

→ Descriptive Analytics

1. Data summarization
2. Visualization
3. Probability

→ Types of Data w.r.t statistics

1. Continuous data
 - There is possibility of decimal values
 - Always be NUMERIC
 - E.g. distance b/ KHGR-CSMT is 85.45metres
2. Discrete data
 - There is no possibility of decimal
 - Only whole numbers allowed, usually counting
 - Always be NUMERIC
 - E.g. count of students in class is 84
3. Attribute data
 - Characteristic data / attributes
 - E.g. gender of Ryan is 'Male'

▼ Descriptive Statistics

- helps to describe, show and summarize data in a meaningful way.
- describing features of a data set by generating summaries
- state facts and proven outcomes from a population
- describe a situation
- achieved with the help of charts, graphs, tables, etc.

▼ NumPy

```

1 import numpy as np
2 import pandas as pd
3 # importing NumPy & pandas

```

```

1 grades = np.array(['A', 'C', 'D', 'A', 'A', 'B', 'C', 'A', 'B', 'B', 'B', 'A', 'C', 'B', 'A', 'A', 'B', 'C'])
2 # creating an ndarray of grades

```

```

1 len(grades)
2 # printing the count of grades

```

```
18
```

▼ np.unique(ndarray)

```

1 np.unique(grades)
2 # printing nd array of only unique values from existing ndarray

```

```
array(['A', 'B', 'C', 'D'], dtype='<U1')
```

▼ np.unique(ndarray, return_counts=True)

```

1 np.unique(grades, return_counts=True)
2 # returns a tuple of two ndarrays, one with unique values and
3 # other with their respective count
4 # unique values are attribute data
5 # count of unique values is discrete data

```

```
(array(['A', 'B', 'C', 'D'], dtype='<U1'), array([7, 6, 4, 1]))
```

```

1 grd_counts = np.unique(grades, return_counts=True)
2 # returns a tuple of two ndarrays, one with unique values and
3 # other with their respective count

```

```

1 grd_counts[0]
2 # printing unique values from tuple grd_counts

```

```
array(['A', 'B', 'C', 'D'], dtype='<U1')
```

```
1 grd_counts[1]
2 # printing count of unique values from tuple grd_counts

array([7, 6, 4, 1])
```

▼ Proportion

- Proportion = count / total number of items

```
1 grd_props = grd_counts[1] / len(grades)
2 grd_props
3 # calculate grade proportion

array([0.38888889, 0.33333333, 0.22222222, 0.05555556])
```

```
1 grd_props = np.round(grd_props, 3)
2 grd_props
3 # rounding off grade proportion

array([0.389, 0.333, 0.222, 0.056])
```

▼ Proportion Percentage

- proportion percentage = Proportion * 100
- proportion percentage = (count/total Number of items)*100

```
1 grd_pct = grd_props*100
2 grd_pct
3 # calculate grade percentage from grade proportion

array([38.9, 33.3, 22.2,  5.6])
```

▼ np.cumsum()

- returns cummulative sum
- keeps summing up values index-wise

```
1 np.cumsum(grd_pct)
2 # calculate summulative sum, means to add grade percentages cummulativly
3 # to tell
```

```
array([ 38.9,  72.2,  94.4, 100. ])
```

▼ Sorting

```
1 x = np.array([45, 78, 23, 56, 12])
2 # creating an ndarray
```

```
1 np.argsort(x)
2 # returns index in order of sorted numbers
```

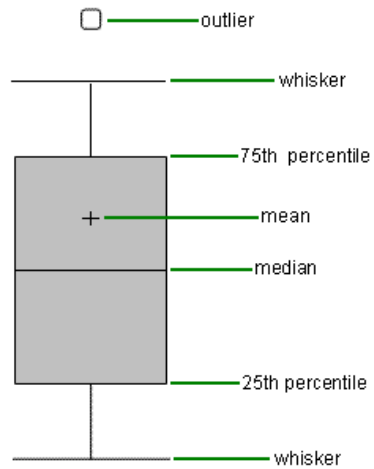
```
array([4, 2, 0, 3, 1])
```

```
1 sorted(x)
2 # returns list in sorted order
```

```
[12, 23, 45, 56, 78]
```

▼ Box Plot

- method for descriptive statistics
- also known as box and whisker plot
- includes two parts, a box and a set of whiskers
- consists of five number summary of data set
 1. Minimum (0th percentile) - lower whisker
 2. first quartile (25th percentile)
 3. Median (50th percentile) - horizontal line in the middle of box
 4. third quartile (75th percentile)
 5. Maximum (100th percentile) - upper whisker
- we draw a box from the first quartile to the third quartile
- A vertical line goes through the box at the median, and



