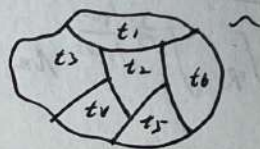


Principles of Data Reduction

- Experimenter uses the information in X_1, X_2, \dots, X_n to infer an unknown parameter θ . If n is large. Then the observed sample is too large to interpret.

\Rightarrow We won't use $T(\vec{X})$ rather than the entire sample \vec{X} .

For $T(\cdot)$, it can partition sample space into different sets $A_t = \{\vec{x} \in \mathcal{X} \mid T(\vec{x}) = t\}$.



\Rightarrow For θ , we're interested in the information relevant θ . Discard information irrelevant θ .

(1) The Sufficient Principle:

- $T(\vec{X})$ is sufficient statistic for θ . Then it captures all information of θ . θ only depends on X through $T(\vec{X})$. If $T(\vec{X}) = T(\vec{y})$, then inference of θ are same.

pf: For $X, Y \sim f(x|\theta)$, $T(\vec{X}) = T(\vec{y})$

$$\text{Then } P_{\theta}(X = \vec{x}) = P_{\theta}(X = \vec{x} \mid T(X) = T(\vec{x})) P_{\theta}(T(X) = T(\vec{x}))$$

$$= P_{\theta}(Y = \vec{x} \mid T(Y) = T(\vec{y})) P_{\theta}(T(Y) = T(\vec{y}))$$

$$= P_{\theta}(Y = \vec{x})$$

① Criterion:

By definition: $P_{\theta}(X=x | T(X)=T(x)) = \frac{P_{\theta}(X=x, T(X)=T(x))}{P_{\theta}(T(X)=T(x))}$

$$= \frac{p(x|\theta)}{q(T(x)|\theta)} = h(x). \quad q(\cdot) \text{ is pdf of } T(X).$$

\Rightarrow Thm. $T(X)$ is sufficient statistic $\Rightarrow \frac{p(x|\theta)}{q(T(x)|\theta)} = h(x)$

★ Factorization Thm:

$X \sim f(x|\theta)$. $T(X)$ is sufficient statistic for θ .

$\Leftrightarrow f(\vec{x}|\theta) = q(T(\vec{x})|\theta) h(\vec{x})$

p.f. (\Rightarrow) . Proved above

(\Leftarrow)
$$\frac{f(x|\theta)}{q(T(x)|\theta)} = \frac{q(T(\vec{x})|\theta) h(\vec{x})}{\int_{A(\theta)} f(x|\theta) dx} = \frac{q(T(\vec{x})|\theta) h(\vec{x})}{\int_{A(\theta)} q(T(\vec{y})|\theta) h(\vec{y}) d\vec{y}}$$

$$= \frac{q(T(\vec{x})|\theta) h(\vec{x})}{q(T(\vec{x})|\theta) \int_{A(\theta)} h(\vec{y}) d\vec{y}} = r(\vec{x})$$

② Sufficient Statistic Vector:

• It's usually for multiple parameters.

Thm. (For exponential family)

$X_k \sim f(x|\vec{\theta})$, $1 \leq k \leq n$, i.i.d. $f(x|\vec{\theta}) = h(x) c(\vec{\theta}) e^{\sum_{k=1}^K \eta_k(\vec{\theta}) t_k(x)}$

$\vec{\theta} = (\theta_1, \dots, \theta_K)$, $1 \leq k$. Then $T(\vec{x}) = (\sum_{j=1}^n t_1(x_j), \dots, \sum_{j=1}^n t_K(x_j))$

is sufficient statistic vector for $\vec{\theta}$.

③ Minimal sufficient Statistic:

• Note that: \vec{X} is the largest sufficient statistic.

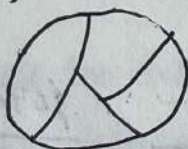
$\forall g, g(x)$ is one-to-one func. Then for $T(x)$ S-S.

$g(T(x))$ is S-S. too.

\Rightarrow We want to find a statistic carrying little information, but retaining all information of θ .

Def: $T(x)$ is minimal sufficient statistic, if any other S-S, $T'(x)$, $T'(x) = T'(y) \Rightarrow T(x) = T(y)$. ($T = f(T')$)

At(x)



$T(x)$ submits a coarsest partition!
 \Rightarrow Achieve greatest reduction.

Remark: i) If MLE is a sufficient statistic. Then it's the minimal.

pf: For \forall sufficient statistic $T(x)$

$$f(x|\theta) = g(T(x)|\theta) h(x)$$

$$\lambda(x) = h(x) \sup_{\theta} g(T(x)|\theta) = h(x) g(T(x)|\hat{\theta})$$

where $\hat{\theta}$ satisfies $g(T(x)|\hat{\theta}) = \sup_{\theta} g(T(x)|\theta)$. $\therefore \hat{\theta} = \hat{\theta}(T(x))$

ii) Most location family $f(x-\theta)$. The most reduction is $(X_{(1)}, X_{(2)}, \dots, X_{(n)})$. So \vec{x} is the minimal S-S.

iii) Minimal S-S isn't unique. Since $\forall g$ one-to-one if T is minimal. Then $g(T)$ is minimal, too.

