

Chapter 6. Principles of Data Reduction

6.1 Introduction

$$X_1, \dots, X_n \stackrel{iid}{\sim} f(x|\theta)$$

추출의 대상.

(x_1, \dots, x_n) : obs., just a long list of numbers that is hard to interpret.

$$\underline{X} = (X_1, \dots, X_n) \xrightarrow{T} T(\underline{X}), \text{ statistic}$$

"data reduction", "data summary".

Three principles

i) sufficiency : data reduction that does not discard information about θ .

ii) likelihood : a function of θ that contains all the information about θ that is available from the sample

iii) equivariance : data reduction that preserves some important features of the model.

6.2 The sufficiency principle.

Sufficiency : If $T(\underline{X})$ is a suff. statistic for θ , then any inference about θ should depend on the sample \underline{X} only through the value $T(\underline{X})$.

That is, if \underline{x} and \underline{y} are two sample points such that $T(\underline{x})=T(\underline{y})$, then the inference about θ should be the same whether $\underline{X}=\underline{x}$ or $\underline{Y}=\underline{y}$ is observed.

Defn: A statistic $T(\underline{X})$ is a SS for θ if the conditional distribution of the sample \underline{X} given the value of $T(\underline{X})$ does not depend on θ .

Th^m: $\underline{X} \sim p(\underline{X}|\theta)$, $T(\underline{X}) \sim f(t|\theta)$ is a SS fr θ , if
 $p(\underline{X}|\theta) / f(T(\underline{X})|\theta)$ is constant as a function of θ .

(e.g.) $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Bernoulli}(\theta), 0 < \theta < 1.$

$T(\underline{X}) = X_1 + \dots + X_n$: SS for θ ?

Sol. To check if $p(\underline{x}|\theta) / f(T(\underline{x})|\theta)$ does not depend on θ .

$$p(\underline{x}|\theta) = \theta^{\sum x_i} (1-\theta)^{n-\sum x_i}$$

$$f(T(\underline{x})|\theta) = \binom{n}{\sum x_i} \theta^{\sum x_i} (1-\theta)^{n-\sum x_i}$$

$$p(\underline{x}|\theta) / f(T(\underline{x})|\theta) = \binom{n}{\sum x_i} : \text{does not depend on } \theta.$$

$\therefore T(\underline{X}) = \sum_{i=1}^n X_i$ is a SS for $\theta.$ //

(e.g.) $X_1, \dots, X_n \stackrel{iid}{\sim} N(\theta, 1), \theta \in \mathbb{R}$

\bar{X} is a ss for θ .

Sol.
$$\frac{f(x|\theta)}{f(T(x)|\theta)} = \frac{(2\pi)^{-n/2} \exp\left(-\frac{1}{2} \sum_{i=1}^n (x_i - \theta)^2\right)}{(2\pi \cdot \frac{1}{n})^{-1/2} \exp\left(-\frac{n}{2} (\bar{x} - \theta)^2\right)}$$
$$= n^{-\frac{1}{2}} \exp\left(-\frac{1}{2} \left\{ \sum x_i^2 - 2n\theta\bar{x} + n\theta^2 \right\} + \frac{n}{2} \left\{ (\bar{x})^2 - 2\bar{x}\cdot\theta + \theta^2 \right\}\right)$$
$$= n^{-\frac{1}{2}} \exp\left(-\frac{1}{2} \sum x_i^2 + \frac{n}{2} (\bar{x})^2\right) \quad : \text{does not depend on } \theta,$$

(e.g.) $X_1, \dots, X_n \stackrel{iid}{\sim} f$. Then $X_{(1)}, \dots, X_{(n)}$ is a sufficient statistic.

