



IIT Madras

ONLINE DEGREE

Statistics for Data Science - 1
Prof. Usha Mohan
Department of Management Studies
Indian Institute of Technology, Madras

Lecture – 3.5

(Refer Slide Time: 00:14)

The next important thing which we are going to discuss is the notion of what we call a Percentile. These days with a lot of competitive examinations, most of these course and competitive examinations are reported as percentiles. Percentile is different from percentage.

What is a percentile? The sample 100 percent percentile is that data value that has the property that at least 100p percent of the data are less than or equal to it and at least 100(1-p) percent of the data are greater than or equal to it.

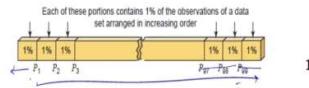
So, now if $p = 1/2$, $100p = 50\%$ of the data, $(1 - p) = 1/2$, $100(1 - p) = 50\%$ of the data and we have already seen a measure which says that 50 percent of the data is less than the value and 50 percent of the data is greater than the value and we call this as the median of a dataset. So, the $100 * \frac{1}{2}$ or the 50th percentile is the median of the dataset.

(Refer Slide Time: 01:59)



Percentiles

- The sample $100p$ percentile is that data value having the property that at least $100p$ percent of the data are less than or equal to it and at least $100(1 - p)$ percent of the data values are greater than or equal to it.



¹Figure source: Mann, P. S. (2007). Introductory statistics. John Wiley & Sons

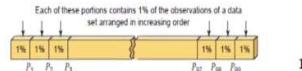
So, to demonstrate it using a figure, you can see that 99th percentile would have $100 * 0.99$ that is 99% of the data is less than it, but 1% is greater than it. Similarly, P_1 says 1% is less than it whereas, 99% is greater than or equal to it. So, the concept of the percentile tells us that value in the dataset below which I have $100 * p$ which are less than or equal and $100 * (1 - p)$ which are greater than or equal to it .

(Refer Slide Time: 02:38)



Percentiles

- The sample $100p$ percentile is that data value having the property that at least $100p$ percent of the data are less than or equal to it and at least $100(1 - p)$ percent of the data values are greater than or equal to it.



¹Figure source: Mann, P. S. (2007). Introductory statistics. John Wiley & Sons

Now, if two data values satisfy the same condition, it is the arithmetic average of these values we have already seen how to compute the median when I have odd or even number of observations.

(Refer Slide Time: 02:55)

The navigation bar at the top left shows the following path: Statistics for Data Science -1 → Numerical summaries → Percentiles. The title 'Computing Percentile' is centered below the path. To the right of the title is the logo of Anna University, which features a lamp and the text 'ANNA UNIVERSITY' and 'TECHNOLOGIES'.

To find the sample $100p$ percentile of a data set of size n

- /1. Arrange the data in increasing order.
2. If np is not an integer, determine the smallest integer greater than np . The data value in that position is the sample $100p$ percentile.



How do I compute a percentile? There are many algorithms to compute a percentile. I now present with a very simple algorithm to computer a percentile which is also very commonly used algorithm.

So, suppose I have a data of size n that is I have n observation, I want to find out what is my sample $100p$ percentile, what I do is I arrange the data in ascending order that is my first step similar to what we did when we computed a median. Once, I arrange my data in ascending order, I find out what $n * p$ is.

So, now let me give you the analogy of what we did with respect to the median. So, I arrange the data in ascending order, this is something which I did for the median also, remember $p = \frac{1}{2}$ so, I check what is $n/2$. If n is even, I know $n/2$ would be an integer. If n is odd, I know $n/2$ is not an integer. So, if np is not an integer, I determine the smallest integer greater than or equal to np that is the value, in that position is the sample $100p$ percentile.

(Refer Slide Time: 04:21)

Statistics for Data Science -1
└ Numerical summaries
└ Percentiles

Computing Percentile

$n=5 \quad p=1/2 \quad np = \frac{5}{2} = 2.5 = 3$

To find the sample 100p percentile of a data set of size n

1. Arrange the data in increasing order.
2. If np is not an integer, determine the smallest integer greater than np . The data value in that position is the sample 100p percentile.
3. If np is an integer, then the average of the values in positions np and $np + 1$ is the sample 100p percentile.

