

Engineering Statistics
Professor Manjesh Hanawal
Industrial Engineering and Operational Research
Indian Institute of Technology, Bombay
Lecture 21
Random Sampling, Sample means and Sample variance

(Refer Slide Time: 0:15)



Previous Lecture:

- ▶ Exponential Family of Distributions

This Lecture:

- ▶ Population and Random Sampling
- ▶ sample mean, variance and standard deviation
- ▶ Sampling from Normal distribution



Now that we have most of the distributions we know, we have written them compactly and put them as an exponential family which has a very compact CPMF or PDF representation. Now, let us start looking into the data itself. How to infer the parameters that we are associated in these distributions? So, for that, if we have to get these underlying parameters, data is fine, but how this, that data should satisfy certain properties or the way the data is generated has to be done in some particular way that is what we are going to know studying sampling.

So, when I say sampling it is about the data points if you let us say there is a black box which is generating data, you do not know what is the underlying distribution according to which it is generating data, you just ask it to give one data point it gave that like a sampling for you. You got one sample like that you can ask more samples and that whole process is called sampling.

(Refer Slide Time: 1:42)

The screenshot shows a presentation slide titled "Random Sampling". The slide content includes:

- ▶ Samples are used to obtain information about large populations by examining only a small fraction. Examples
 - ▶ Who will win the polls?
 - ▶ Will there be demand for a new car
 - ▶ How many pay taxes
 - ▶ Health of people
- ▶ How to sample for better results
- ▶ Do random sampling for unbiased (to be made precise) results

A text box on the slide states: "Random Variables X_1, X_2, \dots, X_n are called random samples of size n from population $f(x)$ if they are i.i.d with common distribution with $f(x)$." The $f(x)$ in this text is circled in red.

Below the text box, it says: "if (x_1, x_2, \dots, x_n) are samples from population $f(x)$ " and the equation $f(x_1, x_2, \dots, x_n | \theta) = \prod_{i=1}^n f(x_i | \theta)$ is shown with red underlines.

The slide footer includes logos for NPTEL, IIT Bombay, and CDEEP, along with the name "Manjesh K. Hanawal" and the page number "3".

Now little bit motivation on the sampling. Sampling, you can just roughly think as data collection mechanism. Now, for time being forget about this underlying distribution and all. Data is collected in various form for various purposes. And nowadays, especially when there is a large population and we want to understand something about this population, but you cannot go and talk everybody, you may only talk to a small fraction and you could just imagine like you go and talking to them and collecting information that itself can be thought as sampling process.

Now, let us see. One way to look into all this is like I am going to let us say X is a random variable, which denotes the score obtained by a particular student in this IE 605, X is that random variable. And all of your realizations of that random variable. I can see what is the score obtained by you, what is the score obtained by you, like that I can think of realization of that random variable. If I got like about 60 students are in this class. If I got the score of the 60 students that I can think of 60 samples of this random variable X .

So, similarly let us think we as a country, we have so many people. Let us look into one particular issue like health status of the Indians, what is the health status of Indians? That is like one value let us say some health has some number which is going to be denoted by X . Now all of us are realization of that. All of us kind of like based on like our health value, like let us say we assuming that like, we have some inherent properties like we as Indians, which is going to go on some medical condition, maybe like all of us are manifestation of that

medical condition, which shows up in that like maybe we can go at each individual ones and ask that the extending this argument suppose that say, this is a classic examples of polls.

We get all the states polls, national level polls, everywhere. And depending on something, maybe let us say all Indians have some tendencies to make, I mean, we have our own way of making opinions about what to do, which party to consider and all. And during these polls, when people do all these opinion polls we have seen like they cannot go and ask everybody that is a massive process that election commission will do. They will only go and talk to certain people that is like a sampling process instead of going and asking everybody they are only going to possibly ask a certain fraction of the people feeling that they are there are like a representative sample of the whole population.

Now, here we can see that, X is a we can think of polling, X is like a random variable, which is going to denote what which party are going to vote. Then I cannot then all of us are like a realizations of that. Now, you can just think of these arguments like suppose there is a new car in the market and all of us have something about like to buy it or not. Now, what, the company who is launching that car want to know that, whether that brand new car will be how much it will be accepted in the market. So, they cannot go and ask everybody. They will only go and ask maybe certain fraction of the population, only collect certain samples.

