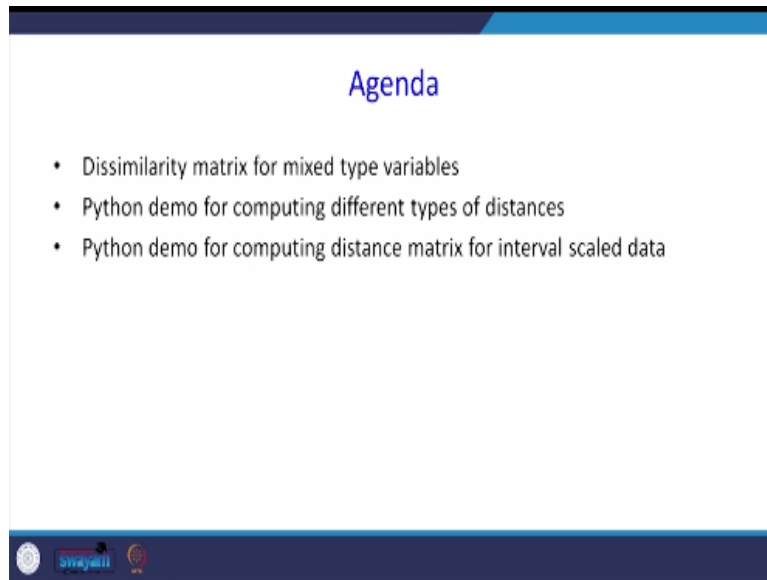


**Data Analytics with Python**  
**Prof. Ramesh Anbanandam**  
**Department of Management Studies**  
**Indian Institute of Technology – Roorkee**

**Lecture – 53**  
**Cluster Analysis - V**

In our previous class, I have explained how to handle different types of data for doing cluster analysis. I have started some theory, there is a mixed data type, how to handle that kind of data.

**(Refer Slide Time: 00:38)**



The agenda for today's lecture is how to find the dissimilarity matrix for mixed type variables. And we will do python demo for computing different types of distances which I have explained theory in my previous lectures and also I will tell you python demo for computing distance matrix for interval scaled data.

**(Refer Slide Time: 01:01)**

## Example

Consider the data given in the following table and compute a dissimilarity matrix for the objects of the table  
Now we will consider all of the variables, which are of different types

object identifier	test-1 (categorical)	test-2 (ordinal)	test-3 (ratio-scaled)
1	code-A	excellent	445
2	code-B	fair	22
3	code-C	good	164
4	code-A	excellent	1,210

Now let us take an example. This is a mixed type dataset. So consider the data given in the following table and compute a dissimilarity matrix for the objects of the table. Now we will consider all of the variables which are different types. So there are three type of dataset one is categorical, ordinal and ratio-scaled. If this kind of mixed data type is there how to use these kind of dataset for finding dissimilarity matrix and giving as an input for the cluster analysis.

(Refer Slide Time: 01:34)

