# Slide 4: More About Sufficient Statistics

## STATS 511: Statistical Inference

Kean Ming Tan

# Curved Exponential Family

Suppose that $X_1, \ldots, X_n \sim N(\mu, \mu^2)$. We have shown that the distribution belongs to an exponential family with

$$T(\mathbf{x}) = \left( \sum_{i=1}^{n} X_i^2, \sum_{i=1}^{n} X_i \right) \qquad \text{and} \qquad w(\boldsymbol{\theta}) = \left( -\frac{1}{2\mu^2}, \frac{1}{\mu} \right).$$

Here, $w(\boldsymbol{\theta})$ forms a curve in the 2-dimensional space. That is, as $\mu$ varies, we get a curve in the "$xy$" plane.

# An Example on Non Exponential Family

Suppose that $X_1, \ldots, X_n$ is a random sample from a "right-half" normal distribution with $\sigma^2 = 1$. Find a sufficient statistic for $\mu$.

# Minimal Sufficient Statistic

As in previous example, there are many choices of sufficient statistics for the parameter of interest, and of course the "smaller" ones are more useful for data reduction. This motivates the concept on minimal sufficient statistic.

**Minimal Sufficient Statistic:** $T(\mathbf{X})$ is a minimal sufficient statistic for $\theta$ if $T(\mathbf{X})$ is sufficient, and is a function of any other sufficient statistic.

# How to Find Minimal Sufficient Statistic (Theorem 6.2.13)

**Theorem:** Let $f(\mathbf{x} \mid \boldsymbol{\theta})$ be the pdf and pmf of a sample $\mathbf{X}$. Suppose there exists a function $T(\mathbf{x})$ such that, for every two sample points $\mathbf{x}$ and $\mathbf{y}$, the ratio $f(\mathbf{x} \mid \boldsymbol{\theta})/f(\mathbf{y} \mid \boldsymbol{\theta})$ is constant as a function of $\boldsymbol{\theta}$ if and only if $T(\mathbf{x}) = T(\mathbf{y})$. Then $T(\mathbf{X})$ is a minimal sufficient statistic for $\boldsymbol{\theta}$.

$$f(x|\theta) = h(x)\, C(\theta) \exp\left\{ \sum_{j=1}^{K} w_j(\theta)\, T_j(x) \right\}$$

let $K = 3$

Proof idea: if $T(x)$ is not linearly independent

then $T_1(x) = a_2 T_2(x) + a_3 T_3(x)$

Ratio: $\dfrac{h(x)}{h(y)} \exp\left\{ w_1(\theta)\big(T_1(x) - T_1(y)\big) \right.$

$+ w_2(\theta)\big(T_2(x) - T_2(y)\big)$

$\left. + w_3(\theta)\big(T_3(x) - T_3(y)\big) \right\}$

$\Rightarrow \dfrac{h(x)}{h(y)} \exp\left\{ w_1(\theta)\Big( a_2\big(T_2(x) - T_2(y)\big) \right.$

$+ a_3\big(T_2(x) - T_2(y)\big)\Big)$

$+ w_2(\theta)\big(T_2(x) - T_2(y)\big)$

$\left. + w_3(\theta)\big(T_3(x) - T_3(y)\big) \right\}$

$= \dfrac{h(x)}{h(y)} \exp\left\{ \big(w_2(\theta) + a_2 w_1(\theta)\big)\big(T_2(x) - T_2(y)\big) \right.$

$\left. + \big(w_3(\theta) + a_3 w_1(\theta)\big)\big(T_3(x) - T_3(y)\big) \right\}$

so, if $w(\theta)$ is not linearly independent,
it does not imply $T(x) = T(y)$

# An Example: Multinomial Distribution

We have three boxes labelled Box 1, Box 2, and Box 3. We toss $n$ balls into Box 1, Box 2, or Box 3 ($n$ is given). Suppose that Box 1, Box 2, and Box 3 each has probability $p_1$, $p_2$, and $p_3$ of a ball landing in their respective box. Let $X_1$, $X_2$, and $X_3$ be the number of balls that land in Box 1, Box 2, and Box 3 respectively. Then,

$$(X_1, X_2, X_3) \sim \textit{Multinomial}(n, p_1, p_2, p_3).$$

*known* ⇓

*unknown* (under $p_1, p_2, p_3$)