So, the point here is like, we can just imagine that all of us are like samples when we are looking into the big population and let population is like a realization. What we want to do is we cannot go and ask everything, we will be only sampling certain fraction. Now the question here now comes is how to do better sampling? So, let us take this voting issue itself. When the election commission conducts voting, everybody is going and waiting, like, I assume that mostly all the populations are voting. But now we are doing the opinion poll, you do not have that much of resources going and asking everybody that is a massive work, you are only going to ask certain fraction.

Now how to choose that fraction itself? That is a question. So, that whatever that fraction tells that is actually going to represent actual outcome. Now, this motivation, the question is how to do better sampling. And now, on that front, we are going to refer to one important term called random sampling. So, that it results in an unbiased result. So, what is unbiased? We will make it precise, what is the mathematical meaning of unbiased. I mean, we use the word biased, unbiased very frequently. When something is not favor, we will say, that guy is biased

like, he is not thinking the way I want, but we will make it precise what we mean by unbiased.

Now, we are going to say that the random variables X_1, X_2 up to X_n , these are called random samples, of a population f of x . Now you see that I am using the word population for f of x , f of x is actually a PDF, but now I am calling it as a population if they are i.i.d with common distribution f of x . So, simply random samples of size n from population f of x , that means there are i.i.d samples drawn from distribution f of x . And here distribution f of x when I say it could be probability mass function or probability density function, depending on you are dealing with a discrete case or continuous case. Is that fine?

And now if you have one particular sample, let us say you have obtained n samples, x_1, x_2, \dots, x_n and if they are indeed random sample coming from population f of x , then what we know, the joint distribution of them should be written as the product of their marginal PDF. So, this is the definition we have set already. This is just I am using independence here. And identically distribution is coming because they are all following the same distribution. So, the same distribution, the common distribution, sometimes I am also now henceforth calling it as a population.

(Refer Slide Time: 10:48)

Sampling with and without replacement

With replacement

- ▶ After sampling, the sample is put back before the next sample is drawn randomly.
- ▶ Each sample comes from a new fresh experiment
- ▶ sampling with replacements gives i.i.d samples (random sample)

Without replacement

- ▶ After sampling, the sample is not put back, before the next sample is drawn randomly.
- ▶ sampling with replacements can give identical samples but not independent.

Handwritten notes:

- $y = x^2$
- $p(y=f(x))$
- $x_1 x_2 \dots x_n$
- $y_1 y_2 \dots y_n$
- $X \in \{1, 2, \dots, 6\}$
- with replacement = $\frac{5}{4}$
- $P(x=1) = \frac{1}{6}$
- without replacement $P(x_2=1) \neq \frac{1}{6}$

Now when I am doing the sampling, I had to be using words with and without replacements. So, what is with and without replacement? In the with replacement, after you do a sampling, the sample outcome, whatever that you are seeing, this will be put back before the next sample is drawn randomly. And each sample comes from a new fresh experiment, like, if you

put it back, all the possibilities are now there, when you want to do the experiment again, all the possible outcomes are there. So, that is why when I put it back, it looks like it is coming from a fresh experiment.

And when I do sampling with replacement, they will give i.i.d samples. On the contrast, if you are going to do without replacement, so, what I will do is after sampling, the sample is not put back. So, you are sampling outcomes get restricted when you do experiment again.

Student: Sir, if y equals to x squares?

Professor: y equals to?

Student: x square than that is also called sample of the population (θ) (12:07) sample of the population?

Professor: So, let us say you have this sample X_1, X_2 up to X_n and from this you got Y_1, Y_2, Y_n using this relation this is your new sample.

Student: Of the same population?

Professor: Yeah because you have this let us say this is coming from some underlying distribution. Now, one you do, y equals to $f(y)$ this will induce another distribution and y , and all these new samples, they are also coming from this distribution only.

Student: But here the distribution changes from x to x square.