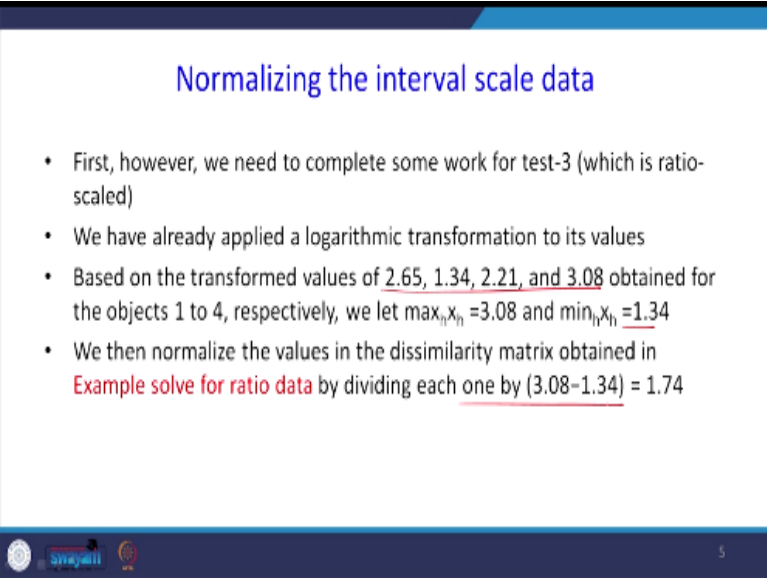
## Example

- The procedures we followed for test-1 (which is categorical) and test-2 (which is ordinal) are the same as outlined above for processing variables of mixed types
- For categorical variable -  $d(i, j) = \frac{p - m}{p}$ ,
- For ordinal variable -  $z_{ij} = \frac{r_{ij} - 1}{M_j - 1}$
- For interval scale variable -  $d_{ij}^{(f)} = \frac{|x_{ij} - x_{ij'}|}{\max_k x_{kj} - \min_k x_{kj}}$ ,

The procedures we followed for the test-1 which is categorical data and test-2 which is ordinal data are the same as outlined above for processing variables of mixed type. So what we have done, if it is a categorical variable type we have used this  $p - m$  divided by  $p$  so where the  $p$  is number of variables where  $m$  is number of matches, that we have discussed in our previous class.

If it is an ordinal data, we have to standardize into 0 to 1 by using this formula  $z_{if} = (r_{if} - 1) / (r_{max} - 1)$ , that is the current rank – 1) divided by (maximum rank – 1). So this after converting into 0 to 1 scale then you can use our simple Euclidean method for finding the dissimilarity matrix. For interval scale variable so  $d_{ij}^{(f)} = \text{modulus of } (x_{if} - x_{jf}) \text{ divided by } (\max_h x_{hf} - \min_h x_{hf})$ .

**(Refer Slide Time: 02:33)**



**Normalizing the interval scale data**

- First, however, we need to complete some work for test-3 (which is ratio-scaled)
- We have already applied a logarithmic transformation to its values
- Based on the transformed values of 2.65, 1.34, 2.21, and 3.08 obtained for the objects 1 to 4, respectively, we let  $\max_h x_{hf} = 3.08$  and  $\min_h x_{hf} = 1.34$
- We then normalize the values in the dissimilarity matrix obtained in **Example solve for ratio data** by dividing each one by  $(3.08 - 1.34) = 1.74$

First normalizing the interval scale data. First, however, we need to complete some work for test-3 which is ratio-scaled. We have already applied a logarithmic transformation to its values. Based on the transformation values we got 2.65, 1.34, 2.21 and 3.08 obtained for the objects 1 to 4, respectively. We let maximum value of  $x_h$  is from this among these values the maximum value is 3.08 and the minimum value is 1.34. The normalized value in the dissimilarity matrix obtained in the example and solve for ratio data by dividing each one by this difference that is  $(3.08 - 1.34)$  that is 1.74.

**(Refer Slide Time: 03:21)**

### Dissimilarity matrix for test-3

- This results in the following dissimilarity matrix for test-3:

Object Identifier	Ratio scaled Data (x)	Log (x)
1	445	2.65
2	22	1.34
3	164	2.21
4	1210	3.08

$$\begin{bmatrix}
 0 & & & \\
 0.75 & 0 & & \\
 0.25 & 0.50 & 0 & \\
 0.25 & 1.00 & 0.50 & 0
 \end{bmatrix}$$

- For 1 and 2 =  $(2.65 - 1.34) / (3.08 - 1.34) = 0.75$

This results in the following dissimilarity matrix for test-3. So what happened there was a object identifier was there, there was a ratio-scale was there, I am calling it as x. So we have taken log of that value 2.65, 1.34, 2.21. So we have to standardize this one for that purpose for example 1 and 2 we use this formula that is this formula to standardize. So here 2.65 this value minus this value divided by the maximum value and minus minimum value.

So when you divide this, this is a 0.75. So 2, 1 this was the distance. Similarly, if you want know 4, 1 so 4, 1 so the difference is  $2.65 - 3.08$  whole divided by this value, that is 1.74 will give you this value. The same way the standardization is done for all the cells.

