



# AI 330: Machine Learning

## Fall 2023

**Dr. Wessam EL-Behaidy**

Associate Professor, Computer Science Department,  
Faculty of Computers and Artificial Intelligence,  
Helwan University.

**Dr. Ensaf Hussein**

Associate Professor, Computer Science Department,  
Faculty of Computers and Artificial Intelligence,  
Helwan University.

# Agenda

**Course Objectives**

**Course Map**

**Course topics**

**References**

**Grading Distribution**

**Introduction to machine learning**

**What is machine learning (ML)?**

**ML Applications**

**ML Life Cycle**

**Types of ML**

# Course Objectives

**Fundamental Understanding:** Provide students with a solid foundation in the basic principles, algorithms, and theories of Machine Learning, ensuring they grasp the core concepts.

*→ How to implement this models?*

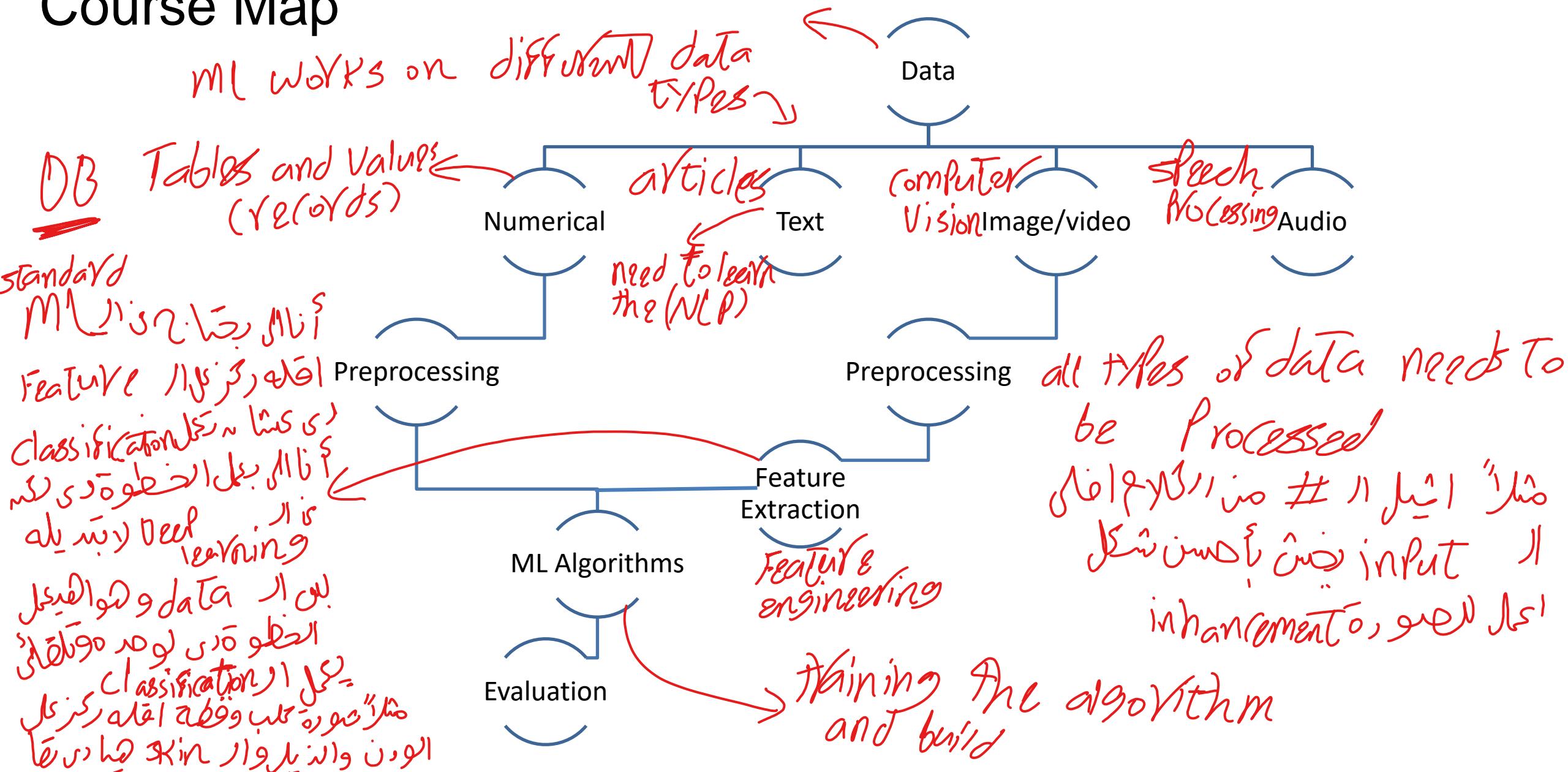
**Practical Application:** Enable students to apply machine learning techniques to real-world problems, fostering their ability to analyze data, build models, and make data-driven decisions.

