

# Introduction to the Theory of Statistics Part 2

## PM522b

Meredith Franklin

Division of Biostatistics, University of Southern California

Slides 2, 2020

# Topics covered

## Data Reduction (CB Ch. 6)

1. Statistics
2. The Sufficiency Principle
3. The Likelihood Principle
4. The Equivariance Principle

# A Statistic

- ▶ The random samples we generated previously are vectors of observations that can be interpreted in statistically meaningful ways
- ▶ We want to use the information contained in our random sample to arrive at conclusions regarding our population
- ▶ A statistic:
  - is a form of data reduction
  - can be thought of as a partition of the sample space
  - is a summary quantity of our random sample
  - is a function of the sample

# A Statistic is a Form of Data Reduction

- ▶ Data reduction means that we use a statistic  $T(\mathbf{x})$  instead of the entire sample  $\mathbf{x} = (x_1, \dots, x_n)$  to make inferences about an unknown parameter  $\theta$ .

## Partitioning the sample space

The sample space  $\mathcal{X}$  can be partitioned and subsequently the observations  $\mathbf{x}$  can be reduced.

Let  $\mathcal{T} = \{t : t = T(\mathbf{x}) \text{ for some } \mathbf{x} \in \mathcal{X}\}$  be the image of  $\mathcal{X}$  under  $T(\mathbf{x})$ .

Then, the statistic  $T(\mathbf{x})$  partitions the sample space  $\mathcal{X}$  into sets  $A_t, t \in \mathcal{T}$  where  $A_t = \{\mathbf{x} : T(\mathbf{x}) = t\}$  for  $t \in \mathcal{T}$

- ▶ So, rather than reporting the whole sample  $\mathbf{x}$  we use  $T(\mathbf{x}) = t$
- ▶ Reporting  $T(\mathbf{x}) = t$  is equivalent to reporting  $\mathbf{x} \in A_t$

# Partitioning the Sample Space

As an example of partitioning the sample space, toss a coin three times and let  $X_1, X_2, X_3$  be the outcome of each toss. So,  $X_1, X_2, X_3 \sim \text{Bernoulli}(\theta)$ . Let  $T = \sum_{i=1}^3 X_i$ . What is the partition of our sample space?

# A statistic is a function of the sample

- ▶ A statistic is formally defined as a function of the observable random variables in a sample and known constants.
- ▶ Functions of observed samples (i.e. data) are used to generate statistics.

## Definition

For an iid sequence of random variables  $X_1, X_2, \dots, X_n$  sampled from our population with distribution function  $f(\mathbf{X}|\theta)$ , the function  $T(\mathbf{X}) = T(X_1, X_2, \dots, X_n)$  which does not contain the unknown parameter  $\theta$  is called a *statistic*.

# A statistic is a function of the sample

- ▶ The statistic  $T(\mathbf{X})$ , a function of random variables, is itself a random variable.
- ▶ When  $T(\mathbf{X})$  is used for inference, two different random samples  $\mathbf{x}$  and  $\mathbf{y}$  that satisfy  $T(\mathbf{x}) = T(\mathbf{y})$  lead to the same inference.
- ▶ The most frequently used statistics are measures of central tendency and measures of concentration or variation of the random sample.

## Simple Examples

Where  $T(X) = T(X_1, X_2, \dots, X_n)$

$$T(X) = \bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \text{ (sample mean)}$$

$$T(X) = S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 \text{ (sample variance)}$$

$$T(X) = M_n \text{ (sample median)}$$

# Statistic

- ▶ A statistic can be basically anything, but the choice of what we use as a statistic depends on the problem at hand.
- ▶ Some important things to note:
  - $T$  need not be a continuous function, but it does need to be measurable, i.e. the mapping  $T : \mathcal{X} \rightarrow \mathcal{T}$  is measurable.
  - By saying it cannot depend on parameter  $\theta$ , that means that the parameter  $\theta$  cannot appear in the formula for  $T$ . However, it is ok if the distribution of  $T$  depends on  $\theta$ .
- ▶ Often we are interested in the distribution of  $T$ . Example from Slides 1:  
 If  $X_1, \dots, X_n \sim N(\mu, \sigma^2)$  then  $\bar{X} \sim N(\mu, \sigma^2/n)$ .  
 Another example, If  $X_1, \dots, X_n \sim \Gamma(\alpha, \beta)$  then  $\bar{X} \sim \Gamma(n\alpha, \beta/n)$ . We derive this result from the mgf:

$$\begin{aligned}
 M_{\bar{X}} &= E[e^{t\bar{X}}] = E[e^{\sum X_i t/n}] = \prod_i E[e^{X_i t/n}] \\
 &= [M_X(t/n)]^n = \left[\left(\frac{1}{1 - \beta t/n}\right)^\alpha\right]^n = \left[\frac{1}{1 - \beta/nt}\right]^{n\alpha}
 \end{aligned}$$

