### STAT 517:Sufficiency

#### Minimal sufficiency and Ancillary Statistics. Sufficiency, Completeness, and Independence

Prof. Michael Levine

March 1, 2016

Levine STAT 517:Sufficiency

- Not all sufficient statistics created equal...
- So far, we have mostly had a single sufficient statistic for one parameter or two for two parameters (with some exceptions)
- Is it possible to find the minimal sufficient statistics when further reduction in their number is impossible?
- Commonly, for k parameters one can get k minimal sufficient statistics

• 
$$X_1, \ldots, X_n \sim Unif(\theta - 1, \theta + 1)$$
 so that

$$f(x;\theta) = \frac{1}{2}I_{(\theta-1,\theta+1)}(x)$$

where  $-\infty < \theta < \infty$ 

The joint pdf is

$$2^{-n}\left\{I_{(\theta-1,\theta+1)}(\min x_i)\right\}\left\{I_{(\theta-1,\theta+1)}(\max x_i)\right\}$$

It is intuitively clear that Y<sub>1</sub> = min x<sub>i</sub> and Y<sub>2</sub> = max x<sub>i</sub> are joint minimal sufficient statistics

→ 御 → → 注 → → 注 →

æ

# Occasional relationship between MLE's and minimal sufficient statistics

- Earlier, we noted that the MLE  $\hat{\theta}$  is a function of one or more sufficient statistics, when the latter exists
- If θ̂ is itself a sufficient statistic, then it is a function of others...and so it may be a sufficient statistic
- ► E.g. the MLE  $\hat{\theta} = \bar{X}$  of  $\theta$  in  $N(\theta, \sigma^2)$ ,  $\sigma^2$  is known, is a minimal sufficient statistic for  $\theta$
- The MLE  $\hat{\theta}$  of  $\theta$  in a  $P(\theta)$  is a minimal sufficient statistic for  $\theta$
- The MLE θ̂ = Y<sub>(n)</sub> = max<sub>1≤i≤n</sub> X<sub>i</sub> of θ in a Unif(0, θ) is a minimal sufficient statistic for θ
- ▶  $\hat{\theta}_1 = \bar{X}$  and  $\hat{\theta}_2 = \frac{n-1}{n}S^2$  of  $\theta_1$  and  $\theta_2$  in  $N(\theta_1, \theta_2)$  are joint minimal sufficient statistics for  $\theta_1$  and  $\theta_2$

□ > < E > < E > < E</p>

► A sufficient statistic T(X<sub>1</sub>,...,X<sub>n</sub>) is called a minimal sufficient statistic if it is a function of any other sufficient statistic

- 4 回 2 - 4 □ 2 - 4 □

æ

# When MLE and minimal sufficient statistics have nothing in common with each other: Example I

- Take again  $X_1, \ldots, X_n \sim Unif(\theta 1, \theta + 1)$
- Clearly,  $\theta 1 < Y_1 < Y_n < \theta + 1$ , or

$$Y_n - 1 < \theta < Y_1 + 1$$

► To achieve the maximum possible value of the likelihood function (<sup>1</sup>/<sub>2</sub>)<sup>n</sup>, choose any θ between Y<sub>n</sub> − 1 and Y<sub>1</sub> + 1; a common choice as MLE is the average of two endpoints

$$\hat{\theta} = \frac{Y_1 + Y_n}{2}$$

▶ Note that the resulting  $\hat{\theta}$  is not even a sufficient statistic...and, therefore, cannot be a minimal sufficient statistic

## A more general location family setting I

- The above example is a location family X<sub>i</sub> = θ + W<sub>i</sub> where W<sub>i</sub> ∼ Unif(-1, 1)
- ► Take a general location family with W<sub>i</sub> having a pdf f(w) and cdf F(w)
- ► We know that the order statistics Y<sub>1</sub> < Y<sub>2</sub> < ··· < Y<sub>n</sub> form a set of sufficient statistics in this case...Can we do better?
- If f(w) is a N(0,1) pdf, X̄ is both the MVUE and MLE of θ; moreover, X̄ is a minimal sufficient statistic
- ► Take f(w) = e<sup>-w</sup> for w > 0 and zero elsewhere; here, Y<sub>1</sub> is a sufficient statistic and the MLE so Y<sub>1</sub> is a minimal sufficient statistic

▲圖 → ▲ 国 → ▲ 国 →

- On the contrary, for the logistic location family, the MLE of θ exists and is easy to compute...nevertheless, the order statistics are *minimal sufficient* in this case
- If f(w) is a Laplace pdf with the location parameter θ, the median Q<sub>2</sub> is an MLE; however, yet again, the order statistics are *minimal sufficient* in this case
- This latter situation is, in general, more common for location models