$n=6 \quad p=1/2 \quad np=3$

If np is an integer, then I look at the average of these values in positions np and $np + 1$. So, what did I say if n/np is not an integer for example, if I look at $n/2$, suppose $n = 5$, my $n/2$ is 2.5, the smallest integer greater than this is 3 so, you saw that in a dataset of 5 points, the third value would give me my median. So, that corresponds to my algorithm here.

I repeat, if np is not an integer, let $n = 5$, $p = 1/2$, $n * p = 5 / 2$ which is 2.5 it is not an integer, the smallest integer greater than this is 3. So, the value or the data which is in the 3rd position because I am already arranging my data in ascending order will give me the sample $100 * p$ or the 50th percentile or the median.

But if np is an integer. So, now, let $n = 6$ and $p = 1/2$, my np is an integer which = 3, then what do I do? I look at the average of the value so, the 3rd + the 4th data value that is what we do and I have six data points x_1, x_2, x_3, x_4, x_5 and x_6 which are arranged in ascending order, then you know that $(x_3 + x_4) / 2$ is my median this is how we define. So, I do the same thing if np is an integer, I look at the average of values in positions np and $np + 1$.

(Refer Slide Time: 06:27)



Example

Let $n = 10$

- Arrange data in ascending order 35, 38, 47, 58, 61, 66, 68,
68, 70, 79
 $\leq 9 \leq 10$

p	np	
0.1	1	$(35+38)/2=36.5$
0.25	2.5	47
0.5	5	$(61+66)/2=63.5$
0.75	7.5	68
1	10	79



So, let us look at $n = 10$. I have my data. I again arrange my data in ascending order this is the same data we have been using. So, I have arranged my data, my $n = 10$. So, for $p = 0.1$, my $n * p$ is $10 * 0.1$ which is 1, it is an integer.

So, going back to your algorithm, I need to look at the value which is np and $np + 1$, the 1st value np is 35, $np + 1$ is 38. So, $(35 + 38) / 2$ is 36.5, this is my 100 so, you look at 0.25, 2.5 it is not an integer so, the 3rd value which is 47 that gives me the 25th percentile.

0.5 we have already seen 63.5, again 61 which is my 5th value 1, 2, 3, 4, 5 $(61 + 66) / 2$ that would give me 63.5. 0.75 is again a fraction so, you can see that when I have 0.75, 7.5 so, the 8th value so, this is the 1st, 2nd, 3rd, 4th, 5th, 6th, 7th and 8th, 9, 10, 8th value is 68, 68 is the 75th percentile, the 100th percentile is the maximum which is 79. So, this is how you compute percentiles and I the example shows how to compute percentile for the given dataset.

(Refer Slide Time: 08:14)

Statistics for Data Science -1
└ Numerical summaries
└ Percentiles



Computing percentile using googlesheets-PERCENTILE function

Step 1 Paste the dataset in a column.

