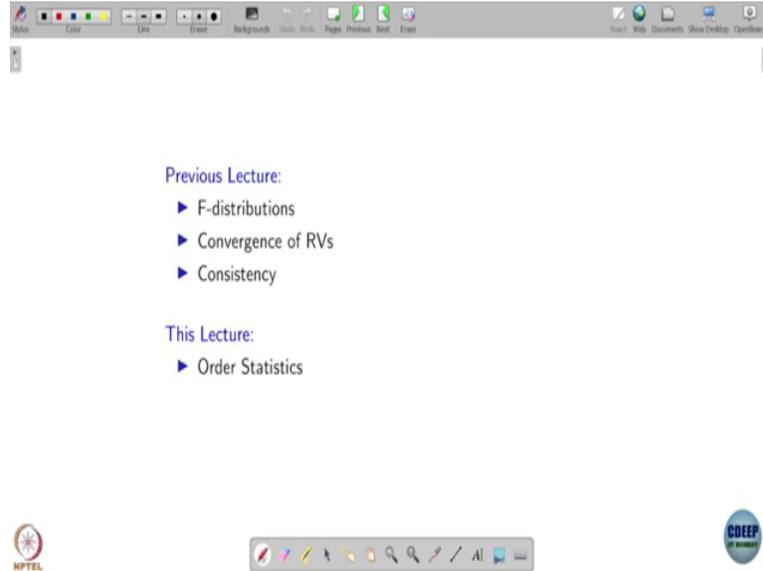


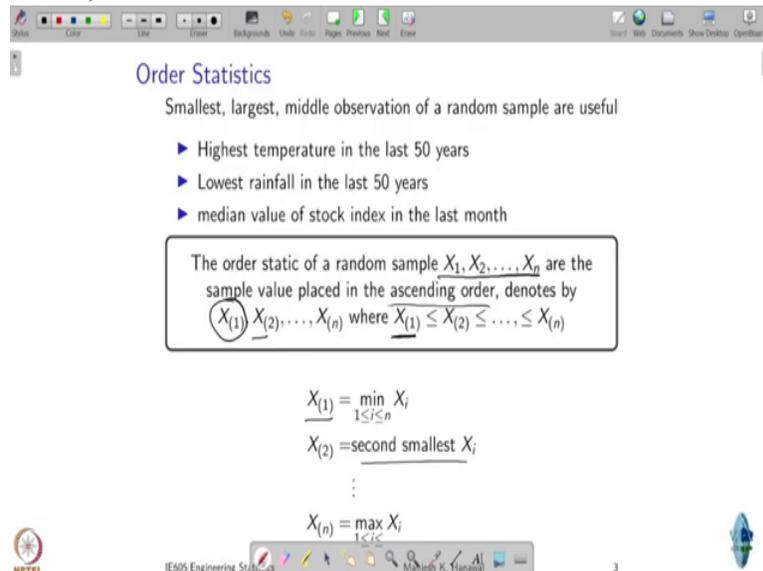
Engineering Statistics
Professor Manjesh Hanawal
Industrial Engineering and Operation Research
Indian Institute of Technology Bombay
Lecture 26
Order Statistics, Median and Percentiles

(Refer Slide Time: 00:15)



Now, I want to talk about something called order statistics.

(Refer Slide Time: 00:19)



So, as the name indicates, here I am interested in ordering of the samples. And look into the statistics of this ordered sample. So, where I should worry about ordering? So, maybe I will be

offered, if you just give me samples, I may want to find out which is the smallest, which is the largest, which is the middle of them. If I have to do this, I need to order them.

If you are just to give me samples, and if I want to find which is the smallest, maybe I want to order them increasing or decreasing way and from that I will find out which is the smallest and which is the largest of them. And this is also important in many applications, for example, which is the hot day in the last 50 years or, which year got the lowest rainfall. So, you may have data like a health, this metrological department may have recorded this data of each day of each month like that, that is like a samples for you.

And from that you want to identify which was the smallest. That identifying the smallest is easy, but what we want to do is understand its distributions. Now, because now, we have to order we will use now one notation to put our samples into orders. What we will do is let us say I have n samples given to me I want to put them in an ascending order. What I will do is now I am going to take the smallest among X_1, X_2, \dots, X_n as $X_{(1)}$.

