

5 Oct - 2022

Data Science

Wednesday!

### Pandas :

- It offers data structures and operations for manipulating numerical tables and time series
- Pandas is used for analyzing the data and its manipulation
- \* Series is a specific row / specific column
- DL does prediction on basis of statistics
- 0.025% of DL-applied
- shape returns a tuple
- iloc: integer location.
- iloc can be used to fetch columns.
- \* `data.iloc[:, [4, 5, 16]]` → all rows with 4, 5, 16 no column

### Filtering DataFrame on a condition

`mask = data['city'] == "Hyderabad"`  
`data[mask]` → DataFrame / mask: boolean series

`check = data[mash1 & mash2]`  
check

`data['season'].value_counts()`

`data['winner'].value_counts().head().plot(kind='bar')`

winning is on horizontal (in histogram)  
frequency on Y-axis

13/07/2022

## Data Science

## Mohit's bday

- `plt.style.use("ggplot")` { built-in option available  
color | shape | background } seaborn and matplotlib folders
- Matplotlib and seaborn are used for data visualization in the pictorial view / form
- import `matplotlib.pyplot as plt`  
main library      sublibrary      alias
- `plt.tight_layout()` → to take graph in center / a little bit alignment from edges
- \* for arrow we write 'v' in marker
- `plt.xkcd()` like graph is sketched
- `plt.bar()`      } bar and plot (line graph)  
`plt.plot()`      } will be mixed
- `arrange(18)` means indexess from 0 to 17
- `Counter()` built in function of pandas  
for different languages
- `most_common()` function  
[name: index[0]] → index[1] → count
- `wedgeprops()` { for boundary specification  
13hr56m → 1 decimal places
- `autopct = '%.1f %%'` → float
- `explode` is for the pie chunk which is out of whole body
- `shadow = True` for shadow of c

# Data Science

→ startangle = 90, cutting pie chart at 90° angle

→ Stackplot helps in visualizing pattern of every user after every minute

for first value it shows the area of first, then add the value with second one and then add it with third one.

But sum of all values must be same in whole list

→ plt.fill\_between (ages>py-salaries, dev\_salaries, where=(py-salaries>dev\_salaries)

interpolate=True, color="red", alpha=0.25, label="Above-Avg")

Transparent/opacity

→ Histogram tells you count of occurrence

→ We take log to make values in a range of small values.

→ log helps out to remove outliers from data

29/08/2022

## 1) Data Science

Thursday

→ Scatter plot is for bivariate analysis (it needs two variables for plotting)

→ sizes(15, 200)  
minimum radius      maximum radius

→ relplot stands for relational plot

→ Scatter plot is for numerical data  $x$  and  $y$  should be numerical. But hue parameter is always categorical. style can or cannot be numeric.

→ catplot (x-axis: categorical, y-axis: numerical)  
categorical plot

→ First quartile means 25% values of your data is less than this value.  
median is 50% of data.

→ Boxplot is used to detect outliers  
→ Outlier disturbs your mean in a drastic way.

→ Median is changed according to confidence interval

→ pdf for discrete data

→ Descriptive Statistics      Inferential Statistics

→ Central Tendency: values of data where values are more

→ std is the under root of variance, so std values are small than var

26 Oct 2022

## Data Science

Wednesday

- Histogram is for single variable
- Scatter plot is for 2 variables
- Distplot :- to do univariate analysis on continuous/numerical data
  - combines histogram, kernel density function, Rugplot
  - x-axis: bins, y-axis: count/frequency → histogram
- By default rug is false in .distplot
- Rugplot tells you the information of count of every value in the bin
- Distplot has probabilities on y-axis
- pdf made from histogram / applying kernel density estimation