**Hands-On Experience:** Offer practical, hands-on experience with machine learning tools and libraries, ensuring students can implement and experiment with various algorithms.

*→ What is the best model for this data to implement?*

**Critical Thinking:** Cultivate critical thinking skills by encouraging students to evaluate the suitability of machine learning solutions for different problems, understand ethical considerations, and explore the limitations of the technology.

# Course Map



# Course Topics

Date	Lecture	Description
Week1	Introduction to Machine Learning	<ul style="list-style-type: none"><li>• Understanding the fundamentals of machine learning.</li><li>• Types of machine learning: supervised, unsupervised, and more.</li><li>• Applications and importance of machine learning.</li></ul>
Week2	Data Preprocessing & Feature Extraction	<ul style="list-style-type: none"><li>• Data cleaning and preprocessing techniques.</li><li>• Feature selection and engineering.</li></ul>
Week3	Supervised Regressors	<ul style="list-style-type: none"><li>• Linear Regression with one variable.</li><li>• Gradient Descent optimization</li></ul>
Week4	Supervised Regressors (cont.)	<ul style="list-style-type: none"><li>• Linear Regression with Multiple Variables.</li><li>• Handling categorical data.</li></ul>
Week5	Supervised classifiers	<ul style="list-style-type: none"><li>• Logistic Regression.</li><li>• Regularization techniques (L1 and L2).</li></ul>
Week6	k-Nearest Neighbors (k-NN) algorithm	<ul style="list-style-type: none"><li>• Understanding the k-NN algorithm.</li><li>• Choosing the optimal k value.</li></ul>

# Course Topics

Date	Lecture	Description
Week8	Model Evaluation and Hyperparameter Tuning	<ul style="list-style-type: none"><li>• Evaluation metrics for regression and classification.</li><li>• Cross-validation and hyperparameter tuning.</li></ul>
Week9	Ensemble Methods	<ul style="list-style-type: none"><li>• Bagging and Random Forests.</li><li>• Boosting: AdaBoost and Gradient Boosting.</li></ul>
Week10	Introduction to unsupervised learning.	<ul style="list-style-type: none"><li>• Introduction to unsupervised learning.</li><li>• Clustering techniques: k-Means.</li></ul>
Week11	Dimensionality Reduction	<ul style="list-style-type: none"><li>• Principal Component Analysis (PCA).</li><li>• Applications and interpretation of PCA</li></ul>
Week12		Project Discussion

This course focus on the ~~Standard~~ ML processing, feature abstraction  
Supervised and unsupervised ML

References Advanced ML → Neural networks, CNN, RNN,  
Transformer.

## Lectures (Course slides) are based on :

- Machine Learning Specialization <https://www.coursera.org/specializations/machine-learning-introduction> at Stanford University (Andrew Ng)
- A Course in Machine Learning by Hal Daumé III

## Practical Labs based on:

Machine learning A to Z: Kirill Eremenko ©superdatascience

# Grading Distribution

- Project: (15%)

Code “github”, pitching video, presentation.

