



Lecture 5: Sufficiency & the Rao-Blackwell Theorem

STA261 – Probability & Statistics II

Ofir Harari

Department of Statistical Sciences

University of Toronto



Outline

Sufficient Statistics

- Definition

- The Fisher–Neyman Factorization Theorem

- The Exponential Family of Distributions

- The Rao–Blackwell Theorem



Sufficient statistics

- We have discussed the likelihood function extensively (and will continue doing so)
- It drives parameter estimation via the Maximum Likelihood principle
- We have studied the many good large-sample properties of MLEs (asymptotic normality & efficiency)
- We shall see that it also drives inference on parameters (i.e. hypothesis testing & confidence intervals)
- All in all, the likelihood is the single most important function in statistics, as summarized in the *Likelihood Principle*:

In the inference about θ , after $\underline{x} = (x_1, \dots, x_n)$ is observed, all relevant experimental information is contained in the likelihood function for the observed \underline{x} .



Sufficient statistics (cont.)

- So, all information about θ is encoded in the likelihood
- But what if the likelihood itself depends on the data through a mere summary (statistic)?

- Consider, for example, a sequence of Bernoulli trials,

$$X_1, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} \text{Binom}(1, p),$$

with the likelihood function

$$\mathcal{L}(p) = p^{\sum_{i=1}^n x_i} (1 - p)^{n - \sum_{i=1}^n x_i}.$$

- If we know the total number of “successful” trials, nothing more can be learned on p from knowing the detailed observations (i.e. the sequences $(1, 0, 1, 0, 1)$ and $(1, 1, 1, 0, 0)$ are the equivalent from a likelihood standpoint).
- Data compression!



Sufficient statistics (cont.)

Definition

A statistic $T(\underline{X}) = T(X_1, \dots, X_n)$ is *sufficient* for an unknown parameter θ if the conditional (joint) distribution of X_1, \dots, X_n given $T(\underline{X})$ does not depend on θ .

- In other words: $T(\underline{X})$ teaches us all we need to know about θ .
- To continue with the Bernoulli trials example, let us now verify that $\sum_{i=1}^n X_i$ is indeed sufficient for p .
- Note that $\sum_{i=1}^n X_i \sim \text{Binom}(n, p)$, thus

$$\begin{aligned} \mathbb{P}\left(\underline{X} = \underline{x} \mid \sum_{i=1}^n X_i = t\right) &= \frac{\mathbb{P}(\underline{X} = \underline{x}, \sum_{i=1}^n X_i = t)}{\mathbb{P}(\sum_{i=1}^n X_i = t)} \\ &= \begin{cases} \frac{\mathbb{P}(\underline{X} = \underline{x})}{\mathbb{P}(\sum_{i=1}^n X_i = t)} & , \quad \sum_{i=1}^n X_i = t \\ 0 & , \quad \text{otherwise} \end{cases} \end{aligned}$$



Sufficient statistics (cont.)

$$\begin{aligned}
 \mathbb{P}\left(\underline{X} = \underline{x} \mid \sum_{i=1}^n X_i = t\right) &= \begin{cases} \frac{\mathbb{P}(\underline{X} = \underline{x})}{\mathbb{P}(\sum_{i=1}^n X_i = t)} & , \quad \sum_{i=1}^n X_i = t \\ 0 & , \quad \text{otherwise} \end{cases} \\
 &= \begin{cases} \frac{p^t(1-p)^{n-t}}{\binom{n}{t} p^t(1-p)^{n-t}} & , \quad \sum_{i=1}^n X_i = t \\ 0 & , \quad \text{otherwise} \end{cases} \\
 &= \begin{cases} \frac{1}{\binom{n}{t}} & , \quad \sum_{i=1}^n X_i = t \\ 0 & , \quad \text{otherwise} \end{cases} .
 \end{aligned}$$

The above does not depend on p , hence $\sum_{i=1}^n X_i$ is sufficient for p .



The Fisher–Neyman Factorization Theorem

- You may have noticed that a direct verification of sufficiency can be messy
- Intuitively, the likelihood depends on the data via $\sum_{i=1}^n X_i$, ergo, it is sufficient
- Our intuition is right this time!

Theorem

A statistic $T(\underline{X})$ is sufficient for $\theta \iff$ for any value of θ we can write

$$\mathcal{L}(\theta) = g(T(\underline{x}), \theta) h(\underline{x}).$$

- Note that in the binary case

$$\mathcal{L}(p) = \underbrace{p^{\sum_{i=1}^n x_i} (1-p)^{n-\sum_{i=1}^n x_i}}_{g(p, \sum_{i=1}^n x_i)},$$

with $h(\underline{x}) = 1$, hence $T(\underline{X}) = \sum_{i=1}^n X_i$ is sufficient for p .



