

Chapter 6: Principles of Data Reduction.

6.1 : Introduction:

$$\bar{X} = (X_1, X_2, X_3, \dots, X_n)$$

random sample

$\tilde{x} = (x_1, x_2, x_3, \dots, x_n)$ is a realization from the random sample.

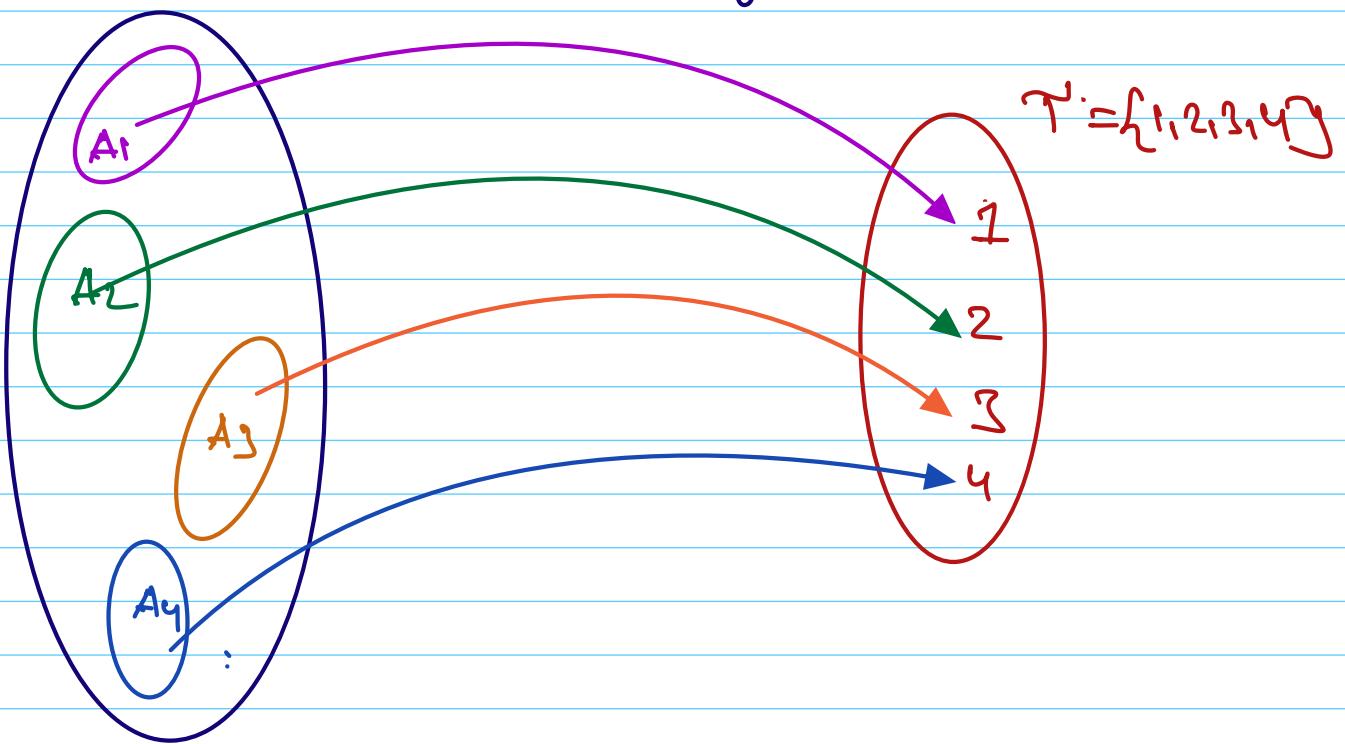
$T(X)$ is a Statistic.

Data reduction in terms of a particular statistic can be thought of as a partition of the sample space.

X : Sample space

$$T = \{t : t = T(x) \text{, for some } x \in X\}$$

Example: X : Sample space of $\bar{X} = (x_1, x_2, \dots, x_n)$



$T(x)$ partitions the sample space

into sets A_t , $t \in T$, defined

$$\text{by } A_t = \{x : T(x) = t\}$$

- * The statistic summarizes the data in that, rather than reporting the entire sample x , it reports only

$$T(x) = t \text{ or equivalently } x \in A_t.$$

\Rightarrow we study three principles of Data reduction.

\Rightarrow interested in data reduction, that do not discard important information about the unknown parameter θ and method's that successfully discard information that is irrelevant or for all gaining knowledge about θ is concerned.