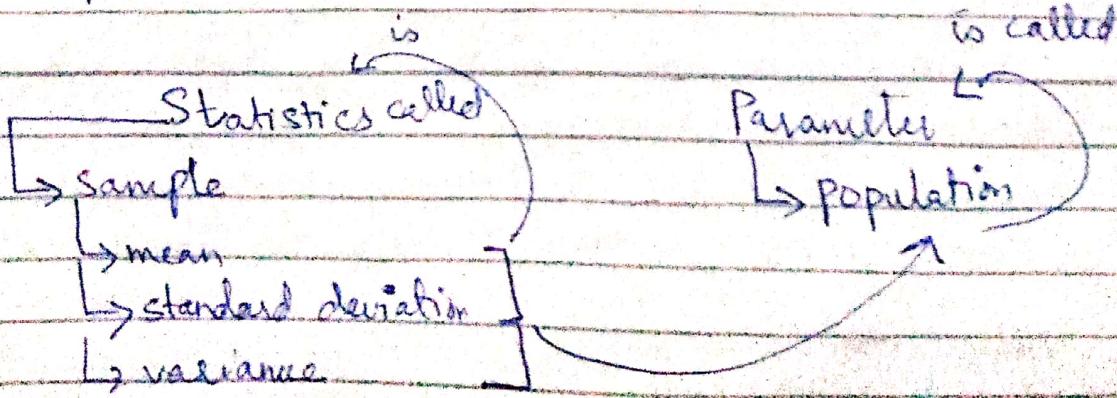
## Statistics

- collect data and draw conclusions from it
- Descriptive and Inferential statistics

Sampling from Population

Drew conclusions from sample data

Sample must be representative of population



26/07/22

DS/

- Mean, std and var of population is called parameter
- Mean, std and var of sample is called statistics

## Descriptive Statistics

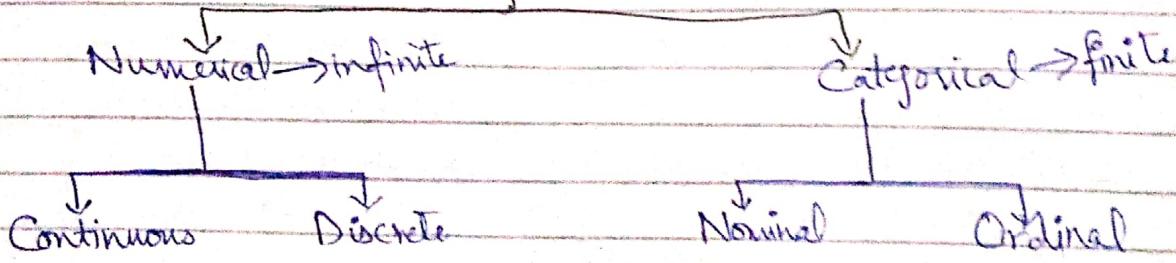
- Mean
- Median, Mode, std deviation
- Variance, MAD (Mean Absolute Deviation or Median Absolute Deviation)

(1)

- ## Inferential Statistics → Is something really different from other or occurred by chance
- Confidence Interval
  - Hypothesis Testing

## Descriptive Statistics

### Types of Data in Statistics



→ Discrete: values are countable / distinct values  
number of mobile apps in phone

→ Nominal: values are independent of each other

→ Ordinal: values have relationship with each other / Order [1 2 3 4 5] 1 and 2, 3 and 4 like have relationship with each other for example, 2 is better than 1 so it's have relationship



26 Oct-2012

Data Science

Wednesday

- Visualization Techniques
- Then Measure of Central Tendency
- Measure of Dispersion

~~mode  
range~~



### Visualization Techniques

- Bar chart or pie chart when data is categorical and nominal
- Graph do not tell no relationship
- Bar and pie chart is used when data is categorical and ordinal

→ Histogram for Numerical data

→ Scatter plot for bivariate numerical data

→ for Numerical and Continuous Data (univariate analysis) you can use boxplot

→ Bar chart when one variable is categorical and one is numerical

### Bivariate Analysis

→ Categorical and Categorical data → Crosstab  
    ↳ bar plot / heatmap / cluster map

Correlation

### Measure of Central Tendency

- where the most values occur in the data
- Median does not change whence you add outlier
- You should use 3 of them (mean, median, mode) to have fair idea of the centrality of data

2 26/ Oct-22

## Data Science

→ Measure of Spread :-

→ Range → max - min

→ Interquartile Range

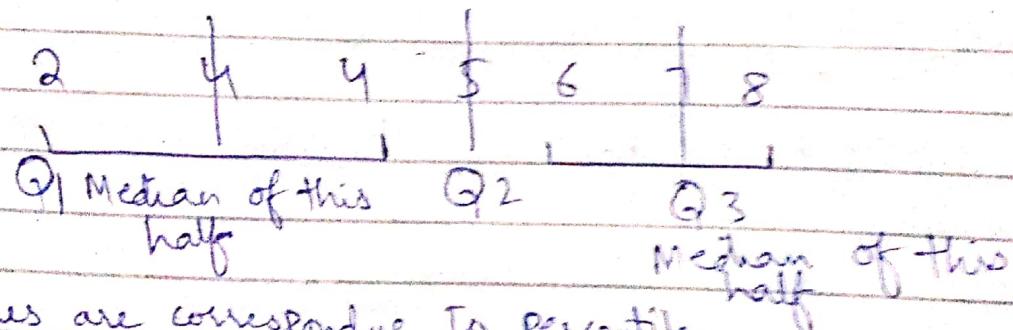
→ Standard deviation

→ Variance

\* Percentile: value at the percentage where percentage of below data is found

(1) Quartiles : Divide the data into quarters

$Q_2$  is the median



Quartiles are corresponding to percentiles

Interquartile Range

Interquartile range is the measure of statistical dispersion between upper (75<sup>th</sup>) and lower (25<sup>th</sup>) quartiles.