IQR

- Inter-Quartile Range
- consists of the central 50% of the data
 - $IQR = Q3 - Q1$

$IQR = \text{third quartile} - \text{first quartile}$
 $IQR = 75\text{th percentile} - 25\text{th percentile}$

Whiskers

- two whiskers represent the maximum & minimum
 - Lower whisker - minimum / 0th Percentile
 - Upper Whisker - maximum / 100th Percentile

Fences

- used to identify outliers

1. lower inner fence: $Q1 - 1.5 \cdot IQR$

2. upper inner fence: $Q3 + 1.5 \cdot IQ$

sometimes we consider $3IQ$ too for fences

3. lower outer fence: $Q1 - 3 \cdot IQ$

4. upper outer fence: $Q3 + 3 \cdot IQ$

Outliers

- a data point that is located outside the fences of the box plot
- represented by dot, a small circle, a star, etc.

▼ Three types of Descriptive summaries

1. Distribution (Probability Distribution)
2. Central Tendency
3. Variability

- Central Tendency

1. mean
2. median
3. mode

- Variation or Dispersion

1. range
2. inter-quartile range
3. variance
4. standard deviation
5. Shape

1. skewness
2. kurtosis

- Probability Distributions

1. Discrete Probability Distributions

1. Binomial Distribution
2. Poisson Distribution

- 3. Hypergeometric Distribution
- 2. Continuous Probability Distributions
 - 1. Exponential Probability Distributions
 - 2. Normal Probability Distributions
 - 3. Student's t Probability Distributions

▼ Central Tendency

- represent the center point or typical value of a dataset
- descriptive "summary" of a data set

- 1. Mean
- 2. Median
- 3. Mode

▼ Mean / Average \bar{x} or μ

- Mean or Average : sum of all elements / count of elements
- Mean or Average : $(x_1 + x_2 + x_3 + \dots + x_n) / n$

1. Sample Mean / Average (\bar{x})

- \bar{x} : Sample Mean of variable x
- $\sum x_n$: sum of n values
- n : number of values in sample
- $\bar{x} = \sum x_n / n$

2. Population Mean / Average (μ)

- μ : Population Mean / Average
- $\sum X_N$: sum of N values
- N : number of values in the population
- $\mu = \sum X_N / N$

Note:

- Inferential statistics use Hypothesis tests such as t-tests, ANOVA tests to determine whether Sample data is different than Population mean

▼ np.mean(ndarray)

```
1 mylist = [3, 6, 5, 4, 7, 6, 8, 9, 1, 10, 13, 4, 7, 11, 4, 7, 9]
2 # create a list of values
3 mylist
4 # printing the list
```

```
[3, 6, 5, 4, 7, 6, 8, 9, 1, 10, 13, 4, 7, 11, 4, 7, 9]
```

```
1 sum(mylist)/len(mylist)
2 # manually calculating mean
```

```
6.705882352941177
```

```
1 np.mean(mylist)
2 # Claculates Population Mean or Average of the the ndarray / list
```

```
6.705882352941177
```

▼ Weighted Mean

- 10 have value as 7
- 15 have value as 8
- 23 have value as 19
- 8 have value as 10
- Assign a weight to each value in the dataset
- frequency is the weight

```
- x1 = 07      w1 = 10
- x2 = 08      w1 = 15
- x3 = 09      w1 = 23
- x4 = 10      w1 = 08
```

- Weighted Mean = $\sum xw / \sum w$
- Weighted Mean = $(x1w1 + x2w2 + x3w3 + x4w4) / (w1 + w2 + w3 + w4)$


```
1 values = np.array([7, 8, 9, 10])
2 # creating ndarray of values / weights to assign weight
```

```
1 freq = np.array([10, 15, 23, 8])
2 # creating ndarray of frequencies
3 freq
```

```
array([10, 15, 23, 8])
```

```
1 freq*values
2 # calculating products of weights & samples
```

```
array([ 70, 120, 207, 80])
```

```
1 freq*values/sum(freq)
2 # average or mean of each weight
```

```
array([1.25      , 2.14285714, 3.69642857, 1.42857143])
```

```
1 sum(freq*values/sum(freq))
2 # weighted mean / average
```

```
8.517857142857142
```

▼ Median

- natural measure of central tendency
- splits the dataset in half
- middle value in sorted list of data
- first sort the data in increasing order, then median is the middle value
 - for odd sample size, Median = middle value

```
Median = middle value
Median = ((n+1)/2)th term
```