Notice that, now I would deliberately put a parentheses on 1. So, this $X_{(1)}$ is minimum of all my samples and the next $X_{(2)}$ is the second smallest among all my samples and that I have calling it as $X_{(2)}$ and like that. And naturally this $X_{(1)}$ with the superscript has to be smaller than $X_{(2)}$ in parentheses like that. So, I have n samples, I have simply took the first smallest one put it first, took the second smallest one put it second like that.

(Refer Slide Time: 03:12)

The slide is titled "Sample mean vs Sample Median". It contains the following content:

- ▶ Sample range: $X_{(n)} - X_{(1)}$
- ▶ Sample median:

$$M = \begin{cases} X_{((n+1)/2)} & \text{if } n \text{ is odd} \\ (X_{(n/2)} + X_{((n/2)+1)})/2 & \text{if } n \text{ is even} \end{cases}$$
- ▶ Example:
 - Random sample: 24, 89, 59, 34, 55, 81, 45, 93, 85, 50
 - Order statistic: 24, 34, 45, 50, 55, 59, 81, 85, 89, 93
 - Sample range: $93 - 24 = 69$
 - Sample mean: 61.5
 - Median: 57
- ▶ Median gives better indication of "typical" values than means!

At the bottom of the slide, there are logos for NPTEL and CDEEP, and the text "IE605 Engineering Statistics" and "M. Jyoti K. Manoj" are visible.

The screenshot shows a presentation slide with the following content:

Order Statistics

Smallest, largest, middle observation of a random sample are useful

- ▶ Highest temperature in the last 50 years
- ▶ Lowest rainfall in the last 50 years
- ▶ median value of stock index in the last month

The order static of a random sample X_1, X_2, \dots, X_n are the sample value placed in the ascending order, denotes by $X_{(1)}, X_{(2)}, \dots, X_{(n)}$ where $X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)}$

$$X_{(1)} = \min_{1 \leq i \leq n} X_i$$

$$X_{(2)} = \text{second smallest } X_i$$

$$\vdots$$

$$X_{(n)} = \max_{1 \leq i \leq n} X_i$$

At the bottom of the slide, there is a navigation bar with a search icon and the text 'Mahesh K. Yanamala'.

Now, putting them in this way will already help me find some information about the data. First thing is sample range or what range my values are being taken. So, if I take my largest value and smallest value, it will give me indication of the range that is right. Largest is this, smallest is this matlab my samples are taking value on this range.

Next something called median. So, suppose let us take n, my n is odd number. So, n plus 1 by 2 is, is going to be an even number. But n plus 1 by 2 is going to be exactly the middle, exactly the middle in my sequence. And that is what is going to be called as median. Like when I in, order my samples in increasing order and one at the middle we are going to call it as median.

And when n is even n plus 1 by 2, I cannot get an integer. So, for that what we will do is we will take the average of the two middle terms. So, when n is even I can take two middle one and then take the average, that is what I am going to call it as median.

So, in a way median gives me information about what is the middle value of my samples, at 50 percent, what is the 50 percent value of my samples. So, I have an example. I will just put some numbers here some random samples, ten samples I have here. And their order statistics, I will just put them in an increasing order, I have just put them in increasing order. And from this it is easy to compute what is the sample range. So, 93 minus 24 and sample median, sample median. So, how did I get 61.5 here? Oh, no, this is sample mean. So, how to get this sample mean?

Student: Add on.

Professor: Just add this 10 numbers and divided by 10 that is 61.5. And what is the median here? So, to apply find out the median I just to first find out whether I have odd number of samples or even number of samples. In this case, I have even number of samples.

So, I am in this category, I will just take the two of the middle ones and their average is going to be 57. And that is why I got median. And usually, this median gives better indication of the typical values of my samples rather than the mean. So, mean here is giving you 61.5 whereas the median is giving you 57. Now, our what we are saying is if I have to take the typical value of this random samples I have, we are going to, we are saying that median is a better representation of that rather than mean value. So, why is that?

Student: Because the last four terms are different than in term (07:28) (07:37) that is why median is the right value because it gives a, it given an information about the central point, whereas the mean (07:50).

Professor: So, is there a case where I should go with mean, rather than taking median? And is there a case where I should prefer median over mean?

Student: In case where the data is not (08:06) you can use in.

Professor: Kyu right like? Let us say, I am interested in the average, I am interested in the performance of this class. And I know that some of you are going to do very well, close to 100 percent. So, you are going to boost the score of this class? And if somebody is going to just ask me tell me, what is the average, then I will just tell the average even though some people failed, but average is very high.

Student: Sir, that is why, that is the reason we cannot gather solution of majority of the class because some people (08:43) majority lies, understood these four in the whole scheme.