$Q_1$	$Q_3$
25.1	25.1

$Q_3 - Q_1$  is Interquartile Range (IQR)  
Boxplot's box difference is IQR

27 Oct 2022

# Data Science

## Mean Absolute Deviation:

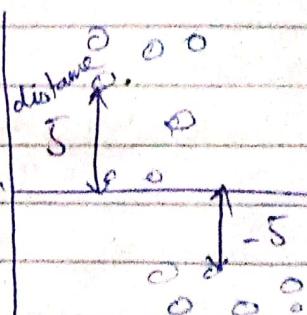
$$MAD = \frac{\sum |x_i - \bar{x}|}{n}$$

$\bar{x} \rightarrow$  mean of your data.  
 $(n) \rightarrow$  individual data point  
 number of points in sample space

## Variance

Means how far the datapoints are spread out from the mean.

$$\sigma^2 = \frac{\sum_{i=1}^N (x_i - \bar{x})^2}{N} \xrightarrow{\text{mean}}$$



It will cancel out each other's mean effect, that's why we use square.

## Problems

Data unit changes from cm to  $\text{cm}^2$

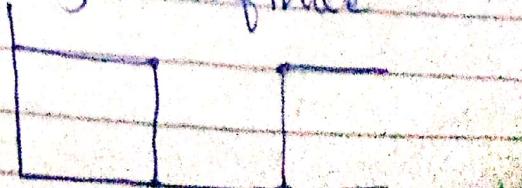
Modulus is not a non-smooth function

The graph shows a V-shaped curve opening upwards, symmetric about the y-axis. The vertex of the V is at the origin (0,0). The right branch of the V passes through points such as (1,1), (2,2), and (3,3). The left branch of the V passes through points such as (-1,1) and (-2,2). The graph is labeled "y = |x|" above the vertex and "graph" below the x-axis.

It just simply  
turns the straight line  
upward or downward

Discontinuity:  $y$  ki aik value k cortex parhig  $x$  ki 2 values ahi hain. Derivative gets infinite so graph gets non-smooth.

$\Rightarrow$  graph gets non-smooth  
if Derivative is 0 or  $\infty$ , this  
means discontinuity



27/Oct/22

## Data Science

### 9. Median Absolute Deviation

→ This function is in  
statsmodels library

→ Gives good results when data has outliers.

### Practical Statistics for Data Scientists

Book

Peter Bruce, Andrew  
Bruce and  
Peter Grodeck

### Probability Density Functions

(PDF)

→ pdf is made from histogram

Kitna Zaidan,  
Kitna Kam

→ kde is applied on histogram

pdf value = 1 occurrence of 1 bin

Total occurrences of all bin values

→ Gives you smooth graph for univariate analysis  
→ for multivariate, you get confused between  
values

But intersection point is taken to calculate high  
probability of species/class/values

### Cumulative Density/Distribution function.

cdf: 1.5 pa 70.% of values <sup>of karn</sup> than

pdf: 1.5 any ki probability 0.28 or 28.1.

b

cumsum in numpy (cumulative sum)

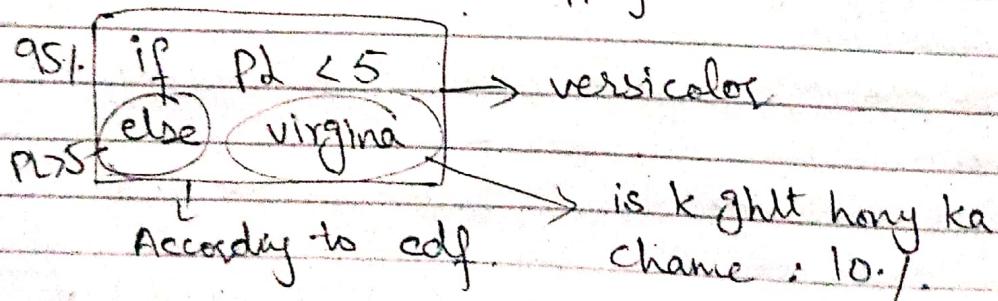
29 Nov. 2011

## Data Science

Wednesday

- KDE is smoothing function
  - PDF (probability ki density show karne ka function)
  - PDF tells you Probability probabilities of occurrence
  - CDF tells you that below this point how much data/people exists
    - ↳ Basically tells you the summary of data till that point
  - PDF is helpful when data is multivariate
    - ↓ occurrence of count of bin
- 1 point: 50  
120 → total count