- e.g. sample size, n=7

```
Median = ((7+1)/2)th term
Median = (8/2)th term
```

Median = 4th term

- for even sample size: Median = Average of two Middle values

Median = Average of two Middle values
 Median = $((n/2)\text{th term} + ((n/2)+1)\text{th term})/2$

- e.g. sample size, $n=8$

Median = $((8/2)\text{th term} + ((8/2)+1)\text{th term})/2$
 Median = $(4\text{th term} + (4+1)\text{th term})/2$
 Median = $(4\text{th term} + 5\text{th term})/2$

```
1 mylist = [3, 6, 5, 4, 7, 6, 8, 9, 1, 10, 13, 4, 7, 11, 4, 7, 9]
2 # create a list of values
3 mylist
4 # printing the list
```

```
[3, 6, 5, 4, 7, 6, 8, 9, 1, 10, 13, 4, 7, 11, 4, 7, 9]
```

```
1 len(mylist)
```

```
17
```

```
1 sorted(mylist)[8]
2 # median: middle value in sorted list
3 #  $17+1/2 = 18/2 = 9\text{th value} = \text{at } 8\text{th index}$ 
```

```
7
```

▼ np.median(ndarray or list)

```
1 np.median(mylist)
2 # median using Numpy median method for the array / list
```

```
7.0
```

▼ Mean vs. Median ... when to use mean and when median

- mean might skew your average value drastically, while median does not skew the middle value

1. mean should be used when

- when data distribution is symmetric
- data doesnot have any extreme values

2. for symmetric distribution, both mean & median are almost same

3. for skewed distribution, median better represents the central tendency

4. median is a robust statistic, mean is sensitive to outliers and skewed distributions

5. use median only for Skewed distribution, Continuous data, Ordinal data

Note

- if data is symmetric, then mean = median
- example of symmetric distribution is normal distribution curve/ bell curve

```
1 sal = [6, 7, 6, 7, 6, 7, 8, 8, 6, 7, 6, 8, 8, 7, 7,8, 6, 7, 8, 6]
2 # list with symmetric/normal distribution
```

```
1 np.mean(sal)
2 # Mean / Average
```

6.95

```
1 np.median(sal)
2 # Median
3
4 # note that, for symmetric data, as there are no extreme values
5 # so mean & median are almost same
```

7.0

```
1 new_sal = [6, 7, 6, 7, 6, 7, 8, 8, 6, 7, 6, 8, 8, 7, 7,8, 6, 7, 8, 6, 54, 50, 52]
2 # list with extreme values, does not follow symmetric distribution
```

```
1 np.mean(new_sal)
2 # as there are extreme values, so mean / average might be misleading
```

12.826086956521738

```

1 np.median(new_sal)
2 # median shows better reality for middle value
3 # note that, as there are extreme values, mean has changed drastically
4 # but median is not impacted

```

7.0

▼ Mode

- most frequent value
- two value can be mode if both hold the highest frequency
- if no value is repeating, (all values have frequency=1), then data does not have a mode
- use mode for categorical , ordinal and discrete data only

Note

- for symmetrical data, mean, median and mode are same

```

1 from collections import Counter

```

```

1 mylist = [3, 6, 5, 4, 7, 6, 8, 9, 1, 10, 13, 4, 7, 11, 4, 7, 9]
2 # create a list of values
3 mylist
4 # printing the list

```

[3, 6, 5, 4, 7, 6, 8, 9, 1, 10, 13, 4, 7, 11, 4, 7, 9]

```

1 Counter(mylist).most_common()
2 # returns unique values and their espective count
3 # Mode: element with most number of occurences
4 # mode is 4, 7 with frequency 3

```

```

[(4, 3),
 (7, 3),
 (6, 2),
 (9, 2),
 (3, 1),
 (5, 1),
 (8, 1),
 (1, 1),
 (10, 1),

```

```
(13, 1),  
(11, 1)]
```

▼ Variation or Dispersion

1. Range
2. Inter-Quartile Range
3. Variance
4. Standard Deviation

▼ Range

- easiest to calculate, simplest measure of dispersion
- denotes the spread of observations
- difference between maximum (0th Percentile) and minimum (100th percentile)
- $\text{range}(X) = \max(x) - \min(x)$