*Ignore*

**Claim:** $(X_1, X_2)$ is the minimal sufficient statistic for $\theta = (p_1, p_2, p_3)$.

$$n = known$$
$$P_1, P_2, P_3 = unknown$$

$$P_1 + P_2 + P_3 = 1$$
$$X_1 + X_2 + X_3 = n$$

$$f(x|\theta) = \frac{n!}{x_1! x_2! x_3!} \, P_1^{x_1} \, P_2^{x_2} \, P_3^{x_3}$$

$$= \frac{n!}{x_1! x_2! x_3!} \exp\left\{ x_1 \log P_1 + x_2 \log P_2 + x_3 \log P_3 \right\}$$

$$= \frac{n!}{x_1! x_2! x_3!} \exp\left\{ x_1 \log P_1 + x_2 \log P_2 + (n - x_1 - x_2) \log(1 - P_1 - P_2) \right\}$$

$T(x) = (x_1, x_2, x_3)$ is not minimal

sufficient statistic because those are

linearly dependent.

Proof: In these situation's we need to

re-write our distribution function's.

# Proof

$$f(x|\theta) = \frac{n!}{x_1! \, x_2! \, (n-x_1-x_2)!} \exp\left\{ x_1 \log p_1 + x_2 \log p_2 + (n-x_1-x_2) \log p_3 \right\}$$

$$= \frac{n!}{x_1! \, x_2! \, (n-x_1-x_2)!} \exp\left\{ x_1 \log \frac{p_1}{p_3} + x_2 \log \frac{p_2}{p_3} + n \log p_3 \right\}$$

$$= \frac{n!}{x_1! \, x_2! \, (n-x_1-x_2)!} \, p_3^{\,n} \, \exp\left\{ x_1 \log \frac{p_1}{p_3} + x_2 \log \frac{p_2}{p_3} \right\}$$

$$T(x) = (x_1, x_2)$$

$$W(\theta) = \left( \log \frac{p_1}{p_3}, \log \frac{p_1}{p_3} \right)$$

$$\Rightarrow \quad T_1(x) = \sum_{i=1}^{n} x_1, \quad T_2(x) = \sum_{i=1}^{n} x_2$$

is Minimal sufficient
Sofistic of $(p_1, p_2, p_3)$

# Ancillary Statistics (Exact opposite of sufficient statistic)

Sufficient statistics contain all the information about $\boldsymbol{\theta}$ that is available in the sample. Now, we consider a statistic $S(\mathbf{X})$ that has no information about $\boldsymbol{\theta}$.

**Definition:** A statistic $S(\mathbf{X})$ whose distribution does not depend on the parameter $\boldsymbol{\theta}$ is called an ancillary statistic.

**Examples:** see 6.2.17–6.2.19 in your textbook

$X_1, X_2, \ldots X_n \sim \text{unif}(\theta, \theta+1)$

$R = X_{(n)} - X_{(1)}$ be the range

$f_R(r) = n(n-1) r^{n-2} (1-r)$ does not depend on $\theta$

# Complete Statistics

**Definition:** Let $f(t \mid \theta)$ be a family of pdfs or pmfs for a statistics $T(\mathbf{X})$. The family of probability distributions is called complete if $E_\theta\{g(T)\} = 0$ for all $\theta$ implies $P_\theta(g(T) = 0) = 1$ for all $\theta$. $T(\mathbf{X})$ is also called a complete statistic.

**Complete Sufficient Statistic:** If $T(\mathbf{X})$ is a complete statistic and a sufficient statistic, then we call $T(\mathbf{X})$ the complete sufficient statistic.

**Theorem 6.2.25:** Let $X_1, \ldots, X_n$ be iid observations from an exponential family. Then the statistic
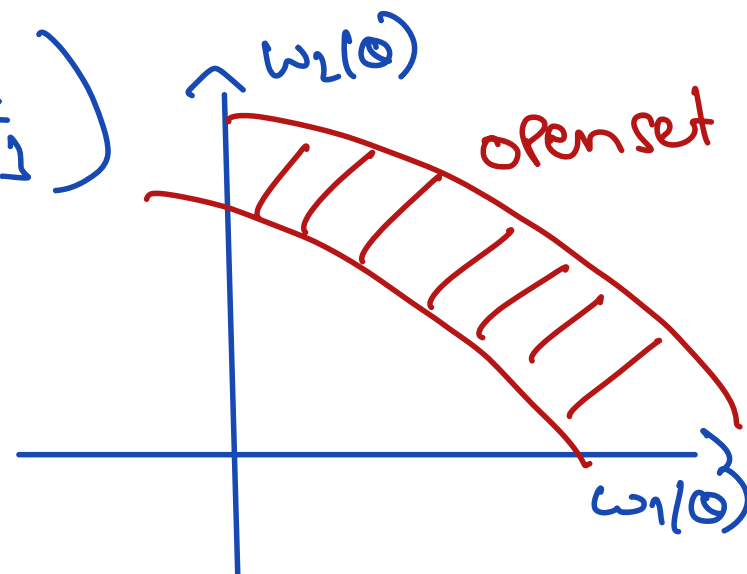
$$T(\mathbf{X}) = \left( \sum_{i=1}^{n} t_1(X_i), \ldots, \sum_{i=1}^{n} t_1(X_k) \right)$$