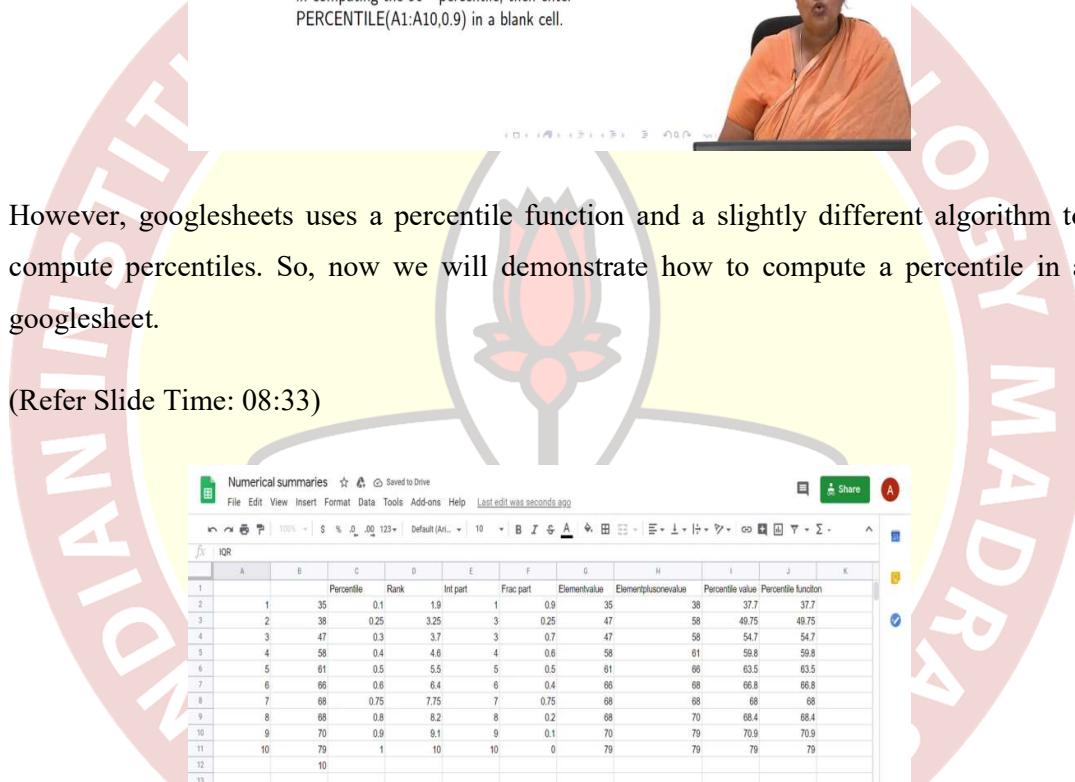
Step 2 In a blank cell enter PERCENTILE(data, percentile), where data indicates the range of data for which percentile needs to be computed, and percentile is the decimal form of the desired percentile.

- For example if the data is in cell A1:A10, and we are interested in computing the 90th percentile, then enter PERCENTILE(A1:A10,0.9) in a blank cell.



However, googlesheets uses a percentile function and a slightly different algorithm to compute percentiles. So, now we will demonstrate how to compute a percentile in a googlesheet.

(Refer Slide Time: 08:33)



A screenshot of a Google Sheets document titled "Numerical summaries". The sheet contains data for calculating quartiles and the interquartile range (IQR). The data is as follows:

	A	B	C	D	E	F	G	H	I	J	K
1		Percentile	Rank	Int part	Frac part	Elementvalue	Elementplusvalue	Percentile value	Percentile function		
2	1	35	0.1	1.9	1	0.9	35	38	37.7	37.7	
3	2	38	0.25	3.25	3	0.25	47	58	49.75	49.75	
4	3	47	0.3	3.7	3	0.7	47	58	54.7	54.7	
5	4	58	0.4	4.6	4	0.6	58	61	59.8	59.8	
6	5	61	0.5	5.5	5	0.5	61	66	63.5	63.5	
7	6	66	0.6	6.4	6	0.4	66	68	68.8	68.8	
8	7	68	0.75	7.75	7	0.75	68	68	68	68	
9	8	68	0.8	8.2	8	0.2	68	70	68.4	68.4	
10	9	70	0.9	9.1	9	0.1	70	79	70.9	70.9	
11	10	79	1	10	10	0	79	79	79	79	
12		10									
13											
14											
15											
16											
17	Q1										
18	Q1	49.75									
19	Q3	68									
20	IQR	18.25									

So, let us go back to this google sheet. I have the following dataset the same dataset which has been given in ascending order. This is the dataset which I have written the dataset in ascending order.

In a blank cell enter percentile and the data. So, I just go here, I type out what is the percentile. So, if I am looking at this dataset, I look at the percentile, so, B2 to B11 that is my dataset with C2, C2 is 0.1 so, this gives me what is a percentile so, for example, if the data we are interested in 90th percentile, I put a 0.9, if I am interested in the 10th percentile, I put a 0.1. So, you can see that 37.7, the 90th percentile I am putting it, here is 70.9, the 100th percentile is 79. So, you can see this is how we have computed or the google sheet computes it.