Fisher–Neyman Factorization Theorem (cont.)



Jerzy Neyman, 1894-1981

Source: statistics.berkeley.edu



Fisher–Neyman Factorization Theorem (cont.)

Proof for the discrete case:

(\implies)

Suppose that $T(\underline{X})$ is sufficient for θ , and let \underline{x} be our sample, such that $T(\underline{x}) = t$.

$$\mathcal{L}(\theta) = \mathbb{P}(\underline{X} = \underline{x} | \theta) = \mathbb{P}(\underline{X} = \underline{x}, T(\underline{X}) = t | \theta)$$

$$= \mathbb{P}(\underline{X} = \underline{x} | T(\underline{X}) = t, \theta) \mathbb{P}(T(\underline{X}) = t | \theta)$$

- $\mathbb{P}(\underline{X} = \underline{x} | T(\underline{X}) = t, \theta)$ does not depend on θ (why?) – call it $h(\underline{x})$
- Denote $g(t, \theta) := \mathbb{P}(T(\underline{X}) = t | \theta)$ – and we're done.



Fisher–Neyman Factorization Theorem (cont.)

Proof (cont.):

(\Leftarrow)

Suppose now that the likelihood can be factorized as

$$\mathcal{L}(\theta) := \mathbb{P}(\underline{X} = \underline{x} | \theta) = g(T(\underline{x}), \theta) h(\underline{x}).$$

Note that

$$\mathbb{P}(T(\underline{X}) = t) = \sum_{T(\underline{x})=t} \mathbb{P}(\underline{X} = \underline{x} | \theta) = g(t, \theta) \sum_{T(\underline{x})=t} h(\underline{x}),$$

hence

$$\mathbb{P}(\underline{X} = \underline{x} | T(\underline{X}) = t) = \begin{cases} \frac{\mathbb{P}(\underline{X} = \underline{x}, T(\underline{X}) = t | \theta)}{\mathbb{P}(T(\underline{X}) = t)} & , \quad T(\underline{x}) = t \\ 0 & , \quad \text{otherwise} \end{cases}$$



Fisher–Neyman Factorization Theorem (cont.)

$$\begin{aligned}
 \mathbb{P}(\underline{X} = \underline{x} | T(\underline{X}) = t) &= \begin{cases} \frac{\mathbb{P}(\underline{X} = \underline{x}, T(\underline{X}) = t | \theta)}{\mathbb{P}(T(\underline{X}) = t)} & , \quad T(\underline{x}) = t \\ 0 & , \quad \text{otherwise} \end{cases} \\
 &= \begin{cases} \frac{\mathbb{P}(\underline{X} = \underline{x} | \theta)}{g(t, \theta) \sum_{T(\underline{x})=t} h(\underline{x})} & , \quad T(\underline{x}) = t \\ 0 & , \quad \text{otherwise} \end{cases} \\
 &= \begin{cases} \frac{g(t, \theta) h(\underline{x})}{g(t, \theta) \sum_{T(\underline{x})=t} h(\underline{x})} & , \quad T(\underline{x}) = t \\ 0 & , \quad \text{otherwise} \end{cases} \\
 &= \begin{cases} \frac{h(\underline{x})}{\sum_{T(\underline{x})=t} h(\underline{x})} & , \quad T(\underline{x}) = t \\ 0 & , \quad \text{otherwise} \end{cases} ,
 \end{aligned}$$

that does not depend on θ .



Example: Poisson distribution

Example

Let $X_1, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} \text{Pois}(\lambda)$. Find a sufficient statistic for λ .

Solution:

- Here $\mathcal{L}(\lambda) = e^{-n\lambda} \frac{\lambda^{\sum_{i=1}^n x_i}}{\prod_{i=1}^n x_i!}$
- Write

$$\mathcal{L}(\lambda) = g\left(\sum_{i=1}^n x_i, \lambda\right) h(\underline{x})$$

where $g\left(\sum_{i=1}^n x_i, \lambda\right) = e^{-n\lambda} \lambda^{\sum_{i=1}^n x_i}$ and $h(\underline{x}) = \frac{1}{\prod_{i=1}^n x_i!}$

- Then, from the factorization Theorem, $\sum_{i=1}^n X_i$ is sufficient for λ .



Example: the Cauchy distribution

Example

Let X_1, \dots, X_n be a random sample from the Cauchy distribution, with pdf

$$f(x|\theta) = \frac{1}{\pi [1 + (x - \theta)^2]}.$$

Does a sufficient statistics for θ exist?