→ Paid) (3 months free)

- Programming Assignment (DataCamp)(10%)

- Practical Assessment (10%)

↳ (Training) واسباب کلی (code Now)  
track within course / design

- Midterm exam (15%)

- Final exam (50%)

# Join Machine Learning(Fall2023) Team Class

To communicate with us and get course materials

Team class code:

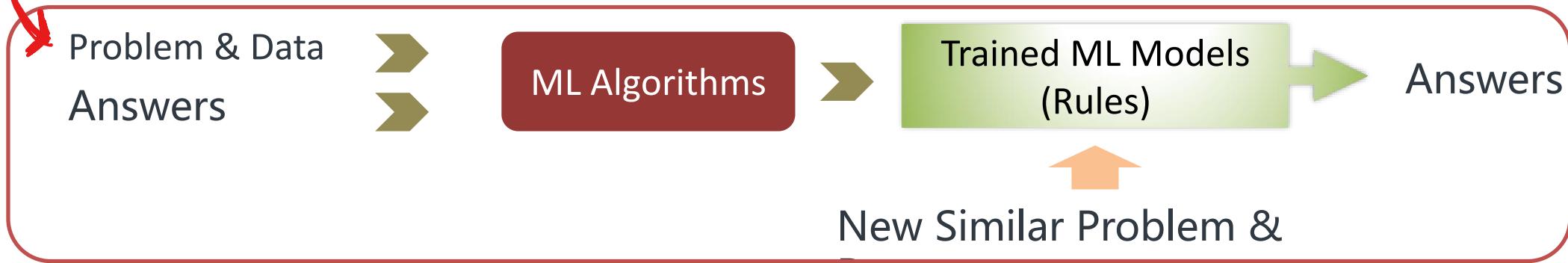
hrglst0

# What is Machine Learning (ML)?

answer \ جواب و رجوع ملحوظ



Rule \ قواعد و معايير ملحوظ

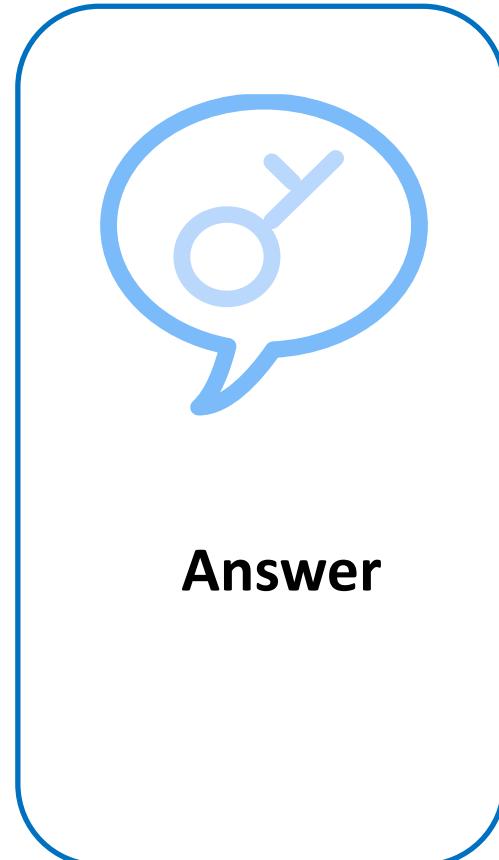
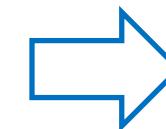
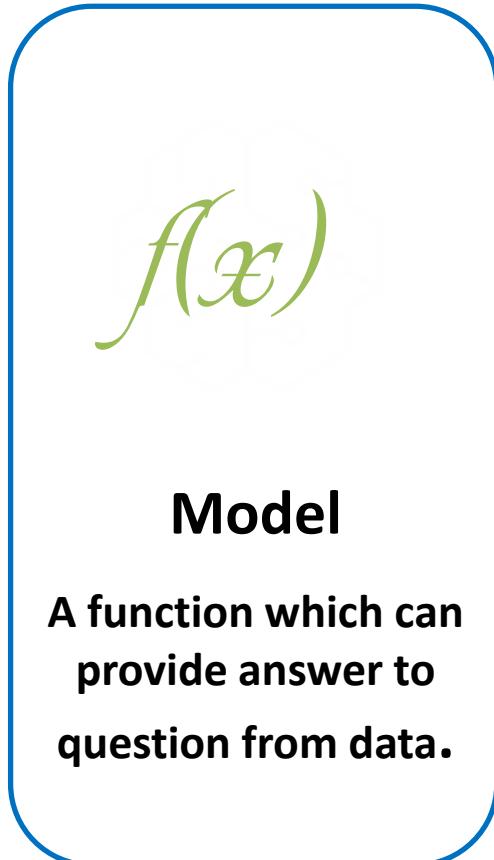
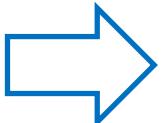