Outside of the exp'l family, it is rare to have a ss of smaller dimension than the sample size, so in many cases it will turn out that the o.s. is the best way that we can do.
(not much of a reduction)

Thm: (factorization theorem) $(X_1, \dots, X_n) \sim f(x|\theta)$

$T(X)$ is a SS for $\theta \Leftrightarrow f(x|\theta) = g(T(x)|\theta) h(x)$

(e.g.) $X_1, \dots, X_n \stackrel{iid}{\sim} N(\theta, 1)$

$$f(x|\theta) = (2\pi)^{-n/2} \exp\left(-\frac{1}{2} \sum_i (x_i - \theta)^2\right)$$

$$= (2\pi)^{-n/2} \exp\left(-\frac{1}{2} \sum x_i^2 + n\theta \cdot \bar{x} - \frac{n}{2}\theta^2\right)$$

Put $g(t|\theta) = \exp\left(n\theta \cdot t - \frac{n}{2}\theta^2\right)$ and $h(x) = (2\pi)^{-n/2} \exp\left(-\frac{1}{2} \sum x_i^2\right)$,

(e.g.) $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Unif}(1, 2, \dots, \theta)$

$$f(x|\theta) = \frac{1}{\theta} I(x_i = 1, 2, \dots, \theta)$$

$$f(x|\theta) = \theta^{-n} [I(x_i = 1, 2, \dots, \theta) \text{ for } i=1, 2, \dots, n] = \theta^{-n} I(\max_{1 \leq i \leq n} x_i \leq \theta)$$

Put $g(t|\theta) = \theta^{-n} I(t \leq \theta)$ and $h(x) = 1$.

(e.g.) $X_1, \dots, X_n \stackrel{\text{ind}}{\sim} N(\mu, \sigma^2), \quad \underline{\theta} = (\mu, \sigma^2) \in \mathbb{R} \times \mathbb{R}_+$

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

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

$\therefore (\sum x_i, \sum x_i^2)$ is a SS for $\underline{\theta} = (\mu, \sigma^2)$

$l-1$

$$(\bar{x}, s^2)$$

$\therefore (\bar{x}, s^2)$ is a SS for (μ, σ^2) . "

Ih^m: $X_1, \dots, X_n \stackrel{\text{Ind}}{\sim} f(x|\underline{\theta}) = h(x) c(\underline{\theta}) \exp\left(\sum_{j=1}^k w_j(\underline{\theta}) t_j(x)\right)$, $\underline{\theta} = (\theta_1, \dots, \theta_d)$, $d \leq k$.
 Then, $T(X) = \left(\sum_{i=1}^n t_1(X_i), \dots, \sum_{i=1}^n t_k(X_i) \right)$ is a SS for $\underline{\theta}$.

Pf: Trivial by the factorization theorem. "

Defn: A suff. stat. $T(\underline{x})$ is called a MSS

if, for any other SS $T'(\underline{x})$, $T(\underline{x})$ is a func. of $T'(\underline{x})$.

"If $T(x) = T(y)$, then $T(x) = T(y)$ ".

- the greatest possible data reduction
for a sufficient statistic.
- any 1-1 func. of a MSS is also a MSS.

Then: $\tilde{X} \sim f(x|\theta)$

$\exists T(\cdot)$ s.t. $f(x|\theta)/f(y|\theta)$ is constant in θ iff $T(x)=T(y)$.

Then, $T(X)$ is a MSS for θ .

(e.g.) $X_1, \dots, X_n \stackrel{iid}{\sim} N(\theta, 1)$

$$\frac{f(x|\theta)}{f(y|\theta)} = \exp\left(-\frac{1}{2}\left\{\sum_{i=1}^n (x_i - \theta)^2 - \sum_{i=1}^n (y_i - \theta)^2\right\}\right) = \exp\left(-\frac{1}{2}\left\{\sum x_i^2 - \sum y_i^2 - \theta \cdot (\sum x_i - \sum y_i)\right\}\right)$$

is const. in θ iff $\sum x_i = \sum y_i$.