Professor: Right. But yes, now you have to say that these samples are coming from this population, so the population is always, what is the PDF of for PMF you are looking at. Now, when you are going to do sampling with replacement, it can give identical samples but need not be independent, so let us quickly look into some example. Let us say X is some random variable, which is taking 6 value, sorry, from that is this experiment is throw of a die, which is going to give you 6 outcomes. Let us say I am going to repeat this experiment 10 times. When I did first time X_1 , let us say I got a value 5. What is the probability of getting 5?

Student: 1 by 6.

Professor: 1 by 6 because all the things are there. Now 5 I got. I took that 5. So, that dice, that the same dice I again used to get the second value and let us say I got a value 4, what is this probability?

Student: 1 by 6.

Professor: This is also 1 by 6. So, I am going to use the same dice and again and again. Now let us now say that x equals to 5 I got. Now I removed that face 5. Now my dice is consisting only of 5 faces and that 5 is removed. Now when I got X equals to 4, what is the probability of getting 4?

Student: 1 by 5.

Professor: 1 by 5. Now, this is a difference. Now, when I get this is let us say this is like a with replacement and now when I am doing with replacement P of X is always 1 by 6 that is uniform, but now I did sorry this is with replacement. And now when I do without replacement for this need not this is let us say t this is this and this is equals to i . This need not be 1 by 6, if you have removed something which have that faces whatever you observed you have not considered again then this need not be all the time 1 by 6. So, that is what we said sampling with replacements this results in i.i.d identical distributions all the time.

But, if you do sampling without replacement that is no more going to be identically distributed. One thing you need to notice is when we are going to do the sampling with replacement also it is not necessarily that it is always going to be i.i.d. In some cases, you may happen that it may be not independent because if you are not going to perform the experiment in an independent fashion that could be dependency there. So, you have to be that is what will careful. Whenever you are doing the sampling you are doing sampling with or without replacement and what going forward will always assume that we are going to do sampling with replacement all the time, so, that we always get i.i.d samples. Is the distinction between with or without sampling clear? Maybe we should say without.

So, when we do, let us read this sampling without replacement can give identical samples. So, this statement is a little ambiguous. Let me see what I wanted to say here.

(Refer Slide Time: 18:05)

Sampling with and without replacement

With replacement

- ▶ After sampling, the sample is put back before the next sample is drawn randomly.
- ▶ Each sample comes from a new fresh experiment
- ▶ sampling with replacements gives i.i.d samples (random sample)

Without replacement

- ▶ After sampling, the sample is not put back, before the next sample is drawn randomly.
- ▶ sampling with replacements can give identical samples but not independent.

Handwritten notes:

- $y = x^2$
- $p(y=i) = \frac{1}{6}$
- $X \in \{1, 2, \dots, 6\}$
- with replacement = $\frac{5}{4}$
- without replacement $p(x_i=1) \neq \frac{1}{6}$

See, sampling with replacement is definitely going to be given identical samples sorry identically distribution, because the distribution remains the same. But now whether it is going to be independent or not, depends on the way you construct your experiment. If you are ensuring independence that is true then it is not only identically distributed this is also going to be independent samples.

Now I am just trying to like this is I was just trying to highlight the fact here that even if you do with replacement, it may do's to identical underlying distribution can remain identical, but it may not remain independent if you are going to bring in some dependency there. So, if I do without that statement breaks, but let us say let us read it in this fashion only. But it is just like the experiments are not done in identical fashion. Even if you do replacement, but if you do not do identically sorry, if you do not do experiments in an independent fashion i.i.d things may not remain. So, that is what now we are going to follow henceforth is whenever I am saying I am doing a random sampling, I mean is it is with replacement and also it is done independently so that I get i.i.d samples.

(Refer Slide Time: 19:40)

Statistic of Random Samples

- ▶ When a sample X_1, X_2, \dots, X_n is drawn, we would be interested in some summary of values
- ▶ Any well defined summary may be expressed as a function $T(X_1, X_2, \dots, X_n)$

The random variable/vector $Y = T(X_1, X_2, \dots, X_n)$ is called statistic. The distribution of the statistic Y is called the sampling distribution of Y .

Often used statistics

- ▶ Sample mean: $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$
- ▶ Sample variance: $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$
- ▶ Sample standard deviation: $S = \sqrt{S^2}$

we will denote the observed values as \bar{x}, s^2, s , respectively.