The Task of ML is generating the rule based on given data

Rule \ قواعد و معايير ملحوظ

# What is a Model?

A Function

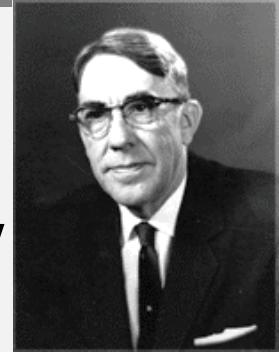


# What is Machine Learning (ML)?

“The field of study that gives computers the ability to learn without being explicitly programmed”

*Arthur Samuel*

Pioneer of AI research



direct knowledge is no longer needed

# ML Terminologies

ML optimizes on **predictive** performance, while statistics places importance on **interpretability** and parsimony/simplicity.

ML / Statistics Jargon	Definition
Label/Target/Output Variable/"y"	The results to predict
Feature/Input Variable/"x"	Input data to help make predictions رسائلات، مدخلات
Feature Engineering / Transformation	Reshaping raw input data to give more insights رسالة، تغيير
Dimensionality / [1 <sup>st</sup> d, 2 <sup>nd</sup> d, ... , n <sup>th</sup> d]	Number of features $\sum_{i=1}^n$ skin, length, width
Model Weights / Parameters	A set of numbers embedded in a model to make predictions رسائل، نتائج
Model Training → To find the optimal value of parameters	Applying optimization techniques to find the “best” set of model weights

→ extract the features from the data such as images,  
(جذب) transform the raw input (data) to features.

$$y = \alpha X + b$$

weight  
parameters

whis a,b نیوں جیسا کسی  
best result کے لئے

Variable  
Feature

z = Prediction decision line

## Machine Learning (Applications) (Techniques)

# ML Applications

like Google , amazon  
جیکس بحث کے لئے اور فوارڈ کرنے والے ملکے

Problem type	Description	Example
Ranking	Helping users find the most relevant <i>items</i>	Ranking algorithm within Amazon Search
Recommendation	Giving users the <i>items</i> they may be most interested in	
Classification	Figuring out what category does an <i>item</i> belongs to	
Regression	Predicting a numerical value of an <i>item</i>	
Clustering	Putting similar <i>items</i> together	
Anomaly Detection	Finding uncommon <i>items</i>	



## ML Applications

وَيُعْرَفُ الْأَنْوَاعُ بِالْمُؤْكِدِاتِ الْمُتَطَابِقَاتِ مُعْلِمَةً لِلْأَنْوَاعِ الْمُتَطَابِقَاتِ وَالْمُؤْكِدِاتِ الْمُتَطَابِقَاتِ

Problem type	Description	Example
Ranking	Helping users find the most relevant <i>items</i>	
Recommendation	Giving users the <i>items</i> they may be most interested in	Recommendations across the website
Classification	Figuring out what category does an <i>item</i> belongs to	Deals recommended for you <a href="#">See all deals</a>
Regression	Predicting a numerical value of an <i>item</i>	\$7.00 - \$147.90 Ends in 03:25:54
Clustering	Putting similar <i>items</i> together	\$79.99 \$139.99 Ends in 03:25:54
Anomaly Detection	Finding uncommon <i>items</i>	\$8.99 - \$37.49 Ends in 03:20:55
		\$4 End
		Amazon's Choice
		Panasonic RP-HJE120-PPK In-Ear Stereo Earpho by Panasonic
		\$8.18  prime   FREE One-Day Get it by Tomorrow, Apr 24 FREE One-Day Shipping on qualifying orders over \$35
		More Buying Choices \$7.99 (37 new offers) <a href="#">See newer model of this item</a>