Solution:

Here

$$\begin{aligned} \mathcal{L}(\theta) &= \prod_{i=1}^n f(x_i|\theta) = \frac{1}{\pi^n \prod_{i=1}^n [1 + (x_i - \theta)^2]} \\ &= \frac{1}{\pi^n} \exp \left\{ - \sum_{i=1}^n \log [1 + (x_i - \theta)^2] \right\} = \dots \end{aligned}$$

- ★ Slice and dice it all you like, the x_i 's and θ cannot be separated
- ★ No sufficient statistic, keep all the data



Example: Uniform distribution

Example

Let $X_1, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} \text{U}[0, \theta]$ (continuous uniform). Find a sufficient statistic for θ .

Solution:

Recall that

$$f(x|\theta) = \begin{cases} \frac{1}{\theta} & , \quad 0 \leq x \leq \theta \\ 0 & , \quad \text{otherwise} \end{cases} = \frac{1}{\theta} \cdot I\{0 \leq x \leq \theta\},$$

for

$$I\{0 \leq x \leq \theta\} = \begin{cases} 1 & , \quad x \in [0, \theta], \\ 0 & , \quad \text{otherwise.} \end{cases}$$

In light of this,

$$\mathcal{L}(\theta) = \prod_{i=1}^n f(x_i|\theta) = \frac{1}{\theta^n} \prod_{i=1}^n I\{0 \leq x_i \leq \theta\} = \frac{1}{\theta^n} I\{x_{\max} \leq \theta\},$$

where $x_{\max} = \max(x_1, \dots, x_n)$.



Example: Uniform distribution (cont.)

Solution (cont):

$$\mathcal{L}(\theta) = \frac{1}{\theta^n} I \{x_{\max} \leq \theta\}$$

★ We can write

$$\mathcal{L}(\theta) = g(x_{\max}, \theta)h(\underline{x})$$

for $g(x_{\max}, \theta) = \frac{1}{\theta^n} I \{x_{\max} \leq \theta\}$ and $h(\underline{x}) = 1$.

★ Then by the factorization Theorem, $X_{\max} = \max(X_1, \dots, X_n)$ is sufficient for θ .



The Exponential family of distributions

The following definition covers a surprisingly large subset of the distributions we have familiarized ourselves with.

Definition

A distribution with cdf/pmf $f(x|\theta)$ is said to belong to a *one parameter exponential family of distributions* if

$$f(x|\theta) = \begin{cases} \exp \{ c(\theta) T(x) + d(\theta) + S(x) \} & , \quad x \in A \\ 0 & , \quad \text{otherwise} \end{cases}$$

where A does not depend on θ .

- Fantastic! What is it good for?
- Sufficiency, among other things



Sufficiency in Exponential families

- Suppose for a moment that we have a random sample from an exponential family, with

$$f(x|\theta) = \exp \{c(\theta) T(x) + d(\theta) + S(x)\}.$$

- Then, the likelihood is given by

$$\begin{aligned} \mathcal{L}(\theta) &= \prod_{i=1}^n \exp \{c(\theta) T(x_i) + d(\theta) + S(x_i)\} \\ &= \exp \left\{ c(\theta) \sum_{i=1}^n T(x_i) + nd(\theta) + \sum_{i=1}^n S(x_i) \right\} \\ &= \underbrace{\exp \left\{ c(\theta) \sum_{i=1}^n T(x_i) + nd(\theta) \right\}}_{g(\sum_{i=1}^n T(x_i), \theta)} \underbrace{\exp \left\{ \sum_{i=1}^n S(x_i) \right\}}_{h(\underline{x})} \end{aligned}$$

- Evidently, $\sum_{i=1}^n T(X_i)$ is sufficient for θ .



Example: Poisson distribution

- Here

$$f(x_i|\lambda) = \frac{e^{-\lambda} \lambda^{x_i}}{x_i!} = \exp \{-\lambda + x_i \log \lambda - \log x_i!\}$$

- Denoting $c(\lambda) = \log \lambda$, $T(x_i) = x_i$, $d(\lambda) = -\lambda$ and $S(x_i) = -\log x_i!$, we have

$$f(x_i|\lambda) = \exp \{c(\lambda) T(x_i) + d(\lambda) + S(x_i)\}.$$

- The Poisson family is an exponential family of distributions then

- And indeed, we have shown that $\sum_{i=1}^n T(X_i) = \sum_{i=1}^n X_i$ is sufficient for λ



Example: Gamma distribution (λ known)

- Here $X_1, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} \text{Gamma}(\alpha, \lambda)$, i.e.

