

Engineering Statistics
Professor Manjesh Hanawal
Industrial Engineering and Operations Research
Indian Institute of Technology, Bombay
Lecture 35
Minimal Sufficient Statistics

(Refer Slide Time: 00:22)

Example 1:

Population distribution $\sim \mathcal{N}(\mu, \sigma^2)$, with μ unknown and σ^2 known

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

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

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

$$= h(x)g(\bar{x}|\mu)$$

Hence $T(x) = \bar{x}$ is a sufficient statistics.

Handwritten notes:
 $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$
 \bar{x} is a sufficient

Functions of Sufficient Statistics

- Can there be one sufficient statistic or multiple?
- The entire sample is always sufficient statistic. Set $T(x) = (x_1, x_2, \dots, x_n)$ and $h(x) = 1$, then

$$f(x|\theta) = f(T(x)|\theta)h(x)$$
- Any one-to-one function of a sufficient statistics is also a sufficient statistics.
- Let $T^*(x) = r(T(x))$ for some invertible r and $T(x)$ is a sufficient statistics

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

$$= g(r^{-1}(T^*(x)|\theta))h(x)$$

$$= g^*(T^*(x)|\theta)h(x)$$

Handwritten notes:
 $g \circ r^{-1} = g^*$

Hence $T^*(x)$ is a sufficient statistics

So, any questions from the last lecture, what are we discussed about this factorization theorem, how to use this factorization theorem to come up with some statistics which are sufficient and then if you know already some function is belongs to exponential family, all we need to do is

write it in its in this form, and they are all the t_i functions will already give me the sufficient statistics for that.

(Refer Slide Time: 00:50)

Functions of Sufficient Statistics

- ▶ Can there be one sufficient statistic or multiple?
- ▶ The entire sample is always sufficient statistic. Set $T(x) = (x_1, x_2, \dots, x_n)$ and $h(x) = 1$, then

$$f(x|\theta) = f(T(x)|\theta)h(x)$$
- ▶ Any one-to-one function of a sufficient statistics is also a sufficient statistics.
- ▶ Let $T^*(X) = r(T(X))$ for some invertible r and $T(X)$ is a sufficient statistics.

$$\begin{aligned} f(x|\theta) &= g(T(x)|\theta)h(x) \\ &= g(r^{-1}(T^*(x)|\theta))h(x) \\ &= g^*(T^*(x)|\theta)h(x) \end{aligned}$$

Hence $T^*(x)$ is a sufficient statistics

Handwritten notes: "sufficient", "g o r^{-1} = g^", "h"*

We also discussed that sufficient statistics need not be unique there could be multiple ways of getting sufficient statistics. And if you give me one sufficient statistics, I may derive multiple sufficient statistics by using some invertible function on it.

(Refer Slide Time: 01:09)

Minimal Sufficient Statistics

A sufficient Statistics $T(X)$ is called a minimal sufficient statistics, if for any sufficient statistics $T'(X)$,

$$T'(x) = T'(y) \implies T(x) = T(y)$$

- ▶ $T = \{t : t = T(x), x \in \mathcal{X}\}$ and $A_t = \{x : T(x) = t\}$
- ▶ $T' = \{t' : t' = T'(x), x \in \mathcal{X}\}$ and $B_{t'} = \{x : T'(x) = t'\}$
- ▶ for any $t' \in T'$, there exists $t \in T$ such that $B_{t'} \subset A_t$.

Then we started talking about minimal sufficient statistics. So, let us repeat some of the part here what you are already done. So, I am slowing myself by repeating things. If you want to further slow, you should interact or ask something related to this, I mean, I am trying to incorporating the bringing the slowness by repeating things, but if you also want to do from your side, you should act do not simply sit.

(Refer Slide Time: 01:41)

Minimal Sufficient Statistics

A sufficient Statistics $T(X)$ is called a minimal sufficient statistics, if for any sufficient statistics $T'(X)$,

$$T'(x) = T'(y) \implies T(x) = T(y)$$

- $T = \{t : t = T(x), x \in \mathcal{X}\}$ and $A_t = \{x : T(x) = t\}$
- $T' = \{t' : t' = T'(x), x \in \mathcal{X}\}$ and $B_{t'} = \{x : T'(x) = t'\}$
- for any $t' \in T'$, there exists $t \in T$ such that $B_{t'} \subset A_t$.