- In general, if the minimal sufficient statistic exists (and it almost always does), any complete sufficient statistic is also a minimal sufficient statistic
- The converse is not true, however; from the uniform example, note that

$$\mathbb{E}\left[\frac{Y_n-Y_1}{2}-\frac{n-1}{n+1}\right]=0$$

for all  $\theta$ 

### Ancillary statistics

- A quick example for X<sub>1</sub>,..., X<sub>n</sub> ~ N(θ, 1) the distribution of S<sup>2</sup> does not depend on θ
- Alternatively, take  $X_1, X_2 \sim \Gamma(\alpha, \theta)$  where  $\alpha > 0$  is known and recall that  $Z = \frac{X_1}{X_1 + X_2}$  has a beta distribution that does not depend on  $\theta$ ...Thus, Z is an **ancillary** statistic for this sample size 2 w.r.t  $\theta$
- In general, select a location family X<sub>i</sub> = θ + W<sub>i</sub>, i = 1,..., n, where −∞ < θ < ∞ is a parameter and W<sub>1</sub>,..., W<sub>n</sub> ~ f(w) that doesn't depend on θ
- The common pdf of X<sub>i</sub> is f(x − θ)...any location-invariant statistic Z = u(X<sub>1</sub>,...,X<sub>n</sub>) s.t. Z = u(W<sub>1</sub> + θ,...,W<sub>n</sub> + θ) = u(W<sub>1</sub>,...,W<sub>n</sub>) for all θ is an ancillary statistic
- Sample variance is one such statistic...Sample range R = max X<sub>i</sub> - min X<sub>i</sub> is another...Finally, the absolute mean deviation from the sample median

- Let X<sub>1</sub>,..., X<sub>n</sub> ~ f(x; θ) where θ ∈ Ω and Ω is an interval.
   Let Y<sub>1</sub> be a complete and sufficient statistic for θ
- $Z = u(X_1, \ldots, X_n)$  is another statistic
- $\blacktriangleright$  Distribution of  $Y_1$  doesn't depend on  $\theta \to Z$  is independent of  $Y_1$

- ▶ If  $Y_1$  is a sufficient statistic and is independent of  $Z \rightarrow$  distribution of Z doesn't depend on  $\theta$
- If Y<sub>1</sub> is a sufficient and a *complete* statistic, distribution of Z doesn't depend on θ → Y<sub>1</sub> and Z are independent
- The second case is very easily satisfied for regular exponential families

- $X_1, \ldots, X_n \sim N(\mu, \sigma^2)$
- X
   is a sufficient statistic for μ while S<sup>2</sup> is ancillary for μ (location-invariant)
- Thus,  $\bar{X}$  and  $S^2$  are independent

→ 御 → → 注 → → 注 →

- $X_1, \ldots, X_n \sim e^{-(x-\theta)}$  for  $\theta < x < \infty$
- $Y_1 = \min_{1 \le i \le n} X_i$  is a complete sufficient statistic for  $\theta$
- ► Any location invariant statistic, e.g. S<sup>2</sup> or the sample range are independent of Y<sub>1</sub>

★御★ ★注★ ★注★

- $X_1, X_2 \sim f(x; \theta) = \frac{1}{\theta} e^{-x/\theta}$
- $Y_1 = X_1 + X_2$  is a complete sufficient statistic for  $\theta$
- ► Y<sub>1</sub> is independent of the scale invariant statistic X<sub>1</sub>/X<sub>1</sub>+X<sub>2</sub> that is beta distributed

▲□ ▶ ▲ □ ▶ ▲ □ ▶

- $X_1,\ldots,X_n \sim N(\theta_1,\theta_2)$
- $\bar{X}$  and  $S^2$  are joint sufficient statistics for  $heta_1$  and  $heta_2$
- Thus, the location and scale invariant statistic

$$Z = \frac{\sum_{i=1}^{n} (X_{i+1} - X_i)^2}{\sum_{i=1}^{n} (X_i - \bar{X})^2}$$

is independent of both  $\bar{X}$  and  $S^2$ 

同 ト イヨ ト イヨト

### When a sufficient statistic is not complete

- An incomplete sufficient statistic may still contain some information about the parameter
- $X_1,\ldots,X_n \sim \frac{1}{2}I_{(\theta-1,\theta+1)}(x)$
- Let  $Y_1 = \min X_i$  and  $Y_n = \max X_i$ ;  $T_1 = \frac{Y_1 + Y_n}{2}$  is an MLE of  $\theta$
- $T_2 = Y_n Y_1$  is an ancillary statistic...however

$$Var(T_1|t_2) = rac{(2-t_2)^2}{12}$$

and so the distribution of  $T_2$  has some information about  $\theta$ 

白 と く ヨ と く ヨ と …