Thm. (Criterion)

$f(x|\theta)$ is pmf or pdf of X . For $T(\vec{x})$, $\forall \vec{x}, \vec{y} \in \mathcal{X}$.

$\frac{f(\vec{x}|\theta)}{f(\vec{y}|\theta)}$ is irrelevant with $\theta \Leftrightarrow T(\vec{x}) = T(\vec{y})$

Then $T(\vec{x})$ is minimal s.s. for θ .

Pf: 1) $T(\vec{x})$ is s.s.:

$$\text{Note that } f(x|\theta) = f(x_{T(\vec{x})}|\theta) \frac{f(x|\theta)}{f(x_{T(\vec{x})}|\theta)}$$

$$\text{where } T(x_{T(\vec{x})}) = T(\vec{x}), \quad f(x_{T(\vec{x})}|\theta) \triangleq g(T(\vec{x})|\theta)$$

$$\therefore f(x|\theta) / f(x_{T(\vec{x})}|\theta) = h(\vec{x}), \text{ By Factorization Thm } \square$$

2) $T(\vec{x})$ is minimal:

By Factorization Thm: $\tilde{T}(\vec{x}) = \tilde{T}(\vec{y})$ for any other s.s. $\tilde{T}(\vec{x})$

$$\Rightarrow \frac{f(x|\theta)}{f(\vec{y}|\theta)} \text{ is irrelevant with } \theta, \therefore T(\vec{x}) = T(\vec{y})$$

Remark: i) \forall s.s. satisfies (\Leftarrow) . But only minimal s.s.

satisfies (\Rightarrow) . Minimal is "necessary".

ii) To find minimal s.s. Firstly calculate

$f(x|\theta) / f(\vec{y}|\theta)$. Find the condition that is necessary and sufficient the equation doesn't contain θ . (i.e. $T(\vec{x}) = T(\vec{y})$)

Thm. For family of density $\{f_k(x)\}_0^k$ has common

$$\text{support. Then: i) } T(\vec{x}) = \left(\frac{f_1(x)}{f_0(x)}, \frac{f_2(x)}{f_0(x)}, \dots, \frac{f_k(x)}{f_0(x)} \right)$$

is minimal sufficient for the family ii) if $T(\vec{x})$

is sufficient for the family. Then $T(x)$ is minimal.

Remark: The sufficient minimal statistic can be extended to nonparametric family:

• If $X_1, X_2, \dots, X_n \sim f$, unknown density. Then the order statistic is minimal sufficient.

Pf: If $f \sim \text{logistic}(\alpha, \beta) \Rightarrow$ By Thm: minimal S.S.

$$T(x) = \left(\frac{\pi f_1(x_1)}{\pi f_0(x_1)}, \dots, \frac{\pi f_n(x_n)}{\pi f_0(x_n)} \right) \xleftrightarrow{\text{one-to-one}} (X_{(1)}, \dots, X_{(n)})$$

If f is nonparametric, by ss, $(X_{(1)}, \dots, X_{(n)})$ is sufficient. So it's the minimal.

④ Ancillary Statistic:

• Def: A statistic $S(x)$ whose dist doesn't depend on θ . It's the ancillary statistic for θ .

Remark: It contains no information about θ .

e.g. $\left\{ \begin{array}{l} \text{i) For location family: } f(x-\theta). \\ \quad R = X_{(n)} - X_{(1)} \text{ is ancillary. (or } X - Y) \\ \text{ii) For scale family: } f\left(\frac{x}{\theta}\right) \\ \quad X/Y \text{ is ancillary statistic.} \end{array} \right.$

Remark: Minimal S.S. isn't indept with ancillary statistic.

It may contain ancillary statistic: $X \sim N(0, \theta^2)$

$X_{(n)}, X_{(1)}$ minimal S.S. But $X_{(n)} - X_{(1)}$ is ancillary.

③ Complete Statistic:

Def: For statistic $T(\vec{X}) \sim f(x|\theta)$. It's called complete if $\forall g, E_\theta(g(T)) = 0, \forall \theta$, then $P_\theta(g(T) = 0) = 1, \forall \theta \in \Theta$.

Remark: i) It means any transformation of T won't contain ancillary statistic. If $E_\theta(g(T))$ is irrelevant with θ , then $g(T) \equiv C$ w.p.1.
So the definition can be refined: $E_\theta(g(T)) \equiv C, \forall \theta$
 $\Rightarrow P_\theta(g(T) \equiv C) = 1, \forall \theta$.