Again I repeat, the way google sheet you do is choose the data that is B2 to B11 that is the data which I want to compute the percentile for and C2, C2 is 0.1, 0.25 will give me the 25th percentile this gives me the 10th percentile, 54.7 will give me the 30th percentile, 0.9 gives me 70.9 is the 90th percentile and 79 is the 100th percentile.

Notice that the percentiles need not be part of the dataset. What you will notice immediately is these percentiles which google sheet percentile function gives us is different from the for the same dataset it is different.

Here I got my 10th percentile to be 36.5 whereas, google sheet gives me 37.7. So, the algorithm that google sheet uses, it is not that this is wrong and that is right, but the algorithms use as I mentioned earlier the algorithms used are different.

(Refer Slide Time: 11:04)

Statistics for Data Science - I

- └ Numerical summaries
- └ Percentiles

Computing percentile using googlesheets-algorithm

Step 1 Arrange data in increasing order.

Order	1	2	3	4	5	6	7	8	9	10
$x_{[i]}$	$x_{[1]}$	$x_{[2]}$	$x_{[3]}$	$x_{[4]}$	$x_{[5]}$	$x_{[6]}$	$x_{[7]}$	$x_{[8]}$	$x_{[9]}$	$x_{[10]}$
Data	35	38	47	58	61	66	68	68	70	79

Let $x_{[i]}$ denote the i^{th} ordered value of the dataset.

Step 2 Find rank using the following formula.

$rank = percentile \times (n - 1) + 1$ where n is total number of observations in the dataset



So, what is the algorithm a google sheet uses to compute percentile? So, first it arranges the data in ascending order. So, I have the data these are $x_1, 2, x_3$ these are the ranks the data is arranged in ascending order that is what we have done here the data is arranged in ascending order.

The second step is for any observation, x_i is denoting the i th order value, find the rank using the following formula. So, the rank is percentile * $(n - 1) + 1$ where n is the total number of observations.

(Refer Slide Time: 11:54)

Statistics for Data Science - I
└ Numerical summaries
└ Percentiles

Computing percentile using googlesheets-algorithm

Step 1 Arrange data in increasing order.

Order	1	2	3	4	5	6	7	8	9	10
$X_{[1]}$	35	38	47	58	61	66	68	68	70	79
Data	35	38	47	58	61	66	68	68	70	79

Let $x_{[i]}$ denote the i^{th} ordered value of the dataset.

Step 2 Find rank using the following formula.
 $rank = percentile \times (n - 1) + 1$ where n is total number of observations in the dataset

- Example: to compute 25 percentile of a set of $n = 10$ observations, $rank = 0.25 \times (10 - 1) + 1 = 3.25$

Step 3 Split the rank into integer part and fractional part.

- Integer part of 3.25 = 3; fractional part is 0.25.

Step 4 Compute the ordered data value $x_{[i]}$ corresponding to the integer part rank.

A photograph of a woman in an orange sari speaking into a microphone is visible in the bottom right corner of the slide.

So, let us look at it n in this case is 10, $10 - 1$ is 9 ok, percentile * $(n - 1)$. So, for example, if I want the 25th percentile, I put $0.25 * 10 - 1$ so, I have $0.25 * 9 + 1$. I get a 3.25, is it clear.

So, you can see that here what I have done is the computed the rank using that same formula. So, what is this formula? It is C2, C2 is your percentile * $(n - 1)$ which is my 9 + 1. So, this is the data. So, percentile into so, I am computing the rank of each of these datasets. So, the rank here is 1.9, 3.25 we have already demonstrated how I got this 3.25.

For each of these ranks, I have an integer part and I have a fractional part. If the rank is 1.9, the integer part is 1, the fractional part is 3.9. For 3.25, the integer part is 3, the fractional part is 0.25 ok.

