

Statistical

Inference

Aakash  
Ghosh

# Contents

1	Introduction	2
1.1	why statistical inference?	2
1.2	Problem of non-identifiability	2
1.2.a	Example	3
1.3	Parameter	3
1.3.a	Examples	3
1.4	Regular families	4
2	Statistic	5
2.1	Introduction	5
2.2	Ancillary Statistics	5
2.2.a	Example 1	5
2.2.b	Example 2	6
2.2.c	Example 3	6
2.3	Scale Family	7
2.4	Location Family	7
2.5	Location-Scale Family	7
2.6	Limiting conditions	8
2.7	Characterization in terms of level sets	8
2.8	Sufficient Statistic	9
2.8.a	Example	9
2.9	Fisher-Neyman factorization	9
2.9.a	Example 1	10
2.9.b	Example 2	10
2.9.c	Example 3	10
2.9.d	Example 4	10

# Chapter 1

## Introduction

### 1.1 why statistical inference?

There are three main questions we would like to answer:

1. Given a random vector  $X$ , determine which  $\theta \in \Theta$  is most compatible with the data given [i.e. which  $\theta \in \Theta$  generates  $X$ ]. This is point estimation.
2. Given an  $X$  we wish to determine if a given value of  $\theta$  or distribution  $f_\theta$  is consistent. This is hypothesis testing.
3. After (1), we would like to find a suitable range of values of  $\theta$  which is compatible. This is interval estimation.

We would also like our answers to be optimal in each of the above cases.

### 1.2 Problem of non-identifiability

This is a problem which arises when more than one value of  $\theta \in \Theta$  can give rise to the same  $f_\theta$ .

$$g : \Theta \rightarrow \mathcal{F}$$

↑  $g$  is onto, but  
not one-one

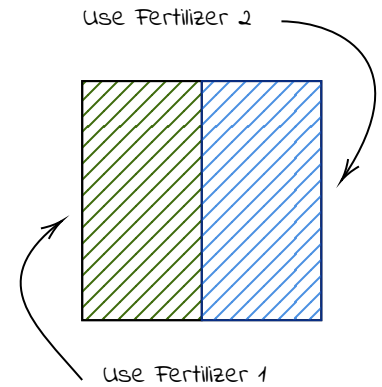
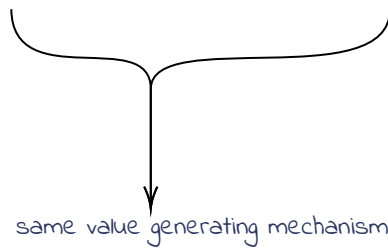
$\exists \theta_1, \theta_2$  such that  $g(\theta_1) = g(\theta_2)$ . The difficulty arising from this is very obvious: with non-identifiability,  $\theta_1$  can be confused with  $\theta_2$  as we are essentially claiming that both  $\theta_1$  and  $\theta_2$  results in the same data generating mechanism. We must understand that this is an inherent problem of our data generating mechanism/experiment and must be addressed before any further analysis is carried out.

### 1.2.a Example

Consider a scenario where we want to test the effectiveness of two fertilizers on a field. what we might do is divide the field in two divisions such, use fertilizer one on the first half and fertilizer two on the second half and then record our crop yield. Assume that the field has yield follows a normal distribution with mean  $\alpha$  with standard deviation 1. Further, assume that each fertilizer  $i$  increases yield by  $\alpha_i$ . we repeat this process several ( $= n$ ) times.

Assume on the  $j^{th}$  iteration with fertilizer  $i$  the yield is  $X_{ij}$ . Define  $X_1 = [X_{11}, X_{12}, X_{13} \dots X_{1n}]$  and  $X_2 = [X_{21}, X_{22}, X_{23} \dots X_{2n}]$ . Then by using our setup, we can use  $X_1$  and  $X_2$  to estimate  $\alpha + \alpha_1$  and  $\alpha + \alpha_2$  but not  $\alpha, \alpha_1, \alpha_2$  individually.

$$\theta = (\alpha, \alpha_1, \alpha_2) \qquad \theta_c = (\alpha - c, \alpha_1 + c, \alpha_2 + c)$$



Again note this inability is an inherent flaw in the experiment itself, and we can't do anything to circumnavigate it. However, we might address this problem by assuming that  $\alpha_1 + \alpha_2 = 0$ . Note, this is definitely not true. But we are claiming something along the lines that the net fertility remains constant, and this helps in computing all the values.