is complete as long as the parameter space $\Theta$ contains an open set in $\mathbb{R}^k$.

# Multinomial and Curved Exponential Family Example

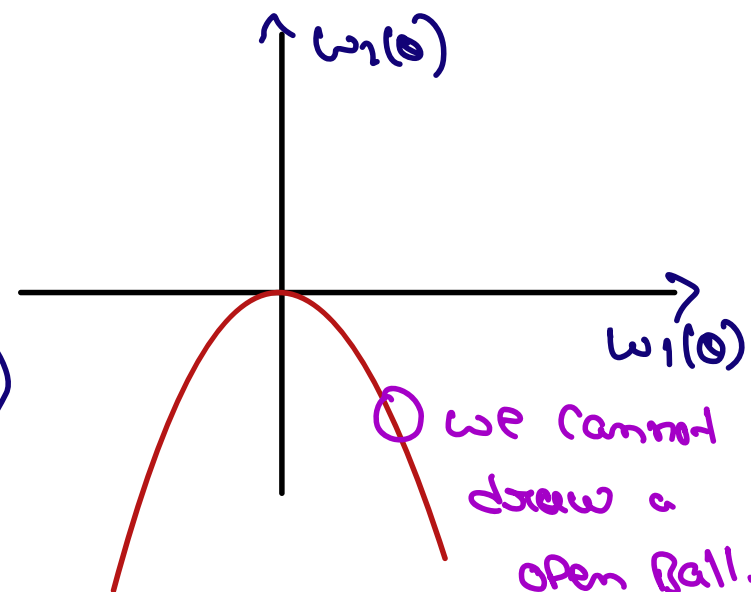① $\omega_1(\theta) = \left( \log \frac{P_1}{P_3}, \log \frac{P_2}{P_3} \right)$

$T(x) = (x_1, x_2)$

open set

$\omega_2(\theta)$

$\omega_1(\theta)$

② $X_1, X_2, \ldots, X_n \sim N(\mu, \sigma^2)$

$W(\theta) = \left( -\frac{1}{2\mu^2}, \frac{1}{\mu} \right)$

$T(x) = \left( \sum x_i^2, \sum x_i \right)$

$\omega_2(\theta)$

$\omega_1(\theta)$

we cannot draw a open Ball.

$$T_1(x) = \sum x_i^2 \qquad T_2(x) = \sum x_i$$

$$\mathbb{E}\left[T_1(x)\right] = \mathbb{E}\left[\sum x_i^2\right]$$

$$g(T) = (n+1)T_1 - 2T_2^2$$

find a $g(t)$ s.t $\quad E_u\left[g(T)\right] = 0$ but

$\quad g(t) \neq 0$

$X \sim N(u, v^2) \qquad E_u\left[T_1\right] = 2nv^2$

$$E_u\left[T_2^2\right] = (n^2+n)v^2$$

$$E_u\left[g(t)\right] = (n+1) E\left[T_1\right] - 2E\left[T_2^2\right]$$

$$= 0$$

But $g(t) \neq 0$ unless all $x_i$s are zero

$\Rightarrow \quad T(x)$ is not a C.S.S for $u$

Quiz type Problem:

$(X_1, X_2, X_3) \sim$ Multinomial $(n, P_1, P_2, P_1 + P_2)$

# Non-Exponential Family Example

Let $X_1, \ldots, X_n \sim Unif(0, \theta)$ be iid Uniform random variables. We know that $T(\mathbf{X}) = X_{(n)}$ is sufficient statistic for $\theta$. Is $T(\mathbf{X})$ complete sufficient statistic?

