

IT 307 ML

UNIT 1

Introduction to ML



Presented by Dr. Ankush Sawarkar
Email: adsawarkar@sggs.ac.in

FACULTY PROFILE

- Qualifications
 - Ph D in Machine Learning & Deep Learning
 - M.Tech in Computer Science & Engineering
 - B.E in Computer Science & Engineering
 - 10 year of experience
 - Data Science Researcher
 - Data Science Consultant
 - Software Developer
 - Assistant Professor
 - Corporate Trainer
 - Area of specialization
 - Data Science
 - AI - Machine Learning, Deep Learning
 - Natural Language Processing
 - Bioinformatics
 - Computer Vision
 - Visualization - Tableau
- Research Publication**
- 15 Research papers
 - 660+ citations
- <https://scholar.google.com/citations?hl=en&user=Y3CY1-wAAAAJ>
- Dr. Ankush D. Sawarkar (Ph.D. VNIT Nagpur)**
Assistant Professor, Dept. Information Technology,
Email| adsawarkar@sggs.ac.in / ankush1sawarkar@gmail.com
<https://orcid.org/my-orcid?orcid=0000-0001-7099-1987>
<https://www.webofscience.com/wos/author/record/ABE-6640-2020>
<https://www.scopus.com/authid/detail.uri?authorId=56669766500>
<https://www.researchgate.net/profile/Ankush-Sawarkar>
<https://sites.google.com/view/ankushsawarkar/home>

Books:

1. Kevin Murphy, **Machine Learning: A Probabilistic Perspective**, MIT Press, 2012.
2. **Machine Learning For Dummies** by John Mueller and Luca Massaron
3. Trevor Hastie, Robert Tibshirani, Jerome Friedman, **The Elements of Statistical Learning**, Springer 2009 (freely available online).
4. Christopher Bishop, **Pattern Recognition and Machine Learning**, Springer, 2007

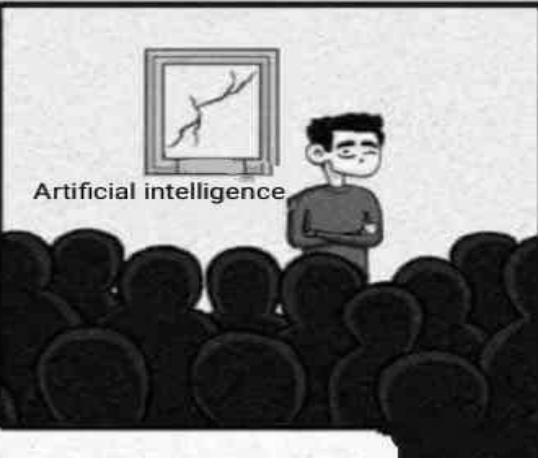
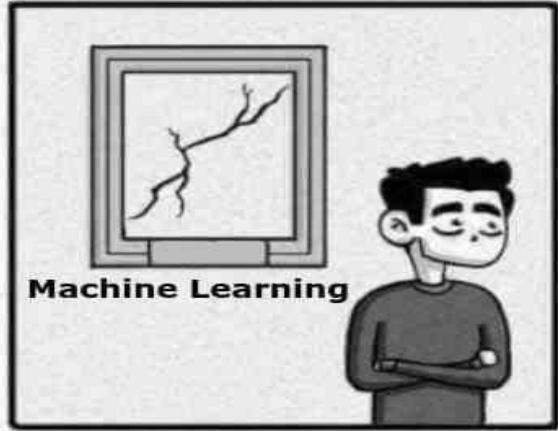
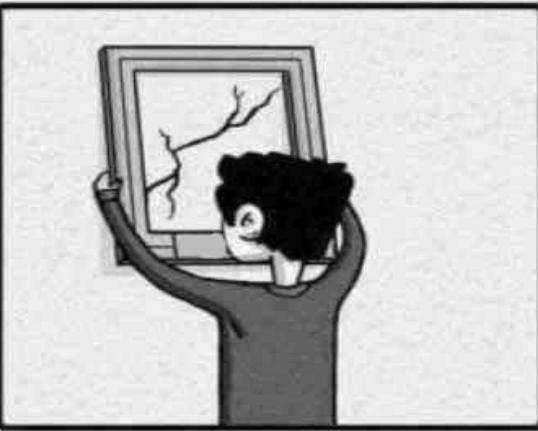
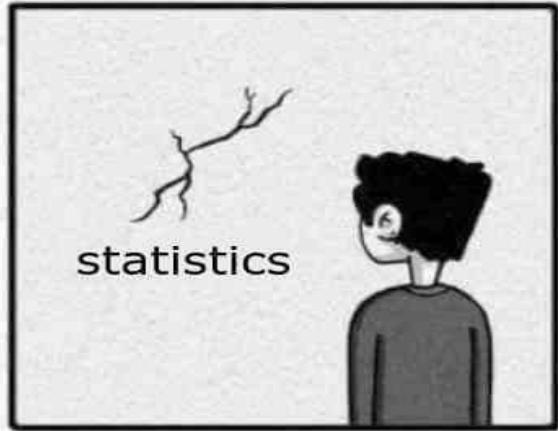
At the end of this course

- Machine Learning, Types of ML
- Regression/ Classification / Clustering
- Overfitting & Underfitting
- Model Selection - MAE/ RMSE/ Regularization
- Evaluation Metrics - Confusion Matrix (TP/TN/FP/FN)/ Accuracy/ Precision/ Recall/
F1Score
- Logistics Regression/ Sigmoid/ Limitation of Logistics Reg
- Decision Tree - Gini Index/ Information Gain
- K Fold Cross Validation/ K- Nearest Neighbors/ Random Forest
- Support Vector Machine (SVM) Tuning Parameter
- Clustering - Kmeans/ Elbow method/ Performance measures

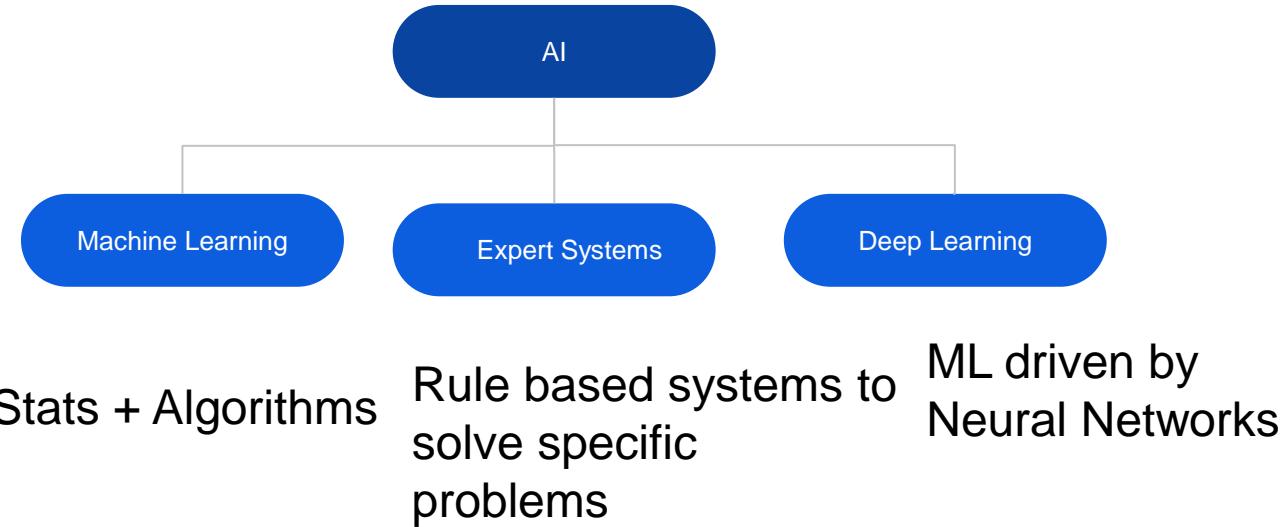
ML Techniques Overview

- “The field of study that gives computers the ability to learn without being explicitly programmed.” - Arthur Samuel
- “A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P if its performance at tasks in T, as measured by P, improves with experience E”. - Tom Mitchell

AI Landscape

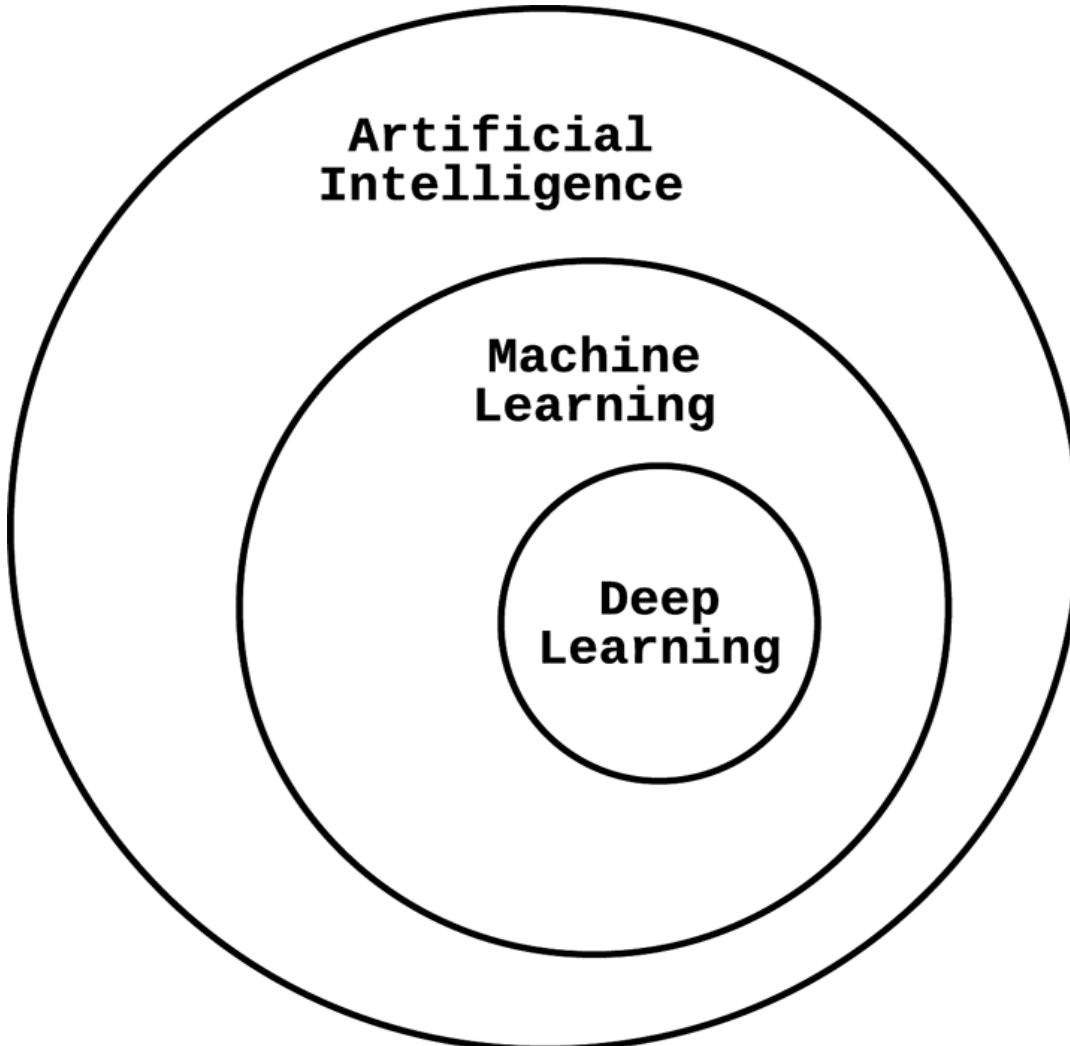


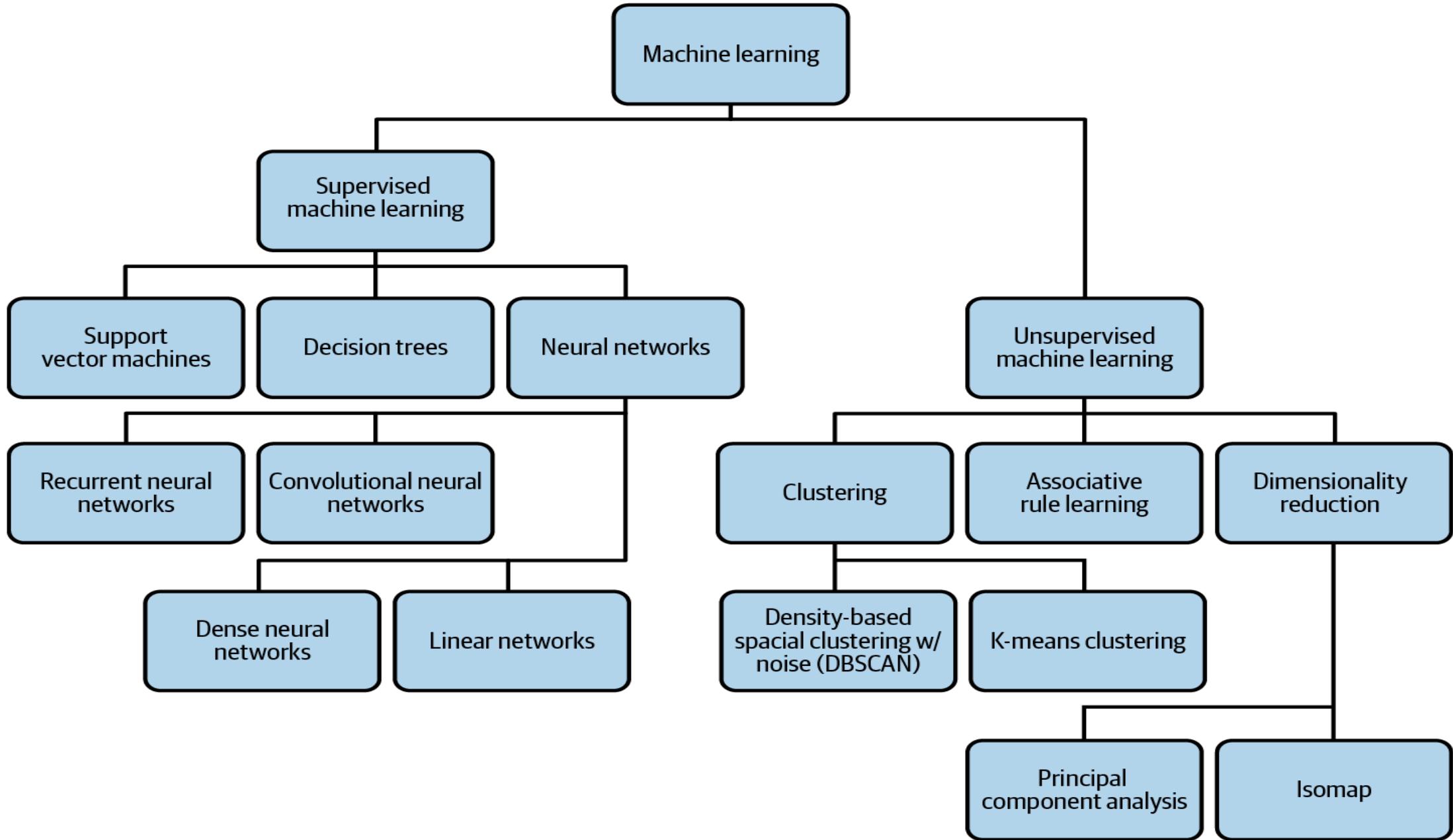
AI Landscape...