(Refer Slide Time: 04:23)

### dissimilarity matrices for the three variables

- We can now use the dissimilarity matrices for the three variables in our computation of Equation  $d_{ij}^{(f)} = \frac{|x_{if} - x_{jf}|}{\max_k x_{ikf} - \min_k x_{ikf}}$ ,
- For example, we get  $d(2,1) = (1(1)+1(1)+1(0.75))/3 = 0.92$

$$\begin{bmatrix}
 0 & & & \\
 1 & 0 & & \\
 1 & 1 & 0 & \\
 0 & 1 & 1 & 0
 \end{bmatrix}$$

Dissimilarity matrix  
for categorical

$$\begin{bmatrix}
 0 & & & \\
 1 & 0 & & \\
 0.5 & 0.5 & 0 & \\
 0 & 1.0 & 0.5 & 0
 \end{bmatrix}$$

Dissimilarity matrix  
for ordinal

$$\begin{bmatrix}
 0 & & & \\
 0.75 & 0 & & \\
 0.25 & 0.50 & 0 & \\
 0.25 & 1.00 & 0.50 & 0
 \end{bmatrix}$$

normalize the values in the  
dissimilarity matrix for ratio data

Now let us consider dissimilarity matrices for all three variables. What are the three variables? One is categorical, ordinal and ratio data. So we can now use the dissimilarity matrices for the three variables in our computation of equation, so  $d_{ij}^{(f)} = \text{modulus of } (x_{if} - x_{jf}) \text{ divided by } (\max_h x_{hf} - \min_h x_{hf})$ .

**(Refer Slide Time: 04:51)**

**Example**

- The resulting dissimilarity matrix obtained for the data described by the three variables of mixed types is:

$$\begin{matrix} (2,1) & \begin{bmatrix} 0 & & & \\ 0.92 & 0 & & \\ 0.58 & 0.67 & 0 & \\ 0.08 & 1.00 & 0.67 & 0 \end{bmatrix} \end{matrix}$$

For example, we got between d 2 1. So this is the location, 2, 1. I will explain how we got 0.92. So how we got this one is, so there are three variables is there; for dissimilarity matrix it is 1, for categorical variable, for ordinal variable, for 2, 1 portion is 1, for ratio data it is a 0.75. So we will find out the weighted mean. So weight is we are giving equal weigh for all kind of dataset. So  $(1 \text{ into } 1) + (1 \text{ into } 1) + (1 \text{ into } 75)$  so total sum of weight is 3, so 0.92 that is why we got this value.

Similarly, each element you can do that one, for example here  $1 \text{ into } 1 + 1 \text{ into } 0.5 + 1 \text{ into } 0.25$  divided by 3 so we will get this value. So this matrix is a resulting dissimilarity matrix attained for data described by the 3 variables for the mixed type. So this dissimilarity matrix is given as a input for doing the cluster analysis. So what happened this is a combined matrix, combined dissimilarity matrix for all three kind of variables, so what are that variables it is for categorical, ordinal and ratio data.

**(Refer Slide Time: 06:12)**

## Interpretation

- If we go back and look at Table of given data, we can intuitively guess that objects 1 and 4 are the most similar, based on their values for test-1 and test-2
- This is confirmed by the dissimilarity matrix, where  $d(4,1)$  is the lowest value for any pair of different objects
- Similarly, the matrix indicates that objects 2 and 4 are the least similar

If you go back and look at the table of the given data we can intuitively guess that the object 1 and 4 are most similar, based on their values for test-1 and test-2. You see that, because if you look at this one, this is 0, this is 0, this is 0.25. In the given matrix, for example, for, say for categorical data it is 0, for ordinal data it is 0, it look like the position of 4, 1 for both the variables is seems to be very similar. But what is happening here if the third where it is ratio data it is 0.25.

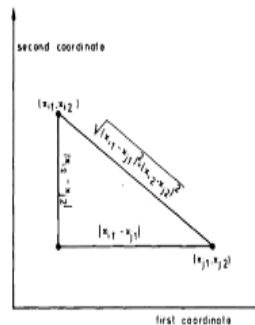
But when you find the average of 3 weighted average by looking at the dissimilarity matrix for categorical and ordinal look in the position of 4, 1. So this seems to be very similar because it is close to 0. But what is happening the 0.25 it is not very close to 0, so this can be verified by when you look at this one, so among the all the dataset in the combined the resulting dissimilarity matrix the value of 0.0 it is very similar, so the position of 4, 1 is very similar to each other, that is the point.

