

Data Science is generation of actionable knowledge directly from huge amount of complex data. The goal of data science is to gain insights and knowledge from any type of data — both structured and unstructured

Structured data:

Rows and columns

Name , age, height, weight, location, education, working, exp, tech skills →

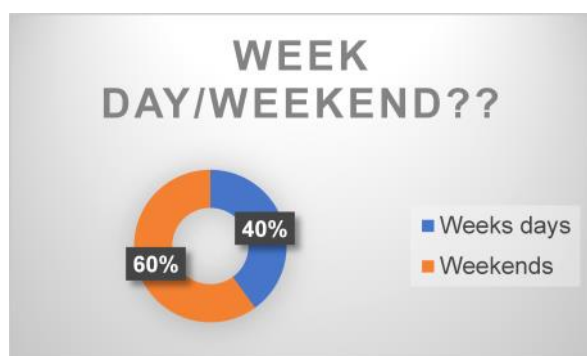
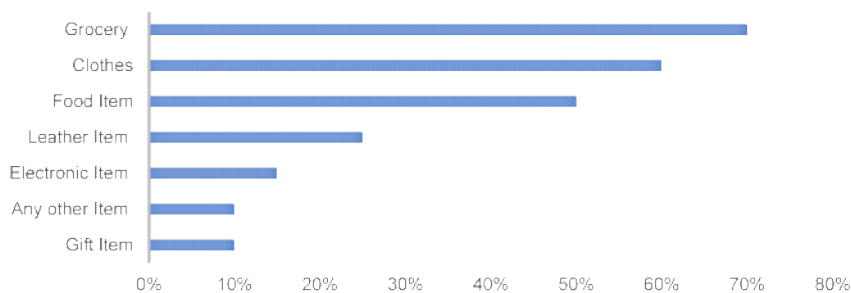
1
2
3
4
5

Banking, hospitals, universities, sales, insurance, telecom, pharamacy, tours and travels, bookmy show,

Unstructured data:

Any text data, google search, images, videos, comments, reviews, blogs, websites , feedback, whatstapp chat, scanned Reports, doctor's prescription,

SALES CATEGORY

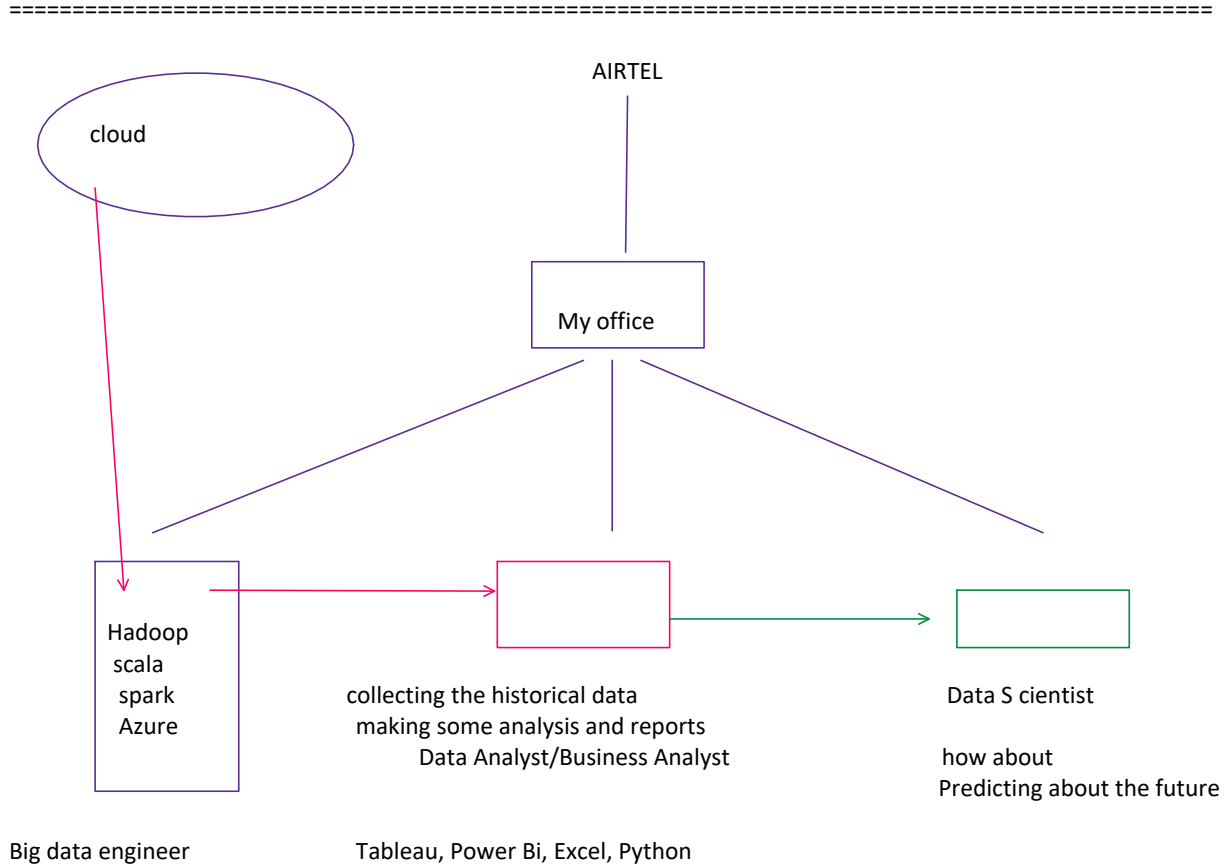
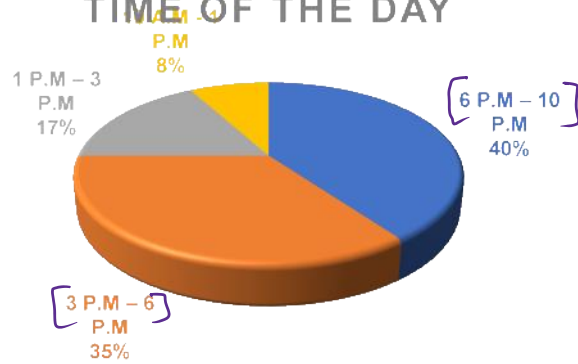