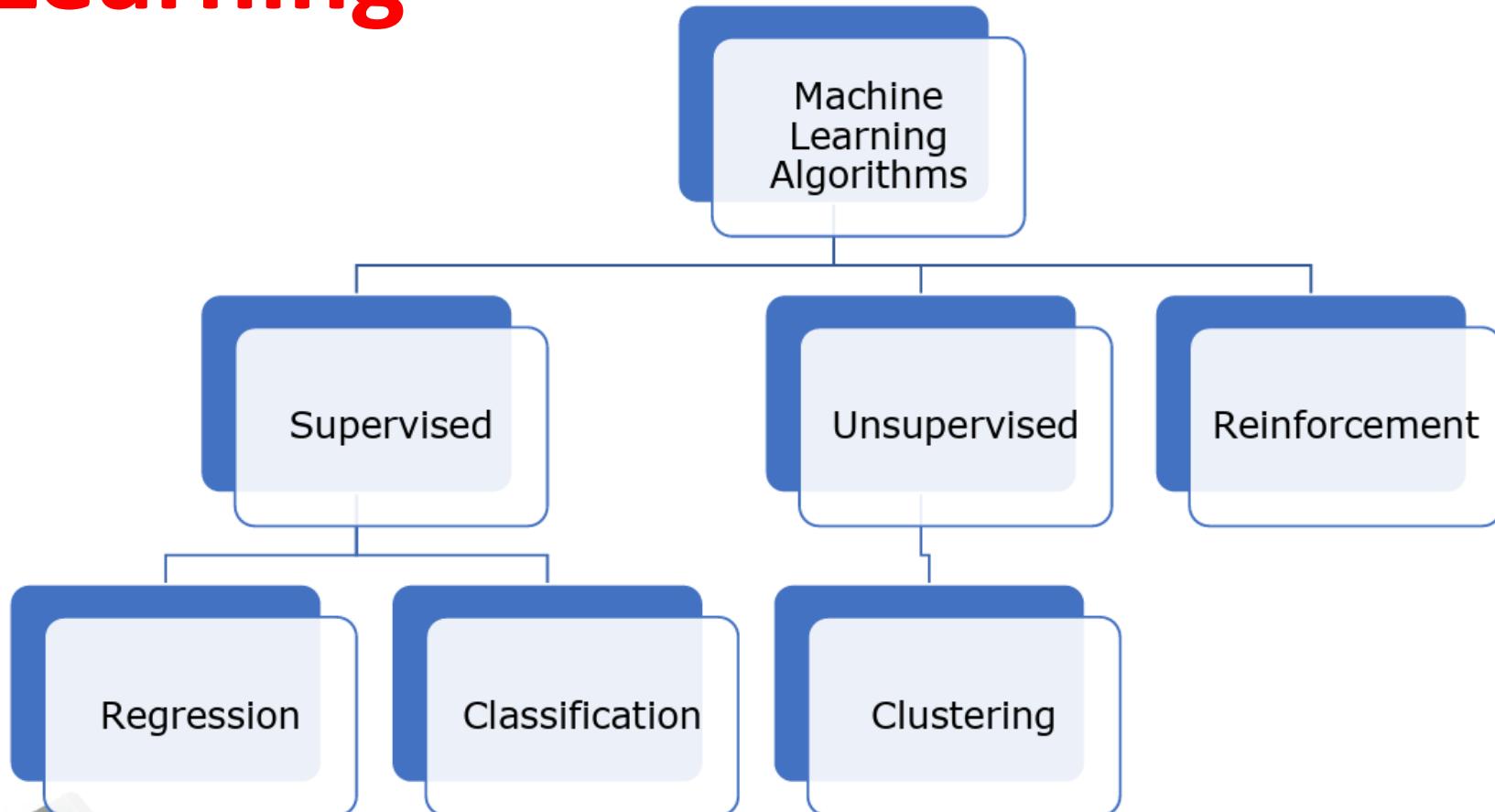
What are the factors driving the development of AI ?

Another way of looking





Machine Learning



Supervised Learning

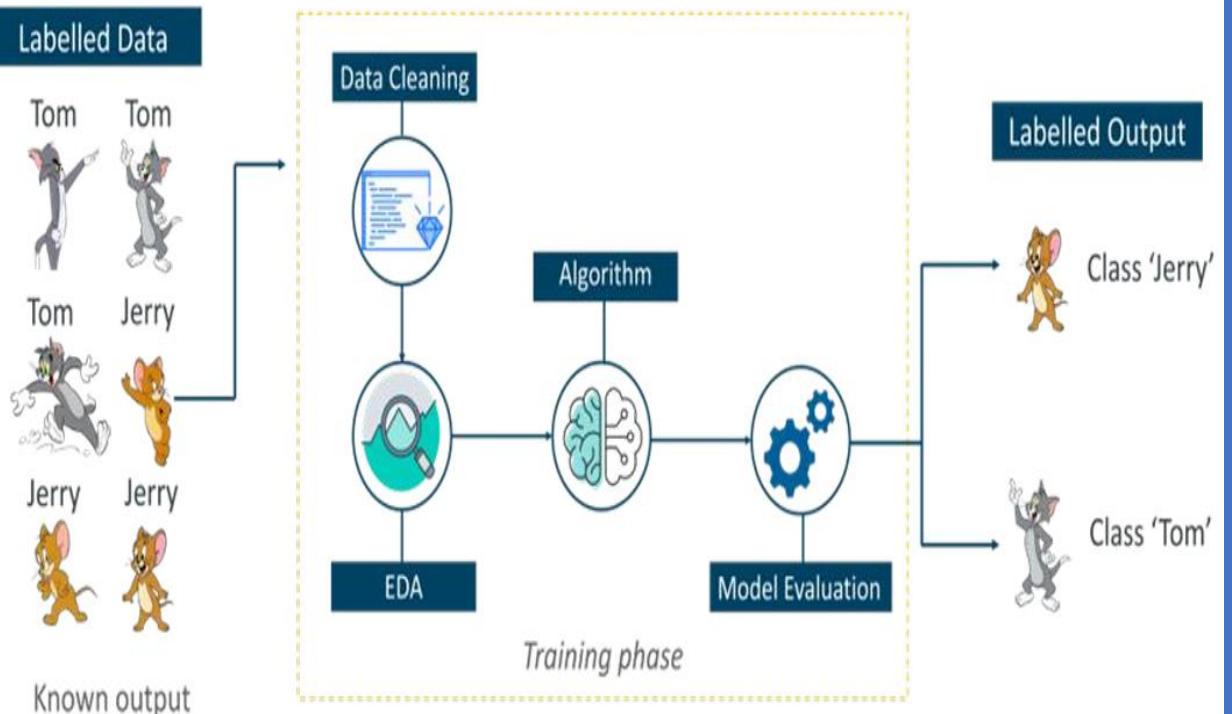


Fig. 2 Supervised Learning REF[1]

Unsupervised Learning

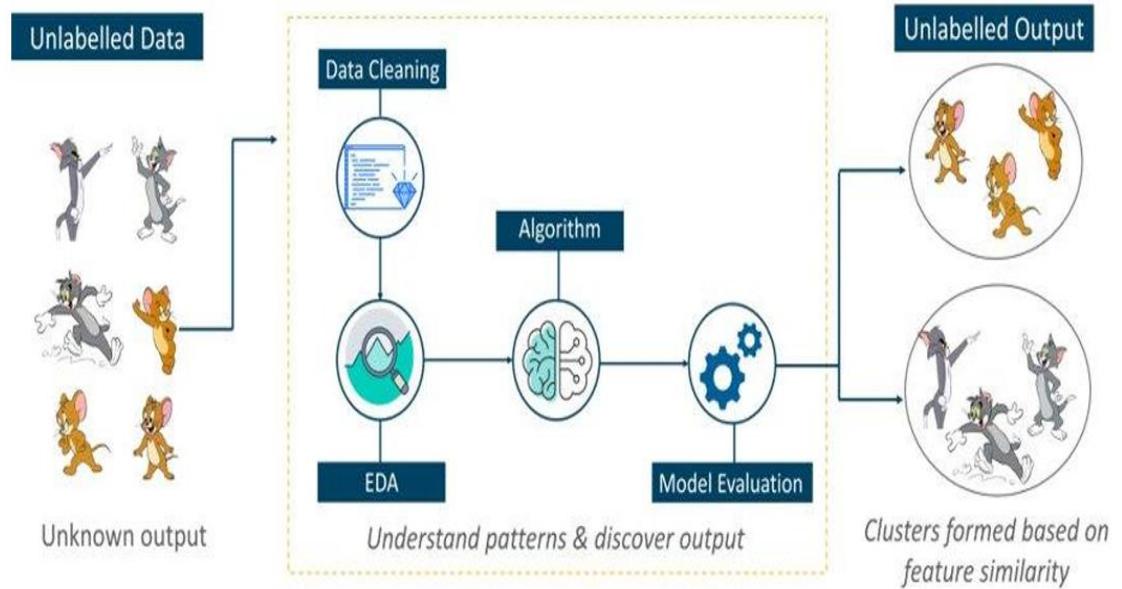


Fig. 3 Unsupervised Learning REF[1]

Supervised Learning

There are two categories of supervised learning:

- Classification task
- Regression task

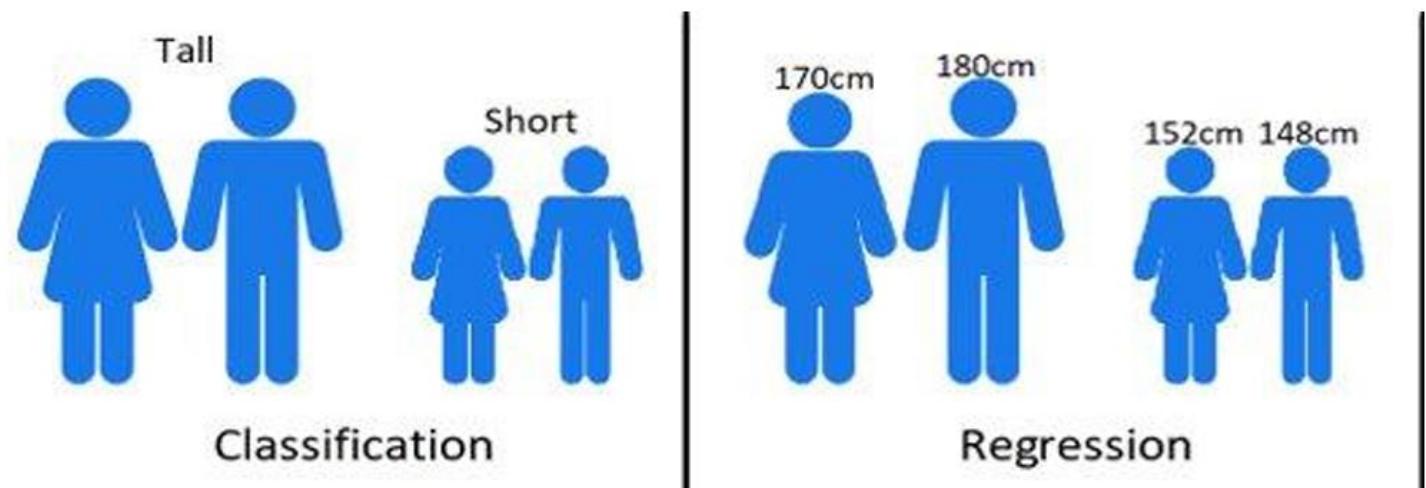


Fig. 5 Supervised Learning- Classification and Regression REF[3]



Unsupervised Learning- Clustering

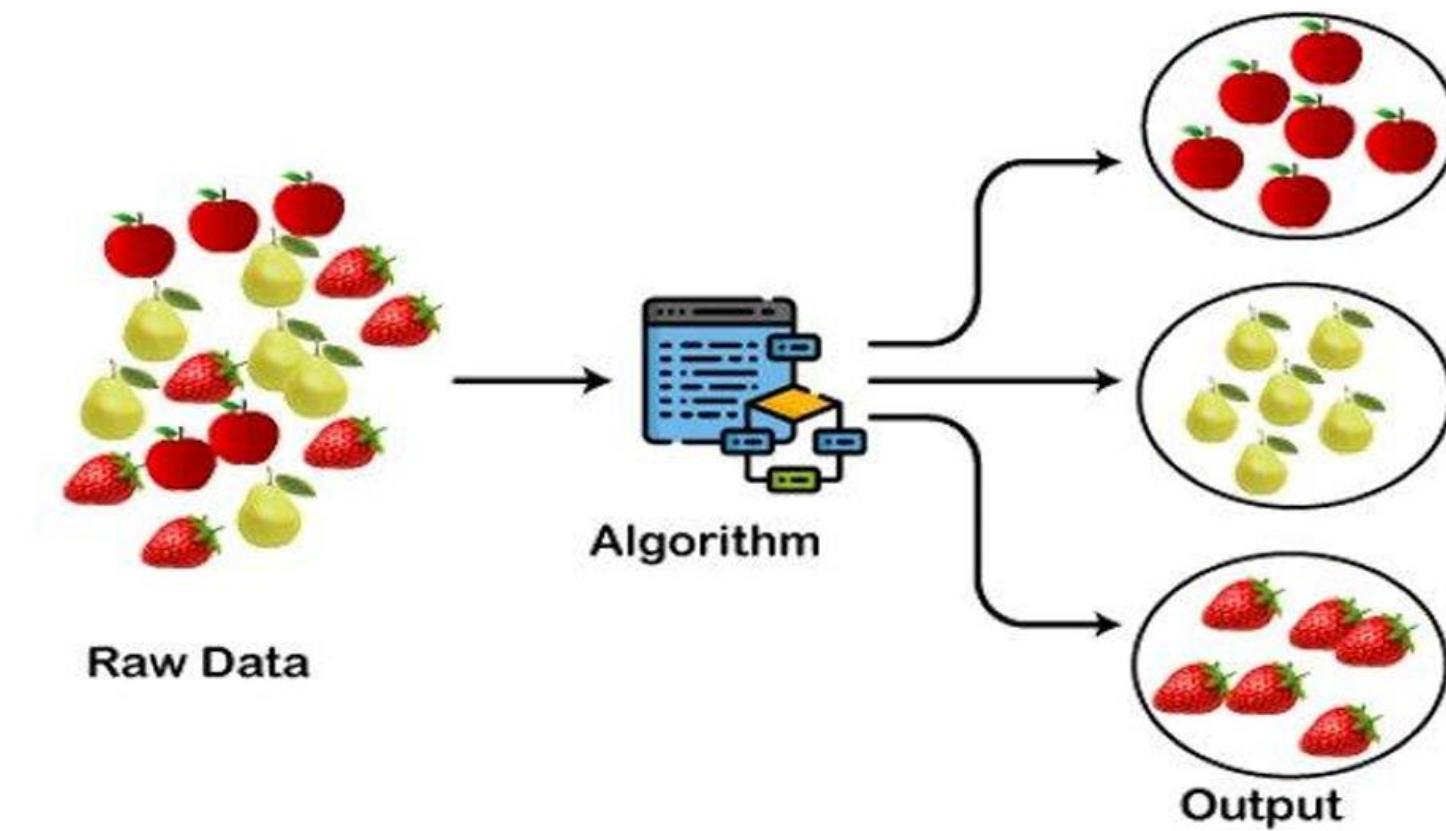


Fig.6 Clustering- Unsupervised Learning REF[2]

Characteristics Comparison

Regression

- Supervised Learning
- Output is a continuous quantity
- Main aim is to forecast or predict
- Eg: Predict stock market price
- Algorithm: Linear Regression

Classification

- Supervised Learning
- Output is a categorical quantity
- Main aim is to compute the category of the data
- Eg: Classify emails as spam or non-spam
- Algorithm: Logistic Regression