## 1.3 Parameter

A parameter is a function

$$\nu : \mathcal{F}_\theta \rightarrow \mathcal{N}$$

where  $\mathcal{N}$  is an arbitrary set. when  $F_\theta$  is an identifiable parameterization, then there exists function  $q$  such that:

$$\nu(F_\theta) = q(\theta)$$

### 1.3.a Examples

1.  $\mu(\theta) = \int x dF_\theta(x)$
2.  $\sigma^2(\theta) = \int x^2 dF_\theta(x) - (\mu(\theta))^2$
3.  $m(\theta) = \text{median of } F_\theta$

## 1.4 Regular families

For our work we will assume that:

1. All  $F_\theta \in \mathcal{F}$  have continuous CDF with density  $f_\theta$  or are discrete with pmf  $f_\theta$
2. There is identifiable parameterization.

## Chapter 2

# Statistic

### 2.1 Introduction

Any measurable function  $T : X \rightarrow \mathbb{R}^k$  is called a statistic.

The idea is to compress the given data without losing valuable information. we look at two extreme cases:

- $T(X) = \text{constant}$ : This tells us nothing and all information is lost,
- $T(X) = X$ : while this retains all information, there is no compression that takes place.

### 2.2 Ancillary Statistics

we wish to find good statistics. And it turns out that to achieve this it is better to sieve out the bad ones and then look further in whatever is left\*.

An ancillary statistic  $T = T(X)$  is such that the distribution of  $T(X)$  is independent of  $\theta$ .

\*Honestly, this feels like something you find on cleanup: You do everything and then understand that doing it this way is smoother.

#### 2.2.1 Example 1

Assume  $X_1, X_2 \dots X_n \sim (\mu, 1), \mu \in \mathbb{R}$ . Let

$$T(X) = \sum_{i=0}^n \sum_{j=0}^n (X_i - X_j)^2$$

Two approaches for solving this are:

1. Define  $\tilde{X}_i = X_i - \bar{X}$ . Note that:

$$\begin{aligned}
 T(X) &= \sum_{i \neq j} (X_i - X_j) \\
 &= \sum_{i \neq j} (\tilde{X}_i - \tilde{X}_j) \\
 &= (n-1) \sum_i \tilde{X}_i^2 - 2 \sum_{i \neq j} \tilde{X}_i \tilde{X}_j \\
 &= (n+1) \sum_i \tilde{X}_i^2 - 2 \sum_{i,j} \tilde{X}_i \tilde{X}_j \\
 &= (n+1) \sum_i \tilde{X}_i^2 - 2 \left( \sum_i \tilde{X}_i \right) \left( \sum_j \tilde{X}_j \right) \\
 &= (n+1) \sum_i (X_i - \bar{X})^2
 \end{aligned}$$

Now distribution of  $X_i - \bar{X} \sim \mathcal{N}(0, 1)$  which is independent of  $\mu$ . Therefore,  $T$  is ancillary.

2. The above argument can be considerably shortened by arguing that the distribution of  $X_i - X_j$  is independent of  $\mu_i$  and therefore the whole sum is independent of  $\mu$ .

### 2.2.b Example 2

Assume  $X_1, X_2 \dots X_n \sim \text{Exp}(\theta)$ . The density function is given by  $f_\theta(x) = \theta e^{-\theta x}$ . Such a distribution often occurs when modelling things like lifetime of a particular make of a bulb, etc. Define:

$$T(X) = \frac{\max_i X_i}{\min_i X_i}$$

Finding the distribution of  $T$  explicitly is a nightmare. Instead, what we do is define  $Y_i = \theta X_i$ . Note that:

$$P(Y_i < y) = P(\theta X_i < y) = \int_0^{y/\theta} \theta e^{-\theta x} dx = \int_0^y e^{-u} du$$

Therefore,  $Y_i \sim \text{Exp}(1)$ . we rewrite  $T$  as:

$$T(X) = \frac{\max_i X_i}{\min_i X_i} = \frac{\max_i Y_i}{\min_i Y_i} = T(Y)$$

But, it is trivial to note that any parameter on  $Y$  is independent of  $\theta$ . Therefore,  $T(X)$  is ancillary.

### 2.2.c Example 3

Assume  $X_1, X_2 \dots X_n \sim \text{Unif}(0, \theta)$ . Define:

$$T(X) = \frac{\max_i X_i}{\min_i X_i}$$

Define  $Y_i = X_i/\theta$ . Note that:

$$P(Y_i < y) = P(X_i/\theta < y) = \int_0^{y\theta} \frac{1}{\theta} dx = \int_0^y du$$

Therefore,  $y \sim \text{Unif}(0, 1)$ . we rewrite  $T$  as:

$$T(X) = \frac{\max_i X_i}{\min_i X_i} = \frac{\max_i Y_i}{\min_i Y_i} = T(Y)$$

Again as any parameter on  $Y$  is independent of  $\theta$ . Therefore,  $T(X)$  is ancillary.