TIME OF THE DAY



12:00 → 10%

TIME OF THE DAY



$$Y = mx + C$$

Basis statistics --> 2 days

Python --> scratch --> 2 days

Statistics + python --> 1 days

Probability , Test of hypothesis --> 10-12 hours

Machine Learning:

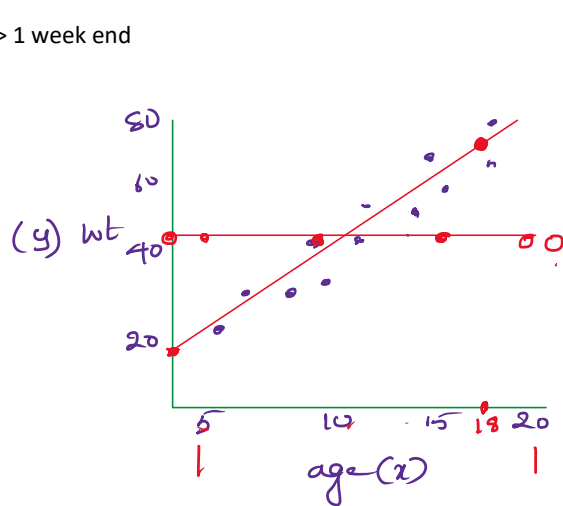
2 months -->

NLP --> 1 weekend

Deep learning --> 1 week end

Equation

Deep learning --> 1 week end



$$y = mx + c$$

$wt = m \cdot age + c$

\downarrow Slope \downarrow Intercept

Model →

$$70kg = 2(25) + 20$$

Value added courses:

Assignments

---> Project

AI classes --> 3 weeks

20 --> assignment team

Hadoop, spark

Azure

Tableau

Powerbi

Core python

SQL --> 10 hrs

Steps behind the project life cycle of data scientist:

1. Framing the problem:

2. SQL -> collect the data from the server to our local system

SQL to Python

3. Exploratory data analysis:

Structured data --> EDA

Unstructured --> preprocessing

Which columns are really important to us.

Dependent ← wt → age, ~~ht~~, ~~Gender~~, ~~Skills~~, ~~Salary~~

Dependent Salary → exp, ~~Skills~~, education, ~~ht~~?, ~~Gender~~

\downarrow

Y

Independent variables (X)

Scatter plot, box plot, histogram, bar graph

Data cleaning:

Data cleaning:



Boxplot

Data Transformation:

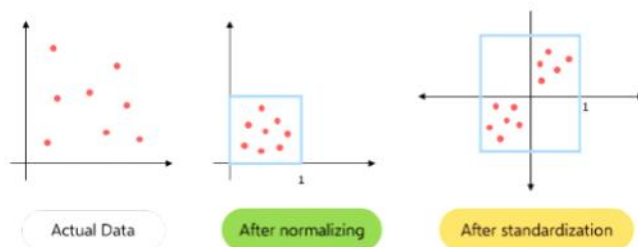
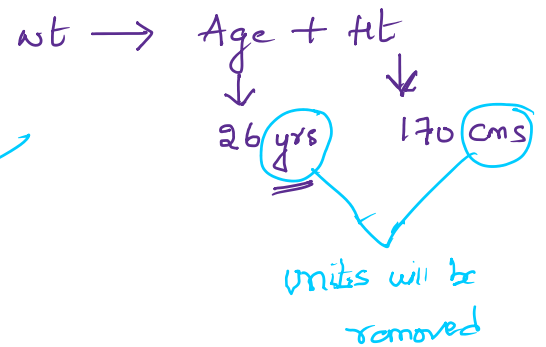
What is transformation?

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

\downarrow \downarrow \downarrow
 WL 1 1

$$= 10 + 26 + 170$$

206 Kgs \approx 206 $\begin{cases} \text{yrs} \\ \text{cms} \end{cases}$



Data partition:



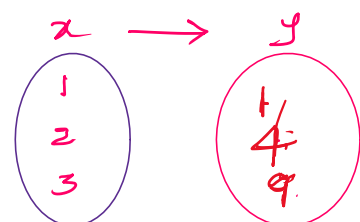
- ① Training $\rightarrow \frac{1000}{600} \rightarrow \underline{\text{model}}$
- ② Test $\rightarrow 400 \rightarrow \underline{\text{tested}}$

Selection of model:

Machine learning

$$f(x) = x^2 \quad \boxed{y = x^2}$$

$$f(x) = 2x$$



Brand is 'Y' $\rightarrow x^2$

$$f(x) = 2x$$

$$f(x) = 2x + 1$$

$$\text{Based of } Y' = \begin{matrix} x^2 \\ 2x \\ 2x+1 \end{matrix}$$

1. Linear Regression $\rightarrow 2x+1$
2. Logistic Regression
3. Support vector machine
4. Decision Trees
5. KNN classifier
6. Naïve Bayes classifier

Probability, matrices, distance method

Cross validation:

Model should be verified with various methods --> wrong results

Evaluation:

Deployment:

=====