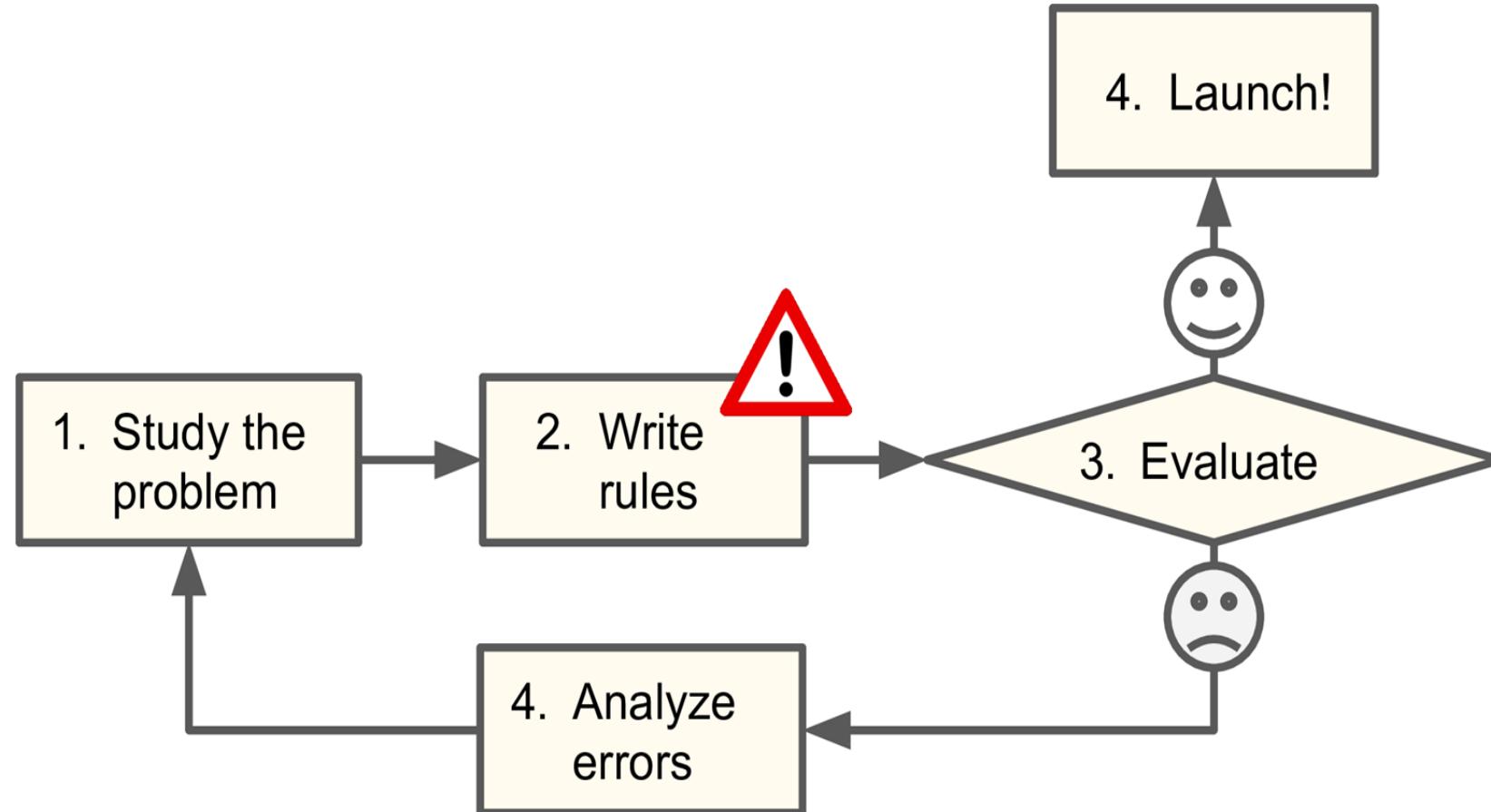
Clustering

- Unsupervised Learning
- Assigns data points into clusters
- Main aim is to group similar items clusters
- Eg: Find all transactions which are fraudulent in nature
- Algorithm: K-means

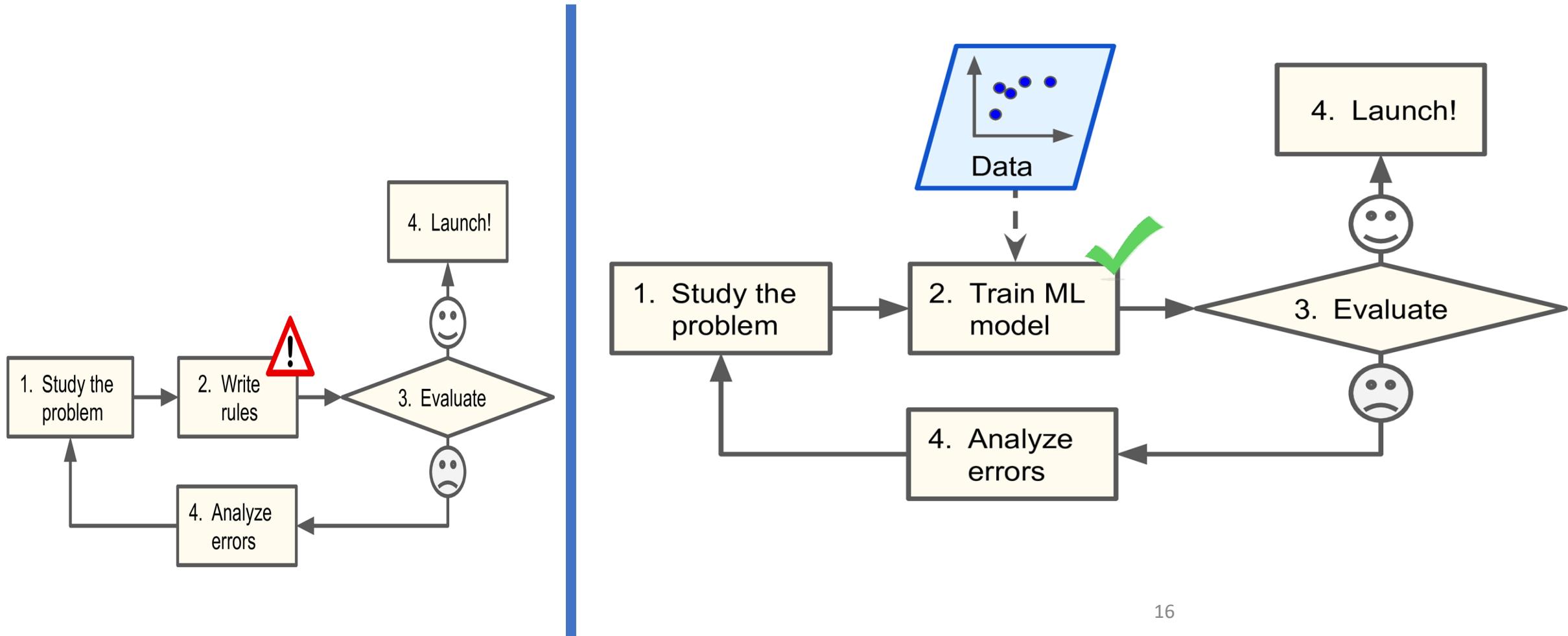


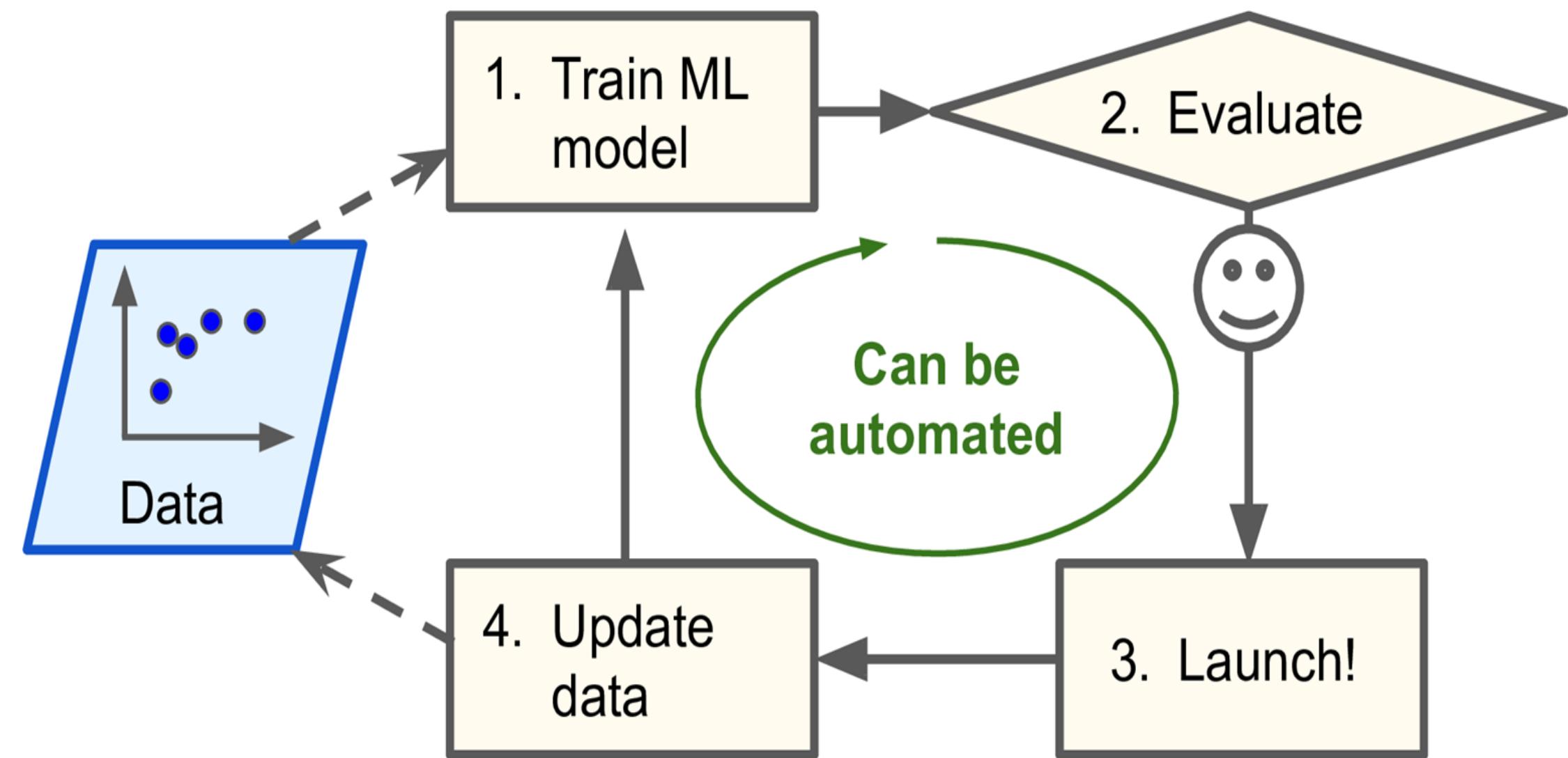
Fig. 7 Characteristic Comparison REF[1]

Why do we need ML? The Email Spam Problem

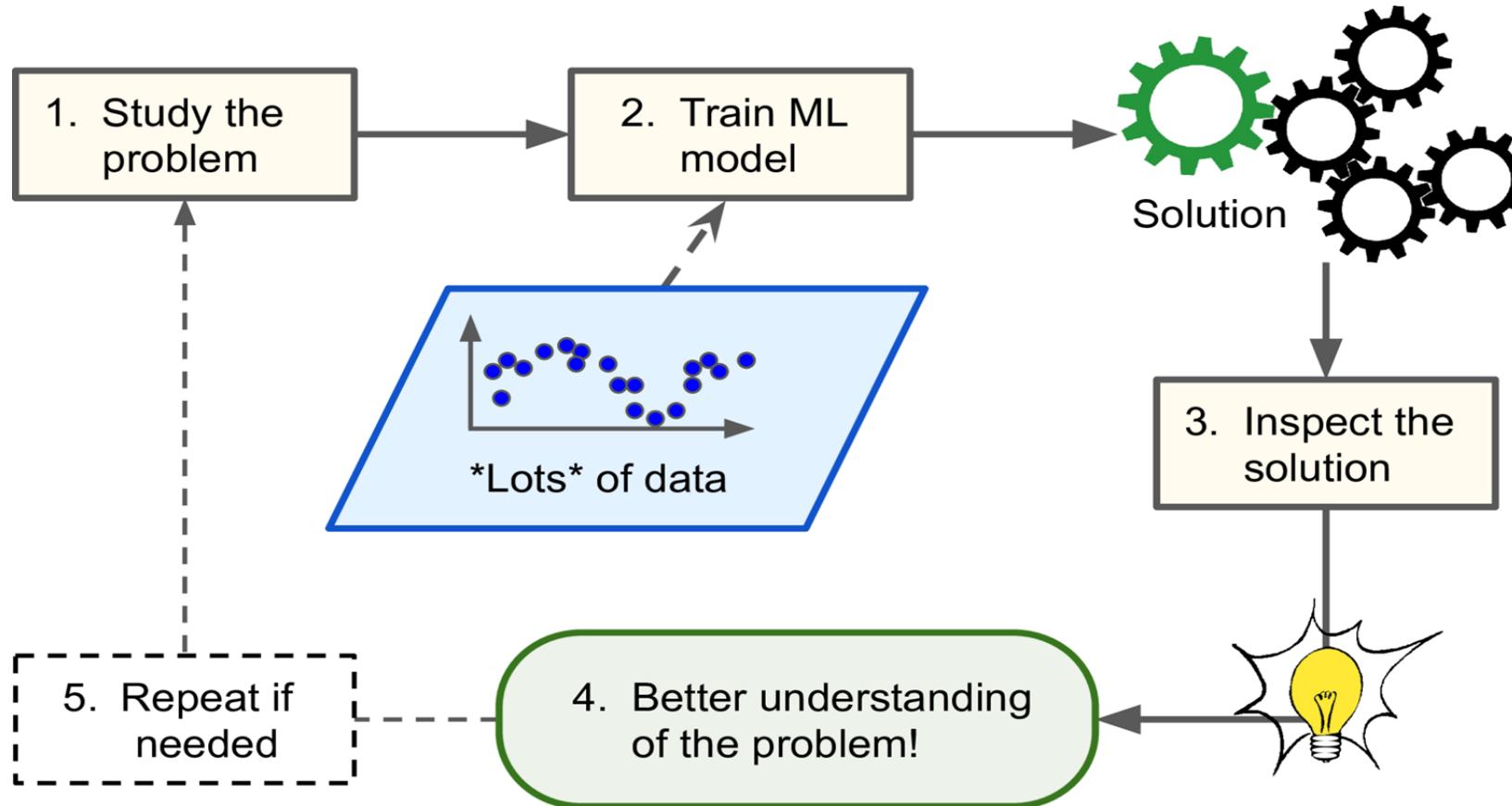


Why do we need ML? The Email Spam Problem





Help Humans Learn



For instance, once a spam filter has been trained on enough spam, it can easily be inspected to reveal the list of words and combinations of words that it believes are the best predictors of spam. Sometimes this will reveal unsuspected correlations or new trends, and thereby lead to a better understanding of the problem.

Where to use ML?

- Problems for which existing solutions require a lot of fine-tuning or long lists of rules: one Machine Learning model can often simplify code and perform better than the traditional approach.
- Complex problems for which using a traditional approach yields no good solution: the best Machine Learning techniques can perhaps find a solution.
- Fluctuating environments: a Machine Learning system can easily be re-trained on new data, always keeping it up to date.
- Getting insights about complex problems and large amounts of data.

Some Typical Applications

- Analyzing images of products on a production line to automatically classify them
- Detecting tumors in brain scans
- Forecasting your company's revenue next year, based on many performance metrics
- Segmenting clients based on their purchases so that you can design a different marketing strategy for each segment
- Recommending a product that a client may be interested in, based on past purchases

Questions to be kept in mind when designing an ML System?