Which is the mgf of  $\Gamma(n\alpha, \beta/n)$



# Principles of Data Reduction

- ▶ We need to evaluate how good a statistic really is, and to do this we rely on the three principles of data reduction:
  - Sufficiency
  - Likelihood
  - Equivariance

# Sufficiency

We assess our statistic for certain properties:

- ▶ Does the statistic retain all of the information about the true population parameters?
- ▶ Has some information about our parameters been lost or obscured through the process of reducing our data?
- ▶ A sufficient statistic for  $\theta$  is one that captures all of the information about  $\theta$  contained in our sample.
- ▶ This leads to the sufficiency principle:

## Sufficiency Principle, CB 6.2

If  $T(\mathbf{X})$  is a sufficient statistic for  $\theta$ , then inference about  $\theta$  should depend on the sample  $\mathbf{X}$  only through the value of the statistic  $T(\mathbf{X})$ . If  $\mathbf{x}$  and  $\mathbf{y}$  are two sample points such that  $T(\mathbf{x}) = T(\mathbf{y})$ , the inference about  $\theta$  should be the same whether  $\mathbf{X} = \mathbf{x}$  or  $\mathbf{Y} = \mathbf{y}$  is observed.

Basically, if we know the value of the sufficient statistic  $T$  we can do just as good of a job estimating  $\theta$  as someone who knows the entire sample.

# Sufficiency

R.A. Fisher published an article in 1922 in Philosophical Transactions of the Royal Society stating that a statistic is sufficient if "no other statistic that can be calculated from the same sample provides any additional information as to the value of the parameter".

# Sufficiency

## Definition: Sufficient statistics

For a random sample  $X_1, X_2, \dots, X_n$  with pdf  $f(\mathbf{x}|\theta)$ , the statistic  $T(\mathbf{X})$  is said to be sufficient if the conditional distribution of  $X_1, X_2, \dots, X_n$  given  $T(\mathbf{X})$  does not depend on  $\theta$ .

- ▶ A statistic  $T(\mathbf{X})$  is sufficient for  $\theta$  if inferences about  $\theta$  depend on  $\mathbf{X}$  only through  $T(\mathbf{X})$ . (Informal definition)
- ▶ A statistic  $T(\mathbf{X})$  is sufficient for  $\theta$  if the conditional distribution of  $\mathbf{X}$  given  $T(\mathbf{X})$  does not depend on  $\theta$ . (Formal definition)

# Sufficiency

## Example

To illustrate sufficiency, we devise a scenario where we have two 522b students A and B. Student A knows the entire random sample  $X_1, \dots, X_n = \mathbf{x}$  and can compute the statistic  $T(\mathbf{X}) = t(\mathbf{x})$ . This student can make inference about the parameter  $\theta$  using this information. On the other hand, student B only knows the value of the statistic  $T(\mathbf{X}) = t(\mathbf{x})$ . Since the conditional distribution of  $X_1, \dots, X_n$  given  $T(\mathbf{X})$  does not depend on  $\theta$ , student B knows  $P(\mathbf{X} = \mathbf{y} | T(\mathbf{X}) = t(\mathbf{x}))$ , which is a probability distribution on  $A_{T(\mathbf{x})} = \{\mathbf{y} : T(\mathbf{y}) = T(\mathbf{x})\}$  that can be calculated without knowledge of the true value of  $\theta$ . So, student B can use this distribution to generate a random sample  $\mathbf{y}$  satisfying  $P(\mathbf{Y} = \mathbf{y} | T(\mathbf{X}) = t(\mathbf{x})) = P(\mathbf{X} = \mathbf{y} | T(\mathbf{X}) = t(\mathbf{x}))$ . This means that for each  $\theta$ ,  $\mathbf{X}$  and  $\mathbf{Y}$  have the same unconditional pdf (shown on next slide). Student B knows just as much about  $\theta$  via  $T(\mathbf{X}) = t(\mathbf{x})$  as student A who knows the entire sample  $\mathbf{X} = \mathbf{x}$ .