التطبيقات الذكاء الاصطناعي في الواقع

Problem type	Description	Example
Ranking	Helping users find the most relevant <i>items</i>	
Recommendation	Giving users the <i>items</i> they may be most interested in	
Classification	Figuring out what category does an <i>item</i> belongs to ☞ <i>تصنيف المنتج</i>	Product classification for our catalog
Regression	Predicting a numerical value of an <i>item</i>	
Clustering	Putting similar <i>items</i> together	
Anomaly Detection	Finding uncommon <i>items</i>	



High-Low Dress



Straight Dress



Striped Skirt



Graphic Shirt

# ITEM HIS RECORD

بيانات وبيانات وبيانات وبيانات

Problem type	Description	Example
Ranking	Helping users find the most relevant <i>items</i>	
Recommendation	Giving users the <i>items</i> they may be most interested in	
Classification	Figuring out what category <u>does an item belong to</u> gpa → A, B, C, D, F	
Regression	Predicting a numerical value of a item gpa 3.3 → <u>within range</u>	
Clustering	Putting similar items together (witht) Outliers (without) (without) outliers (with)	
Anomaly Detection	Finding uncommon <i>items</i>	<p>Predicting sales for specific ASINs</p> <p>Amazon Fashion WOMEN MEN KIDS &amp; BABY LUGGAGE SEARCH</p> <p>elements: SILVER</p> <p>Elamanis Silver 925 Ladies' Heart Tag T-Bar Sterling Silver Necklace of 46 cm</p> <p>Customer reviews: 4.5 out of 5 stars, 10 customer reviews</p> <p>Price: £99.99 &amp; FREE Delivery in the UK. Delivery details</p> <p>In stock.</p> <p>Was it delivered by Monday, 04 April? Order within 11 hrs 38 mins and choose Priority Delivery. Learn more Details</p> <p>Customer reviews: 4.5 out of 5 stars, 10 customer reviews</p> <p>Note: This item is eligible for click and collect. Details</p> <ul style="list-style-type: none"> <li>• Silver with Sterling Silver wire</li> <li>• Presented in a pale blue Elements gift box</li> <li>• The necklace is 46 cm/18 inches in length</li> <li>• Gold silver plate weighing over 40 g</li> <li>• Chain: Fine rolo chain with delicate heart charms</li> <li>• Adjustable length 46-48cm</li> </ul> <p>INPUTS: جملات دخل وبيانات وبيانات جديدة (أو قديمة) → <u>within range</u></p> <p>predictions sample paths</p> <p>zt</p> <p>xt</p> <p>Seasonality   Out of stock   Promotions</p>

# Grouping Similar Items Together

*Categories*

Problem type	Description	Example
Ranking	Helping users find the most relevant <i>items</i>	
Recommendation	Giving users the <i>items</i> they may be most interested in	
Classification	Figuring out what category does an <i>item</i> belongs to	
Regression	Predicting a numerical value of an <i>item</i>	
Clustering	Putting similar <i>items</i> together	<p><b>Close-matching for near-duplicates</b></p>  <p>Sheriff Walt Longmire Robert Taylor Trench Coat by NMFashion \$109.00 - \$160.00</p> <p>Sheriff Walt Longmire Robert Taylor Trench Coat by Spender \$154.00 FREE Shipping on eligible orders</p> <p>Robert Taylor Longmire Sheriff Walt Trench Coat by IMPRESSIONS \$175.00 FREE Shipping on eligible orders</p>
Anomaly Detection	Finding uncommon <i>items</i>	