What question(s) am I trying to answer? Do I think the **data collected** can answer that question?

What is the best way to phrase my question(s) as a machine learning problem?

Have I collected enough data to represent the problem I want to solve?

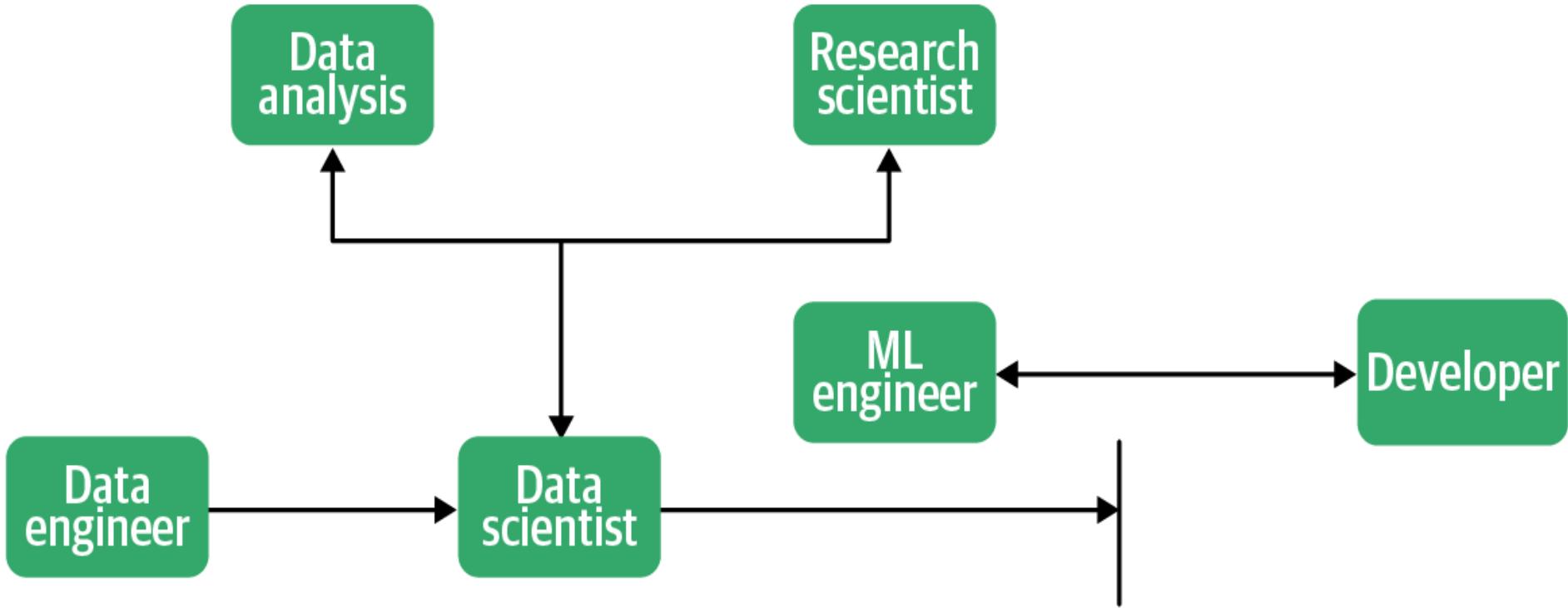
What **features of the data did I extract**, and will these enable the right predictions?

How will I **measure success** in my application?

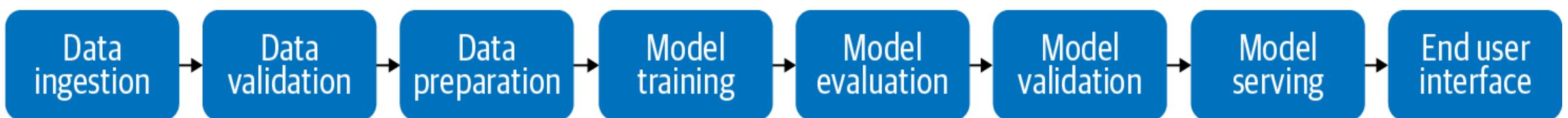
How will the machine learning solution interact with other parts of my research or business product?

Different Roles in an Organization's ML Model Development Process

Business



Engineering



Machine Learning Process

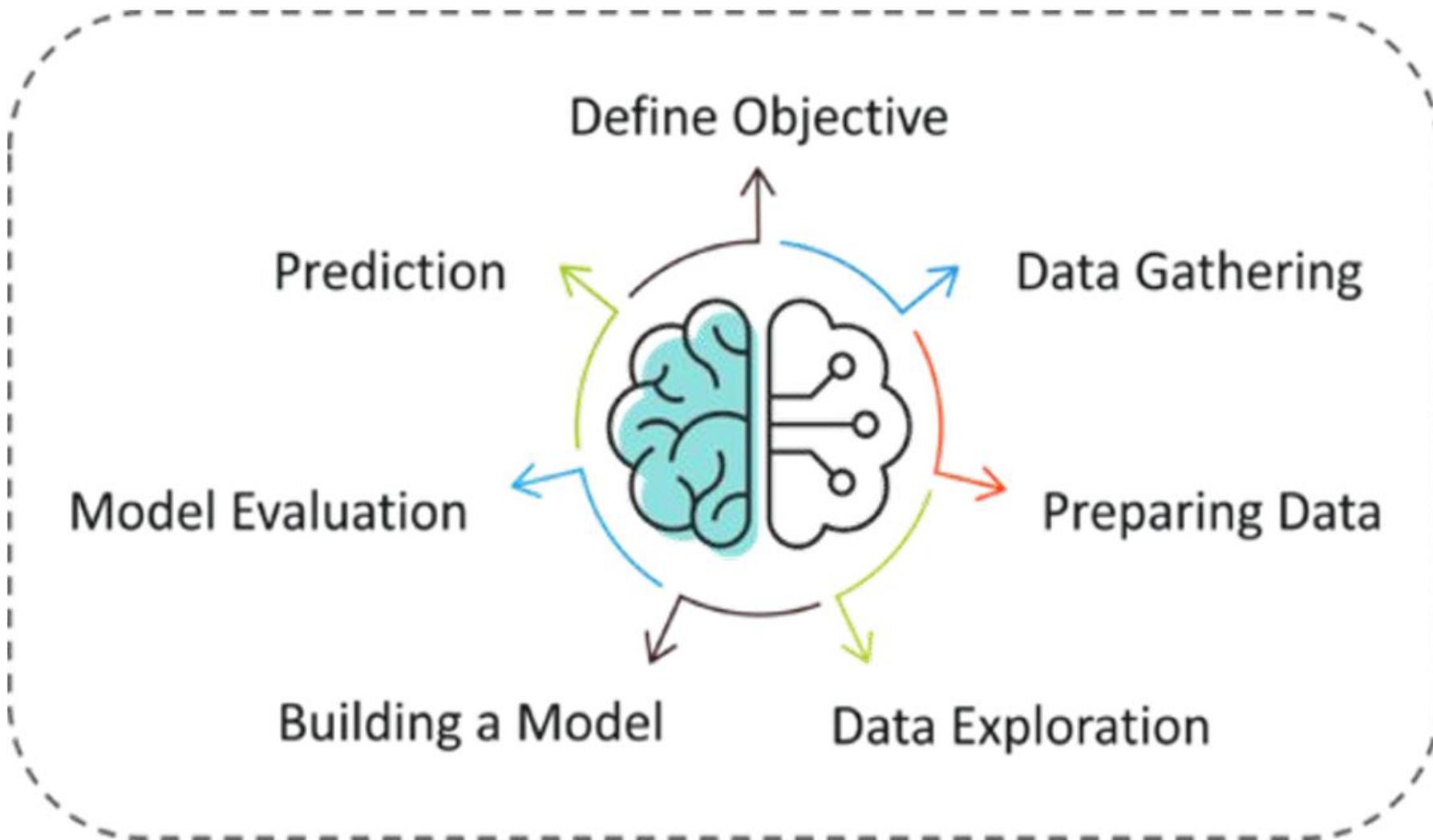


Fig. 1 Machine Learning process REF[1]

Machine Learning Key Terms

Algorithm: (Logic) A Machine Learning algorithm is a set of rules and statistical techniques |

Model: (Main component of Machine Learning) Algorithm maps how the decision is taken |

Machine Learning Key Terms

- **Predictor Variable:** It is a feature(s)/ input of the data that can be used to predict the output.
- **Response Variable:** It is the feature or the output variable that needs to be predicted by using the predictor variable(s).
- **Training Data:** The machine learning model is built using the training data (Learning).
- **Testing Data:** After the model is trained, it must be tested to evaluate how accurately it can predict an outcome.

Machine Learning Key Terms

Hours- Independent Variable → X→ Predictor

Scores – Dependent Variable→ Y→ Response

student_score.csv

Hours	Scores
2.5	21
5.1	47
3.2	27
8.5	75
3.5	30
1.5	20
9.2	88
5.5	60
8.3	81
2.7	25
7.7	85
5.9	62
4.5	41
3.3	42
1.1	17
8.9	95
2.5	30
1.9	24
6.1	67
7.4	69
2.7	30
4.8	54
3.8	35
6.9	76
7.8	86

Machine Learning Key Terms

Country, Age, Salary

Independent Variable → X → Predictor

Purchased

Dependent Variable → Y → Response

Country	Age	Salary	Purchased
France	44	72000	No
Spain	27	48000	Yes
Germany	30	54000	No
Spain	38	61000	No
Germany	40		Yes
France	35	58000	Yes
Spain		52000	No
France	48	79000	Yes
Germany	50	83000	No
France	37	67000	Yes

Machine Learning Key Terms

Algorithm: (Logic) A Machine Learning algorithm is a set of rules and statistical techniques

Model: (Main component of Machine Learning) Algorithm maps how the decision is taken

Machine Learning Key Terms

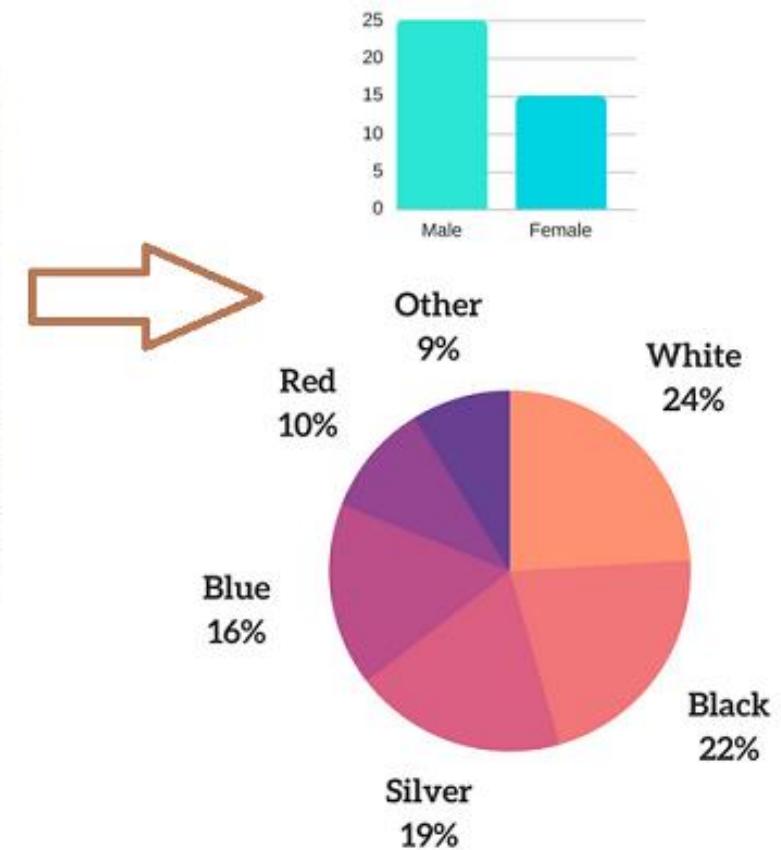
- **Predictor Variable:** It is a feature(s)/ input of the data that can be used to predict the output.
- **Response Variable:** It is the feature or the output variable that needs to be predicted by using the predictor variable(s).
- **Training Data:** The machine learning model is built using the training data (Learning).
- **Testing Data:** After the model is trained, it must be tested to evaluate how accurately it can predict an outcome.

Descriptive Statistics Examples

- Descriptive statistics about a class involve the average score of marks obtained by students for different courses including Maths, Physics, Chemistry.
- Analyzing favorite color of CAR liked by students in a class and representing it in the form of graph or pie chart.

	A	B	C	D
1	Respondent Number	Age	Gender	Favorite Car Color
2	1	22	M	White
3	2	37	F	Silver
4	3	45	F	Black
5	4	62	F	Gray
6	5	28	M	Red
7	6	45	M	Green
8	7	88	F	Brown
9	8	61	M	White
10	9	95	M	Black
11	10	27	M	White
12	11	39	F	Green
13	12	43	M	Brown
14	13	55	F	Black
15	14	59	F	White

RAW DATA



Descriptive Statistics

Descriptive Statistics Types

- Measures of Central Tendency
- Measures of Dispersion or Variation

	A	B	C	D	E
1	Student	Scores		Scores	
2	Student 1	72		Mean	70.2
3	Student 2	91		Standard Error	5.050412524
4	Student 3	77		Median	74.5
5	Student 4	80		Mode	72
6	Student 5	72		Standard Deviation	15.9708067
7	Student 6	46		Sample Variance	255.0666667
8	Student 7	81		Kurtosis	-1.005270166
9	Student 8	54		Skewness	-0.646893849
10	Student 9	83		Range	45
11	Student 10	46		Minimum	46
12				Maximum	91
13				Sum	702
14				Count	10
15					
16					

Descriptive Statistics Types

- Measures of Central Tendency

- Mean:

- Formula 1: Sum of all values/ No. of samples= $\Sigma (x) / N$

Student	Scores
Student 1	72
Student 2	91
Student 3	77
Student 4	80
Student 5	72
Student 6	46
Student 7	81
Student 8	54
Student 9	83
Student 10	46

Sum= 702, N=10, then what is Mean?

Descriptive Statistics Types

- Measures of Central Tendency

- Mean:
- Formula 1: Sum of all values/ No. of samples= $\Sigma (x) / N$
- If values are repetitive, grouped data and getting frequency, Mean is expressed as

$$\text{Formula 2: Mean} = \sum(x_i \cdot f_{x_i}) / N$$

- Notations:

- For Population- μ
- For sample- \bar{x}

Population: The group we are interested in studying

Sample: A subset of Population

x	f _x	x.f _x
72	2	=72*2 = 144
91	1	91
77	1	77
80	1	80
46	2	= 46*2 = 92
81	1	81
54	1	54
83	1	83
	$\sum x_i = 10$	$\sum (x \cdot f_{x_i}) = 720$

Student	Scores
Student 1	72
Student 2	91
Student 3	77
Student 4	80
Student 5	72
Student 6	46
Student 7	81
Student 8	54
Student 9	83
Student 10	46

Descriptive Statistics Types

Student	Scores
Student 1	72
Student 2	91
Student 3	77
Student 4	80
Student 5	72
Student 6	46
Student 7	81
Student 8	54
Student 9	83
Student 10	46

- Measures of Central Tendency

- Median: Middle observation
- Arrange in ascending order: 46, 46, 54, 72, 72, 77, 80, 81, 83, 91
- If N is even- Median is mean of [$N/2$ th and $(N/2+1)$ th] elements
46, 46, 54, 72, 72, 77, 80, 81, 83, 91
- If N is odd- Median is median is $(N+1)/2$ th element
46, 46, 54, 72, 72, 77, 80, 81, 83

- Measures of Central Tendency

- Mode: Most repetitive value
- If there is clear such unique value

Descriptive Statistics Types

• Measures of Central Tendency

- Mode: Most repetitive value

Student	Scores
Student 1	72
Student 2	91
Student 3	77
Student 4	80
Student 5	72
Student 6	46
Student 7	81
Student 8	54
Student 9	83
Student 10	46

If there is no such clear unique value


$$\text{Mode formula} = L + h \frac{(f_m - f_1)}{(f_m - f_1) + (f_m - f_2)}$$

where,

- 'L' is the lower limit of the modal class.
- 'h' is the size of the class interval.
- ' f_m ' is the frequency of the modal class.
- ' f_1 ' is the frequency of the class preceding the modal class.
- ' f_2 ' is the frequency of the class succeeding the modal class.