$\therefore T(X) = \sum_{i=1}^n x_i$ is a MSS for θ . "

(e.g.) $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Unif}(\theta, \theta+1)$.

$$f(x|\theta) = I(\max_{i=1}^n x_i - 1 < \theta < \min_{i=1}^n x_i)$$

$f(x|\theta)/f(y|\theta)$ is const. in θ iff $\min_{i=1}^n x_i = \min_{i=1}^n y_i$ and $\max_{i=1}^n x_i = \max_{i=1}^n y_i$.

$\therefore (X_{(1)}, X_{(n)})$ is a MSS for θ .

Defn: A statistic $S(X)$ whose distribution does not depend on the parameter θ is called an ancillary statistic.

(e.g.) $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Unif}(\theta, \theta+1)$

$R = X_{(n)} - X_{(1)}$ is an ancillary statistic for θ .

(e.g.) (location family ancillary) $X_1, \dots, X_n \stackrel{iid}{\sim} f(x-\theta), \theta \in \mathbb{R}$.

$R = X_{(n)} - X_{(1)}$ is an ancillary statistic for θ .

sol. Consider $Z_1, \dots, Z_n \stackrel{iid}{\sim} f(x)$. " $\theta = 0$ "

$$X_1 \stackrel{d}{=} Z_1 + \theta, X_2 \stackrel{d}{=} Z_2 + \theta, \dots, X_n \stackrel{d}{=} Z_n + \theta.$$

$$F_R(r|\theta) = P_\theta(R \leq r)$$

$$= P_\theta(X_{(n)} - X_{(1)} \leq r)$$

$$= \underbrace{P_\theta(Z_{(n)} - Z_{(1)} \leq r)}$$

does not depend on θ . "

(e.g) (scale family ancillary) $X_1, \dots, X_n \stackrel{iid}{\sim} \frac{1}{\sigma} f(\frac{x}{\sigma}), \sigma > 0$.

$(X_1/X_n, X_2/X_n, \dots, X_{n-1}/X_n)$ is an ancillary stat. for σ .

- A MSS achieves the maximal data reduction possible while retaining "all" the information about θ .

\Rightarrow MSS $\perp\!\!\!\perp$ ancillary stat. ?

NOT necessarily !

BUT in many important cases, Yes !

Defn: $T(\underline{X})$ is called a complete statistic if

$$E_{\theta} g(T) = 0 \quad \forall \theta \quad \text{implies} \quad P_{\theta}(g(T) = 0) = 1 \quad \forall \theta.$$

$$\left(\int g(t) f_T(t|\theta) dt \stackrel{\forall \theta}{=} 0 \iff \int_{g(t)=0} f_T(t|\theta) dt \stackrel{\forall \theta}{=} 1 \right)$$

"Completeness is a property of a family of prob. dist's."

(e.g.) $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Bernoulli}(p)$, $0 < p < 1$

$T = \sum_{i=1}^n X_i \sim \text{Bin}(n, p)$ is a complete statistic.

Sol. Let g be a func. s.t. $E_p g(T) = 0 \quad \forall p \in (0, 1)$.

Then $0 = E_p g(T) = \sum_{t=0}^n g(t) \binom{n}{t} p^t (1-p)^{n-t}$
 $= (1-p)^n \sum_{t=0}^n g(t) \binom{n}{t} \left(\frac{p}{1-p}\right)^t \quad \forall p \in (0, 1)$.

$$\underbrace{\sum_{t=0}^n g(t) \binom{n}{t} r^t}_> 0 \quad \forall r > 0.$$

$\therefore g(t) = 0$ for all $t = 0, 1, 2, \dots, n$.

$$\therefore P_p(g(T) = 0) \stackrel{\forall p}{=} 1 .$$

(e.g.) $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \text{Unif}(0, \theta), \theta > 0$.

$T(X) = X_{(n)}$ is a complete statistic.

Sol. $f_T(t|\theta) = n t^{n-1} \theta^{-n} I(0 < t < \theta)$

Let g be a func. s.t. $E_\theta g(T) = 0 \forall \theta > 0$.

