \cdot (2\pi\sigma^2)^{-n/2} \exp \left\{ -\frac{1}{2} \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{\sigma^2} \right\} \\ &= g(T(\underline{x}) | \mu) \cdot h(\underline{x}) \end{aligned}$$

where $T(\underline{x}) = \bar{x}$,

$$g(t | \mu) = \exp \left\{ -\frac{1}{2} \frac{n(t - \mu)^2}{\sigma^2} \right\},$$

$$h(\underline{x}) = (2\pi\sigma^2)^{-n/2} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{\sigma^2} \right\}$$

Now suppose both μ, σ^2 are unknown. Set $\Theta = (\mu, \sigma^2)$

Our sufficient statistics are $\underline{T}(\underline{X}) = (T_1(\underline{X}), T_2(\underline{X}))$, where
 $T_1(\underline{X}) = \bar{X}$, $T_2(\underline{X}) = S^2 = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n-1}$.

Set $h(\underline{x}) = 1$ and

$$\begin{aligned} g(\underline{t} | \underline{\theta}) &= g(t_1, t_2 | \mu, \sigma^2) \\ &= (2\pi\sigma^2)^{-n/2} \exp \left[-\left\{ n(t_1 - \mu)^2 + (n-1)t_2 \right\} / 2\sigma^2 \right] \end{aligned}$$

Then $f(\underline{x} | \mu, \sigma^2) = g(\underline{t} | \underline{\theta}) \cdot h(\underline{x})$

so $\underline{T}(\underline{X}) = (\bar{X}, S^2)$ is a sufficient statistic for the normal model.

Example: Suppose that X_1, \dots, X_n iid $\text{Binomial}(n, \theta)$, θ unknown,

Show that $\sum_{i=1}^n X_i$ is a sufficient statistic for θ .

Exponential Family (not to be confused with the exponential distribution)

5

A family of probability distributions is called an exponential family if its pdfs/pdfs can be expressed as

$$f(x|\theta) = h(x) \cdot c(\theta) \cdot \exp \left\{ \sum_{i=1}^k w_i(\theta) t_i(x) \right\}$$

Example: Binomial exponential family

Suppose $X \sim \text{Binomial}(n, \theta)$

$$\begin{aligned} f(x|\theta) &= \binom{n}{x} \theta^x (1-\theta)^{n-x} \\ &= \binom{n}{x} (1-\theta)^n \left(\frac{\theta}{1-\theta} \right)^x \\ &= \binom{n}{x} (1-\theta)^n \exp \left\{ \log \left(\frac{\theta}{1-\theta} \right) x \right\} \end{aligned}$$

\therefore this is an exponential family with

$$h(x) = \binom{n}{x}, \quad c(\theta) = (1-\theta)^n, \quad w_i(\theta) = \log \left(\frac{\theta}{1-\theta} \right),$$

$$\text{and } t_i(x) = x.$$

Other distributions that are exponential families:

Normal, exponential, gamma, χ^2 , beta, Dirichlet,

Poisson, geometric

Why do we care?

Thm. (Pitman-Koopman-Darmois)

Suppose that X_1, \dots, X_n are iid r.v.'s with pdf $f_{X_i}(x_i|\theta)$, and the support of $f_{X_i}(x_i|\theta)$ does not depend on θ (Counter-example: Uniform($\frac{1}{2}, 1$, θ)). Only if $f_X(x|\theta)$ is an exponential family is there a sufficient statistic $T(X) = (T_1(X), \dots, T_k(X))$ whose length k does not increase as n increases.

Thm: Let X_1, \dots, X_n be iid where $f_{X_i|\theta}(x_i|\theta)$ is in an exponential family with pdf is

$$f(x|\theta) = h(x) \cdot c(\theta) \cdot \exp \left\{ \sum_{j=1}^k w_j(\theta) t_j(x) \right\},$$

where $\underline{\theta} = (\theta_1, \dots, \theta_d)$ for $d \leq k$.

$$\text{Then } T(X) = \left(\sum_{j=1}^n t_1(x_j), \dots, \sum_{j=1}^n t_k(x_j) \right)$$

is a sufficient statistic for $\underline{\theta}$.

Rao-Blackwell Thm:

Let $\hat{\theta}$ be an estimator of a parameter θ , and $T(X)$ a sufficient statistic for θ .

Define $\tilde{\theta} = E[\hat{\theta} | T(X)]$.

Then $MSE(\tilde{\theta}) \leq MSE(\hat{\theta})$.

We have "Rao-Blackwellized" the original estimator $\hat{\theta}$ to obtain an improved estimator $\tilde{\theta}$.

Note: If $\hat{\theta}$ was unbiased, $\tilde{\theta}$ is still unbiased.

Basic idea: If an estimator $\hat{\theta}$ wasn't based on a sufficient statistic, it can be improved by conditioning on a sufficient statistic.

Suppose $X_1, \dots, X_n \sim \text{Normal}(\theta, \sigma^2)$ with σ^2 known.

Set $\hat{\theta} = X_1$. Not based on the sufficient statistic \bar{X} , clearly suboptimal.
It can be shown that $X_1 | \bar{X} \sim \text{Normal}(\bar{X}, \frac{n-1}{n}\sigma^2)$

$$\tilde{\theta} = E(\hat{\theta} | \bar{X}) = E(X_1 | \bar{X}) = \bar{X}$$

$\tilde{\theta} = \bar{X}$ is the Rao-Blackwellized estimator, and

$$MSE(\tilde{\theta}) = \text{Bias}(\tilde{\theta}) + \text{Var}(\tilde{\theta}) = \frac{\sigma^2}{n} < \sigma^2 = \text{Var}(\hat{\theta}) = MSE(\hat{\theta}),$$

as given by the Rao-Blackwell theorem,