This is confirmed by the dissimilarity matrix, where 4, 1 is the lowest value for any pair of objects. Similarly, the matrix indicates that the object 2, 4 are the least similar. When you look at this 2, 4, so the highest value is 1. When you go there, here also 2, 4 here also it is 1, here also 1, here also 1 so in the resulting dissimilarity matrix also this seems to be 1 so highly dissimilar.

**(Refer Slide Time: 07:55)**

## Distance Measurement using python - Euclidean Distance

$$d(i, j) = \sqrt{(x_{i1} - x_{j1})^2 + (x_{i2} - x_{j2})^2 + \dots + (x_{ip} - x_{jp})^2}$$



Now, we will go to the distance measurement using python. In my previous class, I have explained how to find out Euclidean distance, Manhattan distance, Minkowski distance. So this is an example for Euclidean distance, the formula for finding Euclidean distance is  $(x_{i1} - x_{j1})$  whole square +  $(x_{i2} - x_{j2})$  whole square +  $(x_{ip} - x_{jp})$  whole square, then square root. If there are only two variable the distance formula is  $(x_{i1} - x_{j1})$  whole square +  $(x_{i2} - x_{j2})$  whole square, then square root.

**(Refer Slide Time: 08:30)**

## Python Demo for Euclidean Distance

```
In [1]: import scipy
        from scipy.spatial import distance

        #Euclidean Distance

In [2]: import numpy as np
        a = [1,2,3]
        b = [4,5,6]
        dst = distance.euclidean(a,b)

In [3]: dst
Out[3]: 5.196152422706632
```

But how to use python command, I brought the screenshot of that, for that import scipy, from scipy.spatial import distance, so how to find the Euclidean distance? So import numpy as np, so a and b there are two; a is one point where 1, 2, 3; the b is another point 4, 5, 6. If you want to

know the distance Euclidean distance between a, b so we have to write `dst = distance.euclidean(a, b)`. So when you type `dst`, so this is the our Euclidean distance between point a and b.

**(Refer Slide Time: 09:07)**

Distance Measurement using python – Minkowski Distance :

$$d(i, j) = (|x_{i1} - x_{j1}|^p + |x_{i2} - x_{j2}|^p + \dots + |x_{in} - x_{jn}|^p)^{1/p},$$

- $p=1$  Manhattan distance
- $p=2$  Euclidean distance

The next, the distance is Minkowski distance. So the Minkowski distance is, the combination of both Manhattan distance and Euclidean distance. So  $d_{i,j} = |x_{i1} - x_{j1}|^p + |x_{i2} - x_{j2}|^p + \dots + |x_{in} - x_{jn}|^p$  to the power  $1/p$ . So what is happening, this is small  $p$ . If  $p = 1$ , then it is a Manhattan distance. If the  $p = 2$  it is an Euclidean distance. Let us see how to find out this Minkowski distance using python.

**(Refer Slide Time: 09:46)**

Python Demo for Minkowski Distance

```
#Minkowski Distance

In [4]: distance.minkowski([1, 0, 0], [0, 1, 0], 1) #manhattan distance
Out[4]: 2.0

In [5]: distance.minkowski([1, 0, 0], [0, 1, 0], 2) #Euclidean distance
Out[5]: 1.4142135623730951

In [6]: distance.minkowski([1, 2, 3], [4, 5, 6], 2)
Out[6]: 5.196152422706632

In [7]: distance.minkowski([1, 2, 3], [4, 5, 6], 3)
Out[7]: 4.3267487309222245
```



So I have taken two points, Minkowski distance let us see 1, 0, 0; 0, 1, 0. So the comma 1, this represent that we are finding Manhattan distance. So when you enter this the Manhattan distanced is 2. So this same dataset if you type 2, you will get Euclidean distance because the  $p = 2$  will get you formula for the Euclidean distance that is 1.41. So this was our another example just to verifying this (1, 2, 3); (4, 5, 6), 2 so this will represent our Euclidean distance.