6.2 The Sufficiency Principle:

- * A Sufficient statistic for a parameter θ is a statistic that in a certain sense, captures all the information about θ contained in the sample.
- * Any additional information in the sample, besides the value of the sufficient statistic

does not contain any more information about θ .

Sufficiency Principle:

If $T(x)$ is a sufficient statistic for θ , then any inference about θ should depend on the sample X only through $T(x)$

* If x, y are two sample's such that $T(x) = T(y)$, then inference about θ should be same whether $X=x$ or $X=y$ is observed

Definition 6.2.1 :-

A statistic $T(x)$ is a sufficient statistic for θ if the conditional distribution of the sample X given the value of $T(x)$ does not depend on θ

$$\text{IP}(\bar{X} = \bar{x} | T(\bar{X}) = t)$$

$$= \text{IP}(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n | T(x_1, x_2, \dots, x_n) = t)$$

does not depend on θ .

What happens after conditioning
on $T(X)$?

When we condition the sample X
on the value of $T(X)$:

$$\text{IP}(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n | T(\bar{X}) = t)$$

the Parameter θ is no longer relevant.

This is because:

- ① $T(\bar{X})$ already captured everything about θ .
- ② The remaining data X (after

(knowing $T(x)$) contains no extra information about θ .

\Rightarrow once we know $T(x)$, we have already extracted everything relevant to θ .

\Rightarrow The remaining variability in x (after knowing $T(x)$) is unrelated to θ . It only depends on the internal structure of x , like how the data can be rearranged while keeping $T(x)$ constant.

key takeaway:

The conditional distribution of x , given $T(x)$, does not depend on θ because $T(x)$ already contains all the information about θ .

$$P_0(X=x \mid T(x) = T(\omega))$$

$$= P_0(X_1=x_1, X_2=x_2, \dots, X_n=x_n \mid T(x_1, x_2, \dots, x_n) = T(\omega))$$

$$= \frac{P_0(X=x \text{ and } T(x) = T(\omega))}{P_0(T(x) = T(\omega))}$$

$$= \frac{P_0(x \mid \theta)}{q(T(\omega) \mid \theta)}$$

$$\Rightarrow \frac{P_0(X_1=x_1, X_2=x_2, \dots, X_n=x_n \mid \theta)}{P_0(T(x) = T(\omega) \mid \theta)}$$

* $P(x \mid \theta)$ is the joint pmf of all sample X

* $q(\cdot \mid \theta)$ is the pmf of $T(x)$

$\Rightarrow T(x)$ is a sufficient statistic for θ



$\forall x \in X \in \mathbb{R}^n$, the above ratio of
pmf's is constant as a function of
 θ .

Theorem 6.2.3:

If $P(X|\theta)$ is the joint pdf or
pmf of X and $q(f|\theta)$ is the pdf
or pmf of $T(x)$ is a sufficient
statistic for θ if, for every x in the
sample space, the ratio $\frac{P(X|\theta)}{q(f|\theta)}$
is constant as a function of θ .

Ex 6.2.3: (Binomial sufficient statistic)

$X_1, X_2, \dots, X_n \sim \text{Ber}(\theta)$

we need to show $T(x) = X_1 + X_2 + \dots + X_n$
is a sufficient statistic of θ .

$$\frac{P(x|\theta)}{q(T(x)|\theta)} = \frac{\prod \theta^{x_i} (\theta - \bar{x})^{1-x_i}}{\binom{n}{t} \theta^t (\theta - \bar{x})^{n-t}}$$

$$= \frac{\theta^{\sum x_i} (\theta - \bar{x})^{n-\sum x_i}}{\binom{n}{t} \theta^t (\theta - \bar{x})^{n-t}}$$

$$= \frac{\theta^t \cdot (\theta - \bar{x})^{n-t}}{\binom{n}{t} \theta^t (\theta - \bar{x})^{n-t}}$$

$$= \frac{1}{\binom{n}{\sum x_i}}$$

\Rightarrow does not depend on θ

$\Rightarrow T(x)$ is sufficient statistic of θ .