NPTEL IEO05 Engineering Statistics Manjesh K. Hanawal 5

Statistic of Random Samples

- ▶ When a sample X_1, X_2, \dots, X_n is drawn, we would be interested in some summary of values
- ▶ Any well defined summary may be expressed as a function $T(X_1, X_2, \dots, X_n)$

The random variable/vector $Y = T(X_1, X_2, \dots, X_n)$ is called statistic. The distribution of the statistic Y is called the sampling distribution of Y .

Often used statistics

- ▶ Sample mean: $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$
- ▶ Sample variance: $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$
- ▶ Sample standard deviation: $S = \sqrt{S^2}$

we will denote the observed values as \bar{x}, s^2, s , respectively.

$\bar{X} = \bar{x}$
 $S^2 = s^2$
 $S = s$

NPTEL IEO05 Engineering Statistics Manjesh K. Hanawal 5

Now, you have samples. You have random samples. Next question is what you are going to do with that? Now with this sample, you are actually trying to extract some information. Now you want to get some summary of this. And whenever we are trying to do summary basically we are looking into some function of this T , some function of the samples, which we are going to denote it as T .

Now, let us say I am output some function on the samples and I got a value Y . Usually this we got by processing this data sample that is called statistic. Notice that here you are X_1 to X_n , now you are summarizing them by getting one random variable Y which we are through some function T that is called statistic. And whatever the distribution of resultant Y that is called sampling distribution of Y . Now, what are the possible operations I can do, you can any

operations you want, if you have n samples, the T function can be anything. But we will focus on some often used statistics, one is simple take the average of all the samples you have. And this is called sample mean.

Another thing you do something more you first compute sample means subtracted from the samples and then square them. And divided by anything you want like maybe let us say n minus 1. This we are going to call it a sample variance and denoted by S square and the square root of this sample variance we are going to call it as standard deviation. So, I am highlighting this because these are the most common statistics we are going to use, sample mean sample variance and sample standard deviation.

The realized values of the sample means, sample variance x, we are going to denote it by small letters x bar, s square and s. So, notice that here all I am just representing the samples as random variables, not the particular realizations, but when you plug in the particular realization whatever the value you got, you are going to, let us see, is this sample mean here x bar is a random quantity?

Student: Yes.

Professor: Yes, because that is gotten by averaging random variables. So, this x bar can take some particular value, x bar small x bar, and this s square is also random variable, whatever the realization, it can take, we are going to denote and similarly, s square can also take some particular small value realization or we are going to denote it like s.

(Refer Slide Time: 22:59)

Properties of statistics \bar{X} and S^2

X_1, X_2, \dots, X_n is random sample from a population with mean μ and variance σ^2

▶ $\mathbb{E}(\bar{X}) = \mu$

$$\mathbb{E}\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}(X_i) = \mu$$

▶ $\text{Var}(\bar{X}) = \sigma^2/n$

$$\begin{aligned} \text{Var}(\bar{X}) &= \text{Cov}(\bar{X}, \bar{X}) = \text{Cov}\left(\frac{1}{n} \sum_{i=1}^n X_i, \frac{1}{n} \sum_{j=1}^n X_j\right) \\ &= \mathbb{E}\left(\left(\frac{1}{n} \sum_{i=1}^n (X_i - \mu)\right)\left(\frac{1}{n} \sum_{j=1}^n (X_j - \mu)\right)\right) \quad \text{iid} \\ &= \frac{1}{n^2} \sum_{i=1}^n \mathbb{E}((X_i - \mu)^2) = \frac{\sigma^2}{n} \end{aligned}$$

NPTEL IIT Bombay IIT Madras IIT Kharagpur IIT Gandhinagar IIT Roorkee IIT Guwahati IIT Jammu IIT Patna IIT Varanasi IIT BHU Varanasi IIT Kanpur IIT Roorkee IIT Guwahati IIT Jammu IIT Patna IIT Varanasi IIT BHU Varanasi IIT Kanpur

Manjesh K. Hanawal 6

Now, let us look into some other properties of this sample means statistics and sample variance we have. Let us say these X_1, X_2, X_n they are random samples coming from a population with mean, μ and variance σ^2 . The expected value of this sample mean is always going to be the underlying mean value of that population. So, now, let us say why is that so, \bar{X} is nothing but this quantity. You should take the expectation of this. This is nothing but average of this expectation. Is this step correct from this to this?