## ML Applications

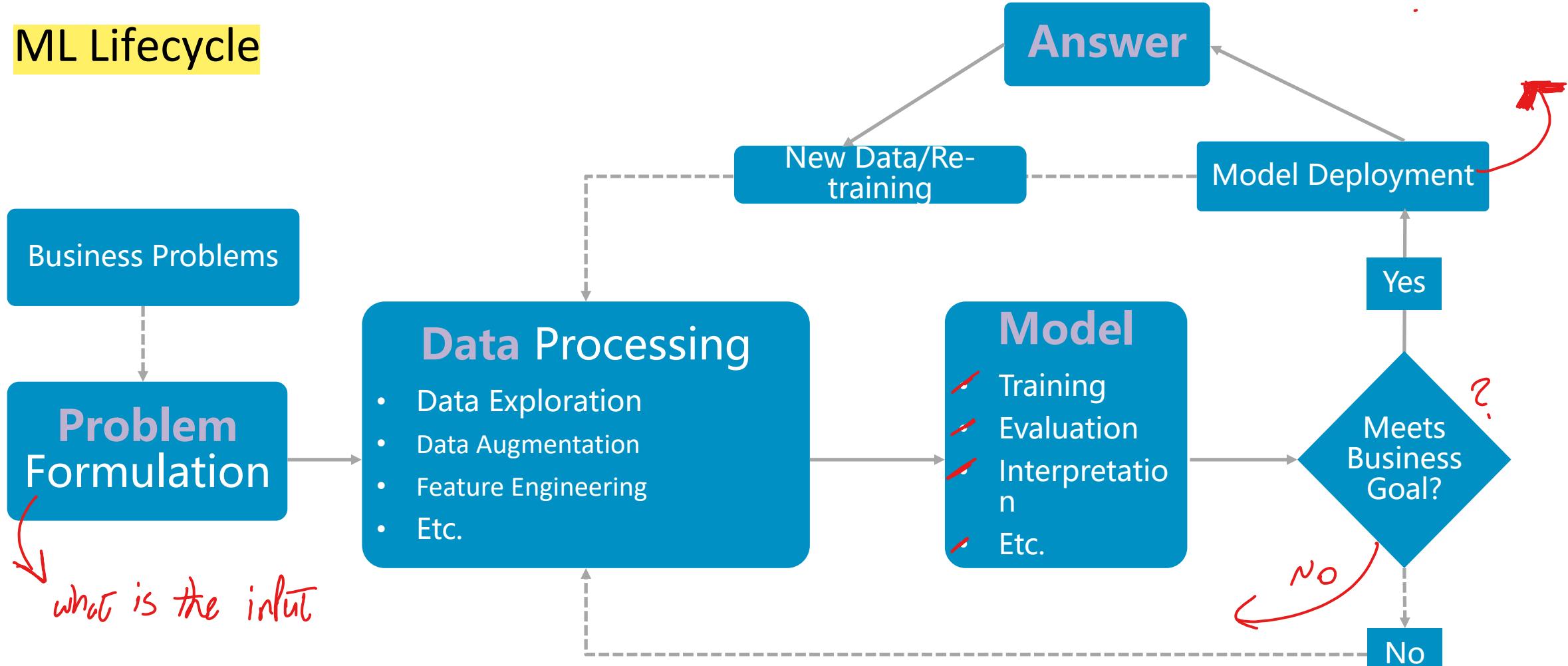
عاجز عن اكتشاف جميع ال\_items لـ detect "شيء"

- - -  
عاجز عن اكتشاف جميع ال\_items لـ detect "شيء"  
عاجز عن اكتشاف جميع ال\_items لـ detect "شيء"  
عاجز عن اكتشاف جميع ال\_items لـ detect "شيء"  
عاجز عن اكتشاف جميع ال\_items لـ detect "شيء"

Problem type	Description	Example				
Ranking	Helping users find the most relevant <i>items</i>					
Recommendation	Giving users the <i>items</i> they may be most interested in					
Classification	Figuring out what category does an <i>item</i> belongs to					
Regression	Predicting a numerical value of an <i>item</i>					
Clustering	Putting similar <i>items</i> together					
Anomaly Detection	Finding uncommon <i>items</i>	<p>Fruit freshness</p> <p>Before                          After</p> <p>amazonfresh</p> <table> <tr> <td>Good</td> </tr> <tr> <td>Damage</td> </tr> <tr> <td>Serious Damage</td> </tr> <tr> <td>Decay</td> </tr> </table>	Good	Damage	Serious Damage	Decay
Good						
Damage						
Serious Damage						
Decay						