- In what kind of scenarios CDF helps.
- When you have overlapping curves



- \* CDF helps in quantifying your decision

## Normal Distribution

- unknown task is done using known patterns (molding unknown in a way so that it can be resolved with known techniques)

2 Nov. 2022

## Data Science

Wednesday

→ fourier: sine and cosine wave combination can transform every signal

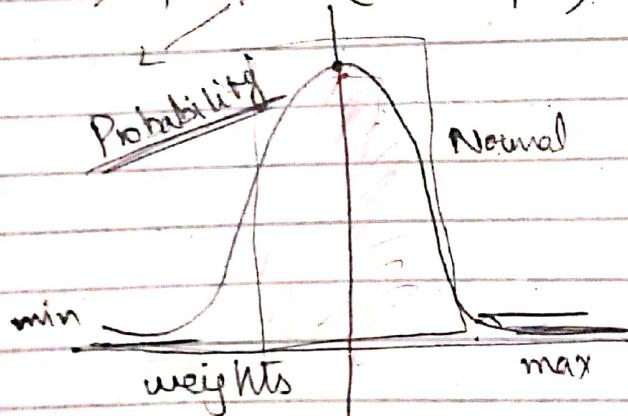
Area under the curve is calculated by (integration)

No sum

### Distribution:

- ① Poisson
- ② Bernoulli
- ③ Binomial
- ④ Normal Distribution

↳ Gaussian (Bell shaped)



- There are less people in world, whose weights are less
- There are less . . . , whose weights are maximum
- There are more people, whose weights are normal

Central Limit Theorem says

\* Every distribution can be converted into normal distribution

## Normal Distribution Curve Equation

$$P(X=z) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left(\frac{z-\mu}{\sigma}\right)^2}$$

$\mu$  = mean

$\sigma$  = standard deviation

$\sigma^2$  = variance

Normal distribution

has same mean,  
median, mode

If  $\mu=0, \sigma=1$

Equation:

$$P = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x^2}$$

constant

Always for Normal curve, variance (std = 1)  
and  $\mu=0$

Standard Normal Variate :  $\boxed{\mu=0, \sigma=1}$

$\boxed{\text{loc} = \mu, \text{scale} = \text{std in formula.}}$

Subtract mean and

Divide by std to make the  
distribution standard normal variate

-deviate np. std (distr)

Z-score :  $\frac{\text{Score} - \text{avg score}}{\text{std}}$

$$\underline{1.66} = \underline{1.6 + 0.06}$$

2-tailed ~~0.9515~~ 0.9515

xply by 100 = 95.15 · P

Student scored more than 95.15% of all

3 Nov 2022

Data Science

Thursday

statistics X

## Data Analysis Process

- Steps

- Asking questions
- Data Wrangling / Munging
- Exploratory Data Analysis
- Drawing conclusions
- Communicating results

Speech therapy  
for  
Autistic patients

### ① Asking questions:

- for better questions:-
- You must have subject/domain expertise
  - Experience

### ② Data Wrangling:

Steps:

- ① Gathering data
- ② Assessing Data
- ③ Cleaning data

#### - Gathering data (API : Application program interface)

\* CSV files

\* API (data in json format)

\* Web Scraping (data in html format)

\* databases

#### - Cleaning data

\* Missing data (mean)

, Remove duplicate data (drop\_duplicates)

, Incorrect data type (astype)

3 Nov 2022

# Data Science

Thursday

## ③ EDA:

### ① Explore

↓  
Visualization

### ② Augment

changing according to the patterns seen in visualization

- Removing Outliers
- Merging Dataframes
- Adding new column

This EDA is also called Feature Engineering.

## ④ Drawing Conclusions:

- Machine Learning
- Inferential Statistics

} Data Scientist

- Descriptive Statistics

} Data Analyst

Till EDA, Data Scientist and Data Analyst have to cover all the steps.

### Note:

All the steps are not linear though!!

Dirty Data: Quality issues (Duplicates, missing, corrupt, incomplete)

Messy Data: Structural issues (each variable forms a column; observation form a row each observational unit forms a table)

Name      First Name  
            Surname  
Contact

John Doe  
Johndoe