

Statistics: numbers that have been collected in order to provide information about something.

Also defined as science of collecting and studying numbers.

In []:

1

Types of Statistics

1) **Descriptive Statistics:** In descriptive statistics the data is described in a summarized way. The summarization is done from the sample of the population using different parameters such as **Mean or Standard Deviation**.

Descriptive statistics are a way of using charts, graphs & summary measures to organize, represent and explain the set of data.

Data is typically arranged and displayed in table or graphs summarizing detail such as histogram, pie chart, bars, scatter plots etc.

Descriptive statistics are just descriptive and thus does not require any generalization beyond the data collected.

In []:

1

2) Inferential Statistics:

In Inferential statistics we try to interpret the outcome. After the data has been collected, analysed & summarized we use inferential statistics to describe the meaning of the collected data.

In Inferential statistics, it is majorly intended to test hypothesis and investigate relationship between variables and can be used to make predictions.

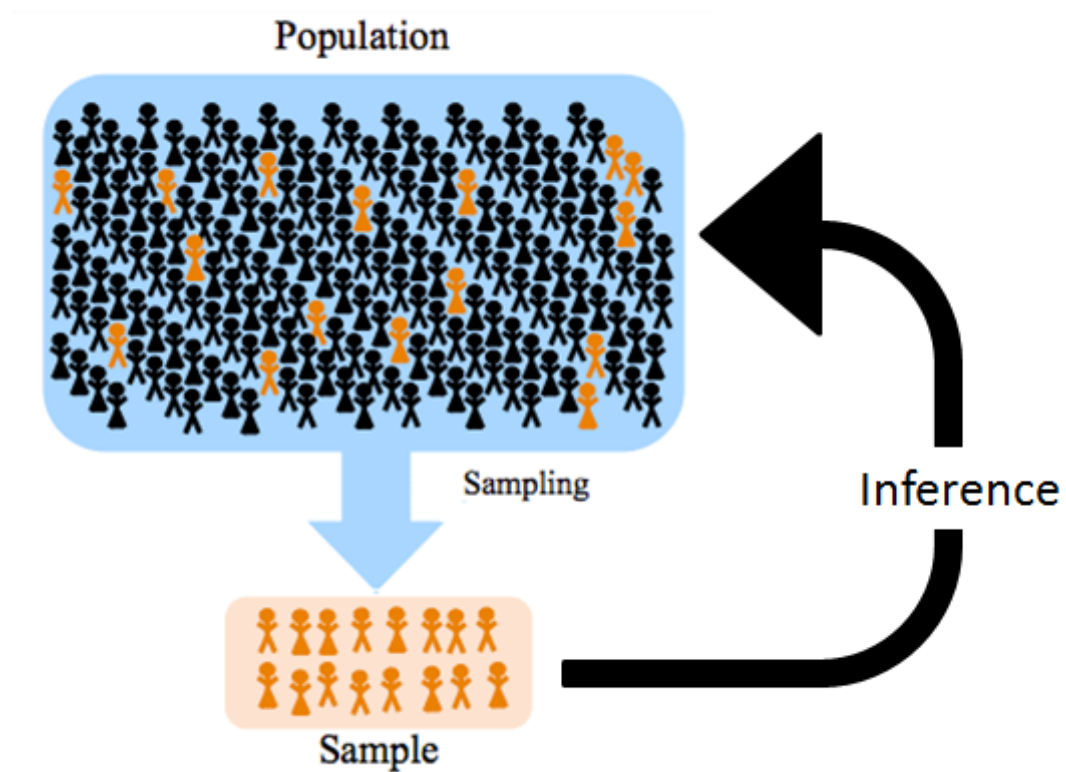
They are typically used to draw conclusion and inference, i.e to make a valid generalization from sample data

In []:

1

Population: It is actually a collection of a set of individual or object or events whose properties are yet to be analyzed.

Sample : It is the subset of the population



In []:

1

Measure of Central Tendency

Measure of Central Tendency is also known as summary in statistics that is used to represent the center point of a particular value of data set or sample set.

1) Mean

2) Median

3) Mode

```

1 1) Mean (Average)
2
3
4 Formula: sum of all the terms / total no of terms
5
6
7 Example:
8
9
10 Cars      Milage      Cylinder
11 Swift     21.3         3
12 BMW       20.8         4
13 Audi      18          5

```

```

1 Milage Mean = 21.3+20.8+18 / 3
2
3           = 20.03 ==> Mean value for milage

```

```

1 Cylinder Mean = 3+4+5/3
2
3           = 4 ==> Mean value for cylinder

```

```
In [1]: 1 21.3+20.8+18
```

```
Out[1]: 60.1
```

```
In [2]: 1
```

```

1 2) Median :
2
3
4 Cars      Milage      Cylinder
5 Swift     21.3         3
6 BMW       20.8         4
7 Audi      18          5
8 Ertiga    15          4