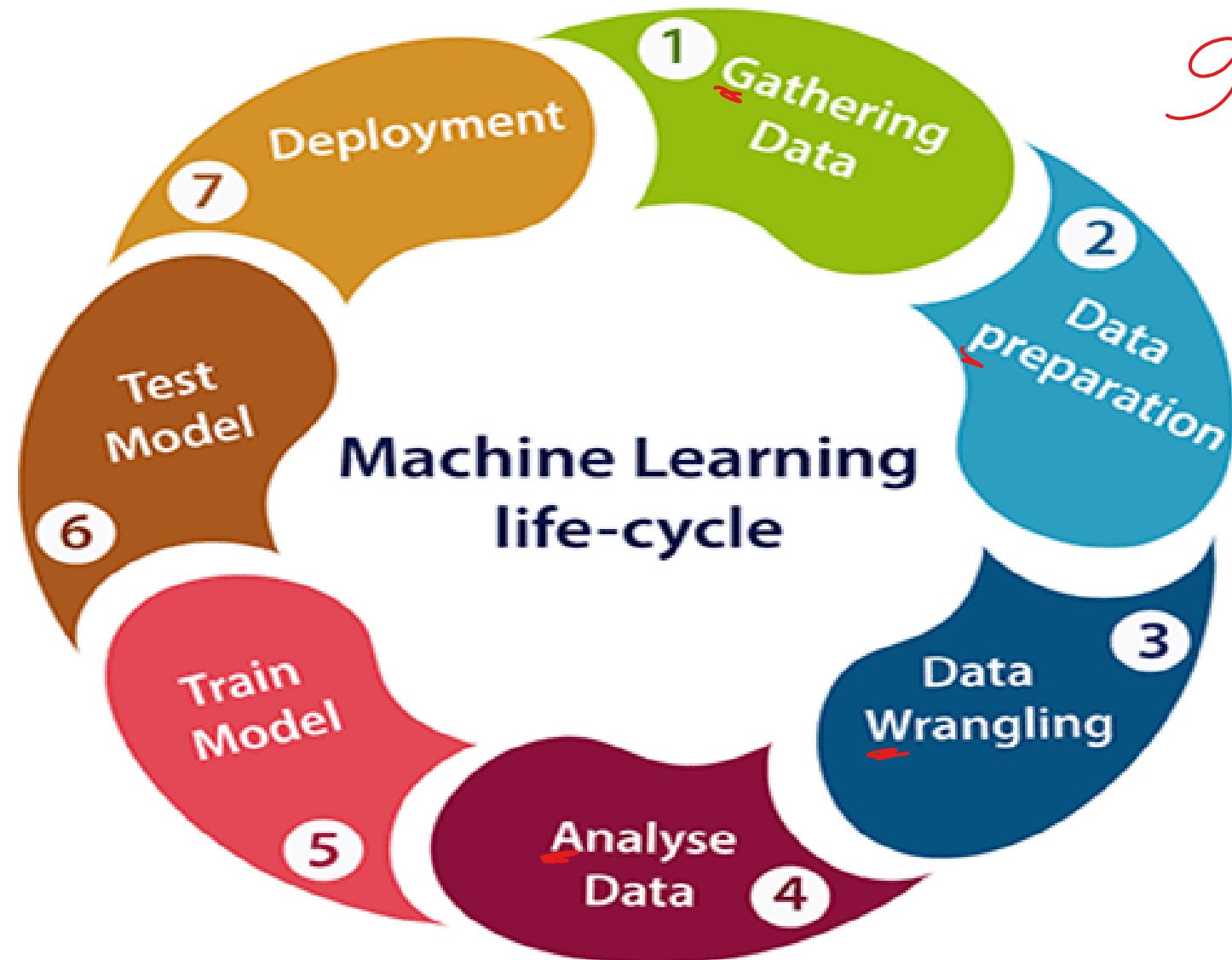
# **Machine Learning Life Cycle**

# ML Lifecycle



**Problem** > **Data** > **Model** > **Answer**

9PWA