# Sufficiency

## Example, con't

For this example to work,  $\mathbf{X}$  and  $\mathbf{Y}$  must have the same unconditional distribution, namely  $P_\theta(\mathbf{X} = \mathbf{x}) = P_\theta(\mathbf{Y} = \mathbf{x}) \forall \mathbf{x}$  and  $\theta$ .

$$\begin{aligned} P_\theta(\mathbf{X} = \mathbf{x}) &= P_\theta(\mathbf{X} = \mathbf{x} \text{ and } T(\mathbf{X}) = T(\mathbf{x})) \\ &= P(\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = T(\mathbf{x})) P_\theta(T(\mathbf{X}) = T(\mathbf{x})) \\ &= P(\mathbf{Y} = \mathbf{x} | T(\mathbf{X}) = T(\mathbf{x})) P_\theta(T(\mathbf{X}) = T(\mathbf{x})) \\ &= P_\theta(\mathbf{Y} = \mathbf{x} \text{ and } T(\mathbf{X}) = T(\mathbf{x})) \\ &= P_\theta(\mathbf{Y} = \mathbf{x}) \end{aligned}$$

# Sufficiency

To verify that a statistic  $T(\mathbf{X})$  is indeed sufficient for parameter  $\theta$ , we must verify that for any fixed values of  $\mathbf{x}$  and  $t$ ,  $P_\theta(\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = T(\mathbf{x}))$ .

## Theorem, CB 6.2.2

$T(\mathbf{X})$  is sufficient for  $\theta$  iff the ratio  $p(\mathbf{x}|\theta)/q(T(\mathbf{x}|\theta))$  is independent of  $\theta$  where  $p(\mathbf{x}|\theta)$  and  $q(T(\mathbf{x}|\theta))$  are the joint pmfs or pdfs of  $\mathbf{X}$  and  $T(\mathbf{X})$ , respectively.

$$\begin{aligned} P_\theta(\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = T(\mathbf{x})) &= \frac{P_\theta(\mathbf{X} = \mathbf{x} \text{ and } T(\mathbf{X}) = T(\mathbf{x}))}{P_\theta(T(\mathbf{X}) = T(\mathbf{x}))} \\ &= \frac{P_\theta(\mathbf{X} = \mathbf{x})}{P_\theta(T(\mathbf{X}) = T(\mathbf{x}))} \\ &= \frac{p(\mathbf{x}|\theta)}{q(T(\mathbf{x}|\theta))} \end{aligned}$$

# Sufficiency

A couple of examples:

## Sufficiency of sample mean for the normal distribution

Given  $X_1, \dots, X_n$  iid  $N(\mu, \sigma^2)$  with  $\sigma^2$  known, is the sample mean,  $\bar{X} = (X_1, \dots, X_n)/n$  a sufficient statistic for  $\mu$ ?

The joint pdf for the sample  $\mathbf{X}$  is

$$\begin{aligned} f_{\mathbf{X}}(\mathbf{x}|\mu) &= \prod_{i=1}^n (2\pi\sigma^2)^{-1/2} \exp(-(x_i - \mu)^2/(2\sigma^2)) \\ &= (2\pi\sigma^2)^{-n/2} \exp\left(-\sum_{i=1}^n (x_i - \mu)^2/(2\sigma^2)\right) \\ &= (2\pi\sigma^2)^{-n/2} \exp\left(-\sum_{i=1}^n (x_i - \bar{x} + \bar{x} - \mu)^2/(2\sigma^2)\right) \text{ (add and subtract } \bar{x}) \\ &= (2\pi\sigma^2)^{-n/2} \exp\left(-\sum_{i=1}^n (x_i - \bar{x})^2 + n(\bar{x} - \mu)^2/(2\sigma^2)\right) \end{aligned}$$



# Sufficiency

## Sufficiency of sample mean for the normal distribution, con't

The joint pdf for the sample mean  $\bar{X}$  which is iid  $N(\mu, \sigma^2/n)$  is:

$$f_{\bar{X}}(\bar{x}|\mu) = (2\pi\sigma^2/n)^{-1/2} \exp(-n(\bar{x} - \mu)^2/(2\sigma^2))$$