Basu's Thm:

If $T(\vec{X})$ is complete sufficient statistic (minimal)
Then $T(\vec{X})$ is indept with every ancillary statistic.

Pf: Suppose $S(\vec{X})$ is ancillary, indept with θ .

Prove: $P(S(\vec{X}) = s | T(\vec{X}) = t) = P(S(\vec{X}) = s)$

$$\begin{aligned} \text{Note that } P(S(\vec{X}) = s) &= \sum_t P(S(\vec{X}) = s, T(\vec{X}) = t) \\ &= \sum_t P(S(\vec{X}) = s | T(\vec{X}) = t) P(T(\vec{X}) = t) \end{aligned}$$

$$g(T(\vec{X})) = P(S(\vec{X}) = s | T(\vec{X}) = t) - P(S(\vec{X}) = s). \quad \square$$

Remark: i) The converse is false.

ii) "sufficient" is necessary. Since there exists complete statistic not being sufficient:

$\{X_k\} \sim \text{Poisson}(\lambda)$. Let $T(\vec{X}) = X_1$. Complete.

iii) The minimal can be omitted:

Prop. If minimal s.s exists. Then any complete sufficient statistic is minimal.

For exponential family:

$X_k \sim f(x|\vec{\theta}) = h(x) c(\vec{\theta}) e^{\sum_{i=1}^k w(\theta_i) t_i(x)}$, i.i.d.

$1 \leq k \leq n$, $\vec{\theta} = (\theta_1, \dots, \theta_k)$. Then $T(X) = (\sum_{i=1}^n t_1(X_i), \dots, \sum_{i=1}^n t_k(X_i))$

is complete if Θ contains an open set in \mathbb{R}^k .

Pf: By n -dimensional Laplace Transforming.

(2) The Likelihood Principle:

① Def: Given $\vec{X} = \vec{x}$ observed. $L(\theta|\vec{x}) = f(\vec{x}|\theta)$ is called likelihood function.

Remark: It can be used to summarize data.

Likelihood Principle:

If \vec{x} and \vec{y} are two sample points from $f(x|\theta)$.

$\forall \theta$, $L(\theta|\vec{x}) = c(x, y) L(\theta|\vec{y})$. Then the conclusion drawn from \vec{x} and \vec{y} concerning θ are same.

② The formal form:

Def: Experiment $E: (\vec{X}, \theta, \{f(x|\theta)\})$, which means

\vec{X} is random vector with pdf $f(x|\theta)$ for some fixed θ .

\Rightarrow After the experiment was performed and having observed sample $\vec{X} = \vec{x}$, we make some conclusion $Ev(E, \vec{x})$ about θ .

Formal Sufficiency Principle:

For $E = (\vec{X}, \theta, \{f(x|\theta)\})$, $T(\vec{x})$ is s.s. for θ .

If $T(\vec{x}) = T(\vec{\eta})$, $\vec{x}, \vec{\eta}$ are 2 sample points.

Then $Ev(E, \vec{x}) = Ev(E, \vec{\eta})$

Remark: The Likelihood Principle can be used to derive it:

$$L(\theta|\vec{x}) = g(T(\vec{x})|\theta)h(\vec{x}) = g(T(\vec{\eta})|\theta)h(\vec{x}) = \frac{h(\vec{x})}{h(\vec{\eta})} L(\theta|\vec{\eta})$$

Conditionality Principle:

$E_1 = (X_1, \theta, \{f_1(x|\theta)\})$, $E_2 = (X_2, \theta, \{f_2(x|\theta)\})$, only

the unknown θ are common. Consider mixed

experiment: J is r.v. observed. Indep't of X_1, X_2, θ .

$P(J=1) = P(J=2) = \frac{1}{2}$. Perform $E_J = E^* = ((j, X_j), \theta, \{f_j^*(x|\theta)\})$

$X^* = (j, X_j)$, $f^*(X^*|\theta) = \frac{1}{2} f_j(X_j|\theta)$, $\Rightarrow Ev(E^*, X^*) = Ev(E_j, X_j)$

Remark: It means the conclusion only depends on which experiment is performed. Indep't with which one we choose.