So, for each one of them, I have an integer part and a fractional part. So, in the 3rd step, I split it into an integer and fractional part and that is what we have done here. For the rank, I have split it into an integer part and the fractional part. For every rank, first we compute the rank, then I split the rank into integer and fractional part.

Once that is done, I look at what is the I compute the ordered data value corresponding to the integer part rank. So, what do I mean by this? For 3.25, the integer part is 3, the fractional part is 0.25.

(Refer Slide Time: 14:05)

Statistics for Data Science - I
└ Numerical summaries
└ Percentiles

Computing percentile using googlesheets-algorithm

Step 1 Arrange data in increasing order.

Order	1	2	3	4	5	6	7	8	9	10
$x_{[i]}$	$x_{[1]}$	$x_{[2]}$	$x_{[3]}$	$x_{[4]}$	$x_{[5]}$	$x_{[6]}$	$x_{[7]}$	$x_{[8]}$	$x_{[9]}$	$x_{[10]}$
Data	35	38	47	58	61	66	68	68	70	79

Let $x_{[i]}$ denote the i^{th} ordered value of the dataset.

Step 2 Find rank using the following formula.
 $rank = percentile \times (n - 1) + 1$ where n is total number of observations in the dataset

- Example: to compute 25 percentile of a set of $n = 10$ observations, $rank = 0.25 \times (10 - 1) + 1 = 3.25$

Step 3 Split the rank into integer part and fractional part.

- Integer part of 3.25 = 3; fractional part is 0.25.

Step 4 Compute the ordered data value $x_{[i]}$ corresponding to the integer part rank.

- The ordered data value corresponding to integer part rank of 3, $x_{[3]}$ is 47.

3 | 2 25
| 0 50

A photograph of a woman in an orange sari sitting at a desk, looking towards the camera. A small video camera icon is visible in the bottom right corner of the slide.

Compute the ordered data value corresponding to the integer part rank. The integer part rank is 3 and you can see that the ordered data corresponding to this integer part rank is 47.

(Refer Slide Time: 14:26)



Computing percentile using googlesheets-algorithm

Step 1 Arrange data in increasing order.

Order	1	2	3	4	5	6	7	8	9	10
$x_{[i]}$	35	38	47	58	61	66	68	68	70	79
Data	35	38	47	58	61	66	68	68	70	79

Let $x_{[i]}$ denote the i^{th} ordered value of the dataset.

Step 2 Find rank using the following formula.

$$\text{rank} = \text{percentile} \times (n - 1) + 1 \text{ where } n \text{ is total number of observations in the dataset}$$

- Example: to compute 25 percentile of a set of $n = 10$ observations, $\text{rank} = 0.25 \times (10 - 1) + 1 = 3.25$

Step 3 Split the rank into integer part and fractional part.

- Integer part of 3.25 = 3; fractional part is 0.25.

Step 4 Compute the ordered data value $x_{[i]}$ corresponding to the integer part rank.

- The ordered data value corresponding to integer part rank of 3, $x_{[3]}$ is 47.

$$3 \frac{0.25}{\parallel} 0.25$$



Again, for this 25th percentile 3.25, integer part is 3, the fractional part is 0.25. So, what I do is I look at computing or I look at that value which is corresponding to the integer part rank which is 47 that is what we have done here ok.

So, I have this element value. So, 1, 35 the integer value corresponding to 1 is 35, corresponding to 3 is 47, 4 is 58, 5 is 61, 6 is 66, 7 is 68 again here, I have the integer part which is 8, again it is 68, with 9 it is 70 and with 10 it is 79. This is the value which corresponds to the integer. So, 1 you can see its 35, 3 it is 47 and so forth.

