

Lecture 5. Sufficient Statistics. Intro to Estimation

Sufficient statistics produce a smaller well-behaved model
imposing assumption.

Parametric family: $X_1, \dots, X_n \sim i.i.d. F(\cdot, \theta)$. θ is the only unknown parameter

\uparrow
don't know

Eg¹) $x_1, \dots, x_n \sim \text{iid } N(\mu, \sigma^2) \quad (\theta = \mu, \sigma^2)$

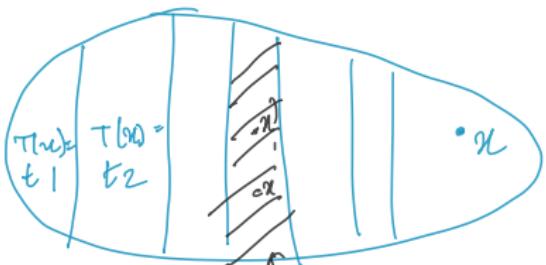
Eg²) $x_1, \dots, x_n \sim \text{iid } X_i = \begin{cases} 1 & P \\ 0 & 1-P \end{cases} \quad \text{so} \quad \theta = p$

$$x = (x_1, \dots, x_n)$$

Definition 1.

Statistic $T(X)$ is sufficient for θ if the conditional distribution of X given $T(X)$ does not depend on θ .

$X | T(x)$ doesn't depend on θ



Just reporting statistics/summary
one slice - diff. value of t

Suppose $T(x) = t$ is given.
What's the dist'n of x — doesn't depend on θ .

→ how x is distributed depends on θ

Sufficient statistics

* Joint distn of $T(n)$ comes from X

Example Let $X = (X_1, \dots, X_n)$ be a random sample from $N(\mu, \sigma^2)$.

Suppose that σ^2 is known.

$T(X) = \bar{X}_n$ is sufficient for μ .

$X_1, \dots, X_n \sim \text{iid } N(\mu, \sigma^2)$

$\theta = \mu$

↓
pdf
known.
is a product of pdf

i) Figure out distn of $X | X_n$

Joint distn of X : $f_n(x_1, \dots, x_n) = \prod_i f_{x_i}(x_i) = \text{Pdf of gaussian}$

$$= \frac{1}{\sqrt{2\pi}\sigma} \exp \left\{ -\frac{1}{2\sigma^2} \sum_i (x_i - \mu)^2 \right\}$$

$$= \frac{1}{(\sqrt{2\pi}\sigma)^n} \exp \left\{ -\frac{1}{2\sigma^2} \sum_i (x_i - \mu)^2 \right\}$$

$$T(X) = \bar{X}_n \sim N(\mu, \sigma^2/n)$$

now figure pdf of suff. stat.

$$f_{T(n)}(t) = \frac{1}{\sqrt{2\pi} \frac{\sigma}{\sqrt{n}}} \exp \left\{ -\frac{n}{2\sigma^2} (t-\mu)^2 \right\}$$

which is ratio of their

$$f_{x|T(x)}(x_1, \dots, x_n | \bar{x}_n) =$$

$$\frac{(\sqrt{2\pi}\frac{\sigma}{\sqrt{n}})^n \exp \left\{ -\frac{1}{2\sigma^2} \sum (x_i - \mu)^2 \right\}}{\sqrt{2\pi} \frac{\sigma}{\sqrt{n}} \exp \left\{ -\frac{n}{2\sigma^2} (\bar{x} - \mu)^2 \right\}}$$

constant

$$= \text{const.} \cdot X \exp \left\{ -\frac{1}{2\sigma^2} \sum x_i^2 + \frac{2\mu}{2\sigma^2} \sum x_i - \frac{n\mu^2}{2\sigma^2} + \frac{n}{2\sigma^2} [\bar{x}^2 - 2\bar{x}\bar{x} + \mu^2] \right\}$$

cancel out

& thus doesn't depend on μ .

Factorization Theorem

Theorem 1 (Factorization Theorem).

Let $f(x|\theta)$ be the pdf of X . Then $T(X)$ is a sufficient statistic if and only if there exist functions $g(t|\theta)$ and $h(x)$ such that $f(x|\theta) = g(T(x)|\theta)h(x)$.