Suppose if you take 1, 2, 3 and 4, 5, 6 the  $p$  value can be 3 also; if it is 3 then Minkowski distance is 4.32. This formula `distance.minkowski (1, 2, 3); (4, 5, 6), 2` this data previously we have used the distance between by using this formula `distance.euclidean` distance a, b we got 5.19. So by using Minkowski formula also the Minkowski formula the same dataset if you type 2 you will get the same answer. So the `distance.minkowski` between the two points (1, 2, 3) and (4, 5, 6), 3 where the  $p = 3$  so this our Minkowski distance.

**(Refer Slide Time: 11:07)**

**Dissimilarity matrix**

```
#dissimilarity or distance matrix

In [9]: import pandas as pd
        from scipy.spatial import distance_matrix

        data = [[1, 0], [0, 1], [0, 0]]
        df = pd.DataFrame(data, columns=['a', 'b'])
        df

Out[9]:
   a b
0  1  0
1  0  1
2  0  0

In [10]: pd.DataFrame(distance_matrix(df.values, df.values))

Out[10]:
      0      1      2
0  0.00000  1.414214  2.828427
1  1.414214  0.000000  1.414214
2  2.828427  1.414214  0.000000
```

Now, can we explain how to find out dissimilarity matrix? So for that you have to import pandas as pd from `scipy.spatial` import `distance_matrix`. So there are 1, 2; there are 3 points; 1, 4; 2, 5; 3, 6. So there are two columns a, b. So `pd.DataFrame` data, columns equal to this one, you will get this kind of output. So if you want to know the distance between a, b so there are identifier name is 0 1 and 2. There are two variables a and b. So if you want to know the distance matrix between different identifier 0 0 is 1; 1 0 is 1.41; 2 0 is 2.84 by using this command `distance_matrix df.values and df.values`.

(Refer Slide Time: 12:05)

### Distance matrix calculation for Interval-Scaled Variables

- For example :
- Take eight people, the weight (in kilograms) and the height (in centimetres)
- In this situation,  $n = 8$  and  $p = 2$ .

Person	Weight(Kg)	Height(cm)
A	15	95
B	49	156
C	13	95
D	45	160
E	85	178
F	66	176
G	12	90
H	10	78



Now, distance matrix calculation for interval-scaled variables. For example, there are 1, 2, 3, 4, 5, 6, 7, 8 there is a 8 persons, we can say object identifier, take 8 people, the weight is given in kilograms and the height is given in centimeters, so  $n = 8$ , so  $p = 8$ . Now how to find out the distance matrix?

(Refer Slide Time: 12:30)

```
#data matrix

In [5]: import pandas as pd
        from scipy.spatial import distance_matrix

        data = [[15, 95], [49, 156], [13, 95], [45, 160], [85, 178], [66, 176], [12, 90], [10, 78]]
        ctys = ['A', 'B', 'C', 'D', 'E', 'F', 'G', 'H']
        df = pd.DataFrame(data, columns=['weight', 'Height'], index=ctys)

In [6]: df

Out[6]:
   weight  Height
A      15      95
B      49     156
C      13      95
D      45     160
E      85     178
F      66     176
G      12      90
H      10      78
```



So I have taken the data 15, 95; 49, 156 and so on. So there is a labels A, B, C, D, E, F, G, H. So the data frame is so Weight and Height. So we got this table. So weight and height is there, these are different identifier.

(Refer Slide Time: 12:54)

In [7]: `Distance_matrix = pd.DataFrame(distance_matrix(df.values, df.values), index=df.index, columns=df.index)`  
Distance\_matrix

Out[7]:

	A	B	C	D	E	F	G	H
A	0.000000	69.835521	2.000000	71.589105	108.577162	95.718337	5.830952	17.720045
B	69.835521	0.000000	70.830784	5.656854	42.190046	26.248809	75.663730	87.206651
C	2.000000	70.830784	0.000000	72.449683	109.877204	96.798760	5.099020	17.262677
D	71.589105	5.656854	72.449683	0.000000	43.863424	26.400758	77.388630	89.157165
E	108.577162	42.190046	109.877204	43.863424	0.000000	19.104973	114.337221	125.000000
F	95.718337	26.248809	96.798760	26.400758	19.104973	0.000000	101.548018	112.871608
G	5.830952	75.663730	5.099020	77.388630	114.337221	101.548018	0.000000	12.165525
H	17.720045	87.206651	17.262677	89.157165	125.000000	112.871608	12.165525	0.000000