(Refer Slide Time: 15:22)



Computing percentile using googlesheets-algorithm-contd

Step 5 The percentile value is given by the formula

$$\text{Percentile} = x_{[i]} + \text{fractional part} \times [x_{[i+1]} - x_{[i]}]$$

- $\text{Percentile} = 47 + 0.25 \times [58 - 47] = 47 + 0.25 \times 11 = 47 + 2.75 = 49.75$

$$3 \frac{0.25}{\parallel} 0.25$$



Once that is done, in the 5th step you find out what is the percentile value. The percentile value is the x_i which is $47 + \text{the fractional part}$, remember the part 25th percentile was 3.25, the fractional part was 0.25, the integer part was 3, the value corresponding or the ordered corresponding to x_3 was 47 I take the fractional point which is 0.25 and I look at $x_i + 1$. What is $x_i + 1$? It is 58.

So, I look at $x_i + 1$ which is 58, $- x_i$ which is - 47. So, $47 + 0.25 * [58 - 47]$ is 49.75 and that is what I have here. In the first thing, I have again element value is 35, element value + 1 is 38 which is x_2 , $x_2 - x_1$ is what I have is 2, I multiply that with my fractional part that is what you do here which is 0.9. I add that and I get a percentile function which is 37.7. So, I have $33 * 0.9$ which is 2.7, $35 + 2.7$ is 37.7.

So, you can see that this column I gives us the percentile value computed using this equation whereas, column J gives the percentile value computing the percentile function of the google sheets and we can see that both of them are exactly the same.

(Refer Slide Time: 17:31)

Statistics for Data Science -1
└ Numerical summaries
└ Percentiles

Quartiles	$Q_1 \rightarrow$ First (LOWER) $Q_2 \rightarrow$ Second MEDIAN $Q_3 \rightarrow$ Third (UPPER)
Definition	The sample 25th percentile is called the first quartile. The sample 50th percentile is called the median or the second quartile. The sample 75th percentile is called the third quartile.

A photograph of a woman in an orange sari sitting at a desk, looking down at a screen.

So, why are percentiles mentioned when we talk about the measures of dispersion? As I earlier mentioned, a very important measure is called the quartile. The sample 25th percentile is called the first quartile, the 50th percentile we already know is the median or the second quartile and the sample 75th percentile is called the third quartile.

So, I have Q1, Q2 and Q3 this is referred to as the first quartile or in some books as the lower quartile, this is the third quartile or referred to as the upper quartile, this is the median or the second quartile which is already we have seen is the median ok.

(Refer Slide Time: 18:22)

The sample 25th percentile is called the first quartile. The sample 50th percentile is called the median or the second quartile. The sample 75th percentile is called the third quartile.

In other words, the quartiles break up a data set into four parts with about 25 percent of the data values being less than the first(lower) quartile, about 25 percent being between the first and second quartiles, about 25 percent being between the second and third(upper) quartiles, and about 25 percent being larger than the third quartile.

Now, when I have the quartiles, the quartiles in fact, break up a dataset into four parts. So, if I have this as my dataset, I know already the median breaks up a dataset into two parts, the first quartile.

The first quartile or the lower quartile I have a upper quartile this is Q2, this breaks up so, I have the least, this is the minimum, I have the maximum, I have the first quartile, the second quartile, the third quartile also referred to as a lower quartile and the upper quartile ok. It breaks up the dataset so, I have part 1, part 2, part 3 and part 4.

So, you can see that the quartiles break up an entire dataset into four parts. 25% of the data lie here, 50% lie here, 75% lie here and entire dataset lies between minimum and maximum.

(Refer Slide Time: 19:34)



The Five Number Summary

- ▶ Minimum ✓
- ▶ Q_1 : First Quartile or lower quartile ✓
- ▶ Q_2 : Second Quartile or Median ✓
- ▶ Q_3 : Third Quartile or upper quartile ✓
- ▶ Maximum ✓



So, why are these again an important measure? If you look at it any descriptive statistics most of the summaries are given / what we refer to as a five-number summary. The five-number summary includes the minimum, the first quartile, the median, the third quartile and the maximum as we have given here. So, this five-number summary is a very good way of summarizing a dataset.