```

If your number of data point is even:

Step 1: Arrange data in ascending order

15 18 20.8 21.3 10 20

=====

10 15 18 20 20.8 21.3

$18+20 \implies 38/2 \implies 19$

Step 2: Check the midpoint : $18 + 20.8 \implies 38.8/2 \implies 19.4$ (median value)

In case of Odd Numbers:

21.3 20.8 18 10 122

Median Value : 18

1	Mode :
2	
3	10
4	2
5	3
6	4
7	2
8	6
9	4
10	2
11	7
12	7
13	10
14	4
15	2

In []:

1

In []:

1

Measure of variability

In []:

1

What is variability : It describes how far apart the data points lie from each other and from a center of distribution

Variability is also referred as spread OR scatter OR dispersion

In []:

1

Why variability matters

- 1) While measure of central tendency tells us where most of the point lie, variance on the other hand tells how far apart the data points are
- 2) This is important because it tells you whether the points tends to be clustered or more widely spread out.
- 3) Low variability is ideal because it means that you can better predict information about sample data based on population.
- 4) High Variability means that the value are less consistent so it is harder to make predictions.

In []:

1

Range:

It shows the spread of the data from the lowest to the highest value in distribution.

It is the easiest way of finding the measure of variability.

```
1 Example:
2
3 Data
4 72
5 110
6 134
7 190
8 238
```

9	287
10	305
11	324

Range = Highest value - lowest value

$$324 - 72$$

$$= 252$$

In []:

1

Inter Quartile Range (IQR) : It gives you the spread of the data from middle of your distribution

```

1 Q1 ==> Quartile 1 ==> 25% of value from data.
2 Q2 ==> (the median) ==> median value (mid value)
3 Q3 ==> (Quartile 3) ==> 75% of the value
4 Q4 ==> (Quartile 4) ==> Highest value
5
6
7 Formula : Q3 - Q1
8
9 Example:
10
11 72
12 110
13 134
14 190
15 238
16 287
17 305
18 324
19
20
21 Q1 ==> 0.25 * 8 ==> 2
22
23     2 ==> (value which is at the second position)
24

```

```
25         value is for q1 ==> 110
26
27
28 Q3 ==> 0.75 * 8 ==> 6
29
30         6 ==> (value which is at the 6th Position)
31
32         value for q3 is 287
33
34
35
36 IQR = Q3 - Q1
37     = 287 - 110
38
39     = 177
40
```

```
1 72
2 110
3 134
4 190
5 238
6 287
7 305
8
9 Q1 = 0.25 * 7
10    = 1.75 ==> 2
11
12    == 110
13
14
15 Q3 = 0.75 * 7
16
17    = 5.25 ==> 5
18
19    = 238
20
21
22
23 IQR = 238 - 110
24
```

```
25     = 128
26
27
```

```
In [5]: 1 import numpy as np
```

```
In [ ]: 1
```

Standard Deviation

It finds out the average of variability in your data set

It tells you on an average how far the score (data points lie) from the mean

The larger the standard deviation is the more variance (spread) the data is

Steps:

- 1) List each score and find out the mean value.
- 2) subtract mean value from each data to get the deviation from the mean
- 3) square each of those deviation.
- 4) Add up all the squared deviation
- 5) Divide the sum of squared deviation by n-1 (where n is the sample data size)
- 6) find the square root of the number

```
1 Example:
2
3 72
4 110
5 134
6 190
7 238
8 287
```



```
9 305
10 324
11
12
13
14 Step 1:
15
16 72+110+134+190+238+287+305+324 / 8
17
18 ==> 207.5
19
20
21
22 Step 2:                Deviation from mean
23
24 72                      72 - 207.5 ==> -135.5
25 110                     110-207.5 ==> -97.5
26 134                     134-207.5 ==> -73.5
27 190                     190-207.5 ==> -17.5
28 238                     238-207.5 ==> 30.5
29 287                     287-207.5 ==> 79.5
30 305                     305-207.5 ==> 97.5
31 324                     324-207.7 ==> 116.5
32
33
34 Step 3
35
36 Square Of Deviation
37
38 18360.25
39 9506.25
40 5402.25
41 306.25
42 930.25
43 6320.25
44 9506.25
45 13572.25
46
47
48 Step 4
49
50 Add up all the squared deviation
```

```
51
52 18360.25
53 9506.25
54 5402.25
55 306.25
56 930.25
57 6320.25
58 9506.25
59 13572.25
60
61 =====
62 63904.0
63
64
65
66 Step 5:
67
68 Divide the sum of squared deviation by n-1
69
70
71 63904.0 / 7
72
73 ==> 9129.14
74
75
76 Step6 : Find square root
77
78 sqrt(9129.14)
79
80 ==> 95.54
81
82
83
84
85 It means that the average sum of deviation from mean is 95.54
86
87
88
89
90
```