Handwritten notes on the right:
 T' is a sufficient statistic
 $T'(x) = t' \implies A_t$

So, by definition, we said that we are going to call a statistic T or rather a sufficient statistic T to be minimal sufficient statistic, if you have any other sufficient statistics T' and that T' take same value on 2 samples x and y. Then it implies that my T also takes the same value on those 2 samples x and y.

That means if two samples are indistinguishable under sufficient statistics t, then they are also indistinguishable under the sufficient statistic, or minimal sufficient statistic T. Now, let us quickly understand this suppose, let us say what I have written in text, let us understand that through diagram.

Let us hypothetically just make a 2 dimensional case x1 and x2. Let us say this is my region for some reason. And hypothetically assume that this region gets partitioned into 3 parts under my hypothesis T. What I mean by that? Take any point here, let us say x1 comma x2, this is going to

map if you apply T on that, that your function t on that it is going to give you a value small t and that value small t is the same in this region.

You understand that part, all the points in this region map to the small t and similarly all the points in this also mapped to $1/T$ and like this I have made it into 3 partition. And let us call this region as A_1 , this region as A_2 and this region as A_3 . So, all the points in the region A_1 , they have the same value. When you apply T on that it they will map to the same point.

Now, let us take another T' . Sorry, let us take another statistics T' this another statistics. So, let us call this is T and let us call let us do another partition under this this is under T' . Not necessary that this T' will also result in the same partition. It may result in some other partition. Let us call that as simply like this, under this let us call this A_1' . And let us call this A_2' . And let us call this A_3' .

In fact, it may result in more than 3 partitions. It may have more than 3 or maybe less than 3, whatever. I am just for representation purposes I am just writing it as 3. Now, suppose T' is a sufficient statistic let us assume that and let us call the point to which all these maps here as called, let us call T_1' , let us call this as T_2' and let us call all the value it maps to all the points in this region maps to as T_3' . And similarly, here that is called T_1 , and this is T_2 , and let us call this T_3 .

Now, from this region, I have all the points gets either mapped to T_1' , T_2' , or T_3' . Now, arbitrarily, let us pick T_2' . And let us take some point. Let us, let us, let us know I know let some point here. Let some point let us call this point some x again, that consists of x_1 and x_2 , let us say this $T_1' x$ maps to T_2' . In fact, all points in this region are going to map to T_2' . Now, this region let us call this whatever this where this T_2' , let us call this A_2' , I know that all these points here in this region, they are all getting mapped to the same value of T_2 point. Now, if I look into this portion here will they also under T will they also map to the same point? Not necessarily not map to the same point same as T_2' , but all their values will be same?

Student: No.

Professor Manjesh Hanawal: If T happens to be minimal sufficient statistic so, all these points are taking same values, let us say any point you take x and y here, they are going to take the

same value here, but if T is a sufficient statistic by definition, they should also be taking the same value. So, all these points that are here they are also taking the same value under T even though I have partitioned incorrectly it is not necessary. So, then there should be region here under T which incorporates all these points some region here, which incorporates all these points that are mapped to the same value under t prime. And that is what we are saying.

So, this region here, which has become a partition, this is going to be a subset of another partition under T maybe ideally, I should have drawn this like this it would have been if I had to draw this first figure with respect to this this should have been this maybe this like this, not even this and this, so how many partitions had got now?

So, this has 123, 123 and here, I wanted to have only maybe I want only this much now it has only 2 partitions in this, actually all these points are also here and that is now incorporated in a bigger partition. We just level a presentation. So, if this happens to be sorry, here it I call him b in this set. I am calling them as b, b_1, b_2 and b_3 . So, this b_2 ' here is going to be part of a_2 a_2 here just by this definition that if T is a sufficient statistic every point that are mapping to the same they should also be mapping to the same value and under my T . So, obviously the number of partition under minimal sufficient statistics is going to be large or small compared to another sufficient statistics?

Student: Small.

Professor Manjesh Hanawal: It is going to be smaller, if this is some sufficient statistics we saw it is going to partition my space into 3, but if T is sufficient statistics the number of partition it is going to have is going to smaller than this. So, in this way among all the statistics with statistics is going to partition your space into a small number of subsets? Minimal sufficient statistic.