(Refer Slide Time: 20:05)



The Interquartile Range (IQR)

$$\text{IQR} = \frac{\text{Range}}{\text{Max} - \text{Min}} = Q_3 - Q_1$$

Definition

The interquartile range, IQR, is the difference between the first and third quartiles; that is,

$$IQR = Q_3 - Q_1$$

- ▶ IQR for the example



The interquartile range is a very important measure of dispersion which you have we have already seen the range is the difference between maximum and minimum. The

interquartile range is the difference between the third quartile and the first quartile and this is referred to as the interquartile range.

(Refer Slide Time: 20:37)

Statistics for Data Science -1
└ Numerical summaries
└ Percentiles

The Interquartile Range (IQR)

Definition
The interquartile range, IQR, is the difference between the first and third quartiles; that is,

$$IQR = Q_3 - Q_1$$

- ▶ IQR for the example
 - ▶ First quartile, $Q_1 = 49.75$
 - ▶ Third quartile, $Q_3 = 68$
 - ▶ $IQR = Q_3 - Q_1 = 18.25$

So, for our example, the first quartile was 49.75, the third quartile was 68, the interquartile range is 18.25. So, the interquartile range is also a measure of dispersion.

(Refer Slide Time: 20:55)

Statistics for Data Science -1
└ Numerical summaries
└ Percentiles

Section summary

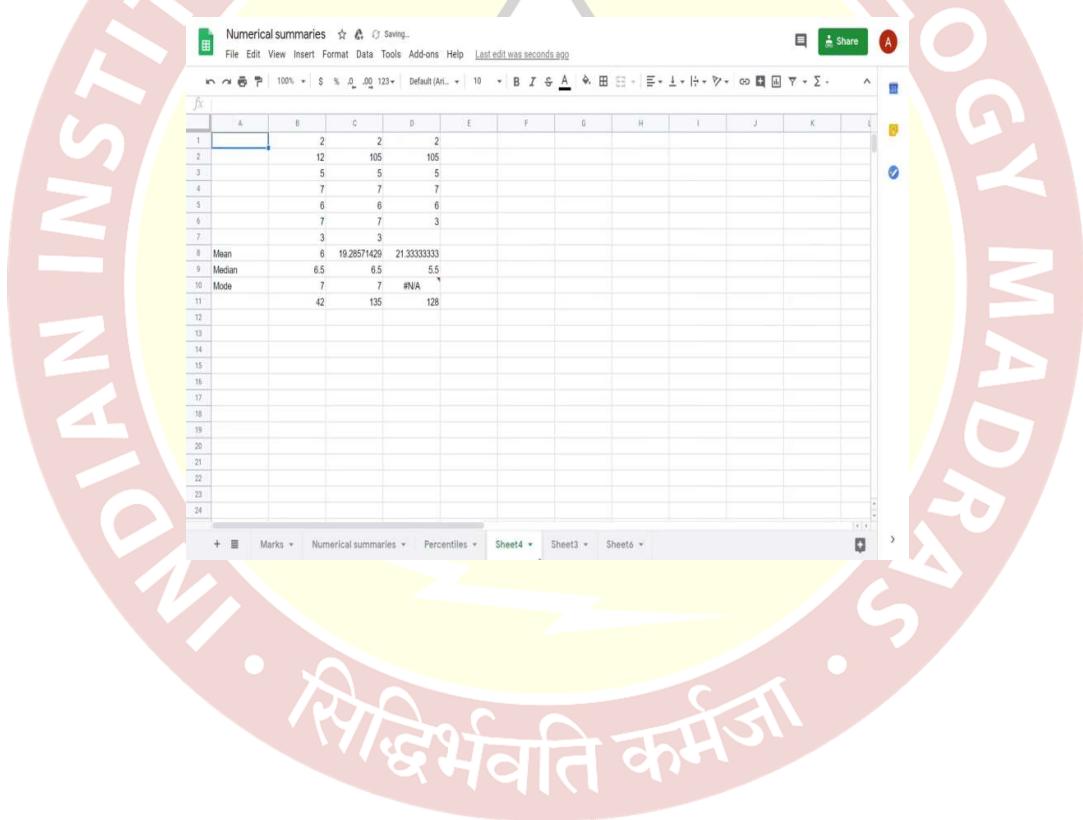
- ▶ Definition of percentiles.
- ▶ How to compute percentiles.
- ▶ Definition of quartile.
- ▶ Five-number summary.
- ▶ Interquartile range as a measure of dispersion.