## location- scale family

① Let $f_0(x)$ be some pdf

$$f_\theta(x) = f_0(x - \theta)$$

↑ location parameter

EX: $X \sim N(\mu, 1)$

$$f_0(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$$

$$f_\mu(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x-\mu)^2}$$

② 

## Scale family

$$f_\sigma(x) = \frac{1}{\sigma} f\left(\frac{x}{\sigma}\right)$$

EX: $X \sim EXP(\theta)$    $f_\theta(x) = \frac{1}{\theta} e^{-\frac{x}{\theta}}$    $0 < x < \infty$

① if $X$ is location family, difference of two $X$'s will remove the location parameter

EX: $X_1, X_2, \ldots, X_n \sim N(\mu, 1)$

$X_1 - X_2$ is ancillary

$R := X_{(n)} - X_{(1)}$ is ancillary

② Scale family

    The ratio of any two $x$ remove the scale parameter

Ex:      $X_1, X_2 \sim N(0, \sigma^2)$      Scale family

     $S(x) = \dfrac{X_1}{X_2}$    is    this   ancillary?

let    $Y = \dfrac{x}{\sigma}$    the      $Y \sim N(0,1)$

So    $S(y) = \dfrac{\sigma y_1}{\sigma y_2} = \dfrac{y_1}{y_2}$

another example: (Both location & scale family)

     $X_1, X_2, \ldots, X_n \sim N(u, \sigma^2)$

ancillary statistic          $\dfrac{X_1 - u}{\sigma} \sim N(0,1)$

      $\dfrac{X_n - X_1}{X_2 - X_1}$  ,    $\dfrac{X_{(n)} - X_{(1)}}{X_1 - \text{median}(x)}$

# Applications of Ancillary and Complete Sufficient Statistic

**Basu's Theorem:** If $T(\mathbf{X})$ is a complete and minimal sufficient statistic, then $T(\mathbf{X})$ is independent of every ancillary statistics $S(\mathbf{X})$.

**Example:** $X_1, \ldots, X_n \sim Uniform(0, \theta]$. Show that $T(\mathbf{X}) = X_{(n)}$ and $S(\mathbf{X}) = (X_{(n)} - X_{(1)})/(X_{[2n/3]} - X_{[n/3]})$.

Comple sufficiet C.S.S Statistic

does not depend on $\theta$

$Y = \frac{X}{\theta} \sim uniform\ (0,1)$

$$S(y) = \frac{\frac{Y_{(n)}}{\theta} - \frac{Y_{(1)}}{\theta}}{\frac{Y_{[2n/3]}}{\theta} - \frac{Y_{[\frac{n}{3}]}}{\theta}} = \frac{Y_{(n)} - Y_{(1)}}{Y_{(\frac{2n}{3})} - Y_{(\frac{n}{3})}}$$

(ancillary)

$T(x)$ is C.S.S
$S(x)$ is ancillary

$\Longrightarrow$

BASU'S theorem

$T(x), S(x)$ are independent

## Quiz Type Problem:

$(X_1, X_2, X_3) \sim$ Multinomial $(n, P_1, P_2, P_1+P_2)$

(a) what's the orange of $P_1$

(b) Find a M.S.S for $(P_1, P_2)$

(c) Is $T$ in Complete

---

① $P_1 + P_2 + P_1 + P_2 = 1 \implies P_1 + P_2 = \frac{1}{2}$

$$\implies P_2 = \frac{1}{2} - P_1 \qquad 0 < P_1 < \frac{1}{2}$$

$$f(\vec{x}) = h(\vec{x}) \, P_1^{x_1} \left(\frac{1}{2} - P_1\right)^{x_2} \left(\frac{1}{2}\right)^{n - x_1 - x_2}$$

$$= \tilde{h}(x) \, P_1^{x_1} \left(\frac{1}{2} - P_1\right)^{x_2}$$

$$= \tilde{h}(\alpha) \exp\left\{ x_1 \log P_1 + x_2 \log\left(\frac{1}{2} - P_2\right)\right\}$$

$T(x) = (x_1, x_2) \qquad W(\theta) = \left( \underset{\theta_1}{\log P_1} , \underset{\theta_2}{\log\left(\frac{1}{2} - P_1\right)} \right)$

Both $\log P_1$ and $\log\left(\frac{1}{2} - P_1\right)$ in not linearly independent. (NONlinear dependent)