Professor: So, for that, why should I go with median?

Student: Because that gives a more, accurate idea about the most of the trans. That is why, so we should go with the mean, because the mean will give the false (09:06) find that (9:08)now.

Professor: But let us suppose now. So, maybe like that is what like if I want to just to make a binary decision, if somebody who is getting above mean value pass, below mean is fail, then you will fail?

And if I say median then you will, you will fail or pass? So, that is how like come in if I want to that is what like in a classes like this, where I have to look into the overall. I mean, where I care about majority, not like only few, maybe median could be a better score. But on the other hand, let us take an another example. Instead this class let us say I put money, I may stock, I am interested in investing in stocks, and I invest in three, four stocks. All I care is on an average how much money I make. In that I should prefer mean or median?

Student: Mean.

Professor: In that case mean reward, I am interested in just reward like how the average, I will just put everybody I do not care like how they are performing individually, what matters is together how much return they give me. In that case, maybe I will go with mean. In this case, like where I have to worry about social things like, I mean, whether most of them will pass or not, I have to worry such kind of maybe median and I will go with.

So, depending on your application, mean may be suitable or median and may be suitable. So, you can choose as per your applications. But why to focus on 50 percent? We can focus on anything you want. So, that is where the sample percentile comes into picture.

(Refer Slide Time: 11:11)

Sample Percentile

For any $p \in [0, 1]$, the $(100p)$ th percentile is the observation such that approximately np of the observations are less than this observation and $n(1-p)$ of the observation are greater.

- ▶ For $p = 0.5$, 50th percentile gives median
- ▶ For any $b \in \mathbb{R}_+$, define

$$\{b\} = \begin{cases} \lfloor b \rfloor & \text{if } \lceil b \rceil \leq b + 0.5 \\ \lfloor b \rfloor & \text{if } b - 0.5 < \lceil b \rceil \end{cases}$$

Handwritten examples:

- $b = 5.6$ → $\{b\} = 6$
- $b = 5.4$ → $\{b\} = 5$
- $b = 5$ → $\{b\} = 5$

Additional handwritten note: $\frac{1}{2} < np < n - \frac{1}{2} \Rightarrow \frac{1}{2n} < p < 1 - \frac{1}{2n}$

Now, for any p between 0, 1, 100 pth percentile is observation such that approximately, np of the observations are less than this observation and remaining n 1 minus p of the observations are greater. And it is just like a generalization of the median we focused about. Median actually

looked into the 50th percentile. But I should, I may be worried about 40 percentile, 30 percentile, 99 percentile, or 98th percentile, and all of that.

So, that is what this is just like a generalization of that. And to just make this more formal, just like additional definition, we are just introducing this another term floor bracket b , this is going to take the value of c if it is b plus 0.5 . So, we will just it is an indication that it is like suppose, let us apply a definition. Let us take b equals to some 5.6 . So, what will be the value of b ? As per this definition? What is a b plus 0.5 ? B plus 0.5 is?

Student: 6.1 .

Professor: 6.1 , and b minus 0.5 is 5.1 . And what is a floor of b ? No sorry, ceil, ceil is going to be 6 , is 6 less than 6.1 ?

Student: Yes.

Professor: So, I am going to take this to be?

Student: 6 .

Professor: 6 , so, it is just like if it is more close, this is more closer towards 6 . So, I will take it as 6 and if this is 6 and if this is 0.4 , then I would have taken b to be 5 , it is simple notation for that. And now, I want to identify suppose you have some p given and you have n samples. If np has to be between half and n minus half then you can find out that that P is simply $\frac{1}{2n} \frac{1}{1 - 1/2n}$. If you want the p if you are np to be between half and 1 minus half. Let us say I do not know why we are all using this.

(Refer Slide Time: 14:07)

Lower and Upper Quartile

$$(100p)\text{th sample percentile is } = \begin{cases} X_{(\lfloor np \rfloor)} & \text{if } p < 0.5 \\ X_{(n+1-\lfloor n(1-p) \rfloor)} & \text{if } p > 0.5 \end{cases}$$

Example 1: $n = 50, p = .35, np = 17.5, \lfloor np \rfloor = 18$. 35th sample percentile is $X_{(18)}$

Example 2: $n = 50, p = .65, n(1-p) = 17.5, \lfloor n(1-p) \rfloor = 18$
 $n + 1 - \lfloor n(1-p) \rfloor = 50 + 1 - 18 = 33$. 65th sample percentile is $X_{(33)}$

- ▶ For $p < 0.5$ and $p > 0.5$ sample percentiles exhibit symmetry
- ▶ if $(100p)$ th sample percentile is i th smallest observation, then $100(1-p)$ th sample percentile is the i th largest observation
- ▶ 25th sample percentile is called lower quartile
- ▶ 75th sample percentile is called upper quartile

$p = 0.2$
 0.8
 $p = 0.75$
 $p = 0.75$

Sample Percentile

For any $p \in [0, 1]$, the $(100p)$ th percentile is the observation such that approximately np of the observations are less than this observation and $n(1-p)$ of the observations are greater.