So, what we have seen so far was how to define percentiles, how to compute percentile, what is the definition of a quartile, the five-number summary which is very very

important in many descriptive statistics and interquartile range which is a measure of dispersion.

The interquartile range so, this if you look at googlesheet, the quartile function with the array and 1 the lower quartile is the first quartile gives me the lower quartile. The quartile function with 3 gives me the upper quartile and the difference between the upper quartile and the lower quartile gives me the interquartile range and we can see that the quartile function if I put 2 gives me the median which you have already seen is 63.5, 3 gives me the upper quartile. So, this is my Q1, this is my Q3 and this is my IQR or my interquartile range.

(Refer Slide Time: 22:08)



(Refer Slide Time: 22:10)

Statistics for Data Science -1

- └ Numerical summaries
- └ Percentiles

Summary

1. Frequency tables
 - 1.1 Frequency table for discrete data. → class
 - 1.2 Frequency table for continuous data. → class intervals
2. Graphical summaries
 - 2.1 Histograms. → class intervals of equal length
 - 2.2 Stem-and-leaf plot. → Mean sensitive outliers
3. Numerical summaries
 - 3.1 Measures of central tendency →
 - 3.1.1 Mean, Median, Mode
 - 3.2 Measures of dispersion
 - 3.2.1 Range, Variance, Standard deviation
 - 3.3 Percentiles
 - 3.3.1 Interquartile range as a measure of dispersion.

So, a summary of the this module where we summarize numerical data. What you should be knowing at the end of this module is we first started by looking at frequency tables, we looked at discrete data where each data point was considered as a category or a class, then we looked at group data, we defined what were class intervals, we defined what was lower class limit, upper class limit and we saw how to construct the frequency tables for both discrete and grouped data.

We followed it with histograms. Again, here we assumed class intervals were of equal length and we saw how to construct histograms using google sheets. The stem-and-leaf plot, we also demonstrated how to come up with a stem-and-leaf plot for an example data.

One then, we moved on to numerical summaries. We started with the measure of central tendency although we have looked at mode and median in the earlier categorical data module, but then now we introduce a very important measure called the mean.

What we observed was mean was very sensitive to outliers is something which we saw and then, we saw what would happen to each of these measures of central tendency if you add a constant and if you multiply with the constant. The reason is it is always helps us to know what would happen to the measures whenever we manipulate the data and the way we manipulated the data here was to add a same value or multiply it with the same value.

Then, we moved on to look at measures of variability or dispersion or spread. We started with the measure range, again we saw that the range is extremely sensitive to the outliers we showed this through illustrated it through an example, then we talked about variance.

We also defined the population variance and the sample variance, but we stuck on to the sample variance. At this point of time, we said that all these measures take the same units of your original data whereas, variance takes the units of squared units of the original data hence, we define a measure which is standard deviation which is square root of the variance which takes the units of your original data.

Then, we went on to percentiles. When we discussed about percentile, we introduced what were percentiles, then we introduced important percentiles namely the 25th percentile, the 50th percentile and the 75th percentile.

We call this the first quartile or the lower quartile, the 75th percentile is the 3rd quartile or the upper quartile and the 50th percentile is what we already have seen as the median. The difference between the third and the first quartile is what we refer to as the interquartile range and we see that interquartile range is a measure of dispersion.

So, with this, we come to the end of the discussion on how to come up with numerical summaries for a single variable. What we are going to see next is the association between two variables. So, far we have looked both in the categorical case and the numerical case, measures of graphical summaries and numerical summaries of a single variable.

What would happen when I have more than one variable. So, we look at having two variables, we start with having two categorical variables look at how we, look at the association between two categorical variables, we look at association between 2 numerical variables and then, we looked at association between a categorical and a numerical variable. So, this is what we are going to do next.

Thank you.