Ex 6.2.4

(Normal statistic)

$X_1, X_2, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$

σ^2 is known.

$$T(x) = \bar{x} = \frac{1}{n}(X_1 + X_2 + \dots + X_n)$$

is sufficient statistic of μ .

$$f(X_1, X_2, \dots, X_n | T(x) = t)$$

$$= f(X_1, X_2, \dots, X_n, T(x) = t | \mu)$$

$$\frac{f(T(x) = t | \mu)}{f(T(x) = t | \mu)}$$

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

$$\bar{x} = \frac{1}{n}(X_1 + X_2 + \dots + X_n) = T(x)$$

$$T(x) \sim N(\mu, \frac{\sigma^2}{n})$$

$$q(T(x)|\mu) = \frac{1}{(2\pi\sigma^2)^{n/2}} \cdot \exp\left(-\frac{1}{2\sigma^2} (\bar{x} - \mu)^2\right)$$

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

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

$$\frac{f(x|\mu)}{q(T(x)|\mu)} = \frac{\frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\frac{1}{2\sigma^2} (\sum x_i^2)\right) \exp\left(-\frac{1}{2\sigma^2} \left(\mu^2 - 2\mu\bar{x} \right)\right)}{\frac{1}{(2\pi\frac{\sigma^2}{n})^{n/2}} \exp\left(-\frac{1}{2\sigma^2} (\bar{x}^2 + \mu^2 - 2\mu\bar{x})\right)}$$

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

$$= \frac{1}{\sqrt{n}} \cdot \frac{1}{(2\pi\sigma^2)^{\frac{n-1}{2}}} \exp\left(-\frac{1}{2\sigma^2} \left(\sum_{i=1}^n (x_i^2 - \bar{x}^2)\right)\right)$$

this does not depend on θ

$\Rightarrow \bar{x}$ is sufficient statistic.

To use the definition of sufficient statistic, we must guess a statistic $T(x)$ to be sufficient, find the pdf or pmf of $T(x)$ and then check the ratio of pd's or pm's does not depend on θ .

\Rightarrow 1st step requires a good deal of intuition

\Rightarrow 2nd step requires a tedious analysis.

Theorem 6.2.6 (Factorization theorem)

Let $f(x|\theta)$ denote the joint pdf or pmf of a sample x .

A statistic $T(x)$ is a sufficient statistic for θ

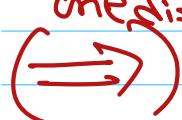


\exists functions $g(T|\theta)$ and $h(x)$

such that, \forall sample points x and all parameter points θ

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

Proof:



Proof only for discrete distributions.

SUPPOSE $T(x)$ is sufficient statistic.

Choose $g(t|\theta) = \Pr_{\theta}(T(x)=t)$

$$h(x) = \underbrace{\Pr(x=x | T(x)=T(x))}_{\text{because } T(x) \text{ is sufficient}}$$

$h(x)$ is independent of θ .

$$\begin{aligned}
 f(x=s) &= P_{\theta}(x=s) \\
 &= P_{\theta}(X=s \text{ and } T(X)=s) \\
 &= P_{\theta}(X=s | T(s)=T(s)) P(T(X)=T(s)) \\
 &= h(s) g(T(s) | \theta)
 \end{aligned}$$

if $T(X)$ is sufficient statistic
 \Rightarrow factorization exists

\iff other direction
 suppose the $f(x|\theta) = g(T(s)|\theta) h(x)$
 factorization exists.

Let $q(t|\theta)$ is pmf of $T(x)$
 To show that $T(x)$ is sufficient
 statistic we examine the ratio

$$\frac{f(x|\theta)}{q(T(x)|\theta)}$$

suppose $A_{T(x)} = \{y : T(y)=T(x)\}$
 all the elements in X
 s.t $T(y) = T(x)$

$$\frac{f(x|\theta)}{g(\tau(x)|\theta)} = \frac{g(\tau(x)|\theta) h(x)}{g(\tau(x)|\theta)}$$

$$= \frac{g(\tau(x)|\theta) h(x)}{\sum_{\tau(x)} f(y|\theta)}$$

$$= \frac{g(\tau(x)|\theta) h(x)}{\sum_{\tau(x)} g(\tau(y)|\theta) h(y)}$$

$$= \frac{g(\tau(x)|\theta) h(x)}{g(\tau(x)|\theta) \sum_{\tau(x)} h(y)}$$

$$= \frac{h(x)}{\sum_{\tau(x)} h(y)}$$

The ratio does not depend on

$\theta \Rightarrow$ factorization exists $\Rightarrow \tau(x)$ is sufficient

EX 6.2.7 (Continuation of Example 6.2.4)

for the normal distribution

$$X_1, X_2, \dots, X_n \text{ iid } N(\mu, \sigma^2)$$

$$\sigma^2 = \text{known}$$

$$\mu = \text{unknown}$$

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

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

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

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

$$= h(x) \circ g(\sum x_i | \mu)$$

$$\Rightarrow T(x) = \sum x_i \text{ is a sufficient statistic.}$$

Ex : 6.2.8 (uniform sufficient statistic)

X_1, X_2, \dots, X_n is Discrete uniform
(1, 2, ..., θ)

where θ is unknown parameter.