The idea in the examples above is that we are "scaling" the function to fit our needs.

## 2.3 Scale Family

A scale family is a family  $\mathcal{F}$  of distribution given by:

$$\left\{ f_\theta | f_\theta(x) = \frac{1}{\theta} g\left(\frac{x}{\theta}\right) \right\}$$

where  $g$  is a known distribution. Set  $Y_i = X_i/\theta$ . Any  $T(X)$  which can be written as  $S(Y)$  is ancillary.

## 2.4 Location Family

A scale family is a family  $\mathcal{F}$  of distribution given by:

$$\{f_\theta | f_\theta(x) = g(x - \theta)\}$$

where  $g$  is a known distribution. Set  $Y = X_i - \theta$ . Any  $T(X)$  which can be written as  $S(y)$  is ancillary.

## 2.5 Location-Scale Family

A scale family is a family  $\mathcal{F}$  of distribution given by:

$$\left\{ f_\theta | f_\theta(x) = \frac{1}{\sigma} g\left(\frac{x - \mu}{\sigma}\right) \right\}$$

where  $g$  is a known distribution. Set  $Y_i = (X_i - \mu)/\sigma$ . Any  $T(X)$  which can be written as  $S(Y)$  is ancillary.



## 2.6 Limiting conditions

other ways of identifying a "bad-statistic" is to look at limiting conditions. For example, let  $X_i \sim \text{Unif}(0, \theta)$ . Consider two candidates:

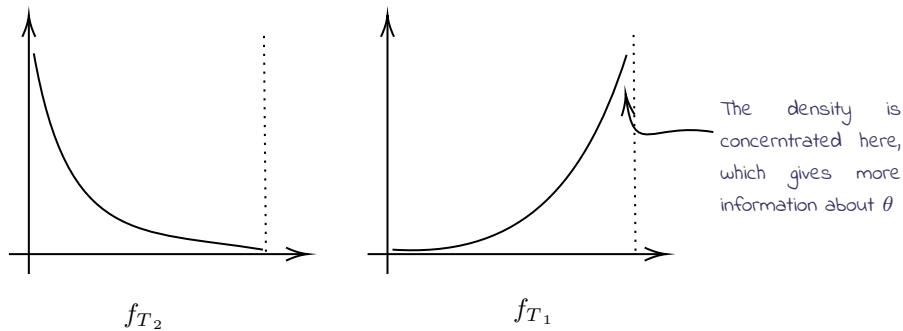
$$T_1(X) = \max_i X_i$$

$$T_2(X) = \min_i X_i$$

It is easy to compute that:

$$f_{T_2}(t) = \frac{n}{\theta} \left(1 - \frac{t}{\theta}\right)^{n-1}, 0 < t < \theta$$

$$f_{T_1}(t) = \frac{n}{\theta} \left(\frac{t}{\theta}\right)^{n-1}, 0 < t < \theta$$

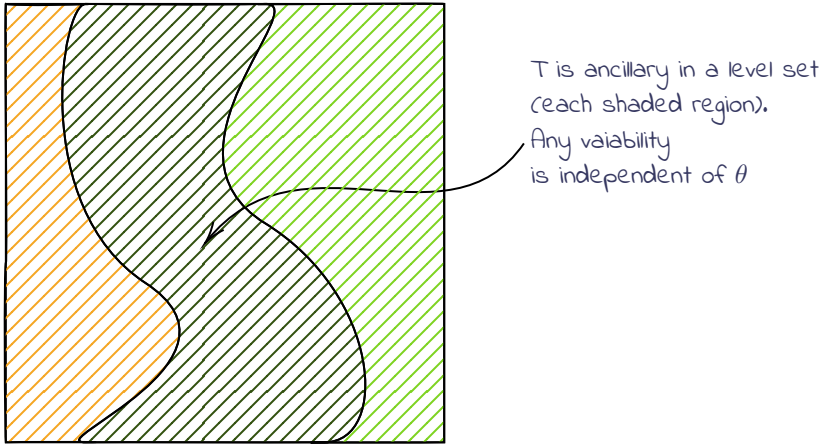


## 2.7 Characterization in terms of level sets

Consider a level set of  $T$

$$A_t = \{x \in \mathbb{R}^n | T(x) = t\}$$

Ideally what we want is  $T$  to do is compress all information about  $\theta$  and leave out all the rest. This essentially translates to arguing that any variability of  $x$  in  $A$  is independent of  $\theta$ , or more formally, conditional distribution of  $X$  given  $T(X) = t$  is independent of  $\theta$ . But that implies that for any  $t$ ,  $T$  restricted to  $T^{-1}(t)$  is ancillary in nature.