Drawbacks of range

1. it uses only two elements of multiple elements, which does not describe intermediate values, making it unreliable if intermediate values change
2. calculation of range can be affected by presence of extreme values
3. does tell anything about variability of other data

```
1 mylist = [3, 6, 5, 4, 7, 6, 8, 9, 1, 10, 13, 4, 7, 11, 4, 7, 9]  
2 # create a list of values  
3 mylist  
4 # printing the list
```

```
[3, 6, 5, 4, 7, 6, 8, 9, 1, 10, 13, 4, 7, 11, 4, 7, 9]
```

```
1 max(mylist) - min(mylist)  
2 # range: max - min  
3 # printing the range for list 'mylist'
```

```
12
```

```
1 x = [1, 4, 4, 7, 5, 8, 9, 3, 10]
2 max(x)-min(x)
3 # printing the range for list 'x'
```

9

```
1 y = [1, 4, 4, 4, 4, 4, 4, 4, 10]
2 max(y)-min(y)
3 # printing the range for list 'y'
```

9

▼ Inter-Quartile Range (IQR)

- Inter-Quartile Range
- consists of the central 50% of the data
 - $IQR = Q3 - Q1$

IQR = third quartile - first quartile
IQR = 75th percentile - 25th percentile

- cannot calculate IQR directly in python, first calculate 75th percentile & 25th percentile using `np.percentile(ndarray, percentile)`, and then find the difference

```
1 mylist = [3, 6, 5, 4, 7, 6, 8, 9, 1, 10, 13, 4, 7, 11, 4, 7, 9]
2 # create a list of values
3 mylist
4 # printing the list
```

[3, 6, 5, 4, 7, 6, 8, 9, 1, 10, 13, 4, 7, 11, 4, 7, 9]

▼ np.percentile(ndarray, percentile)

```
1 np.percentile(mylist, 75)
2 # calculates 75th percentile
```

9.0

```
1 np.percentile(mylist, 75) - np.percentile(mylist, 25)
2 # IQR = 75th percentile - 25th percentile
3 # calculates IQR
```

5.0

```
1 x = [1, 4, 4, 7, 5, 8, 9, 3, 10]
2 # asymmetric data - has extreme values/ skewed distribution
3 y = [1, 4, 4, 4, 4, 4, 4, 4, 10]
4 # symmetric data - has normal distribution
```

```
1 np.percentile(x, 75) - np.percentile(x, 25)
2 # for asymmetric data, there will be some value of IQR
```

4.0

```
1 np.percentile(y, 75) - np.percentile(y, 25)
2 # for symmetric data, IQR is zero
3 # it has same value on both sides of median, so IQR = zero
```

0.0

Deviation, D

- calculated by taking the differences between each number in the data set and the mean
- deviation of each observation (x) w.r.t. Sample Mean (\bar{x}), which can be +ve or -ve
- sum of all deviations is zero
- Deviation, $D = \text{Observation} - \text{Sample Mean}$
- Deviation, $D = X_n - \bar{x}$

▼ Variance S^2 or σ^2

- measures how far each number in the set is from the mean
- measurement of the spread between numbers in a data set
- square of deviation for each observation
- unit is also squared, as it is squared to remove negative values

1. Sample Variance, S^2

- Sample Variance, $S^2 = (\text{Observation} - \text{Sample Mean})^2 / \text{Number Of Values}$

- Sample Variance, $S^2 = \frac{\sum (X_n - \bar{x})^2}{n-1}$

2. Population Variance, σ^2

- Population Variance, $\sigma^2 = \frac{(\text{Observation} - \text{Population Mean})^2}{\text{Number Of Values}}$
- Population Variance, $\sigma^2 = \frac{\sum (X_n - \mu)^2}{N}$

▼ np.var(ndarray, ddof=0)

```
1 x = [4, 6, 7, 9, 10, 3, 5, 8, 3, 7]
```

```
1 np.var(x, ddof=0)
2 # when ddof=0, calculates Population Variance
```

5.36

```
1 np.var(x, ddof=1)
2 # when ddof=0, calculates Sample Variance
```

5.9555555555555556

```
1 np.var(x, ddof=0)**0.5
2 # population standard deviation is square root of population variance
```

2.3151673805580453

▼ Standard Deviation, σ or S.D.