So the ratio  $p(\mathbf{x}|\theta)/q(T(\mathbf{x}|\theta))$  is  $f_X(\mathbf{x}|\mu)/f_{\bar{X}}(\bar{x}|\mu)$  which expands to:

$$\begin{aligned} \frac{f_X(\mathbf{x}|\mu)}{f_{\bar{X}}(\bar{x}|\mu)} &= \frac{(2\pi\sigma^2)^{-n/2} \exp(-\sum_{i=1}^n (x_i - \bar{x})^2 + n(\bar{x} - \mu)^2/(2\sigma^2))}{(2\pi\sigma^2/n)^{-1/2} \exp(-n(\bar{x} - \mu)^2/(2\sigma^2))} \\ &= n^{-1/2} (2\pi\sigma^2)^{-(n-1)/2} \exp(-\sum_{i=1}^n (x_i - \bar{x})^2/(2\sigma^2)) \end{aligned}$$

and does not depend on  $\mu$ . Thus the sample mean is a sufficient statistic for the parameter  $\mu$ .

# Sufficiency with order statistics

Sometimes we can't reduce the sample and have to resort to other means for determining sufficiency.

## Sufficiency when density is unknown

Let  $X_1, \dots, X_n$  be iid with pdf  $f(x)$  which is unknown. The best we can do in this case is show that the order statistics  $X_{(1)}, \dots, X_{(n)}$  are sufficient for  $f(x)$ .  
(Example in class)

# Factorization Theorem

## Theorem

A statistic  $T(\mathbf{X})$  is sufficient for  $\theta$  iff there exists functions  $g(t|\theta)$  and  $h(\mathbf{x})$  such that the joint pdf or pmf,  $f(\mathbf{x}|\theta)$  can be written as:

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

Proof (discrete case):

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

# Factorization Theorem and Exponential Families

- ▶ It is easy to find sufficient statistics for exponential family distributions using the Factorization Theorem.
- ▶ Exponential families are described in CB 3.4. They include many of the most common distributions (both discrete and continuous): normal, exponential, gamma, chi-squared, beta, Dirichlet, binomial, Bernoulli, negative binomial, Poisson, Wishart, Inverse Wishart.

## Exponential Families

Distributions belonging to the exponential family can be expressed as:

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

Where  $h(x) \geq 0$  and  $t_1(x), \dots, t_k(x)$  are real valued functions of the observations  $x$  (they cannot depend on  $\theta$ ),  $c(\theta) \geq 0$ , and  $w_1(\theta), \dots, w_k(\theta)$  are real valued functions of the parameter(s)  $\theta$  (they cannot depend on  $x$ ).

# Factorization Theorem and Exponential Families

- ▶ The important thing to notice is what characterizes the exponential family distributions—the parameter(s) and observation variable(s) must factorize.
- ▶ This means the distribution can be separated into products that each involve either the parameters or the observations.
- ▶ To verify that a pdf or pmf belongs to the exponential family, the functions  $h(x)$ ,  $c(\theta)$ ,  $w_i(\theta)$  and  $t_i(\theta)$  must be identified and shown to have the form shown above.
- ▶ Example in class of exponential family  $N(\mu, \sigma^2)$

# Factorization Theorem and Exponential Families

## Sufficiency, Factorization Theorem and Exponential Families

Let  $X_1, \dots, X_n$  be iid observations from a pdf or pmf  $f(x|\theta)$  that belongs to an exponential family:

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

where  $\theta = (\theta_1, \dots, \theta_d)$ ,  $d \leq k$ .

Then

$$T(\mathbf{X}) = \left(\sum_{j=1}^n t_1(X_j), \dots, \sum_{j=1}^n t_k(X_j)\right)$$

is sufficient for  $\theta$ .

Examples in class of Poisson and normal exponential family factorization for finding sufficient statistics.

# Minimal Sufficient Statistics

- ▶ As we have seen, there are cases where there are many sufficient statistics for a particular model.
- ▶ Sufficient statistics are not unique. If  $T(X)$  is sufficient and  $T^*(X)$  is another statistic such that  $T(X) = g_1(T^*(X))$  for some function  $g_1$  then  $T^*(X)$  is also sufficient.

$$\begin{aligned}f(x|\theta) &= g(T(x)|\theta)h(x) \\ &= g(g_1(T^*(x))|\theta)h(x) \\ &= g^*(T^*(x)|\theta)h(x)\end{aligned}$$

- ▶ So if  $T(X)$  is sufficient, so is  $T^*(X) = (T(X), T_1(X))$  where  $T_1(X)$  is any other statistic.
- ▶ If  $T(X) = g_1(T^*(X))$  then the partition of  $\mathcal{X}$  defined by  $T(x)$  is coarser than that defined by  $T^*(x)$ .