$\boxed{\text{Bot we only check for linear dependence}}$

$\implies W(\theta) \in \mathbb{R}^2$

can't apply our exponential theorem because we only have one free parameter, (we cannot have a open subset $W(\theta) \implies T(x)$ is not C.S.S

TO Prove T is not complete, we WTS

$\exists$ a $g(\cdot)$ s.t $\mathbb{E}[g(T)]$ but $g(t) \neq 0$

$T(x) = (x_1, x_2)$     $(x_1, x_2, x_3) \sim$ multinomial $\left(n, P_1, \frac{1}{2} P_1, \frac{1}{2}\right)$

$\left. \begin{array}{l} E[x_1] = nP_1 \\ E[x_2] = \frac{n}{2} - nP_1 \end{array} \right\} \Rightarrow$ Let

$$g(t) = x_1 + x_2 - \frac{n}{2}$$

$$\Rightarrow E[g(t)] = 0$$

But $g(t) \neq 0$   $\forall t \in T$

$\Rightarrow$    Therefore the $T(x)$ is not complete.

## Another Example:

$$x_1, x_2, \ldots, x_n \sim N(u, \sigma^2)$$

                             $\uparrow$
                          known

Show that

$$\bar{x} \perp S_n^2 \qquad S_n^2 = \frac{1}{n-1} \sum (x_1 - \bar{x})^2$$

We show that $T(x) = \bar{x}$ is C.S.S for

$\underline{u}$. By Basu's theorem ; $\bar{x} \perp\!\!\!\perp$ of any ancillary statistics.

Check if $S_n^2$ is ancillary

$\frac{1}{n} \sum (x_i - \bar{x})^2$ does not depend on

$u \implies S_n^2$ is ancillary

$\implies \underset{(C.S.S)}{\bar{x}} \perp\!\!\!\perp \underset{(ancillary)}{S_n^2}$

Because of Basu's theorem.

## A non-homogenous Poission Process

$\{ x(t) \mid 0 < t < \infty \}$  $x(0) = 0$  $x(t) = \lambda t$
$\lambda > 0$

$x(t_{i+1}) - x(t_i) \sim Poission \left( \int_{t_i}^{t_{i+1}} \lambda t \, dt \right)$

$x(t_1) \perp\!\!\!\perp [x(t_2) - x(t_1)] \perp\!\!\!\perp \cdots \perp\!\!\!\perp x(t_n) - x(t_{n-1})$

Let say we observe $x(1), x(2), \ldots, x(n)$
Find a $C.S.S$ for $\lambda$

$$Y \sim \text{Poisson} \qquad f(y) = \frac{e^{-\lambda} \lambda^y}{y!}$$

Define
$$Y_1 = X(1)$$
$$Y_2 = X(2)$$
$$\vdots$$
$$Y_n = X(n) - X(n-1)$$

By *  $\qquad Y_1, Y_2, \ldots, Y_n$ are independent

By **  $\qquad Y_i \sim \text{Poisson} \left( \int_i^{i+1} \lambda t \, dt \right)$
$$\sim \text{Poisson} \left( \lambda(i - \tfrac{1}{2}) \right)$$

$$f(Y_1, Y_2, \ldots Y_n \mid \lambda) = \prod_{i=1}^{n} \left[ \frac{\left( e^{-\lambda(i-\frac{1}{2})} \right) \left[ \lambda(i-\frac{1}{2}) \right]^{y_i}}{y_i!} \right]$$

$$= \frac{e^{\frac{n}{2}}}{y_1! \, y_2! \ldots y_n!} \cdot e^{-\lambda \frac{i(i-1)}{2}} \cdot \prod_{i=1}^{n} \left( \lambda(i-\tfrac{1}{2}) \right)^{y_i}$$

$$= h(y) \cdot c(\lambda) \prod_{i=1}^{n} \lambda^{y_i} (i-\tfrac{1}{2})^{y_i}$$

$$= h(y) \cdot c(\lambda) \prod_{i=1}^{n} \exp \left\{ y_i \log \lambda + y_i \log(i-\tfrac{1}{2}) \right\}$$

$$= \hat{h}(y) \cdot C(\lambda) \; \exp\left\{ \log\lambda \sum y_i \right\}$$

$$\Rightarrow \quad W(\lambda) = \log\lambda \qquad T(Y) = \sum y_i$$

$$\|$$

we can draw openset

$$\Rightarrow T(Y) \text{ is } C.S.S$$

# Up Next – Methods for Estimation