- square root of variance
- unit is same as observation data, as it is square root of variance
- represents the average distance from the mean
- S.D. reflects the extent of variation in data
- As variation increases, value of S.D also increases

1. Sample S.D.

- Sample S.D., $\sigma = \sqrt{\frac{\sum (\text{Observation} - \text{Sample Mean})^2}{\text{Number Of Values}}}$
- Sample S.D., $\sigma = \sqrt{\frac{\sum (X_n - \bar{x})^2}{n-1}}$

2. Population S.D.

- Population S.D., $\sigma = \sqrt{\frac{\sum (\text{Observation} - \text{Population Mean})^2}{\text{Number Of Values}}}$

- Population S.D., $\sigma = \sqrt{\sum(X_n - \mu)^2 / N}$

- Sample S.D. > Population S.D. because there is uncertainty in sample S.D. as dof is more
- degree of freedom (dof) = Sample Size (N) - 1

▼ np.std(ndarray, ddof=0)

```
1 np.std(x, ddof=0)
2 # when ddof=0, calculates Population standard deviation
3 # ddof : delta degree of freedom
```

2.3151673805580453

```
1 np.std(x)
2 # when ddof is not specified, ddof=0, calculates Population standard deviation
```

2.3151673805580453

```
1 np.std(x, ddof=1)
2 # sample standard deviation
3 # when ddof=1, calculates Sample standard deviation
```

2.440400695696417

```
1 x = [4, 7, 6, 8, 9, 2, 3, 5, 4, 7]
```

```
1 np.std(x, ddof=0)
2 # population standard deviation
```

2.1563858652847827

```
1 np.std(x, ddof=1)
2 # sample standard deviation
```

2.273030282830976

▼ Shape

1. skewness - for location of peak
2. kurtosis - for nature of peak

▼ Skewness, Skew[X]

- measure of the lack of symmetry

1. Symmetric

- A distribution, or data set, is symmetric if it looks the same to the left and right of the center point
- for symmetric curve
- normal distribution
- $\text{Skew}[X] = 0$
- e.g. bell curve / normal distribution

2. peak at left

- when peak is at left, it is Right skewed
- positively skewed, $\text{Skew}[X] > 0$ (+ve)
- e.g. salary distribution of employees, seniors have more salary but have less frequency,

3. peak at right

- when peak is at right, it is Left skewed
- negatively skewed, $\text{Skew}[X] < 0$ (-ve)
- e.g. age distribution of members of parliament, young people are very

Note:

- Using skewness, we can predict if calculating mean would be risky or not, because mean for data which is not normally distributed, would be incorrect

```
1 import scipy
2 from scipy import stats
```

```
1 x = [1,1,2,2,2,2,2,2,2,2,2,3,3,3,4,4,5,6]
2 y = [1,2,2,3,3,4,4,4,5,5,5,5,5,5,5,5,6,6,7]
3 z = [1,1,2,2,2,3,3,3,3,4,4,4,4,4,4,4,4,4,4,5,5,5,5,6,6,6,7,7]
```

```
1 scipy.stats.skew(x)
2 # Positive skewness means Right Skewed or peak at left
```

```
1.1839325019418967
```

```
1 scipy.stats.skew(y)
2 # Negative skewness means Left Skewed or peak at right

-0.6877802002110032
```

```
1 scipy.stats.skew(z)
2 # Zero skewness means Normal Distribution

0.0
```

▼ Kurtosis, K

- is a unitless measure that describes sharpness of the peak
- describes range of data, how much of variation is there in data
- degree to which data values are concentrated around the mean
- measure of whether the data are heavy-tailed or light-tailed relative to a normal distribution

1. mesokurtic

- for normal distribution curve or bell curve
- `kurtosis = 0`

2. leptokurtic if curve is sharper than normal distribution curve

- heavy tails, `kurtosis > 0 (+ve)`
- more peaked as compared to normal distribution curve
- peak is prominent, so many outliers

3. platykurtic

- light tails, `kurtosis < 0 (-ve)`
- if curve is flat than normal distribution curve
- less peaked than normal distribution curve
- peak is less prominent, so fewer outliers

Note:

- higher the Kurtosis (+ve or -ve), means some elements have very high frequency

```
1 x = [1,1,1,2,2,2,3,3,3,3,4,4,4,4,4,5,5,5,5,6,6,6,7,7]
2 y = [1,1,2,2,2,3,3,3,3,4,4,4,4,4,4,4,4,4,4,4,4,4,4,4,4,4,5,5,5,5,6,6,6,7,7]
```

```
1 scipy.stats.kurtosis(y)
2 # positive kurtosis means sharp curve
```

```
0.70242214532872
```

```
1 scipy.stats.kurtosis(x)
2 # negative kurtosis means flat curve
```

```
-0.9451498127839182
```

Laptop Project

Refer: use data.xlsx file

- Q1. How many laptops did I sell overall? How many of those were Dell Laptops(absolute and percentage)?
- Q2. What was my total revenue? How much of that can be contributed to Dell(absolute and percentage)?
- Q4. How is the revenue distributed over laptop categories?
- Q5. What is the best selling continent, country and city?