Class L to U	f_{Class}
30-40	0
40-50	2
50-60	1
60-70	0
70-80	4
80-90	2
90-100	1
	$\sum x_i = 10$

Central Tendency Question

- The marks scored by students are given here:

Marks scored X_i	No. of people scoring it Freq of X_i , as F_{X_i}
1	10
2	12
3	20
4	15
5	10

- Compute Mean, Median and mode F_{X_i}

Descriptive Statistics Types

- Measures of Central Tendency
- Measures of Dispersion or Variation

	A	B	C	D	E
1	Student	Scores		Scores	
2	Student 1	72		Mean	70.2
3	Student 2	91		Standard Error	5.050412524
4	Student 3	77		Median	74.5
5	Student 4	80		Mode	72
6	Student 5	72		Standard Deviation	15.9708067
7	Student 6	46		Sample Variance	255.0666667
8	Student 7	81		Kurtosis	-1.005270166
9	Student 8	54		Skewness	-0.646893849
10	Student 9	83		Range	45
11	Student 10	46		Minimum	46
12				Maximum	91
13				Sum	702
14				Count	10
15					
16					

Measures of Dispersion or Variation

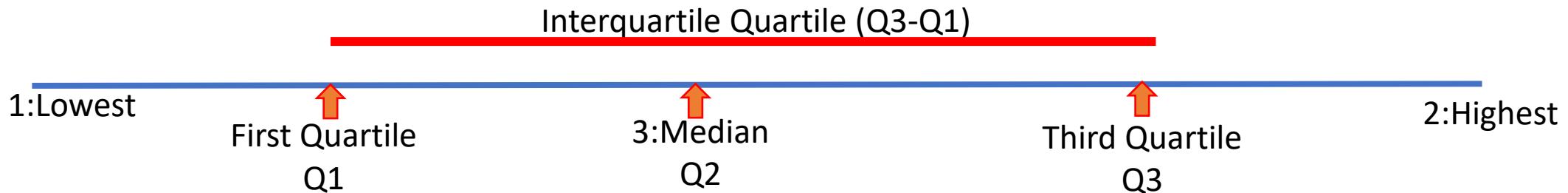
- Dispersion in statistics describes the **spread** of the data values in a given dataset.
 - It reveals the **extent to which the values** of the individual items **differ** in a data set.
- Minimum : 46**
- Maximum:91**
- Sum:702**
- Count:10**
- Range: Highest-lowest : 45**

Sample standard error = $\frac{\sigma}{\sqrt{n}}$, if σ is known

A	B	C	D	E
1	Student	Scores	Scores	
2	Student 1	72		
3	Student 2	91		
4	Student 3	77		
5	Student 4	80		
6	Student 5	72		
7	Student 6	46		
8	Student 7	81		
9	Student 8	54		
10	Student 9	83		
11	Student 10	46		
12				
13				
14				
15				
16				

Measures of Dispersion or Variation

- **Quartiles:** The quartile measures the **spread of values** above and below the median by dividing the distribution into **four groups**.



- It is a measure of variability around the median.
- A quartile divides data into three points—a lower quartile (Q1), median (Q2), and upper quartile (Q3)—to form four groups of the dataset.
- **Interquartile range:** $Q3 - Q1$

Find out Q1, Q2 and Q3 for given dataset

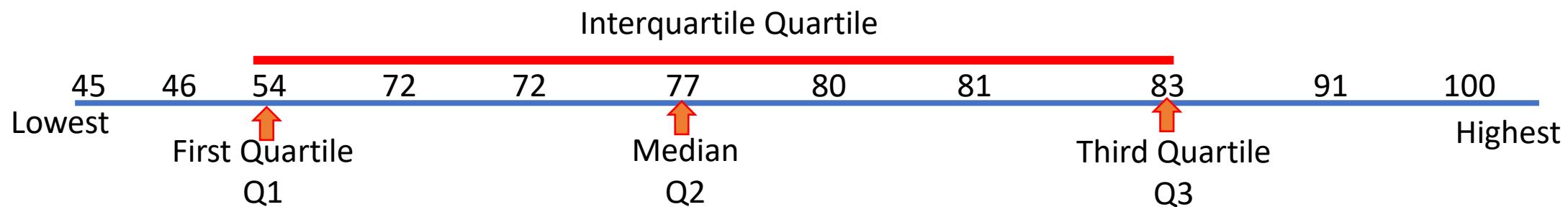
45 100 54 72 91 80 77 81 46 83 72

Measures of Dispersion or Variation

- Steps to find Quartile: Find out Q1, Q2 and Q3 for given dataset

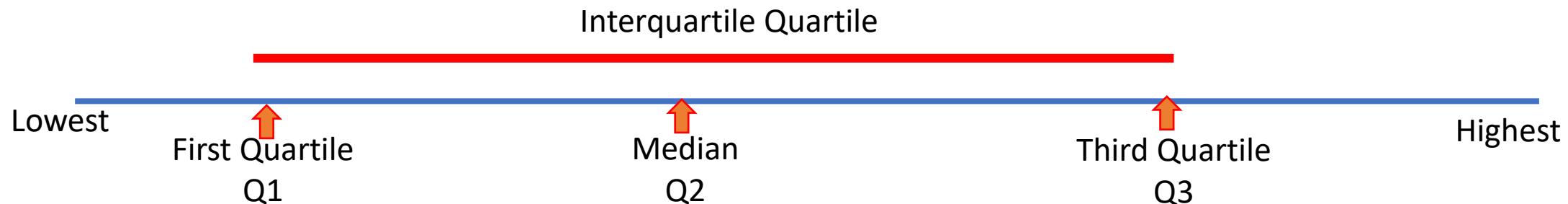
45 100 54 72 91 80 77 81 46 83 72

- 1. Arrange data in order from lowest to highest.
- 2. Find Median (Q2) = $(n+1)*50\% = (n+1)*0.5 = (n+1)/2 = (11+1)/2 = 6\text{th}$
- 3. Find the median of the data values that fall below Q2, that gives Q1= $(n+1)*25\% = (n+1)*0.25 = (n+1)/4 = (11+1)/4 = 3\text{rd}$
- 4. Find the median of the data values that fall above Q2, that gives Q3 = $(n+1)*75\% = (n+1)*0.75 = (n+1)*(3/4) = (11+1)*(3/4) = 9^{\text{th}}$



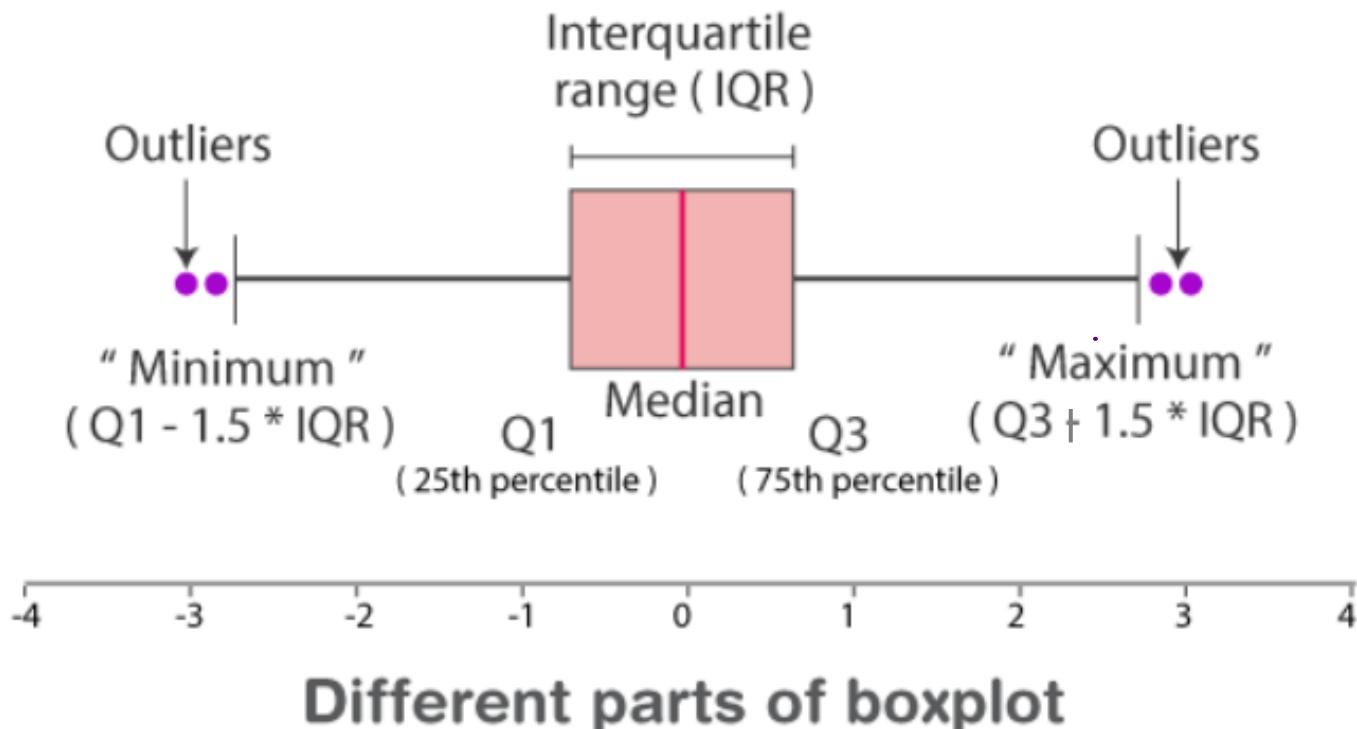
Measures of Dispersion or Variation

- Steps to find Quartile:
- 1. Arrange data in order from lowest to highest.
- 2. Find Median (Q2)
- 3. Find the median of the data values that fall below Q2, that gives Q1
- 4. Find the median of the data values that fall above Q2, that gives Q3
- [1,2,3,4,4,5,6,8,9,10,12]



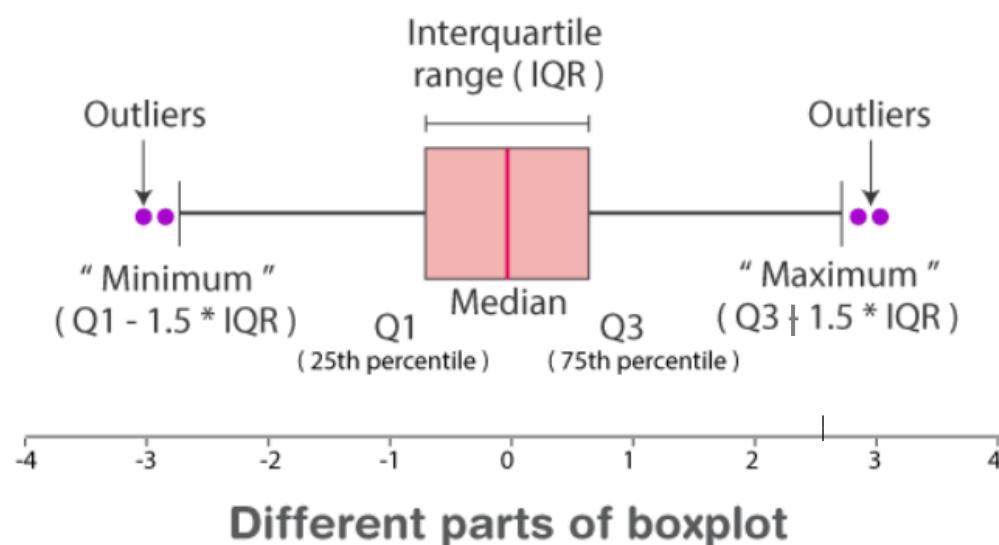
Boxplot

- It is a method for graphically demonstrating the locality, spread and skewness groups of numerical data through their quartiles.
- Box plot (called whiskers) indicating variability outside the upper and lower quartiles.



Boxplot for detecting Outliers

- **Minimum:** The minimum value in the given dataset
- **First Quartile (Q1):** The first quartile is **the median of the lower half** of the data set.
- **Median:** The median is the middle value of the dataset, which divides the given **dataset** into two equal parts. The median is considered as the second quartile.
- **Third Quartile (Q3):** The third quartile is **the median of the upper half of the data**.
- **Maximum:** The maximum value in the given dataset.



- **Interquartile Range (IQR):** The difference between the third quartile and first quartile is known as the interquartile range. (i.e.)
$$IQR = Q3 - Q1$$
- **Outlier:** The data that **falls on the extreme left or right side** of the ordered data is called **outliers**.
- Generally, the outliers fall more than the specified distance from the first and third quartile.
- Outliers are greater than $Q3 + (1.5 \times IQR)$ or less than $Q1 - (1.5 \times IQR)$.

Measures of Dispersion or Variation

Measures of Dispersion with RANGE

- **Standard Deviation:**
 - It provides information on how much variation from the mean exists.
 - However, the standard deviation shows how each value in a dataset varies from the mean.
- **Variance:** Square of standard deviation

Sample Standard Deviation, s :

$$s = \sqrt{\frac{\sum(x_i - \bar{x})^2}{n-1}}$$

Population Standard Deviation, σ :

$$\sigma = \sqrt{\frac{\sum(x_i - \mu)^2}{n}}$$

Sample Variance, s^2 :

$$s^2 = \frac{\sum(x_i - \bar{x})^2}{n-1}$$

Population Variance, σ^2 :

$$\sigma^2 = \frac{\sum(x_i - \mu)^2}{n}$$

Variance and Standard Deviation

Find the population variance for the following set of numbers: 26, 29, 30, 38, 32.

1. Find the Mean $\mu = (26 + 29 + 30 + 38 + 32) / 5.0 = 155/5 = 31$
2. Complete the table.

x_i	$(x_i - \mu)$	$(x_i - \mu)^2$
26	-5	25
29	-2	4
30	-1	1
38	7	49
32	1	1

3. Add all numbers of column 3: $25 + 4 + 1 + 49 + 1 = 80$
4. Divide it by the number of items in your data set: $80 / 5 = 16$. Thus 16 is the population variance, σ^2 for this set of data.
5. Standard Deviation, σ is square root of variance: $\sigma = \sqrt{16} = 4.0$

Population Variance, σ^2 :
$$\sigma^2 = \frac{\sum(x_i - \mu)^2}{n}$$

Population Standard Deviation, σ :

$$\sigma = \sqrt{\frac{\sum(x_i - \mu)^2}{n}}$$

Covariance and Correlation

- Covariance: Shows how the two variables differ
- Correlation: Shows how the two variables are related

Covariance and Correlation

- Covariance: Shows how the two variables differ

$$s_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n - 1}$$

Range: -Inf to +Inf

Covariance of x with itself is
Variance of x

- The numerical value of covariance does not have any significance however if it is positive then both variables vary in the same direction else if it is negative then they vary in the opposite direction.

- Correlation: Shows how the two variables are related

$$\text{Corr}(x, y) = \frac{\text{Cov}(x, y)}{\sigma_x \cdot \sigma_y} = \frac{s_{xy}}{\sigma_x \cdot \sigma_y}$$

Range: -1 to +1

If value is -1/+1: Linear relationship
If value is 0: No linear relationship

Covariance Computation

$$s_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n - 1}$$

For given data, let us find
Covariance.