Eg 1: $x_1, \dots, x_n \sim \text{iid } N(\mu, \sigma^2) \quad \theta = (\mu, \sigma^2)$

(pdf) of gaussian

$$f_{\mu, \sigma^2}(x) = \left(\frac{1}{\sqrt{2\pi}\sigma} \right)^n \exp \left\{ -\frac{1}{2\sigma^2} \sum_i (x_i - \mu)^2 \right\} =$$

$$= \left(\frac{1}{\sqrt{2\pi}\sigma} \right)^n \exp \left\{ -\frac{1}{2\sigma^2} \sum x_i^2 + \frac{2\mu}{2\sigma^2} \sum x_i - \frac{n\mu^2}{2\sigma^2} \right\}$$

parameter

data

$$h(x) = 1$$

suff. st
 $T(x) = (\sum x_i^2, \sum x_i)$

- 1 to 1 correspondence - s^2 defined before in class/ex

eg say: (s^2, \bar{x}_n)

$$s^2 = \frac{1}{n-1} (\sum x_i^2 - n(\bar{x})^2)$$

Eg²) Bernoulli(p)

$$X = (x_1, \dots, x_n) \sim \text{iid } \{ \begin{matrix} 1 & p \\ 0 & 1-p \end{matrix}$$

y Q - what is joint statis suff.

Step 1:
* write down joint distn'

$$f_X(x_1, \dots, x_n) = \prod_{i=1}^n f_{X_i}(x_i)$$

$$= \prod p^{x_i} (1-p)^{1-x_i}$$

$$= p^{\sum x_i} (1-p)^{\sum (1-x_i)}$$

$$f_{X_i}(x_i) = p^{x_i} (1-p)^{1-x_i}$$

$$x_i = 1 : p^1 (1-p)^0 \\ p^1 \cdot 1 = p$$

$$x_i = 0 : p^0 (1-p)^1 \\ 1 \cdot 1$$

Step 2: Sufficient statistics

$$= p^{\sum x_i} (1-p)^{n - \sum x_i}$$

data shows so is suff. st.

$\hookrightarrow T(x) = \sum x_i$

$$g(T, \theta) \cdot h(n) = 1$$

Factorization Theorem: examples

drew from
unit interval

Ex3) $x_1, \dots, x_n \sim \text{iid } U[\theta, 1+\theta]$

* what is pdf for one obs $f_{x_i}(x_i) = \begin{cases} 1 & x_i \in [\theta, 1+\theta] \\ 0 & \text{otherwise} \end{cases}$

= $\underbrace{\prod_{i=1}^n}_{[\theta, 1+\theta]}(1)$ indicator fn.

* their joint distn:

(pdf) $f_{\mathbf{x}}(\mathbf{x}) = \prod_i f_{x_i}(x_i) = \begin{cases} 1 & \theta \leq \min_i x_i \leq \max_i x_i \leq 1+\theta \\ 0 & \text{otherwise} \end{cases}$

sufficient statistic

$$T(\mathbf{x}) = (\min_i x_i, \max_i x_i)$$

points in between not informative

whose pdf is $g(T; \theta) \cdot h(n) = 1$

Factorization Theorem: examples

Estimators. Unbiased estimators

calculated from data

unbiased: if on avg correct

We say that $\hat{\theta} = T(X)$ is unbiased for θ if $E_{\theta}[T(X)] = \theta$ for all possible values of θ where E_{θ} is the expectation when θ is the true parameter value.

$$\text{Bias}(\hat{\theta}) = E_{\theta}[\hat{\theta}] - \theta$$

$$E[T(X)] = E[\hat{\theta}] = \theta$$

$$\text{Bias} = E[\hat{\theta}] - \theta$$

Eg 1) $X_1, \dots, X_n \sim \text{Poisson}(\lambda)$ is a distn that takes values $n \in \{0, 1, 2, 3, \dots\}$ and probability

unknown parameter

values $n \in \{0, 1, 2, 3, \dots\}$

$$P(X_i = k) = \frac{e^{-\lambda} \lambda^k}{k!} \quad \left. \begin{array}{l} \text{all prob.} \\ \text{are the} \end{array} \right\}$$