Suppose if you want to know the distance matrix, so for that `pd.DataFrame distance_matrix` you take these values. So we are getting the distance matrix between A and A, 0, B and A. So you see that the diagonal value 0 because it is a Replica of because the distance between F and F is 0; G and G 0, H and H is 0. So what is happening between A and; B and A it is 69.83; A and B also 69.83 just a mirror value.

(Refer Slide Time: 13:29)

### Distance matrix calculation using Python

In [8]: `Distance_matrix.round(decimals=1, out=None)`

Out[8]:

	A	B	C	D	E	F	G	H
A	0.0	69.8	2.0	71.6	108.6	95.7	5.8	17.7
B	69.8	0.0	70.8	5.7	42.2	26.2	75.7	87.2
C	2.0	70.8	0.0	72.4	109.9	96.8	5.1	17.3
D	71.6	5.7	72.4	0.0	43.9	26.4	77.4	89.2
E	108.6	42.2	109.9	43.9	0.0	19.1	114.3	125.0
F	95.7	26.2	96.8	26.4	19.1	0.0	101.5	112.9
G	5.8	75.7	5.1	77.4	114.3	101.5	0.0	12.2
H	17.7	87.2	17.3	89.2	125.0	112.9	12.2	0.0

We will round into 1 decimal. So by using `matrix; distance_matrix.round (decimals = 1, out = none)` so we got this matrix. This will run with the help of python.

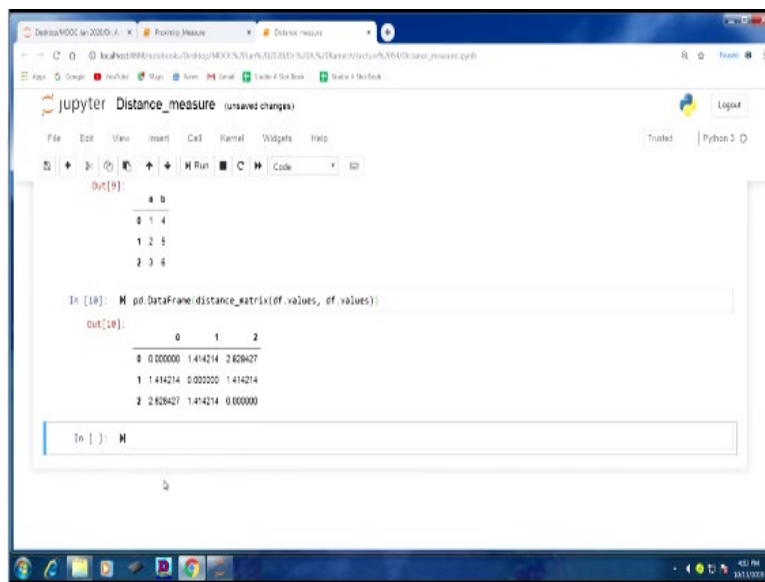
(Refer Slide Time: 13:41)



Instead of 1 if use 2 you will get here. Instead of 1 if you get 2 you will get a Euclidean distance. This is 1.41. We already got 5.19 as Euclidean distance even the Minkowski function you can get; you can verify that answer.

What happened here I gave distance.minkowski, I taken the same point that is 1, 2, 3 and 4, 5, 6 this one I used 2, so number 2 is used to get the Euclidean distance. Number 1 is use to get Manhattan distance. So we got to see that this also 5.19 when in the function Minkowski for use p, otherwise by using distance.euclidean function we got 5.19, so both are same.

**(Refer Slide Time: 15:43)**



The screenshot shows a Jupyter Notebook interface with a browser window. The notebook has a menu bar (File, Edit, View, Insert, Cell, Format, Widgets, Help) and a toolbar with icons for running, saving, and other actions. The code cell contains the following:

```
Out[9]:
```

	a	b
0	1	4
1	2	5
2	3	6

```
In [10]: M = pd.DataFrame(distance_matrix(df.values, df.values))
```

```
Out[10]:
```

	0	1	2
0	0.000000	1.414214	2.628427
1	1.414214	0.000000	1.414214
2	2.628427	1.414214	0.000000

The bottom of the notebook shows an input prompt `In [ ]:` and a status bar at the bottom of the browser window.

