

Lecture 33: Sufficiency and Rao-Blackwell

Ciaran Evans

Sufficient statistics

Question: Given an unbiased estimator, can I improve its variance?

- ▶ Answering this requires us to introduce a new concept:
sufficient statistics

Definition (sufficient statistic): Let X_1, \dots, X_n be a sample from a distribution $f(x|\theta)$. Let $T = T(X_1, \dots, X_n)$ be a statistic. We say that T is a sufficient statistic for θ if the conditional distribution of $X_1, \dots, X_n | T$ does not depend on θ .

Example

Suppose $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Poisson}(\lambda)$

$$\text{Let } T = \sum_{i=1}^n X_i$$

$$T \sim \text{Poisson}(n\lambda)$$

$$f(X_1, \dots, X_n | T) = \frac{f(X_1, \dots, X_n, T)}{f(T)} = \frac{f(X_1, \dots, X_n)}{f(T)}$$

$$\begin{aligned} f(X_1, \dots, X_n, T) &: P(X_1=x_1, X_2=x_2, \dots, X_n=x_n, T=\sum_i x_i) \\ &= P(X_1=x_1, \dots, X_n=x_n) \end{aligned}$$

$$\begin{aligned} \Rightarrow \frac{f(X_1, \dots, X_n)}{f(T)} &= \frac{\prod_{i=1}^n e^{-\lambda} \lambda^{x_i} / x_i!}{e^{-n\lambda} (n\lambda)^T / T!} = \frac{\cancel{e^{-n\lambda}} \lambda^{\sum_i x_i} / \prod_i x_i!}{\cancel{e^{-n\lambda}} \lambda^{n^T} / T!} \\ &= \frac{T!}{n^T \prod_i x_i!} \quad \leftarrow \text{does not depend on } \lambda \\ &\Rightarrow T = \sum_i X_i \text{ is a sufficient statistic} \end{aligned}$$

Example

Suppose $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Bernoulli}(p)$

$$T = \sum_i X_i \sim \text{Binomial}(n, p)$$

$$f(X_1, \dots, X_n | T) = \frac{f(X_1, \dots, X_n)}{f(T)}$$

$$= \frac{\prod_i p^{X_i} (1-p)^{1-X_i}}{\binom{n}{T} p^T (1-p)^{n-T}}$$

$$= \frac{p^{\sum_i X_i} (1-p)^{n - \sum_i X_i}}{\binom{n}{T} p^T (1-p)^{n-T}}$$

$$= \frac{1}{\binom{n}{T}}$$

does not depend on p

$\Rightarrow T$ is sufficient for p

Rao-Blackwell

Let θ be a parameter of interest, and

let $\gamma(\theta)$ be some function of θ .

Let $\hat{\gamma}$ be an unbiased estimator of $\gamma(\theta)$,

and T be a sufficient statistic for θ .

Let $\gamma^* = E[\hat{\gamma} | T]$. Then:

$$\textcircled{1} E[\gamma^*] = \gamma \quad (\text{unbiased})$$

$$\textcircled{2} \text{Var}(\gamma^*) \leq \text{Var}(\hat{\gamma})$$

Example

Suppose $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Poisson}(\lambda)$

sufficient statistic: $T = \sum_i X_i$

consider $\hat{\lambda} = X_1$ $\mathbb{E}[X_1] = \lambda$ $\text{Var}(X_1) = \lambda$

$$\lambda^* = \mathbb{E}[\hat{\lambda} | T]$$

want: $\mathbb{E}[X_1 | T=t]$ i.e. $\mathbb{E}[X_1 | \sum_i X_i = t]$

$$\mathbb{E}[T | T=t] = t$$

$$\Rightarrow \mathbb{E}[\sum_i X_i | T=t] = t$$

$$\Rightarrow \sum_i \mathbb{E}[X_i | T=t] = t$$

$$\Rightarrow n \mathbb{E}[X_1 | T=t] = t \Rightarrow \mathbb{E}[X_1 | T=t] = \frac{t}{n}$$

$$\Rightarrow \lambda^* = \mathbb{E}[\hat{\lambda} | T] = \frac{T}{n} = \frac{1}{n} \sum_i X_i$$

(Rao-Blackwellized estimator) $\text{Var}(\lambda^*) = \frac{\lambda}{n} < \text{Var}(X_1)$

Factorization theorem

Let x_1, \dots, x_n be a sample with joint probability function $f(x_1, \dots, x_n | \theta)$.

A statistic $T = T(x_1, \dots, x_n)$ is sufficient for θ if and only if there exist functions $g(t | \theta)$ and $h(x_1, \dots, x_n)$ such that for any possible x_1, \dots, x_n and any possible θ ,

$$f(x_1, \dots, x_n | \theta) = \underbrace{g(T(x_1, \dots, x_n) | \theta)}_{\text{joint distribution only depends on } \theta \text{ through the sufficient statistic}} h(x_1, \dots, x_n)$$

joint distribution only depends on θ through the sufficient statistic

Ex: $x_1, \dots, x_n \stackrel{iid}{\sim} \text{Bernoulli}(p)$

$$f(x_1, \dots, x_n | p) = p^{\sum x_i} (1-p)^{n - \sum x_i}$$

$$T = \sum x_i \quad g(T | p) = p^T (1-p)^{n-T} \quad h(x) = 1$$

$$\Rightarrow f(x_1, \dots, x_n | p) = g(\sum x_i | p) h(x_1, \dots, x_n)$$

$X_1, \dots, X_n \stackrel{\text{iid}}{\sim} N(\mu, \sigma^2)$ μ, σ^2 are unknown

$$f(x_1, \dots, x_n | \mu, \sigma^2) = \left(\frac{1}{2\pi\sigma^2} \right)^{\frac{n}{2}} \exp \left\{ -\frac{1}{2\sigma^2} \sum_i (x_i - \mu)^2 \right\}$$

$$\begin{aligned} \sum_i (x_i - \mu)^2 &= \sum_i (x_i^2 - 2\mu x_i + \mu^2) \\ &= \sum_i x_i^2 - 2\mu \sum_i x_i + n\mu^2 \end{aligned}$$

$$\begin{aligned} \Rightarrow f(x_1, \dots, x_n | \mu, \sigma^2) &= \left(\frac{1}{2\pi\sigma^2} \right)^{\frac{n}{2}} \exp \left\{ -\frac{1}{2\sigma^2} \left(\sum_i x_i^2 - 2\mu \sum_i x_i + n\mu^2 \right) \right\} \\ &= g(T_1, T_2 | \mu, \sigma^2) \end{aligned}$$

$$T_1 = \sum_i x_i$$

$$T_2 = \sum_i x_i^2$$

$$T = (T_1, T_2) \in \mathbb{R}^2$$

(sufficient statistic)

equivalently:

$$(\bar{X}, \frac{1}{n-1} \sum_i (x_i - \bar{X})^2)$$