⇒ Formal Likelihood Principle:

$E_1 = (X_1, \theta, \{f_1(x_1|\theta)\})$, $E_2 = (X_2, \theta, \{f_2(x_2|\theta)\})$, only the unknown θ are common. x_1^* , x_2^* are sample point from E_1 , E_2 resp. st. $L(\theta|x_1^*) = C L(\theta|x_2^*)$, $C = C(x_1^*, x_2^*)$

Then $E_V(E_1, x_1^*) = E_V(E_2, x_2^*)$

Cor. $E = (\bar{X}, \theta, \{f(x|\theta)\})$ is an experiment. Then :

$E_V(E, \bar{x})$ only depends on \bar{x} and E through $L(\theta|\bar{x})$

Thm. (Birnbawm's Thm)

The formal likelihood principle \Leftrightarrow Formal Sufficient Principle and Conditionality Principle

Def. (\Leftrightarrow) $T(j, x_j) = \begin{cases} (1, x_1^*), & \text{if } j=1, x_1=x_1^* \text{ or } j=2, x_2=x_2^* \\ (j, x_j), & \text{otherwise} \end{cases}$

The key: T

takes on same value

for sample point

$(1, x_1^*)$, $(2, x_2^*)$ from

each experiment

It's on sample space of E^* mixed experiment, $\{(j, x_j)\}$.

(claim: $T(j, x_j)$ is s.s for θ in $E^* \cap f^*$)

$$f^*(j, x_j|\theta) = g(T(j, x_j)|\theta) \quad j \neq 2.$$

$$f^*(2, x_2^*|\theta) = f(x_2^*|\theta) p(J=2) = C(x_1^*, x_2^*) f(x_1^*|\theta) p(J=2)$$

$$= f^*(1, x_1^*|\theta) C(x_1^*, x_2^*) = g(T(2, x_2^*)|\theta) C(x_1^*, x_2^*)$$

$\therefore T(j, X_j)$ is s.s. for $f^*(j, X_j | \theta)$ about θ .

Since $T(1, X_1^*) = T(2, X_2^*)$ By FSP:

$$E_V(E^*(1, X_1^*)) = E_V(E^*(2, X_2^*))$$

$$\text{By CP} = E_V(E^*(1, X_1^*)) = E_V(E_1, X_1^*) = E_V(E_2, X_2^*)$$

(\Rightarrow) Since $f^*(j, X_j | \theta) = p(J=j) f(X_j | \theta)$ By FLP:

$$\therefore L(\theta | (j, X_j)) = C L(\theta | X_j) \Rightarrow E_V(E^*(j, X_j)) = E_V(E_j, X_j)$$

If $T(\vec{x}) = T(\vec{\eta})$. We also have $L(\theta | \vec{x}) = L(\theta | \vec{\eta})$

Remark: i) FSP + CP \Rightarrow Likelihood Principle. (By FLP \Rightarrow LP)

Since $L(\theta | \vec{x}) = C(\vec{x}, \vec{\eta}) L(\theta | \vec{\eta})$. Only one E in LP!

ii) Pf of Cor:

$$\text{For } E = E_1 = E_2, L(\theta | \vec{x}) = L(\theta | \vec{\eta}), C(\vec{x}, \vec{\eta}),$$

$$E_V(E_1, \vec{x}) = E_V(E, \vec{x}) = E_V(E, \vec{\eta}) = E_V(E_2, \vec{\eta})$$

iii) Many common statistical model can't be applied in FLP. Since there may exist some information on θ doesn't base on sufficient statistic. (Violate FSP!)

iv) The pf of Thm isn't compelling. Since before CP coming, we should define s.s for each experiment (Then "the key" won't happen on separate s.s for each experiment; Sample space may be different for E_1, E_2 !)

(3) The equivariance Principle:

It states: If $T(\vec{x}) = T(\vec{\eta})$. Since $E(\vec{E}, \vec{x})$ may be different from $E(\vec{E}, \vec{\eta})$. But there's a certain relationship between the two inferences.

- Type
- i) Measurement Equivariance:
the inference shouldn't depend on measurement scale.
 - ii) Formal Invariance:
two inference problems has same formal structure of model (Dist's $f, f_0 \in \{f(x|0) | 0 \in \Theta\}$ and set of allowable or wrong inferences e.g. $\theta > 0$)
 \Rightarrow Then the same inference procedure should be used.

Equivariance Principle:

$Y = g(X)$ is change of measurement scale, st.

Y has the same formal structure model as X .

\Rightarrow Then the inference procedure should be measure equivariant and formally equivariant.

Remark: The transformations of scale measurement form a "Group". Denote G . For $g \in G$, then

g satisfies

- i) Measurement equivalent: $\exists T(x)$ estimate θ , $Y = g(x)$
 $T(y)$ estimate $g(\theta)$
- ii) Formal invariant: $\exists Y = g(x)$. Then we have
 $T(g(x)) = g(T(x))$
i.e. $T^*(x) = T(x)$

Def: $F = \{f(x|\theta) | \theta \in \Theta\}$. G is a group of transform on sample space X . Then F is invariant under G $\exists \forall \theta \in \Theta, g \in G, \exists \theta' \in \Theta$ st. $Y = g(x) \sim f(x|\theta')$
 $\exists X \sim f(x|\theta)$, initially.