Now I will explain how to find out the dissimilarity or distance matrix; import pandas as pd from scipy.spatial import distance\_matrix. The data equal to 1, 4; 2, 5; 3, 6 so there are three dataset for two variables. If you want to know the, the distance between a and b; now we have three dataset for two variables. If you want to know the distance matrix, so pd.DataFrame(distance\_matrix( df.values, df.values)), so you will get this is distance matrix. So 3 variables and 3 dataset. So this, the distance matrix between 0 and 0 is 0; between 1 and 0 is 1.41, between 2 and 0 is 2.82.

**(Refer Slide Time: 16:42)**

```
In [4]: Distance_matrix = pd.DataFrame(distance_matrix(df.values, df.values), index=df.index, columns=df.index)
Distance_matrix
```

```
Out[4]:
```

	A	B	C	D	E	F	G	H
A	0.000000	69.836521	2.000000	71.589126	108.577162	95.718337	5.832082	17.720345
B	69.836521	0.000000	70.832784	5.658954	42.190248	26.248809	75.863730	87.209551
C	2.000000	70.832784	0.000000	72.449883	109.877254	96.706790	5.096020	17.262677
D	71.589126	5.658954	72.449883	0.000000	43.963424	26.440708	77.388630	88.187185
E	108.577162	42.190248	109.877254	43.963424	0.000000	19.104973	114.337221	125.000200
F	95.718337	26.248809	96.706790	26.440708	19.104973	0.000000	101.548518	112.871828
G	5.832082	75.863730	5.096020	77.388630	114.337221	101.548518	0.000000	12.166525
H	17.720345	87.209551	17.262677	88.187185	125.000200	112.871828	12.166525	0.000000

```
In [ ]: Distance_matrix.round(decimals=1, out=None)
```

```
In [ ]: M
```

Now, we will go for a distance calculation for that import pandas as pd import numpy as np, I will run it, so data matrix is given like this. So there are two variables. There are 8 persons. A, B, C, D, E, F, G, H. Suppose if you want to know the distance matrix so pd.DataFrame distance\_matrix you follow this command, the distance matrix is this one, right between A and A is 0; between B and A is 69.83; between C and A it is 2.00.

(Refer Slide Time: 17:22)

```
In [5]: Distance_matrix.round(decimals=1, out=None)
```

```
Out[5]:
```

	A	B	C	D	E	F	G	H
A	0.0	69.8	2.0	71.6	108.6	95.7	5.8	17.7
B	69.8	0.0	70.8	5.7	42.2	26.2	75.9	87.2
C	2.0	70.8	0.0	72.4	109.9	96.8	5.1	17.3
D	71.6	5.7	72.4	0.0	43.9	26.4	77.4	88.2
E	108.6	42.2	109.9	43.9	0.0	19.1	114.3	125.0
F	95.7	26.2	96.8	26.4	19.1	0.0	101.5	112.9
G	5.8	75.9	5.1	77.4	114.3	101.5	0.0	12.2
H	17.7	87.2	17.3	88.2	125.0	112.9	12.2	0.0

```
In [ ]: M
```

So if you want to have 1 decimal with the rounding of 1 decimal we got this distance matrix. So this distance matrix is given as input for cluster analysis. In this lecture I have explained how to find out the dissimilarity matrix for mixed type of variables. What is the meaning of mixed type

of variables? When we got for a cluster analysis the dataset may be combination ordinal, interval and ratio data.

When these three types of data come together how to find out the dissimilarity matrix? That I have explained. Then, by using python I have explained how to find out the distance. So I have explained how to find out the Euclidean distance. Then I have explained how to find out the Manhattan distance and then I have explained how to find the Minkowski distance. At the end I have explained with the help of python how to find out the distance matrix for the interval-scaled data, because that distance matrix can be used as a input for our cluster analysis. Thank you.