$$f(x|\theta) = \begin{cases} \frac{1}{\theta} & \text{if } x \in \{1, 2, 3, \dots, \theta\} \\ 0 & \text{otherwise} \end{cases}$$

$$f(x_1, x_2, \dots, x_n | \theta) = \begin{cases} \theta^{-n} & x_i \in \{1, 2, \dots, \theta\} \text{ for } i=1, 2, \dots, n \\ 0 & \text{otherwise} \end{cases}$$

$$= \begin{cases} \theta^{-n} & x_i \in \{1, 2, \dots, \theta\} \text{ for } i=1, 2, \dots, n \\ 0 & \max_i x_i > \theta \end{cases}$$

$$T(X) = \max_i x_i$$

$$h(x) = \begin{cases} 1 & x_i \in \{1, 2, \dots, \theta\} \text{ for } i=1, 2, \dots, n \\ 0 & \text{otherwise} \end{cases}$$

$$g(t|\theta) = \begin{cases} \frac{1}{\theta^n} & t \leq \theta \\ 0 & \text{o.w} \end{cases}$$

$$\Rightarrow f(x|\theta) = g(\tau(x)|\theta) h(x)$$

another method using Indicator functions.

Indicator function:

$$I_A(x) = \begin{cases} 1 & \text{if } x \in A \\ 0 & \text{o.w} \end{cases}$$

$$P_{I_A}(x=i) = P(A)$$

$$P_{I_A^c}(x=0) = P(\Omega \setminus A)$$

$$\text{Let } N = \{1, 2, \dots\}$$

$$N_0 = \{1, 2, \dots, 0\}$$

Joint Pdf of x_1, x_2, \dots, x_n

$$f(x|\theta) = \prod_{i=1}^n \frac{1}{\theta} I_{N_\theta}(x_i)$$

$$= \theta^{-n} \prod_{i=1}^n I_{N_\theta}(x_i)$$

define $T(x) = \max_i x_i$

$$\prod_{i=1}^n I_{N_\theta}(x_i) = \left[\prod_{i=1}^n I_N(x_i) \right] I_{N_\theta}(T(x))$$

$$f(x|\theta) = \theta^{-n} I_{N_\theta}(T(x)) \left(\prod_{i=1}^n I_N(x_i) \right)$$

The sufficient statistic can be vector also

$$T(x) = (T_1(x), T_2(x), \dots, T_n(x))$$

when parameter θ is also vector

$$\theta = (\theta_1, \theta_2, \dots, \theta_n)$$

Ex: 6.2.9 (Normal Sufficient statistic,
both Parameter's unknown)

$$X_1, X_2, X_3, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$$

Both μ, σ^2 are unknown

$$\Rightarrow \Theta = (\mu, \sigma^2)$$

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

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

$$\sum_{i=1}^n x_i^2 + n\mu^2 - 2\mu n\bar{x}$$

$$= (n-1)s^2 + n\bar{x}^2 + n\mu^2 - 2\mu n\bar{x}$$

$$= (n-1)s^2 + n(\bar{x}^2 + \mu^2 - 2\mu\bar{x})$$

$$= (n-1) s^2 + n(\bar{x} - \mu)^2$$

$$\Rightarrow f(x|\mu, \sigma^2)$$

$$= \frac{1}{(2\pi\sigma^2)^{n/2}} \cdot \exp\left(-\frac{1}{2\sigma^2} ((n-1)s^2 + n(\bar{x} - \mu)^2)\right)$$

$$= T_1(x) = \bar{x}, T_2(x) = s^2$$

$$\Rightarrow f(x|\mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{n/2}} \cdot \exp\left(-\frac{1}{2\sigma^2} ((n-1)t_2 + n(t_1 - \mu)^2)\right)$$

$$= g(T_1(x), T_2(x) | \mu, \sigma^2) \cdot h(x)$$

$$h(x) = 1$$

therefore (\bar{x}, s^2) is a sufficient statistic for (μ, σ^2)

Theorem 6.2.10: Let x_1, x_2, \dots, x_n iid $f(x|\theta)$

such that

$$f(x|\theta) = h(x) c(\theta) \exp\left(\sum_{i=1}^k \omega_i(\theta) t_i(x)\right)$$

where $\Theta = (\theta_1, \theta_2, \dots, \theta_k)$ $1 \leq i \leq k$ Then

$$T(x) = \left(\sum_{j=1}^n t_1(x_j), \dots, \sum_{j=1}^n t_k(x_j) \right)$$

is a sufficient statistic for Θ