# Minimal Sufficient Statistics

- ▶ Given many possible sufficient statistics, are some better than others?
- ▶ Recall we want a statistic that provides data reduction without loss of information about the parameter  $\theta$ . Thus, a statistic that achieves the most data reduction while retaining all the information about  $\theta$  is preferable. Such a statistic is called a minimal sufficient statistic.

## Minimal Sufficient Statistic

$T(X)$  is a minimal sufficient statistic if it is sufficient and for any other sufficient statistic  $T^*(X)$ ,  $T(X)$  is a function of  $T^*(X)$ .

- ▶ However, this definition of minimal sufficient statistics does not often help identify which of a group of sufficient statistics is actually minimal (normal model example).



# Minimal Sufficient Statistics

## Theorem for Minimal Sufficient Statistics

If  $T(X)$  has the property that the ratio  $f(x|\theta)/f(y|\theta)$  does not depend on  $\theta$  iff  $T(x) = T(y)$  then  $T(X)$  is a minimal sufficient statistic for  $\theta$ .

Proof:

Let  $T(X)$  satisfy the condition of the theorem. We show that  $T(X)$  is sufficient and that it is minimally sufficient.

Let  $x_t$  denote an element of  $A_t$ . Recall  $A_t = \{x : T(x) = t\}$ . So,  $T(x_t) = t$  and  $T(x_{T(x)}) = T(x)$

from the theorem, we have:

$$\frac{f(x|\theta)}{f(x_{T(x)}|\theta)} = h(x)$$

where  $h(x)$  is some function that does not depend on  $\theta$ .

$$f(x|\theta) = f(x_{T(x)}|\theta)h(x)$$

So by the factorization theorem,  $T(X)$  is sufficient.

# Minimal Sufficient Statistics

## Theorem for Minimal Sufficient Statistics, con't

Let  $T^*(X)$  be another sufficient statistic. By the factorization theorem we have

$$\frac{f(x|\theta)}{f(y|\theta)} = \frac{h^*(x)g^*(T^*(x)|\theta)}{h^*(y)g^*(T^*(y)|\theta)} = \frac{h^*(x)}{h^*(y)}$$

Thus  $T^*(x) = T^*(y)$  implies that  $f(x|\theta)/f(y|\theta)$  does not depend on  $\theta$ . From the assumption that  $T(x) = T(y)$  it follows that the partition of  $\mathcal{X}$  induced by  $T^*(x)$  is finer than that induced by  $T(X)$ . This implies that  $T(X)$  is a minimal sufficient statistic.

# Minimal Sufficient Statistics

Some general notes about sufficiency and minimal sufficiency:

- ▶ In terms of partitioning the sample space, any sufficient statistic introduces a partition of the sample space.
- ▶ The partition of the minimal sufficient statistic is the **coarsest** so that it achieves the greatest possible data reduction for a sufficient statistic.
- ▶ A minimal sufficient statistic eliminates all of the extra information in the sample and leaves only that which contains information about  $\theta$ .

# Ancillary Statistics

An ancillary statistic:

- ▶ contains no information about parameter  $\theta$ ; however, it provides a complimentary purpose to a sufficient statistic.
- ▶ An ancillary statistic by itself does not provide any information about a parameter, but in conjunction with another statistic it can (R.A. Fisher).
- ▶ is an observation on a random variable whose distribution is fixed and known, but unrelated to  $\theta$ .
- ▶ is denoted as  $S(X)$

The range statistic  $R = X_{(n)} - X_{(1)}$  is a common example of an ancillary statistic because it does not depend on the distribution of the sample  $\mathbf{x}$  but rather on the parameter of the distribution that relates to *location*.

Other examples include ancillary statistics belonging to the scale family, or a mixture of scale and location.