$$\begin{aligned} f(x_i | \alpha) &= \frac{\lambda^\alpha}{\Gamma(\alpha)} x_i^{\alpha-1} e^{-\lambda x_i} = \exp \{ \alpha \log \lambda - \log \Gamma(\alpha) + (\alpha - 1) \log x_i - \lambda x_i \} \\ &= \exp \{ \alpha \log x_i + \alpha \log \lambda - \log \Gamma(\alpha) - \log x_i - \lambda x_i \}, \quad x \geq 0 \end{aligned}$$

- Set $c(\alpha) = \alpha$, $T(x_i) = \log x_i$, $d(\alpha) = \alpha \log \lambda - \log \Gamma(\alpha)$ and $S(x_i) = -\log x_i - \lambda x_i$, we can write

$$f(x_i | \alpha) = \exp \{ c(\alpha) T(x_i) + d(\alpha) + S(x_i) \},$$

hence the Gamma family is an exponential family of distributions.

- Moreover, $\sum_{i=1}^n T(X_i) = \sum_{i=1}^n \log X_i$ is sufficient for α .



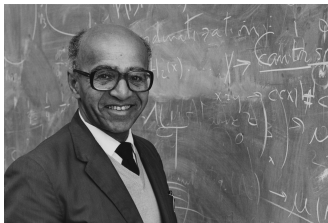
The Rao–Blackwell Theorem

Theorem

Let $\hat{\theta}$ be an estimator of θ with a finite variance. Suppose that T is sufficient for θ , and let $\hat{\theta}^* = \mathbb{E}[\hat{\theta} | T]$. Then for all θ

$$\text{MSE}(\hat{\theta}^*, \theta) \leq \text{MSE}(\hat{\theta}, \theta),$$

where equality holds $\iff \hat{\theta}^* = \hat{\theta}$.



David Blackwell, 1919-2010

Source: nationalmedals.org



The Rao–Blackwell Theorem (cont.)

Proof: Recall that from the law of total expectation

$$\mathbb{E}[\hat{\theta}^*] = \mathbb{E} \left\{ \mathbb{E}[\hat{\theta} | T] \right\} = \mathbb{E}[\hat{\theta}],$$

thus $\hat{\theta}^*$ and $\hat{\theta}$ have the same bias, and to compare their MSEs we need to compare their variances.

★ **In particular, if $\hat{\theta}$ is unbiased then so is $\hat{\theta}^*$.**

Now, applying the law of total variance we have –

$$\text{Var}[\hat{\theta}^*] = \text{Var} \left\{ \mathbb{E}[\hat{\theta} | T] \right\} + \mathbb{E} \left\{ \text{Var}[\hat{\theta} | T] \right\} = \text{Var}[\hat{\theta}] + \mathbb{E} \left\{ \text{Var}[\hat{\theta} | T] \right\} \geq \text{Var}[\hat{\theta}],$$

where

$$\begin{array}{ccccc} \text{equality} & & \hat{\theta} \text{ is a constant} & & \hat{\theta} \text{ is a} \\ \text{holds} & \implies \text{Var}[\hat{\theta} | T] = 0 \implies & \text{w.r.t. } \underline{X} \text{ when} & \implies & \text{function of } T, \\ & & T \text{ is given} & & \text{say, } \hat{\theta} = g(T) \end{array}$$

$$\implies \hat{\theta}^* = \mathbb{E}[\hat{\theta} | T] = \mathbb{E}[g(T) | T] = g(T) = \hat{\theta}$$



Comments on the Rao–Blackwell Theorem

- Where in the proof did we use the fact that $T(\underline{X})$ was sufficient for θ ?
 - ★ Implicitly: we called $\hat{\theta}^* = \mathbb{E}[\hat{\theta} | T]$ an “estimator”, but that is only true because the distribution of $\hat{\theta}^*$ does not depend on θ .
- It is tempting to re-apply Rao–Blackwellization to the resultant estimator – could it be further improved?

— ★ Remember: $\hat{\theta}_{\text{RB}} = \mathbb{E}[\hat{\theta}_0 | T] = g(T)$ (a function of T), therefore

$$\mathbb{E}[\hat{\theta}_{\text{RB}} | T] = \mathbb{E}[g(T) | T] = g(T) = \hat{\theta}_{\text{RB}},$$

suggesting that the process stops after one stage.

- A follow-up result, the **Lehmann–Scheffé Theorem**:
if, in addition to sufficiency, T has a property called *completeness*,
Rao–Blackwellizing an unbiased estimator would yield the unique optimal unbiased estimator.