- ▶ For $p = 0.5$, 50th percentile gives median
- ▶ For any $b \in \mathbb{R}_+$, define

$$\{b\} = \begin{cases} \lfloor b \rfloor & \text{if } \lfloor b \rfloor \leq b + 0.5 \\ \lfloor b \rfloor & \text{if } b - 0.5 < \lfloor b \rfloor \end{cases}$$

$\{b\} = 6$

 $b = 5.4$
 $\{b\} = 5$

▶ $\frac{1}{2} < np < n - \frac{1}{2} \implies \frac{1}{2n} < p < 1 - \frac{1}{2n}$

Now using this, if I am interested in the 100th p sample, our percentile or pth percentile, it depends on whether p is between 0.5 or less than 0.5 or greater than 0.5. If your p is less than 0.5, p over value is going to be exactly this order statistics. And you notice already I have defined what is flower bracket of np? It is approximation to the next immediate integer. And if it is greater than 0.5, again I have to go with this adjustments, to make sure that it, it is more towards the right hand side.

And if it is p is less than 0.5 it is more towards left hand side and if p is greater than 0.5 more towards the right hand side. So, I have to make up for great corrections here. So, this is just a

definitions formal way of putting it. But you I hope you understand how to compute this 100p percentile now.

Now just simply put p equals to, now, this is for $p = 0.5$ is strictly less than 0.5 and p greater than 0.5. And what how to find it for $p = 0.5$? Just use the definition of median, median is exactly for $p = 0.5$. And here I have just given one example. Suppose, n equals to 15 and you want to compute 35th percentile, then first we will compute what is np . And if you compute this floor bracket of np even though notice that it is in the middle, it is going to take the next 18, next value 17.5 is approximated to 18. And the 35th sample is simply x of 18th. The 18th sample when you have ordered them 1 to 50.

And similarly, you can also compute what is the 65th percentile of this sample. In this case It so, happens that 65th percentile is going to be x of 33, 33rd element in my order statistics. So, this $p < 0.5$, $p > 0.5$ samples percentiles exhibit symmetry, I hope you understand what is symmetry. By symmetry I mean suppose you take 0.2 and another value 0.8. So, 0.2 sample is certain numbers away from your left side. The same number of samples it will be away from right hand side for the 80th percentile.

So, if you just take this or let us say 1, 2, 3, up to let us say some n and this is some here, if this is at 80th percent and it is m sample away and if it is a 20th percentile I mean at this point that means that will be also m sample away. In that way they exhibit symmetry.

And if you want to generalize it, if suppose, for a given p 100pth percentile happens to the i th smallest observation like 100pth sample happens to be the i th, let us say if you have this and this is like 100pth sample percentile sample and this happens to be the i th smallest. Then if you look into 100 into $1 - p$ here, assuming that the p is less than 0.5 here then this is going to be what? This is the i th largest element. In that way, they are also showing some symmetry here.

And now the special cases for us are something called lower quartile and upper quartile. So, when you are going to take p equals to 0.25, that number is called lower quartile. And when you take p equals to 0.75 this is called upper quartile. And notice that by our definition, they are symmetric. If p equals to 0.25 corresponds to the i th smallest, then p equals to 0.75 corresponds to the i th largest element in my order statistics.

And after maybe we will see later something called boxplot. When you have given a bunch of samples, you want to see that what samples lies within 25 times 75th percentile. And that means there is the ones where majority of the samples I am interested in. Maybe others are either too low or too large. So, I may be interested in just between something between 25 to 75. So, that is why this lower quartile and upper quartile comes. And maybe in Python, there will be already ready function. If you give a bunch of samples and tell identify 25th quartile, or 75th quartile, they will immediately find the sample value for you.