# Ancillary Statistics

Location family  $\theta$  is the location parameter:

$$\{F(x - \theta) : -\infty < \theta < \infty\}$$

Scale family  $\theta$  is the scale parameter:

$$\{(1/\theta)F(x/\theta) : \theta > 0\}$$

Scale-Location family  $\theta_1$  is the scale parameter and  $\theta_2$  is the location parameter:

$$\{(1/\theta_1)F((x - \theta_2)/\theta_1) : \theta_1 > 0, -\infty < \theta_2 < \infty\}$$

# Ancillary Statistics

## Location Ancillary Statistics

We use the CDF to show how location ancillary statistics do not depend on the parameter  $\theta$ . Let  $X_1, \dots, X_n$  be iid observation from a location parameter family with cdf  $F(x - \theta)$ . Let  $Z_1, \dots, Z_n$  be iid observations with cdf  $F(x)$  (i.e.  $\theta=0$ ) with  $X_1 = Z_1 + \theta, \dots, X_n = Z_n + \theta$ . Show the range  $R = X_{(n)} - X_{(1)}$  is an ancillary statistic.

The cdf of R is

$$\begin{aligned} F(r|\theta) &= P(R \leq r) \\ &= P(\max X_i - \min X_i \leq r) \\ &= P(\max(Z_i + \theta) - \min(Z_i + \theta) \leq r) \\ &= P(\max Z_i - \min Z_i + \theta - \theta \leq r) \\ &= P(\max Z_i - \min Z_i \leq r) \end{aligned}$$

Which does not depend on  $\theta$  because the distribution of  $Z_1, \dots, Z_n$  does not depend on  $\theta$ .

In class example showing the range for  $\text{Uniform}(\theta, \theta + 1)$  is an ancillary statistic.

# Ancillary Statistics

## Scale Ancillary Statistics

Again we use the CDF to show how scale ancillary statistics do not depend on the parameter  $\theta$ . Let  $X_1, \dots, X_n$  be iid observation from a scale parameter family with cdf  $F(x/\theta)$ . Any statistic that depends on the sample through its  $n-1$  values  $X_1/X_n, \dots, X_{n-1}/X_n$  is an ancillary statistic.

Let  $Z_1, \dots, Z_n$  be iid observations with cdf  $F(x)$  (i.e.  $\theta=1$ ) with  $X_i = \theta Z_i$ . The joint CDF of  $X_1/X_n, \dots, X_{n-1}/X_n$  is:

$$\begin{aligned} F(y_1, \dots, y_{n-1} | \theta) &= P(X_1/X_n \leq y_1, \dots, X_{n-1}/X_n \leq y_{n-1}) \\ &= P(\theta Z_1/(\theta Z_n) \leq y_1, \dots, \theta Z_{n-1}/(\theta Z_n) \leq y_{n-1}) \\ &= P(Z_1/Z_n \leq y_1, \dots, Z_{n-1}/Z_n \leq y_{n-1}) \end{aligned}$$

Which does not depend on  $\theta$  because the distribution of  $Z_1, \dots, Z_n$  does not depend on  $\theta$ .

# Location and Scale Families

Examples of Location families:

- ▶ Uniform  $(\theta, \theta + 1)$  with pdf  $f(x|\theta) = I(\theta \leq x \leq \theta + 1)$
- ▶ Cauchy location family with pdf  $f(x|\theta) = \frac{1}{\pi(1+(x-\theta)^2)}$
- ▶  $N(\mu, \sigma^2)$  with  $\mu$  unknown and  $\sigma^2 > 0$  known

Examples of Scale families

- ▶ Uniform  $(0, \theta)$  with pdf  $f(x|\theta) = \theta^{-1}I(0 \leq x \leq \theta)$
- ▶ Cauchy location family with pdf  $f(x|\theta) = \frac{1}{\theta\pi(1+(x/\theta)^2)}$
- ▶  $N(0, \sigma^2)$  with  $\sigma^2 > 0$  unknown
- ▶  $\text{Exp}(\beta)$  with pdf  $f(x|\beta) = \beta^{-1}e^{-x/\beta}I(x \geq 0)$

Examples of Location-Scale families:

- ▶ Uniform  $(\alpha, \beta)$
- ▶ Normal  $N(\mu, \sigma^2)$  both parameters unknown



# Complete Statistics

Ancillary statistics in conjunction with sufficient statistics provides us with a definition for complete statistics.

Basically, if we have a sufficient statistic that optimally summarizes the observations, then there should not be an ancillary statistic that is a function of that statistic.