Since $E_\theta g(T)$ is a const. in θ , $\frac{d}{d\theta} E_\theta g(T) = 0, \forall \theta > 0$.

$$0 = \frac{d}{d\theta} \int_0^\theta q(t) \cdot n t^{n-1} \theta^{-n} dt = \frac{d}{d\theta} \left[\theta^{-n} \int_0^\theta q(t) n t^{n-1} dt \right]$$

$$= -n \theta^{-n-1} \int_0^\theta q(t) n t^{n-1} dt + \theta^{-n} \cdot q(\theta) n \theta^{n-1}$$

$$= -n \theta^{-1} \int_0^\theta q(t) n t^{n-1} \theta^{-n} dt + n q(\theta) \theta^{-1}$$

$$= -n \theta^{-1} \cdot 0 + n q(\theta) \theta^{-1} = \underbrace{n \theta^{-1}}_{> 0, \forall \theta > 0} \cdot q(\theta) \quad \therefore q(t) = 0 \forall t > 0. //$$

Th^m: (Basu) $T(X)$: a complete and (minimal) suff. stat. for θ .
 Then, $T(X)$ is indep. of every ancillary statistic.

pf. (only for discrete distⁿ's)

$S(X)$: an ancillary stat., i.e. $P(S(X)=s)$ does not depend on θ .

$P(S(X)=s | T(X)=t)$ does not depend on θ .

WTS : $P(S(X)=s | T(X)=t) = P(S(X)=s), \forall s, t.$

$$P(S(X)=s) = \sum_t P(S(X)=s, T(X)=t)$$

$$= \sum_t P(S(X)=s | T(X)=t) \cdot P_\theta(T(X)=t)$$

$$\sum_t \underbrace{\left[P(S(X)=s) - P(S(X)=s | T(X)=t) \right]}_{g(T)} \cdot P_\theta(T(X)=t) = 0$$

//

$$\text{Def}: X_1, \dots, X_n \stackrel{\text{iid}}{\sim} f(x|\theta) = h(x) c(\theta) \exp\left(w_1(\theta)t_1(x) + w_2(\theta)t_2(x) + \dots + w_k(\theta)t_k(x)\right)$$

$$\text{where } \underline{\theta} = (\theta_1, \theta_2, \dots, \theta_k).$$

Then, $T(\underline{x}) = \left(\sum_{i=1}^n t_1(X_i), \sum_{i=1}^n t_2(X_i), \dots, \sum_{i=1}^n t_k(X_i) \right)$ is complete

if $\{(\omega_1(\theta), \dots, \omega_k(\theta)) : \theta \in \mathbb{H}\}$ contains an open set in \mathbb{R}^k .

$$(\text{e.g.}) X_1, \dots, X_n \stackrel{\text{iid}}{\sim} f(x|\theta) = \frac{1}{\theta} e^{-x/\theta} I(x>0), \theta > 0.$$

$$E\left(\frac{X_1}{X_1 + \dots + X_n}\right) = ?$$

Sol. $S(\underline{x}) = \frac{X_1}{X_1 + \dots + X_n}$ is an ancillary stat. (\because scale family)

$f(x|\theta)$ forms an exp'l family with $t_1(x)=x$ and $w_1(\theta)=\frac{1}{\theta}$. $\Rightarrow T(\underline{x}) = \sum_i^n X_i$
 $\left\{ \frac{1}{\theta} : \theta > 0 \right\} = (0, \infty)$ contains an open set in \mathbb{R}^1 . is complete.

\therefore By Basu, $S(\underline{x})$ and $T(\underline{x})$ are indep.

$$\underbrace{E(X_1)}_{\theta} = E(S(\underline{x}) \cdot T(\underline{x})) = \underbrace{E(S(\underline{x}))}_{?} \cdot \underbrace{E(T(\underline{x}))}_{n\theta}$$

$$\therefore \frac{1}{n}, //$$

(e.g.) $X_1, \dots, X_n \stackrel{\text{ iid }}{\sim} N(\mu, 1)$.
 \bar{X} and S^2 are indep. !

Th^m: If a MSS exists, then any complete statistic is also a MSS.