In [6]: 1 1

Out[6]: -97.5

In [8]: 1 np.sum([18360.25,
2 9506.25,
3 5402.25,
4 306.25,
5 930.25,
6 6320.25,
7 9506.25,
8 13572.25])

Out[8]: 63904.0

In [10]: 1 import math
2 math.sqrt(9129.14)

Out[10]: 95.54653316578262

In []: 1

Covariance

Covariance provides a measure of strength of the relation between two or more sets of random variates.

In covariance it tells how two variables vary together.

```
1 Cov(X,Y) = SUM E((X-U) E(Y-v)) / n-1
2
3
4
5 where
6
7 X is random variable
8 E(X) = U is the expected value (the mean) of the random variable x
9 E(Y) = v is the expected value (th mean) of the random variable Y
```

```

10 n = number of item in the data set
11
12 SUM ==> Summation Sign

```

$$\text{Cov}(X, Y) = \frac{\sum (X_i - \bar{X})(Y_j - \bar{Y})}{n}$$

```

1 X      Y
2
3 2.1     8
4 2.5    10
5 3.6    12
6 4.0    14
7
8
9
10 mean of X = 2.1+2.5+3.6+4.0 / 4 ==> 3.05
11
12
13 mean of Y = 8+10+12+14/4 ==> 11
14
15
16 COV(X,Y) = (2.1-3.05)(8-11) + (2.5-3.05)(10-11) + (3.6-3.05)(12-11) + (4.0-3.05)(14-11) / 4-1
17
18
19 = (-1)(-3) + (-0.6)(-1) + (0.5)(1) + (1)(3) / 3
20
21 = 3 + 0.6 + 0.5 + 3 / 3
22
23 = 7.1/3
24
25 = 2.36
26
27
28
29

```

```
In [17]: 1 2.5-3.05
```

```
Out[17]: -0.5499999999999998
```

```
In [16]: 1 (-1)*(-3)
```

```
Out[16]: 3
```

```
In [18]: 1 (-0.6)*(-1)
```

```
Out[18]: 0.6
```

```
In [ ]: 1
```

corelation

Corelation is a statistical technique that measure the degree to which two or more variables are related to each other.

The value of corelation co-efficient must fall between -1.0 to 1.0

It can only measure association, but does not show whether X can cause effect on Y OR Vice versa

```
In [30]: 1 57*57
```

```
Out[30]: 3249
```

1	AGE(x)	SUGAR LEVEL(y)
2		
3	43	99
4	21	65
5	25	79
6	42	75
7	57	87
8	59	81

```
1 Step 1 is to calcualte X * Y
```

```
2
```

```
-  
3  
4 XY  
5  
6 4257  
7 1365  
8 1975  
9 3150  
10 4959  
11 4779  
12
```

```
1 Step 2 is to calcualte (X)**2  
2  
3  
4 X**2  
5  
6 1849  
7 441  
8 625  
9 1764  
10 3249  
11 3481
```

```
1 Step 3 is to calculate (Y)**2  
2  
3  
4 Y**2  
5  
6 9801  
7 4225  
8 6241  
9 5625  
10 7569  
11 6561
```

```
In [35]: 1 np.sum([9801,
2 4225,
3 6241,
4 5625,
5 7569,
6 6561])
```

Out[35]: 40022

```
In [ ]: 1
```

```
In [ ]: 1
```

```
1 Step 4 is to do the sum of all columns
2
3
4 SUM X      SUM Y      SUM XY      SUM X**2    SUM Y**2
5
6 247        486        20485      11409      40022
```

```
In [ ]: 1
```

Formula:

$$r = \frac{n(\sum xy) - (\sum x)(\sum y)}{\sqrt{[n\sum x^2 - (\sum x)^2][n\sum y^2 - (\sum y)^2]}}$$

```
In [ ]: 1 6(20485) - (247*486)
        2 -----
        3 sqrt([6(11409) - (247)**2]*[6(40022) - (486)**2])
```

```
In [ ]: 1
```

The corelation coefficient is : **0.5298091541711386**

```
In [ ]: 1
```

```
In [ ]: 1
```



```
In [38]: 1 6*(20485) - (247*486)
```

```
Out[38]: 2868
```

```
In [40]: 1 math.sqrt((6*(11409) - (247)**2) * (6*(40022) - (486)**2))
```

```
Out[40]: 5413.272577655775
```

```
In [41]: 1 2868/5413.27
```

```
Out[41]: 0.5298091541711386
```

```
In [ ]: 1
```