1. Calculate the mean (average) prices for each asset.

Mean(S&P), $m_1 = 2044.80$,

Mean(ABC Corp.) $m_2 = 109.20$,

2. For each security, find the difference between each value and mean price.

Each entry in **column a** is entry in (column S&P 500) – m_1 and

Each entry in **column b** is entry in (column S&P 500) – m_2

	S&P 500	ABC Corp.
2013	1,692	68
2014	1,978	102
2015	1,884	110
2016	2,151	112
2017	2,519	154

	S&P 500	ABC Corp.	a	b	a x b
2013	1,692	68	-352.80	-41.20	14,535.36
2014	1,978	102	-66.80	-7.20	480.96
2015	1,884	110	-160.80	0.80	-128.64
2016	2,151	112	106.20	2.80	297.36
2017	2,519	154	474.20	44.80	21,244.16
Mean	2,044.80	109.20	Sum		36,429.20

Covariance Computation

$$s_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n - 1}$$

For given data, let us find Covariance.

3. Multiply the results obtained in the previous step as $a \times b$.
4. Take the sum of all entries in column $a \times b$.
5. Using the number calculated in step 4, find the covariance.

$$\text{Cov}(\text{S\&P 500}, \text{ABC Corp.}) = \frac{36,429.20}{5 - 1} = 9,107.30$$

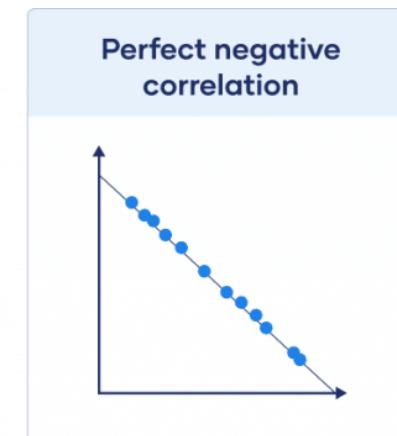
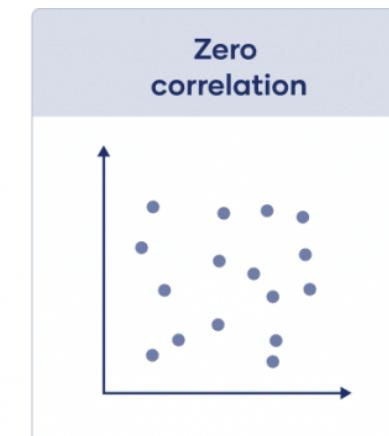
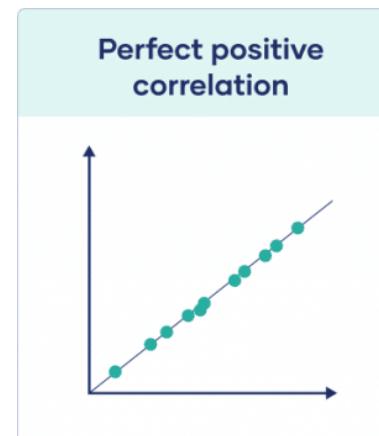
The **positive value** of covariance indicates that the price of the stock and the S&P 500 tend to move in the **same direction**.

	S&P 500	ABC Corp.
2013	1,692	68
2014	1,978	102
2015	1,884	110
2016	2,151	112
2017	2,519	154

	S&P 500	ABC Corp.	a	b	a x b
2013	1,692	68	-352.80	-41.20	14,535.36
2014	1,978	102	-66.80	-7.20	480.96
2015	1,884	110	-160.80	0.80	-128.64
2016	2,151	112	106.20	2.80	297.36
2017	2,519	154	474.20	44.80	21,244.16
Mean	2,044.80	109.20	Sum		36,429.20

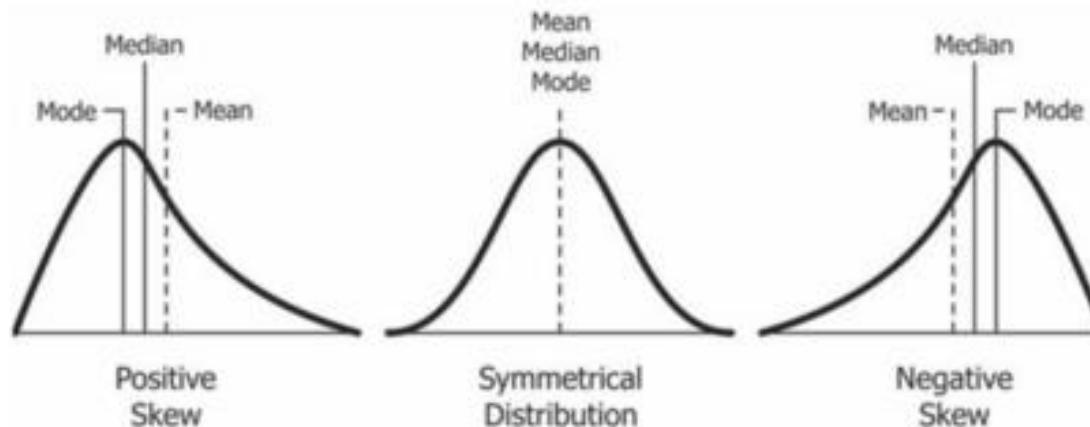
Correlation Positive/Negative

- The more time a student spends watching TV, the lower their exam scores tend to be. In other words, the variable time spent watching TV and the variable exam score have a negative correlation. As time spent watching TV increases, exam scores decrease.
- The more time an individual spends running, the lower their body fat tends to be. In other words, the variable running time and the variable body fat have a negative correlation. As time spent running increases, body fat decreases.
- The height and weight of people have positive correlation with each other.
- The correlation between the temperature and total ice cream sales is positive. In other words, when it's hotter outside the total ice cream sales of companies tends to be higher since more people buy ice cream when it's hot out.



Skewness

- It is a measure of the asymmetry of the probability distribution of a real-valued random variable about its mean.
- The skewness value can be positive, zero, negative, or undefined.
- For Unimodal Distribution:
 - **Symmetric:** Even on either side of mean ($\text{mean} = \text{median} = \text{mode}$)
 - **Positive skew:** tail is on the right side of the distribution ($\text{mean} > \text{median}$)
 - **Negative skew:** tail is on the left side of the distribution ($\text{mean} < \text{median}$)



```
YearsExperience    Salary
0                  1.1    39343
1                  1.3    46205
2                  1.5    37731
3                  2.0    43525
4                  2.2    39891
5                  2.9    56642
6                  3.0    60150
7                  3.2    54445
8                  3.2    64445
9                  3.7    57189
10                 3.9    63218
11                 4.0    55794
12                 4.0    56957
13                 4.1    57081
14                 4.5    61111
15                 4.9    67938
16                 5.1    66029
17                 5.3    83088
18                 5.9    81363
19                 6.0    93940
20                 6.8    91738
21                 7.1    98273
22                 7.9    101302
23                 8.2    113812
24                 8.7    109431
25                 9.0    105582
26                 9.5    116969
27                 9.6    112635
28                10.3   122391
29                10.5   121872
```

Questions:

For the following given data set (Only for Y_test), you have to calculate mean, std deviation, variance, MAE, MSE, RMSE and R2 Score.

```
df=pd.DataFrame({'Actual':y_test,'Predicted':y_pred})
print(df)
```

	Actual	Predicted
0	37731	40748.961841
1	122391	122699.622956
2	57081	64961.657170
3	63218	63099.142145
4	116969	115249.562855
5	109431	107799.502753

Linear Regression

- Linearity
- Linear Regression estimates real values (cost of houses, number of calls, total sales , etc.)
- Base is continuous variables.
- Establish a relationship between the **independent** and **dependent** variables by fitting the best line.
- This best-fit line is known as the regression line and is represented by a linear equation **$Y= b_1X + b_0$** .

Linear Regression :Example

- Linear equation $Y= b_1X + b_0$.
- Remember your childhood
- Teacher assigns a task- Arrange all the children in the class height-wise
- What you as a child would have done?
- Look (visually analyze) at the height and build of people and arrange them using a combination of these visible parameters

Linear Regression :Example

- This is a linear regression in real life!
- The child has figured out that height and build would be correlated to the weight by a relationship, which looks like the equation above.
- Y – Dependent Variable
- a – Slope
- X – Independent variable
- b – Intercept (the predicted value of Y when the X is 0)

Linear Regression :Example

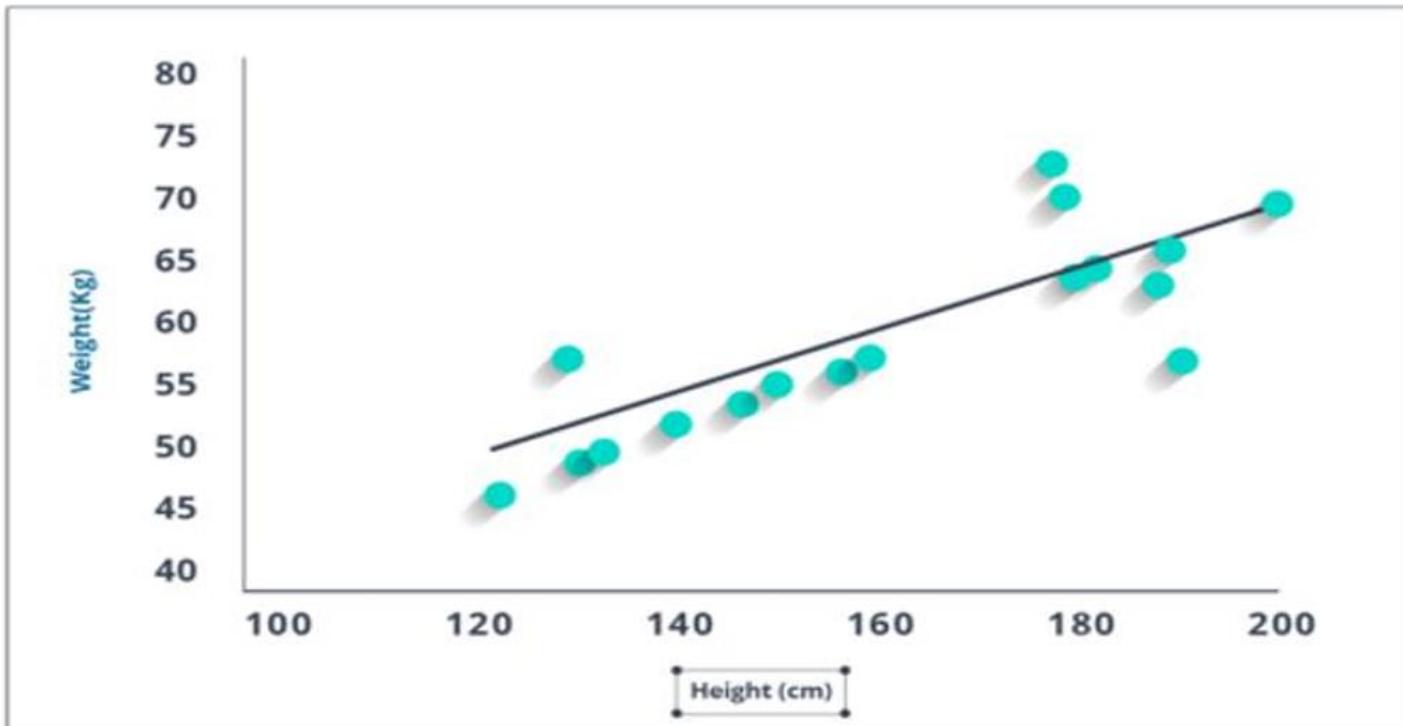


Fig. 8 Linear Regression Plot REF[4]

Use Case: Score Prediction

Problem Statement

Build a model to predict salary based on the number of years of experience.

Data

Use the Salary Data dataset and analyze the relationship between YearsExperience and Salary variables using a linear regression

Machine Learning Key Terms

Hours- Independent Variable → X→ Predictor

Scores – Dependent Variable→ Y→ Response

student_score.csv

Hours	Scores
2.5	21
5.1	47
3.2	27
8.5	75
3.5	30
1.5	20
9.2	88
5.5	60
8.3	81
2.7	25
7.7	85
5.9	62
4.5	41
3.3	42
1.1	17
8.9	95
2.5	30
1.9	24
6.1	67
7.4	69
2.7	30
4.8	54
3.8	35
6.9	76
7.8	86

Linear Regression : Hands On

The screenshot shows a Jupyter Notebook interface with the following details:

- Title Bar:** jupyter Linear_Regression_1 Last Checkpoint: 08/04/2023 (autosaved)
- Toolbar:** File, Edit, View, Insert, Cell, Kernel, Widgets, Help, Not Trusted, Python 3 (ipykernel) ○
- Cells:**
 - Cell 1 (Use Case):** **Linear Regression**

Problem Statement:
Build a model to predict salary based on the number years of experience.
 - Cell 2 (Data):** **Use Case: Score Prediction**

Data:
Use the Salary_Data dataset and analyse the relationship between YearsExperience and Salary variables using a linear regression.
 - Cell 3 (Import Libraries):** **Import Libraries**

In [1]:

```
import numpy as np
import pandas as pd
from matplotlib import pyplot as plt
```

Simple Linear Regression

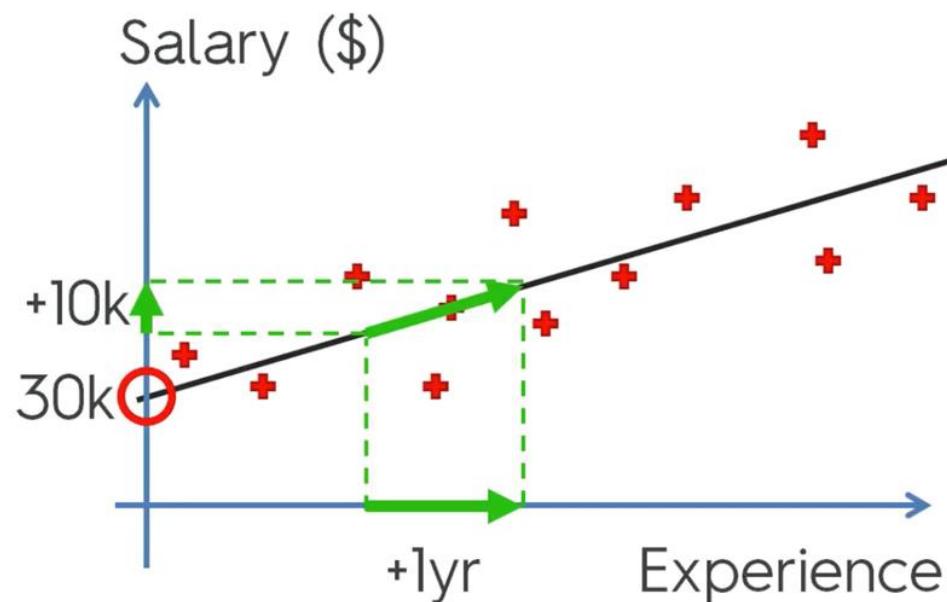
$$y = b_0 + b_1 * x_1$$

Constant Coefficient

Dependent variable (DV) Independent variable (IV)

The diagram illustrates the components of the simple linear regression equation. The equation is $y = b_0 + b_1 * x_1$. A green arrow points from the label "Constant" to the term b_0 . Another green arrow points from the label "Coefficient" to the term b_1 . A third green arrow points from the label "Dependent variable (DV)" to the term y . A fourth green arrow points from the label "Independent variable (IV)" to the term x_1 .