(Refer Slide Time: 20:25)

Distribution of Order Statistics

$X \in \{x_1, x_2, \dots, x_n\}$
 p_1, p_2, \dots, p_n

Discrete Case:
Random sample X_1, X_2, \dots, X_n come from a discrete distributions with pmf $P_X(x_i) = p_i$, where $x_1 < x_2 < \dots$ are the possible realizations in ascending order. For any x_i , what is $P(X_{(j)} \leq x_i)$? $X_{(j)} = j^{\text{th}}$ smallest value

$P_0 = 0$
 $P(X \leq x_1) = P_1 = p_1$
 $P(X \leq x_2) = P_2 = p_1 + p_2$
 \vdots
 $P(X \leq x_i) = P_i = p_1 + p_2 + \dots + p_i$
 \vdots

CDF of j^{th} order statistic $P(X_{(j)} \leq x_i)$

Order Statistics

Smallest, largest, middle observation of a random sample are useful

- ▶ Highest temperature in the last 50 years
- ▶ Lowest rainfall in the last 50 years
- ▶ median value of stock index in the last month

The order static of a random sample X_1, X_2, \dots, X_n are the sample value placed in the ascending order, denotes by $X_{(1)}, X_{(2)}, \dots, X_{(n)}$ where $X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)}$

$\checkmark X_{(1)} = \min_{1 \leq i \leq n} X_i$
 $\checkmark X_{(2)} = \text{second smallest } X_i$
 \vdots
 $\checkmark X_{(n)} = \max_{1 \leq i \leq n} X_i$

In the next five minutes, we will just discuss how to find the distribution of this order statistics. Now let us focus on this some calculations. So, what we want to now do is suppose, consider my discrete random variable. And I am interested in n samples which are coming from the same population mean. And assume that those are going to follow some common probability mass function.

So, all my X are taking value $X_1 X_2$, let us say some X_n value n samples they are taking. And I am also assuming that this value is taken by my random variables are themselves ordered. So, this is the smallest value, it can take this is the next smallest value and like this. Now I am interested in what is the probability that my j th order statistics is going to be taking value less than or equal to X_i . So, I hope you understand it what I am trying to say here.

So, X , all my $X_1 X_2 X_3$, they are taking values realization $X_1 X_2 X_3$, and they it taking that probability that X is taking value X_i , I am going to denote it as P_i . Now, by definition, X of j is j th smallest value. Now I am interested in finding this j th smallest value is less than or equals to some particular value X_i . This X_i is one of these elements in this. Now what I am basically trying to find is I am trying to find cdf of j th order statistics.

So, notice that here, this is called first order statistics, second order statistics. And this is the n th order statistics. And now I am interested in finding the distribution of j th order statistics, how to do that. To do that, I am going to start taking P_0 equals to 0, I am start doing iteration. And now, probability that X is going to be less than or equal to X_1 is P_1 , everybody agree with this? So, notice that probability that $X \leq X_1$ is P_1 ?

So, maybe I will just write this. So, these are ordered. And I am going to take this value with probability P_1 , this value with probability P_1, P_2 , and this value equals to probability P_2 . So, it is like $X_1 X_2$, like X_1 , this is like P_1 and this is like P_2 like this. Just try to follow. So, probability X is less than or equal to X_1 , nothing is there below. So, I takes 1 it is going to get a jump of P_1 . And this is like a discrete case we have.

A probability let us tell X , X is less than or X_2 this point and this point here, this is P_1 plus P_2 , and probability that X is going to be less than equal to to till X_i is like if you are going to take this as X_i . So, this is going to be all the way P_1 plus P_2 all the way up to P_i .

(Refer Slide Time: 24:39)

Discrete Case contd..

$x_i \in \{x_1, x_2, \dots, x_n\}$ $Y_j = \mathbb{1}_{\{X_j \leq x_i\}}$
 $P(Y_j = 1) = P(X_j \leq x_i)$

Binary

- ▶ Fix some x_i . Define $Y_j = \mathbb{1}_{\{X_j \leq x_i\}}$ for all $j = 1, 2, \dots, n$
- ▶ $P(Y_j = 1) = P_i$ for all $j = 1, 2, \dots, n$
- ▶ $Y = \sum_{j=1}^n Y_j$, $Y \in \{0, 1, 2, \dots, n\}$
- ▶ As X_j s are i.i.d, Y_j s are i.i.d. $Y_j \sim \text{Ber}(P_i)$.
- ▶ $Y \sim \text{Bin}(n, P_i)$. Y is sum of n $\text{Ber}(P_i)$ RVs
- ▶ $\{X_{(j)} \leq x_i\} = \{Y \geq j\}$. Hence $P(X_{(j)} \leq x_i) = P(Y \geq j)$

So, now, take some X_i . So, X_i is from one of my outcomes X_1 X_2 up to X_n . And now Y_j is indicator function that my j th random variable is less than or equal to X_i , I hope all of you understand this indicator. So, what are the possible values of Y_j ?