Such a  $T$  is called a sufficient statistic.

## 2.8 Sufficient Statistic

$T = T(x)$  is said to be sufficient for  $\theta$  if  $P(X = B|T = t)$  is independent of  $\theta$  for all  $B \in \mathcal{B}_{\mathbb{R}^n}, t \in \mathbb{R}^k$ . Note that any one-one map of a sufficient statistic is also sufficient.

### 2.8.a Example

Let  $X_1, X_2 \dots X_i$  be Bernoulli iid with parameter  $\theta$ , i.e.

$$P(X_i = 1) = \theta \quad P(X_i = 0) = 1 - \theta$$

Take  $T = \sum X_i$ , we calculate  $P(X = x_i|T = t)$ .

$$P(X = x_i|T = t) = \frac{P(X = x, T(X) = t)}{P(T = t)}$$

the denominator easily comes out to be  $\binom{n}{t}\theta^t(1-\theta)^{n-t}$ . For the numerator, consider two cases:

1.  $T(x) = t$ : The event  $[X = x]$  is a subset of  $[T(x) = t]$ . Therefore, the numerator simplifies to  $P(X = x) = \theta^t(1-\theta)^{n-t}$ , the whole expression simplifies to  $\binom{n}{t}^{-1}$ .
2.  $T(x) \neq t$ : The events  $[X = x]$  and  $[T = t]$  are disjoint. Therefore, the numerator simplifies to 0.

In either case, we see that the value of the expression is independent of  $\theta$ .

## 2.9 Fisher-Neyman factorization

Calculating conditional probabilities is hard for continuous distributions. Finding candidate  $T$  for checking is even harder. A nice characterization for overcoming these difficulties is Fisher-Neyman

factorization.  $T(X)$  is sufficient for  $X$  if and only if  $f_\theta(x)$  can be written as

$$f_\theta(x) = g(T(x), \theta)h(x)$$

### 2.9.a Example 1

Let  $X_i$ s be iid Bernoulli with parameter  $\theta$ .

$$\begin{aligned} f_\theta(x) &= P_\theta(X = x) && \text{[Good practice to include } \theta \text{ to show dependence]} \\ &= \prod_{i=1}^n \theta^{x_i} (1 - \theta)^{1-x_i} 1_{x_i \in \{0,1\}} && \text{[Make sure to characterize the domain]} \\ &= \underbrace{\theta^{\sum x_i} (1 - \theta)^{n - \sum x_i}}_{\substack{\text{Set } T(x) = \sum x_i \\ \text{This part becomes } g(T(x), \theta)}} \underbrace{1_{x_i \in \{0,1\}}}_{\substack{\text{This becomes } h(x)}} \end{aligned}$$

which is exactly what we had concluded before.

### 2.9.b Example 2

Let  $X_i$ s be iid exponentials with parameter  $\theta$ .

$$f_\theta(x) = \prod \theta e^{-\theta x_i} 1_{x_i > 0} = \theta^n e^{-\theta \sum x_i} \prod 1_{x_i > 0}$$

we have a similar decomposition as before with  $T(x) = \sum x_i$

### 2.9.c Example 3

Let  $X_i$ s be iid uniform with parameter  $\theta$ .

$$\begin{aligned} f_\theta(x) &= \prod_i \theta^{-1} 1_{x_i \in (0, \theta)} && \text{[Note how the domain is parameter dependent]} \\ &= \theta^{-n} 1_{0 < \min x_i \leq \max x_i \leq 1} \\ &= \theta^{-n} 1_{0 < \min x_i, \max x_i \leq 1} \\ &= \underbrace{\theta^{-n} 1_{\max x_i \leq 1}}_{\substack{\text{Set } T(x) = \max x_i \\ \text{and get } g(T(x), \theta)}} \underbrace{1_{0 < \min x_i}}_{\substack{\text{This becomes } h(x)}} \end{aligned}$$

### 2.9.d Example 4

Let  $X_i$ s be iid Cauchy with parameter  $\theta$  (pmf  $\frac{1}{\pi} \frac{1}{(x - \theta)^2}$ ).

$$f_\theta(x) = \frac{1}{\pi^n} \prod \frac{1}{1 + (x_i - \theta)^2}$$

Note there is not much to exploit over here. The best (to be proven later) that we can do is to order them in increasing order. This does not provide a reduction in dimensionality but compresses some information as there are multiple combinations of  $x$  which results in the same order.