Simple Linear Regression:



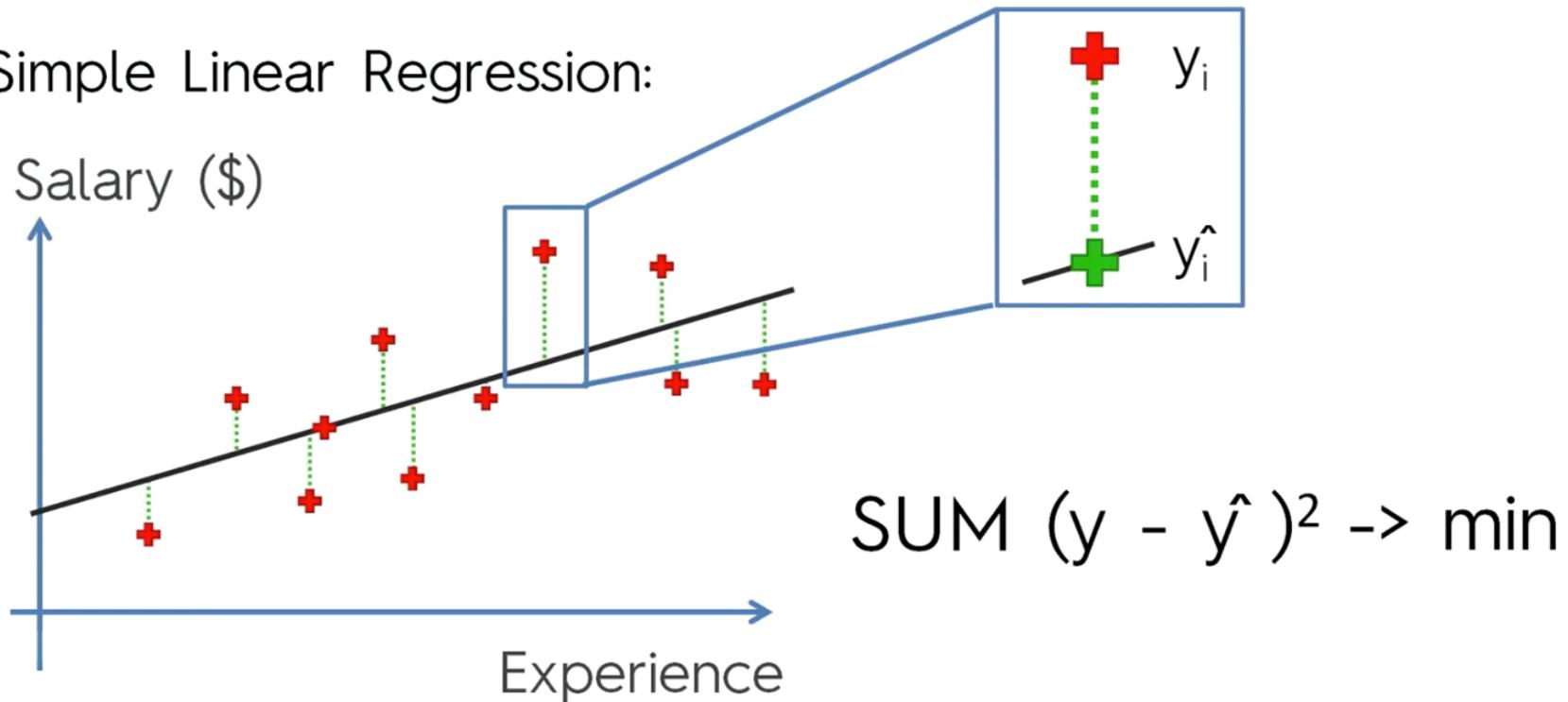
$$y = b_0 + b_1 * x$$



$$\text{Salary} = b_0 + b_1 * \text{Experience}$$

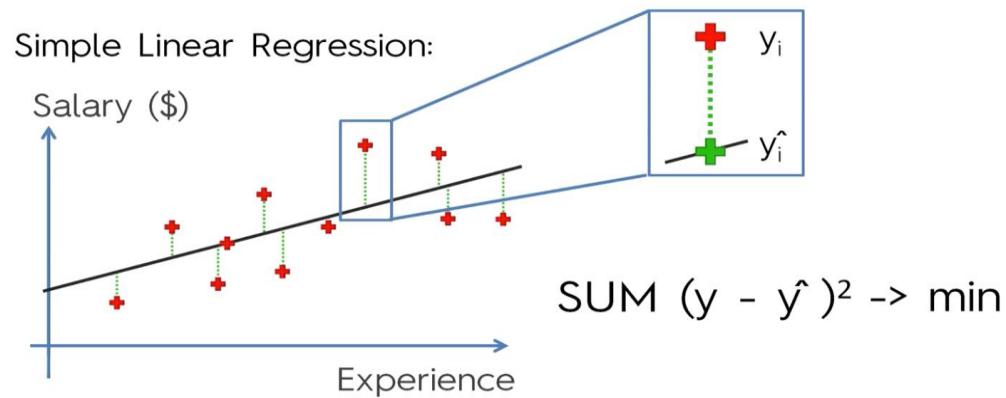
Ordinary Least Squares

Simple Linear Regression:



Evaluation Metrics

- Mean Absolute Error (MAE)
- Mean Squared Error (MSE)
- Root Mean Squared Error (RMSE)



$$\text{MAE} = [\sum \text{abs(actual_value - predicted_value)}] / n$$

$$m_1 = \text{abs}(\text{actual_value} - \text{predicted_value})$$

$$m_1 = \text{abs}(37731.0 - 40748.961841) = 3017.961841$$

$$m_2 = \text{abs}(122391.0 - 122699.622956) = 308.6229559$$

$$m_3 = \text{abs}(57081.0 - 64961.657170) = 7880.65717$$

$$m_4 = \text{abs}(63218.0 - 63099.142145) = 118.8578551$$

$$m_5 = \text{abs}(116969.0 - 115249.562855) = 1719.437145$$

$$m_6 = \text{abs}(109431.0 - 107799.502753) = 1631.497247$$

$$\text{MAE} = (m_1 + m_2 + m_3 + m_4 + m_5 + m_6) / 6$$

$$\text{MAE} = (3017.961841 + 308.6229559 + 7880.65717 + 118.8578551 + 1719.437145 + 1631.497247) / 6$$

$$\text{MAE} = (14677.03421) / 6$$

$$\text{MAE} = 2446.172369$$

$$\text{MSE} = [\sum (\text{actual_value} - \text{predicted_value})^2] / n$$

$m_i = (\text{actual_value} - \text{predicted_value})^2$

$m_1 = (37731.0 - 40748.961841)^2 = 9108093.672$

$m_2 = (122391.0 - 122699.622956)^2 = 95248.12893$

$m_3 = (57081.0 - 64961.657170)^2 = 62104757.43$

$m_4 = (63218.0 - 63099.142145)^2 = 14127.18973$

$m_5 = (116969.0 - 115249.562855)^2 = 2956464.097$

$m_6 = (109431.0 - 107799.502753)^2 = 2661783.266$

$\text{MSE} = (m_1 + m_2 + m_3 + m_4 + m_5 + m_6) / 6$

$\text{MSE} = (9108093.672 + 95248.12893 + 62104757.43 + 14127.18973 + 2956464.097 + 2661783.266) / 6$

$\text{MSE} = (76940473.79) / 6$

$\text{MSE} = 12823412.3$

Root Mean Squared Error (RMSE)

$\text{RMSE} = \text{SQRT}(\text{MSE})$ $\text{RMSE} = \text{SQRT}(12823412.3)$ $\text{RMSE} = 3580.979237$

R-squared?

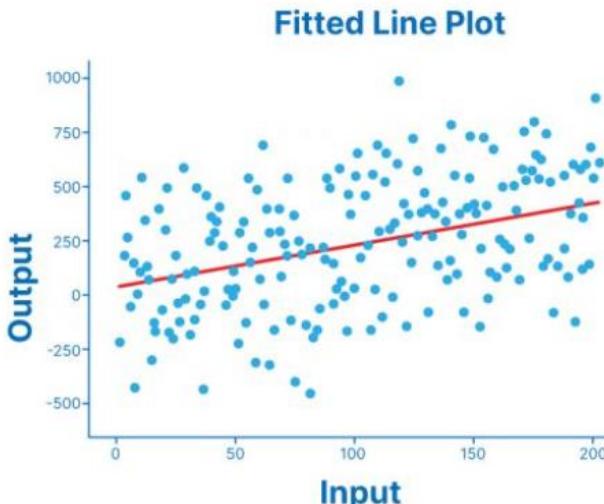
What is R-squared?

R squared or Coefficient of determination, or R^2 is a measure that provides information about the goodness of fit of the regression model. In simple terms, it is a statistical measure that tells how well the plotted regression line fits the actual data. R squared measures how much the variation is there in predicted and actual values in the regression model.

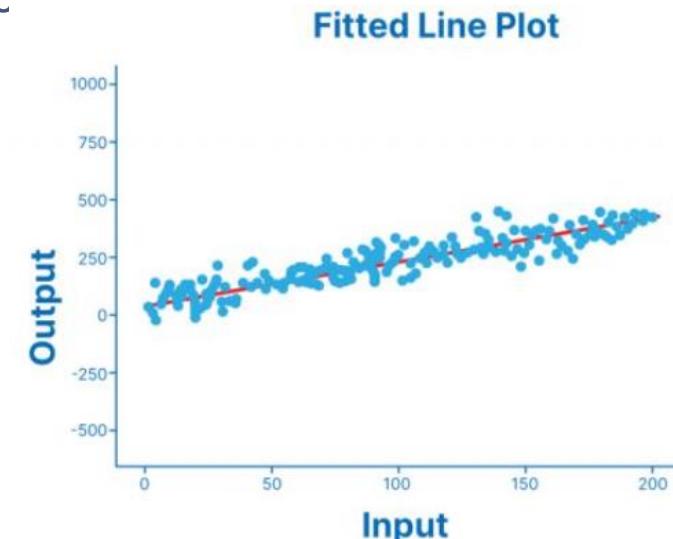
- R-squared values range from 0 to 1, usually expressed as a percentage from 0% to 100%.

$$R\text{-Squared} = 1 - \left(\frac{SSR}{SST} \right) = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

Predicated
where:
SSR is the sum of squared residuals (i.e., the sum of squared errors)
SST is the total sum of squares (i.e., the sum of squared deviations)

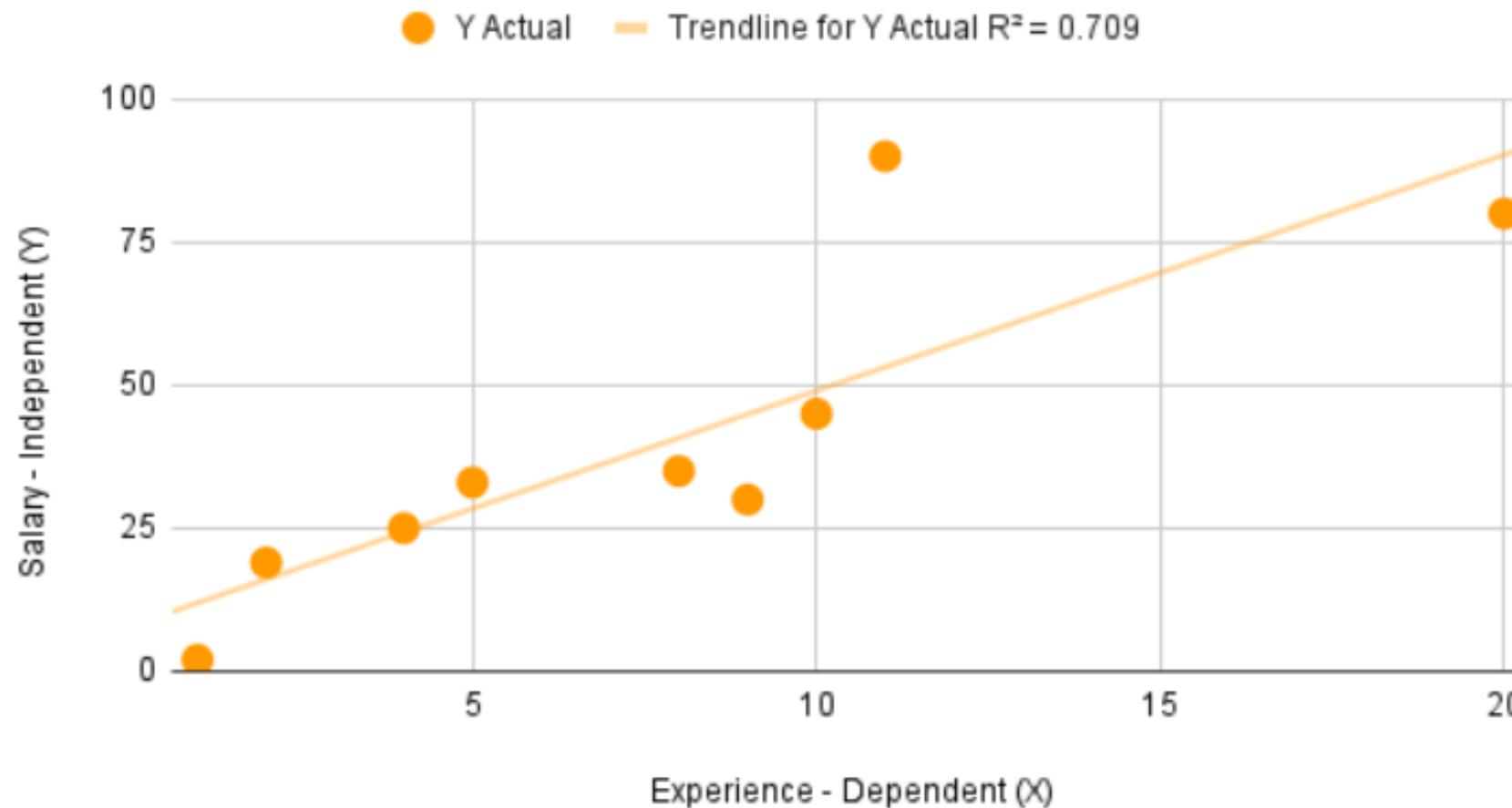


R squared value will be less than 0.5 or near 0.



R squared will be close to 1

Salary vs Experience



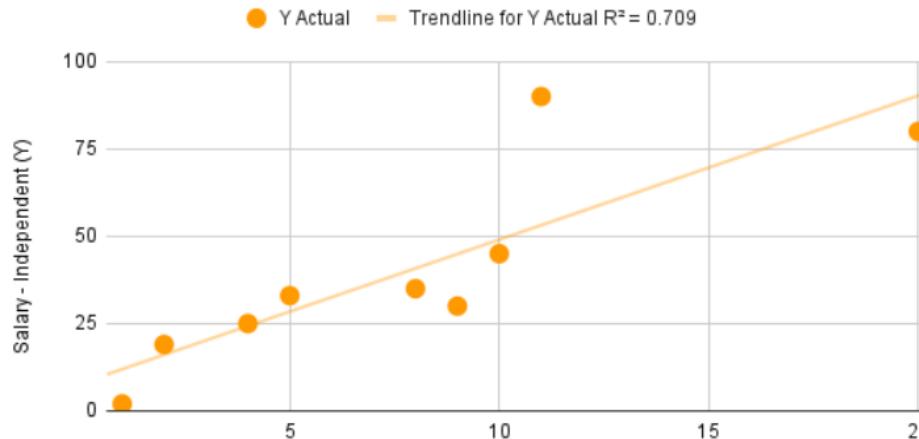
SSR

$$\sum_{i=1}^n (y_i - \hat{y}_i)^2$$

SST

$$\sum_{i=1}^n (y_i - \bar{y})^2$$

Salary vs Experience



x _i	y _i	x _i - x̄ _i	y _i - ȳ _i	(x _i - x̄ _i) ²	(x _i - x̄ _i) (y _i - ȳ _i)
11	90	3.22	50.11	10.38	161.47
10	45	2.22	5.11	4.94	11.36
2	19	-5.78	-20.89	33.38	120.69
8	35	0.22	-4.89	0.05	-1.09
4	25	-3.78	-14.89	14.27	56.25
20	80	12.22	40.11	149.38	490.25
1	2	-6.78	-37.89	45.94	256.80
9	30	1.22	-9.89	1.49	-12.09
5	33	-2.78	-6.89	7.72	19.14

R-Squared =
1-(SSR/SST)

Suppose we have a data set having values of X and Y.