Student: 0, 1.

Professor: 0, 1 and I am doing it for my all end random variables which I have X_2 j equals to 1 to n . Now, what is the probability that Y_j equals to 1. So, now, since Y_j is indicator that X_j is less than or equals to X_i , now probability that y_j is equals to 1 what is this probability? Probability that X_j is less than or equal to X_i agree? So, this is going to happen when only when this event is happening X_j is less than or equal to X_i .

Now, let us focus on this Y , which is now defined as summation of this Y_j s. Y_j is a binary random variable. It is going to take either 1 or 0. Now, then what is capital Y is going to take values? It is sum of n ?

Student: Bernoulli random variables.

Professor: Bernoulli random variables, so that is going to take either 0 or 1 all the way up to n . Now, if all these X_i are i i d you can if X_i is iid you can say notice that Y_j s are also iid. And I just argued that Y_j is nothing but Bernoulli with value P_i . Now, this Y here is sum of Bernoulli each with the same parameter P_i , because I have fixed this X_i . So, now, because of this since y is

sum of n Bernoulli random variables with parameter p_i , I know that Y is binomial with parameter?

Student: $n p_i$.

Professor: $n p_i$. Now, this is true. Now, this is the last one crucial step. Suppose, X_j that is my j th order statistics is less than or equal to X_i , the claim is Y has to be greater than or equal to j that means, X_j , j th order statistics is less than or equals to X_i . That means, at least in this sum Y there has to be j once that is why Y is going to be greater than or equals to j .

And then, if Y has to be greater or equals to j , you can verify that that actually means that X of j has to be less than or equal to X , you just verify that. I mean you need to think about this, these two events are same, if that is the case, probability that X_j j th order statistics is being less than or equal to X is nothing but probability that my Bernoulli random variable is greater than or equal to i .

(Refer Slide Time: 28:27)

Discrete Case contd..

$$P(Y \geq j) = \sum_{k=j}^n \binom{n}{k} p_i^k (1-p_i)^{n-k}$$

$$P(X_{(j)} \leq x_i) = \sum_{k=j}^n \binom{n}{k} p_i^k (1-p_i)^{n-k}$$

$$P(X_{(j)} = x_i) = P(X_{(j)} \leq x_i) - P(X_{(j)} \leq x_{i-1})$$

$$= \sum_{k=j}^n \binom{n}{k} (p_i^k (1-p_i)^{n-k} - p_{i-1}^k (1-p_{i-1})^{n-k})$$

NPTEL IE605 Engineering Statistics Mahesh K. Manasa CDEEP

$$x_i \in \{x_1, x_2, \dots, x_n\}$$

$$Y_j = \mathbb{1}_{\{X_j \leq x_i\}}$$

$$P(Y_j = 1) = P(X_j \leq x_i)$$

- ▶ Fix some x_i . Define $Y_j = \mathbb{1}_{\{X_j \leq x_i\}}$ for all $j = 1, 2, \dots, n$
- ▶ $P(Y_j = 1) = P_i$ for all $j = 1, 2, \dots, n$
- ▶ $Y = \sum_{j=1}^n Y_j$, $Y \in \{0, 1, 2, \dots, n\}$
- ▶ As X_j s are i.i.d, Y_j s are i.i.d. $Y_j \sim \text{Ber}(P_i)$.
- ▶ $Y \sim \text{Bin}(n, P_i)$. Y is sum of n $\text{Ber}(P_i)$ RVs
- ▶ $\{X_{(j)} \leq x_i\} = \{Y \geq j\}$. Hence $P(X_{(j)} \leq x_i) = P(Y \geq j)$



Now, you know, what is why is Bernoulli is greater than or equal to j you know how to compute its probability. And now, X_j is a discrete now, I am trying to find what is the probability that X_j is equals to i. I know this is nothing but X_j is less than or equal to X_i minus X_j is going to be less than X_i minus 1, this is again because the discrete one.

And now, using this relation, like I know that this is nothing but in terms of probability Y greater than I, I use this property and now I can compute this, this and plugin to get what is the probability that my jth order statistics is X_i ? Verify this I know that you need to peacefully sit and read this and you will see this all this calculations come out. Fine, let us stop it.