“Rao–Blackwellization” example

Example

Suppose that the annual number of earthquakes in a certain seismic region follows a $\text{Pois}(\lambda)$ distribution, where different years are assumed to be independent. We wish to estimate the probability that there will be no earthquakes next year, based on a sample X_1, \dots, X_n . “Rao–Blackwellize” the naive estimator

$$\hat{\theta}_0 = \begin{cases} 1 & , \quad X_1 = 0 \\ 0 & , \quad \text{otherwise} \end{cases}$$

to obtain an improved estimator.

Solution:

- First, note that the parameter we wish to estimate is $\theta = \mathbb{P}(X = 0) = e^{-\lambda}$
- Also note that

$$\mathbb{E}[\hat{\theta}_0] = \mathbb{P}(X_1 = 0) = e^{-\lambda} = \theta,$$

hence $\hat{\theta}_0$ (naive as it may be) is unbiased.



“Rao–Blackwellization” example (cont.)

Solution (cont.):

- We have verified that $T = \sum_{i=1}^n X_i$ is sufficient for λ , thus it is also sufficient for $\theta = e^{-\lambda}$ (or any other monotonic transformation of λ)
- Keep in mind that $\sum_{i=1}^n X_i \sim \text{Pois}(n\lambda)$. Likewise, $\sum_{i=2}^n X_i \sim \text{Pois}((n-1)\lambda)$.
- Just like $\hat{\theta}_0, \hat{\theta}_0 | T$ is binary (returns either 0 or 1), hence

$$\begin{aligned} \hat{\theta}_{\text{RB}} &:= \mathbb{E} \left[\hat{\theta}_0 | T \right] = \mathbb{P} \left(\hat{\theta}_0 = 1 \left| \sum_{i=1}^n X_i = T \right. \right) = \mathbb{P} \left(X_1 = 0 \left| \sum_{i=1}^n X_i = T \right. \right) \\ &= \frac{\mathbb{P} \left(X_1 = 0, \sum_{i=1}^n X_i = T \right)}{\mathbb{P} \left(\sum_{i=1}^n X_i = T \right)} = \frac{\mathbb{P} \left(X_1 = 0, \sum_{i=2}^n X_i = T \right)}{\mathbb{P} \left(\sum_{i=1}^n X_i = T \right)} \end{aligned}$$



“Rao–Blackwellization” example (cont.)

$$\begin{aligned}\hat{\theta}_{\text{RB}} &= \frac{\mathbb{P}(X_1 = 0, \sum_{i=2}^n X_i = T)}{\mathbb{P}(\sum_{i=1}^n X_i = T)} = \frac{\mathbb{P}(X_1 = 0) \mathbb{P}(\sum_{i=2}^n X_i = T)}{\mathbb{P}(\sum_{i=1}^n X_i = T)} \\ &= \frac{e^{-\lambda} \cdot e^{-(n-1)\lambda} \frac{[(n-1)\lambda]^T}{T!}}{e^{-n\lambda} \frac{[n\lambda]^T}{T!}} = \left(1 - \frac{1}{n}\right)^T = \left(1 - \frac{1}{n}\right)^{\sum_{i=1}^n X_i}\end{aligned}$$

- Remember: $\hat{\theta}_0$ was unbiased— thus so is $\hat{\theta}_{\text{RB}}$
- It can be shown that it is the best unbiased estimator of $\theta = e^{-\lambda}$, for all λ
- Not the best estimator of θ overall, though
- For large n

$$\hat{\theta}_{\text{RB}} = \left(1 - \frac{1}{n}\right)^{\sum_{i=1}^n X_i} = \left(1 - \frac{1}{n}\right)^{n\bar{X}} \approx e^{-\bar{X}} = e^{-\hat{\lambda}_{\text{MLE}}} = \hat{\theta}_{\text{MLE}}$$



Comparison of $\hat{\theta}_{RB}$ and $\hat{\theta}_{MLE}$

#Calculating MSE by Monte-Carlo Simulation

```
MSEs <- function(lambda, n){
  theta <- exp(-lambda)
  samp <- matrix(rpois(n*1e5, lambda),
                ncol=n)
  xBar <- apply(samp, 1, mean)
  thetaHat1 <- exp(-xBar)
  thetaHat2 <- (1-1/n)^(n*xBar)
  MSE1 <- mean((thetaHat1-theta)^2)
  MSE2 <- mean((thetaHat2-theta)^2)
  return(c(MSE1,MSE2))
}
>
> n <- 5
> Vals1 <- t(sapply(.1*c(1:40), MSEs, n=n))
> plot(.1*c(1:40), Vals1[,1], type='l')
> lines(.1*c(1:40), Vals1[,2], lty=2, col=2)
>
```