Student: Yes.

Professor: Why?

Student: Linearity operation.

Professor: Linearity operation and now, expectation of X_i is all μ . And if I add μ n times and divide by n this is going to be simply μ . Now, similarly, you can also compute the variance of the sample mean. Actually, it is always easier to compute variance thinking this as a covariance, like we already discussed that variance of X is nothing but covariance of that random variable with itself. If you just use this and expand and after doing some simplification, and using the fact that these are all i.i.d samples, you will see that this is σ^2 squared by n .

(Refer Slide Time: 24:38)

The slide displays the following derivation for the variance of the sample mean \bar{X} :

$$\begin{aligned} \mathbb{E}(S^2) &= \mathbb{E}\left(\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2\right) \\ &= \frac{1}{n-1} \mathbb{E}\left(\sum_{i=1}^n (X_i + \mu - \mu - \bar{X})^2\right) \\ &= \frac{1}{n-1} \mathbb{E}\left(\sum_i ((X_i - \mu)^2 + (\bar{X} - \mu)^2 - 2(X_i - \mu)(\bar{X} - \mu))\right) \\ &= \frac{1}{n-1} \left(\sum_i \text{Var}(X_i) + \sum_i \text{Var}(\bar{X}) - \frac{2}{n} \sum_i \mathbb{E}((X_i - \mu)^2)\right) \\ &= \frac{1}{n-1} \left(n\sigma^2 + n\frac{\sigma^2}{n} - \frac{2}{n}n\sigma^2\right) = \frac{1}{n-1}(n\sigma^2 - \sigma^2) = \sigma^2 \end{aligned}$$

Handwritten red notes on the slide include $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ and a circled σ^2 at the end of the final equation.

Bullet points at the bottom of the slide:

- ▶ $\mathbb{E}(\bar{X}) = \mu$: Statistic \bar{X} is unbiased estimator of μ
- ▶ $\mathbb{E}(S^2) = \sigma^2$: Statistic S^2 is unbiased estimator of σ^2

Page number: 7

Statistic of Random Samples

- ▶ When a sample X_1, X_2, \dots, X_n is drawn, we would be interested in some summary of values
- ▶ Any well defined summary may be expressed as a function $T(X_1, X_2, \dots, X_n)$

The random variable/vector $Y = T(X_1, X_2, \dots, X_n)$ is called statistic. The distribution of the statistic Y is called the sampling distribution of Y .

Often used statistics

- ✓ ▶ Sample mean: $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$
- ✓ ▶ Sample variance: $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$
- ✓ ▶ Sample standard deviation: $S = \sqrt{S^2}$

$\bar{X} = \bar{x}$
 $S^2 = s^2$
 $S = s$

we will denote the observed values as \bar{x}, s^2, s , respectively.



Manjesh K. Hanawal


Now, let us look into the next one sample variance. If you look into the sample variance, and again compute its expected value, this needs to be done certain manipulation because \bar{X} here is nothing but average of so, \bar{X} here is nothing but X_i by n . So, if you plug in that, appropriately, you will see that and whatever the definition we will get, we are going to see that this value is exactly equals to sigma square.

So, notice that if you people noticed, when I took a sample mean, I divided by n , but when I took sample variance, I did not divide it by n even though I did n terms here, I divided it by n minus 1. Only when I did that n minus 1, you will see that expected value is going to be sigma square. Had I not divided by n , this would have not been sigma square. This will be some extra factor would have come here.

So, now the definition is, whenever expected value of this sample mean is μ that is the underlying mean value, this statistics \bar{X} is called unbiased estimator of μ . Notice that now I am calling \bar{X} as an estimator. And whenever this expectation of \bar{X} is μ , we are calling it as unbiased estimator. And similarly, expectation of sample variance is also sigma square. And in this case, the statistics S^2 is called unbiased estimator of your variance term. So, sample mean and sample variance, they are the estimators for mean and variance further they are unbiased. So, let us stop here.