(Refer Slide Time: 10:53)

Example: Minimal Sufficient Statistics

(x_1, x_2, \dots, x_n)
 \bar{x}, s^2

Samples are drawn from $N(\mu, \sigma^2)$ with unknown μ and known σ^2

- ▶ $T_1(x) = \bar{x}$ is a sufficient statistics for μ
- ▶ $T_2(x) = (\bar{x}, s^2)$ is also a sufficient statistics for μ (verify!)
- ▶ $T_1(x) = \bar{x} = r(\bar{x}, s^2) = r(T_2(x))$
- ▶ As both T_1 and T_2 are sufficient statistic they contain same knowledge about μ
- ▶ Additional knowledge of s^2 does not add any information about μ .
- ▶ T_1 gives better data reductions!

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

NPTEL
IE605 Engineering Statistics
CDEEP

So, this example we also discussed last time suppose, we have a samples which are drawn from a Gaussian distribution with parameter mu and sigma squared with unknown mean mu and known parameter sigma square. So, only thing that is unknown to me is mu now, I know that sample mean is going to be a sufficient statistic, now if I am going to take another sufficient statistic another statistic will actually come consists of 2 components sample mean as well as sample variance.

Now, actually this is also be sufficient statistics, because every information all the information about my parameter mu is contained in this and in fact I can ignore x square and S bar I can just written which already you know is a sufficient statistics for mu but I know that T1 is a function of T2 if you give me T2 by just dropping the second component I can recover T1.

Now, even though both this statistics or sufficient statistics contain same amount of information about my unknown parameter mu T1 is a offers may a better reduction. So, because of that it compresses information better that should also give a kind of intuition that why the number of partitions is less under minimal sufficient statistics compared to a just an ordinary sufficient statistic.

(Refer Slide Time: 12:56)

Test for Minimal Sufficient Statistics

Let $f(x|\theta)$ be the pmf/pdf of a sample. Suppose there exists a function $T(X)$ such that, for every sample pair (x, y) , the ratio $f(x|\theta)/f(y|\theta)$ is a constant iff $T(x) = T(y)$. Then $T(X)$ is a minimal sufficient statistics.

- ▶ Image of \mathcal{X} under T :
- ▶ Define partition set \mathcal{X}
- ▶ Select one element in each partition
- ▶ Pair each $x \in \mathcal{X}$ with partitions
- ▶ Argue T is sufficient statistics using factorization theorem

NPTEL IE605: Engineering Statistics Manjesh K. Hanawal 10 CDEEP

So, now, how to find minimal sufficient statistics how to test a sufficient statistics is a minimal sufficient statistic. So, first of all we started statistics is any function that is fine check whether something is sufficient statistic that was not obvious, but we come up with a method, what is our method to check whether a statistic is sufficient statistic?

Student: (())(13:28)

Professor Manjesh Hanawal: So, one thing was factorization theorem which is give a four sufficient necessary condition for something to be sufficient statistic and if you have to just to do with simply verify whether something is a sufficient statistics we had some other thing also that ratio test we had, see whether the ratio of unconditional distribution and of your samples and unconditional distribution of your statistics is independent of data.

So, now we have minimum sufficient statistics something more than sufficient statistics how to test that? So, for that we have this result which again gives some complete characterization because it is giving both sufficient and necessary conditions. So, what it says is suppose, let us say f of x by θ is your population density.

Now, if there exists a statistic T , such that you take any pair any pair x and y two samples random samples and if it so happens that the ratio of have those pdfs at those two points x and y happens to be constant, when I say constant here, independent of x sorry independent of θ

here the parameter there are only 2, 3 things that are x , y and θ here, when I say it is a constant it is independent of θ .

This is going to be constant if and only if, and if those two points x and y are such that my statistics use the same value on those 2 points, if this is the case, then my T is going to be a minimal sufficient statistics. First before we will just briefly discuss proof anybody has any difficulty in understanding this statement or interpreting this statement? What does this mean?

What it is saying is you give me a statistic. And what I am going to do is I take two points, two random points and check the ratio of these two points check the PDF of these two points. If this is independent of θ , provided on these two points, T of x equals to t of y only then it should happen. If this is going to be independent, or θ on points x and y where T of x is not equal to t of y , then you are failing this statement. If that is the case, then T is going to be a minimal sufficient statistic.