1. We have to find X̄(mean) and Ȳ(mean).

2. Calculate X_i-X̄ and Y_i-Ȳ and then do (X_i-X̄)²

3. Now calculate (X_i-X̄)(Y_i-Ȳ)

Now we have to calculate R squared so let's try to calculate it

SSR

$$n \sum_{i=1}^{n} (y_i - \hat{y}_i)^2$$

SST

$$n \sum_{i=1}^{n} (y_i - \bar{y})^2$$

y _i [^]	y _i - y _i [^]	SSR	y _i - ȳ _i	SST
53.17	36.83	1356.46	50.11	2511.12
49.05	-4.05	16.39	5.11	26.12
16.07	2.93	8.56	-20.89	436.35
37.10	-2.10	4.39	-4.89	23.90
20.61	4.39	19.29	-14.89	221.68
86.56	-6.56	42.97	40.11	1608.90
8.24	-6.24	38.98	-37.89	1435.57
41.22	-11.22	125.82	-9.89	97.79
24.73	8.27	68.39	-6.89	47.46

$$\text{R-Squared} = 1 - \left(\frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \right)$$

where:

SSR is the sum of squared residuals (i.e., the sum of squared errors)

SST is the total sum of squares (i.e., the sum of squared deviations from the mean)

```
actual_minus_predicted = sum((y_test - y_pred)**2)
actual_minus_actual_mean = sum((y_test - y_test.mean())**2)
r2 = 1 - actual_minus_predicted/actual_minus_actual_mean
print('R2:', r2)
```

R²: 0.988169515729126

Salary Dataset write a program for simple linear regression in python

Simple Linear Regression

$$y = b_0 + b_1 * x_1$$

Multiple Linear Regression

Dependent variable (DV) Independent variables (IVs)

$$y = b_0 + b_1 * x_1 + b_2 * x_2 + \dots + b_n * x_n$$

Constant Coefficients

```
graph TD; DV[Dependent variable DV] --> y; IVs[Independent variables IVs] --> terms; C[Constant] --> b0; Coefficients[Coefficients] --> b1x1; Coefficients --> b2x2; Coefficients --> dots; Coefficients --> bnxn;
```

Multiple Linear Regression

Advertising_sales.csv



Independent Variables:

- TV
- Radio
- Newspaper

Dependent Variable:

- Sales

	TV	Radio	Newspaper	Sales
1	230.1	37.8	69.2	22.1
2	44.5	39.3	45.1	10.4
3	17.2	45.9	69.3	9.3
4	151.5	41.3	58.5	18.5
5	180.8	10.8	58.4	12.9
6	8.7	48.9	75	7.2
7	57.5	32.8	23.5	11.8
8	120.2	19.6	11.6	13.2
9	8.6	2.1	1	4.8
10	199.8	2.6	21.2	10.6
11	66.1	5.8	24.2	8.6
12	214.7	24	4	17.4
13	23.8	35.1	65.9	9.2
14	97.5	7.6	7.2	9.7
15	204.1	32.9	46	19
16	195.4	47.7	52.9	22.4
17	67.8	36.6	114	12.5
18	281.4	39.6	55.8	24.4
19	69.2	20.5	18.3	11.3
20	147.3	23.9	19.1	14.6
21	218.4	27.7	53.4	18
22	237.4	5.1	23.5	12.5
23	13.2	15.9	49.6	5.6
24	228.3	16.9	26.2	15.5
25	62.3	12.6	18.3	9.7
26	262.9	3.5	19.5	12
27	142.9	29.3	12.6	15
28	240.1	16.7	22.9	15.9
29	248.8	27.1	22.9	18.9

Multiple Linear Regression

$$Y = b_1 * X_1 + b_2 * X_2 + b_3 * X_3 + b_0$$

X₁ = TV

X₂ = Radio

X₃ = Newspaper

and

Y = Sales

$$\text{Sales} = b_1 * (\text{TV}) + b_2 * (\text{Radio}) + b_3 * (\text{Newspaper}) + b_0$$

Multiple Linear Regression

Problem Statement :-

Data of 50 companies.

R&D Spends/Administration/Marketing Spend/State (Independent Variables)	Profit (Dependent Variable)
--	--------------------------------

A VC Fund is interested in investing in these companies. But has questions like :-
Where companies perform better? Are these companies are those who spend more money on R&D Spend or on Marketing Spend ?
Help VC Fund to build a model

Dummy Variable/Categorical Variable/One hot encoding

Profit	R&D Spend	Admin	Marketing	State	Dummy Variables
192,261.83	165,349.20	136,897.80	471,784.10	New York	
191,792.06	162,597.70	151,377.59	443,898.53	California	
191,050.39	153,441.51	101,145.55	407,934.54	California	
182,901.99	144,372.41	118,671.85	383,199.62	New York	
166,187.94	142,107.34	91,391.77	366,168.42	California	

$$y = b_0 + b_1 * x_1 + b_2 * x_2 + b_3 * x_3 + b_4 * D_1$$



Dummy Variable Trap

Not Truly Independent Variables

Multicollinearity = High correlation between 2 or more independent variables

Profit	R&D Spend	Admin	Marketing	State
192,261.83	165,349.20	136,897.80	471,784.10	New York
191,792.06	162,597.70			California
191,050.39	153,441.51			California
182,901.99	144,372.41			New York
166,187.94	142,107.34			California

$$D_2 = 1 - D_1$$

Dummy Variables	
New York	California
1	0
0	1
0	1
1	0
0	1

$$y = b_0 + b_1 * x_1 + b_2 * x_2 + b_3 * x_3 + b_4 * D_1 + \underline{b_5 * D_2}$$

Dummy Variable Trap

Profit	R&D Spend	Admin	Marketing	State
192,261.83	165,349.20	136,897.80	471,784.10	New York
191,792.06	162,597.70	151,377.59	443,898.53	California
191,050.39	153,441.51	101,145.55	407,934.54	California
182,901.99	144,372.41	118,671.85	383,199.62	New York
166,187.94	142,107.34	91,391.77	366,168.42	California

Dummy Variables

New York	California
1	0
0	1
0	1
1	0
0	1

$$y = b_0 + b_1 * x_1 + b_2 * x_2 + b_3 * x_3 + b_4 * D_1 + \cancel{b_5 * D_2}$$

Always omit one
dummy variable

Multiple Linear Regression: Hands-On

jupyter Multiple Linear Regression Last Checkpoint: 08/04/2023 (autosaved)

File Edit View Insert Cell Kernel Widgets Help

In [1]: `import pandas as pd
import numpy as np
from matplotlib import pyplot as plt`

In [2]: `advert=pd.read_csv('Advertising_sales.csv')`

In [3]: `advert.describe()`

Out[3]:

	Unnamed: 0	TV	Radio	Newspaper	Sales
count	200.000000	200.000000	200.000000	200.000000	200.000000
mean	100.500000	147.042500	23.264000	30.554000	14.022500
std	57.879185	85.854236	14.846809	21.778621	5.217457
min	1.000000	0.700000	0.000000	0.300000	1.600000
25%	50.750000	74.375000	9.975000	12.750000	10.375000
50%	100.500000	149.750000	22.900000	25.750000	12.900000
75%	150.250000	218.825000	36.525000	45.100000	17.400000
max	200.000000	296.400000	49.600000	114.000000	27.000000

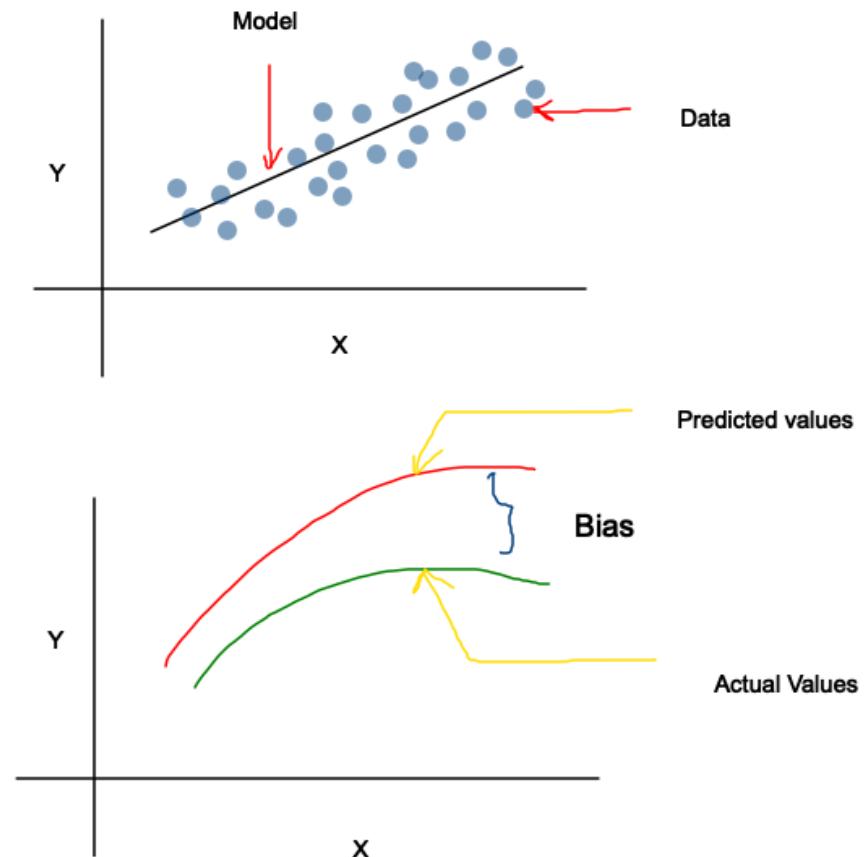
Bias and Variance

Bias-

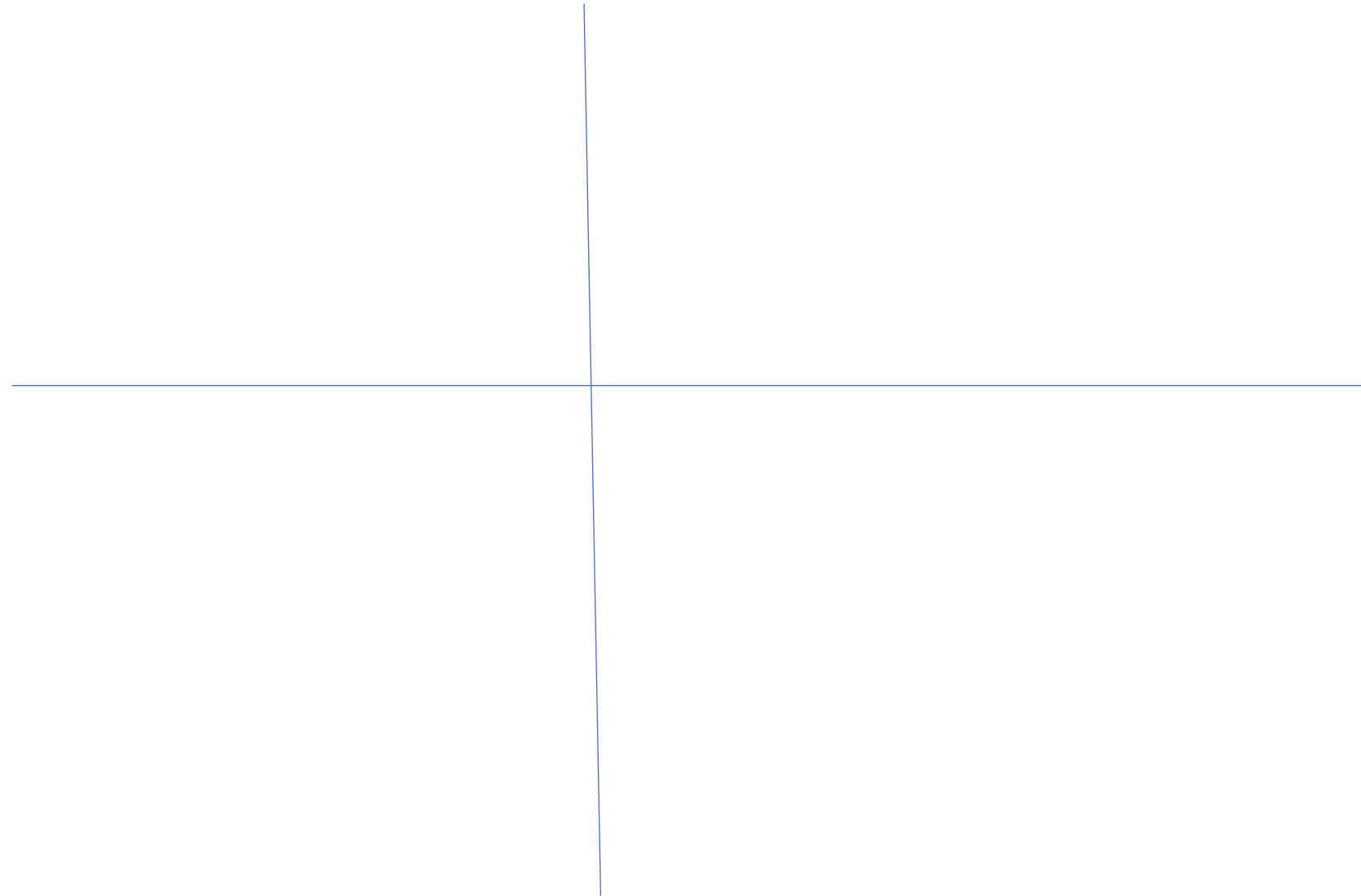
- Prediction Error, introduced in the model
- It is the difference between predicted value and actual value.
- Taring Data.

Variance-

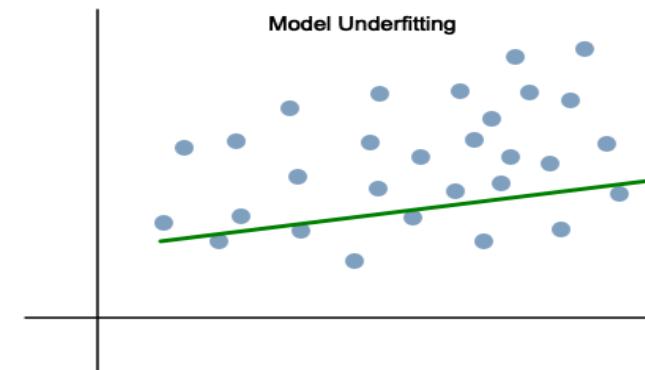
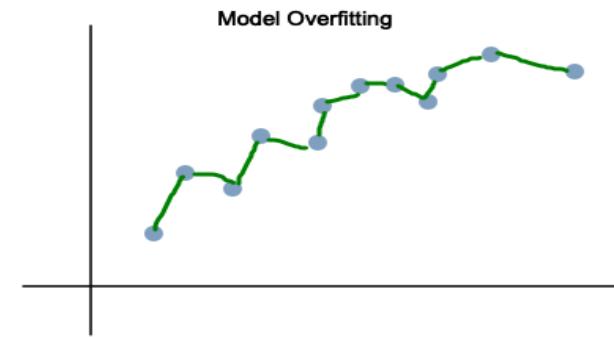
- Test Data.
- Error differences



Model Overfitting and Underfitting



Model Overfitting and Underfitting



How to avoid Underfitting in Model

Increase number of features

Increase training time of model