2 sum
to
1

$$\sum_{k=0}^{\infty} e^{-\lambda} \frac{\lambda^k}{k!} = e^{-\lambda} \sum_{k=0}^{\infty} \frac{\lambda^k}{k!}$$

$$= e^{-\lambda} \cdot e^{\lambda} = 1$$

$$E(X_i) = \sum_{k=0}^{\infty} k \cdot e^{-\lambda} \frac{\lambda^k}{k!} = \sum_{k=1}^{\infty} k e^{-\lambda} \frac{\lambda^k}{k!} = \sum_{k=1}^{\infty} e^{-\lambda} \frac{\lambda^k}{(k-1)!}$$

$$= e^{-\lambda} \cdot \lambda \sum_{k=1}^{\infty} \frac{\lambda^{k-1}}{(k-1)!}$$

$$\hat{\lambda} = \sum_{j=0}^{\infty} \frac{\lambda^j}{j!}$$

$$= e^{-\lambda} \cdot \lambda e^{\lambda} = \lambda$$

$$\hat{\lambda} = \bar{x}_n$$

\uparrow unbiased

$$\hat{x}_i = x_S \text{ (or } x_i)$$

\uparrow unbiased est.

can also do weighted avg.

$$\text{Var}(x_i) = \lambda \text{ also } \underline{\text{unbiased.}}$$

\uparrow

Thus, is a property.

Ex2)

$x_i \sim \text{iid Bernoulli}(p)$

$$E x_i = p$$

$$\hat{p} = \frac{\sum_{i=1}^n x_i}{n}$$

unbiased.

Unbiasness

yes, many cases.

Is it easy to get an unbiased estimator? An example.

e.g.) $x_1, \dots, x_n \sim \text{iid Bernoulli}(p)$

change parameter to be $\theta = \frac{1}{p}$ (not just p as before)

If I put $\hat{\theta} = \frac{1}{\hat{p}}$ will not unbiased bec

$$E\hat{\theta} = E\frac{1}{\hat{p}} \neq \frac{1}{Ep} \quad \begin{array}{|l} \text{face a non-linear} \\ \text{relation.} \end{array}$$

Q) Why not have UB. est.?

Suppose $\hat{\theta} = T(x_1, \dots, x_n)$

$$E(\hat{\theta}) = \sum_{(x_1, \dots, x_n) \in \{0,1\}^n} T(x_1, \dots, x_n) \cdot p^{x_1} (1-p)^{n-x_1} = \frac{1}{p}$$

T
polynomial

↑
 w cannot be
written as
polynomial

Thus, no unbiased est.

be equal to $1/p$.

Non-uniform trans introduce bias — and can correct
by introducing Boot strap bias corr

Bootstrap bias correction

$X_1, \dots, X_n \sim \text{iid } \mu = E[X_i] \hat{\mu} = \bar{X}$
but observed in θ i.e. $= g(\bar{X})$

Setup

Parameter of interest $\theta = g(\mu)$, $\mu = E[X_i]$

Natural estimate $\hat{\theta} = g(\bar{X}_n)$ is biased if g is non-linear.

Bootstrap bias-correction

- (1) For each $b = 1, \dots, B$ generate $\{X_{ib}^*\}$ from $\{X_1, \dots, X_n\}$ with replacement;
of samples generates your sample sample of size n
- (2) Calculate $\bar{X}_b^* = \frac{1}{n} \sum_{i=1}^n X_{ib}^*$;
- (3) Estimate $\theta_b^* = g(\bar{X}_b^*)$; $\theta_b^* = \frac{1}{B} \sum_{b=1}^B \theta_b^*$
- (4) Bias* $= \frac{1}{B} \sum_{b=1}^B \theta_b^* - \hat{\theta} \approx \text{Bias}$.
- (5) Use $\tilde{\theta} = \hat{\theta} - \text{Bias}^*$ as your estimate.

$$\text{Bias} = E[\hat{\theta}] - \theta$$

$\downarrow \Sigma_B$
 $\hat{\theta}$ of the orig. sample

Bootstrap bias correction

Why it works?