## Basu's Theorem

If  $T(X)$  is complete and a minimal sufficient statistic then  $T(X)$  is independent of every ancillary statistic.

(i.e. A complete sufficient statistic is independent of every ancillary statistic.)

(Proof in class)

# Complete Statistics

Complete statistics apply to families of distributions, most importantly the exponential family.

## Complete Statistics in the Exponential Family

Let  $X_1, X_2, \dots, X_n$  be observations from an exponential family with pdf (or pmf) that has the form

$$f(x|\theta) = h(x)c(\theta) \exp\left(\sum_{j=1}^k w(\theta_j)t_j(x)\right)$$

where  $\theta = (\theta_1, \theta_2, \dots, \theta_k)$ . Then the statistic

$$T(X) = \left(\sum_{i=1}^n t_1(X_i), \sum_{i=1}^n t_2(X_i), \dots, \sum_{i=1}^n t_k(X_i)\right)$$

is complete as long as the parameter space  $\Theta$  contains an open set in  $\Re^k$

# Sufficient, Minimal, and Complete Statistics

Sufficiency and completeness is used in constructing estimators of parameters  $\theta$  such as the maximum likelihood estimator and Bayes estimators.

- ▶ Suppose that  $T$  is sufficient for  $\theta$  and there exists a maximum likelihood estimator of  $\theta$ . Then, there exists a maximum likelihood estimator  $U$  that is a function of  $T$ .
- ▶ Rao-Blackwell Theorem: Suppose that  $T$  is sufficient for  $\theta$  and that  $W$  is an unbiased estimator of  $\tau(\theta)$ . Then  $E_{\theta}(W|T)$  is a statistic and is uniformly better than  $W$  as an estimator of  $\tau(\theta)$ .
- ▶ Lehmann-Scheffe Theorem: Suppose that  $T$  is sufficient and complete for  $\theta$  and that  $W = \phi(T)$  is an unbiased estimator of  $\tau(\theta)$ . Then  $W$  is a uniformly minimum variance unbiased estimator (UMVUE) of  $\tau(\theta)$ .

# Likelihood

First we need to go over the likelihood principle:

- ▶ Likelihoods relate data to a population.
- ▶ They arise from a probability distribution function  $f(x|\theta)$  connecting data  $x$  to a population.
- ▶ Used as a data reduction technique.
- ▶ We assume data come from a family of distributions with unknown parameters.
- ▶ We use the data to estimate these unknown parameters.

# Likelihood

## Definition: The Likelihood Principle

The likelihood principle states that given a pdf  $f(\mathbf{x}|\theta)$  and observed data  $\mathbf{x}$ , all of the relevant information regarding the unknown parameter(s)  $\theta$  is contained in the likelihood function for the observed  $\mathbf{x}$

Two likelihood functions contain the same information about  $\theta$  if they are proportional to each other.

# Likelihood

- ▶ In the context of random variables,  $X_1, X_2, \dots, X_n$  are an iid sample from a population with pdf  $f(x|\theta_1, \theta_2, \dots, \theta_k)$  where  $\theta_i, i = 1, \dots, k$  are unknown parameters.
- ▶ The likelihood is  $f$  viewed as a function of  $\theta_i$  for fixed observed values of  $x$ .
- ▶ The joint density of the data evaluated as a function of the parameters with the data fixed.

# Likelihood Function

## Definition

Let  $f(\mathbf{x}|\theta)$  be the joint pmf or pdf of  $\mathbf{X} = X_1, X_2, \dots, X_n$  (iid). Given  $\mathbf{X}=\mathbf{x}$  is observed, then the likelihood function of  $\theta$  defined is by:

$$L(\theta|\mathbf{x}) = f(\mathbf{x}|\theta)$$

for continuous  $\mathbf{X}$ , or for discrete  $\mathbf{X}$ :

$$L(\theta|\mathbf{x}) = P_{\theta}(\mathbf{X} = \mathbf{x})$$

# Equivariance

- ▶ Given some data, statistical decisions should not be affected by simple transformations or reordering of the data.
- ▶ For example, the value of a point estimate will be affected by a transformation, but it should be *equivariant* in the sense that it reflects the transformation in a meaningful way.
- ▶ This is formalized by the equivariance principle through which appropriate classes of transformations are defined and rules that statistical decisions must satisfy are specified.