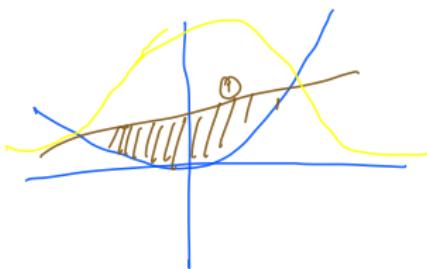
g have $\hat{\mu} = \bar{x} \xrightarrow{P} u$ &
know $\hat{\theta} = g(\hat{\mu}) \xrightarrow{P} g(u) = \theta$

$$E(\hat{\theta} - \theta) = E(g(\hat{\mu}) - g(u)) = g'(u) \underbrace{E(\bar{x} - u)}_0 + \underbrace{\frac{E(g''(u^*) (\bar{x} - u)^2)}{2}}$$

so Bias = $\frac{g''(u)}{2} \cdot \frac{\sigma^2}{n}$

how non linear
①

how spread out



treat as constant, close to u
ss

$$\Rightarrow \frac{g''(u)}{2} E(\bar{x} - u)^2$$

$\text{Var}(x) = \frac{\sigma^2}{n}$

my

And, Bias* $\cong \frac{g''(\hat{u})}{2} \cdot \left(\frac{\sigma^*}{n} \right)^2$

Many unbiased est.

Estimators. Efficiency

how close the realists are

$$\underline{\text{MSE}}(\hat{\theta}) = \mathbb{E}_{\theta}[(\hat{\theta} - \theta)^2]$$

↑ from the truth.

Theorem 2. But don't know so we leave

$$\underline{\text{MSE}(\hat{\theta}) = \text{Bias}^2(\hat{\theta}) + \text{Var}(\hat{\theta})}$$

←

$$\begin{aligned}\mathbb{E}[(\hat{\theta} - \theta)^2] &= \mathbb{E}[(\hat{\theta} - \mathbb{E}\hat{\theta} + \mathbb{E}\hat{\theta} - \theta)^2] \\ &= \mathbb{E}[(\hat{\theta} - \mathbb{E}\hat{\theta})^2 + (\mathbb{E}\hat{\theta} - \theta)^2 + 2(\hat{\theta} - \mathbb{E}\hat{\theta})(\mathbb{E}\hat{\theta} - \theta)] \\ &\quad \text{random} \\ &= \text{Var}(\hat{\theta}) + \text{Bias}^2 + \mathbb{E}[(\hat{\theta} - \mathbb{E}\hat{\theta}) \cdot (\mathbb{E}\hat{\theta} - \theta)] \\ &\quad \uparrow \\ &\quad \theta : \text{pop} \\ &\quad \text{constant,} \\ &\quad \text{exists}\end{aligned}$$

Trade-off between bias and variance

* $\hat{\theta}$ is some unbiased est. of θ — i.e. $E\hat{\theta} = \theta$

Introduce new est. $\tilde{\theta} = \underbrace{(1-c)\hat{\theta} + c \cdot 0}_{\text{shrink it to close 0.}}$ c is small constant
mean of $\tilde{\theta}$:

$$E(\tilde{\theta}) = (1-c)\hat{\theta} + c \cdot 0 = (1-c)\hat{\theta}$$

$$\text{so bias } (\tilde{\theta}) = E\tilde{\theta} - \hat{\theta} = -c\hat{\theta}$$

Now, var of $\tilde{\theta}$

$$\text{var}(\tilde{\theta}) = \text{var}((1-c)\hat{\theta}) = (1-c)^2 \text{var}(\hat{\theta})$$

$$\text{MSE}(\tilde{\theta}) = (1-c)^2 \text{var}(\hat{\theta}) + c^2 \theta^2 = g(c) \text{ trying to minimize it}$$

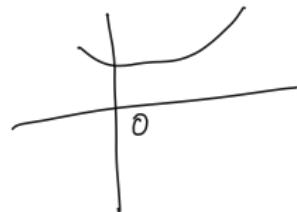
$\approx \sigma^2$

Trade-off between bias and variance

$$g'(c) = -2\sigma^2(1-c) + 2\theta^2 c$$

$$\underline{g'(0)} = -2\sigma^2(1-\theta) \stackrel{< 0}{=} \text{is -ve so } 0 \text{ not an optimal point.}$$

so need more shrinkage.



Connection between efficiency and sufficient statistics

Theorem 3 (Rao-Blackwell).

Assume that $T(X)$ is a sufficient statistic for parameter θ and $\hat{\theta} = \delta(X)$ is an estimator for θ . Define $\hat{\theta}_2 = \mathbb{E}[\delta(X)|T(X)]$ is an estimator for θ as well and $MSE(\hat{\theta}_2) \leq MSE(\hat{\theta})$. In addition, if $\hat{\theta}$ is unbiased, then $\hat{\theta}_2$ is unbiased as well.