# 1- Gathering data (Acquiring Dataset)

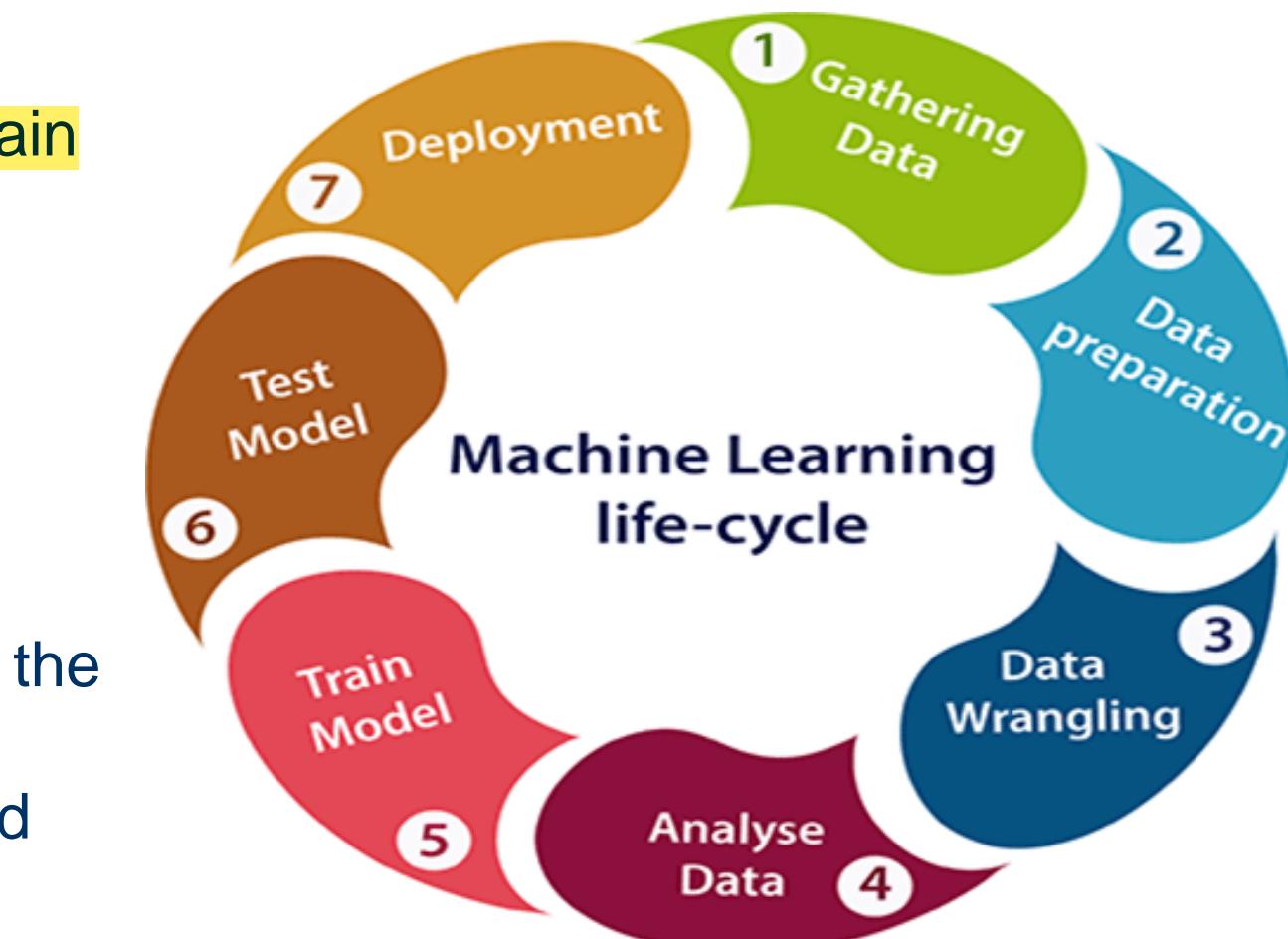
*Data Collection*

The goal of this step is to identify and obtain **all data-related problems**.

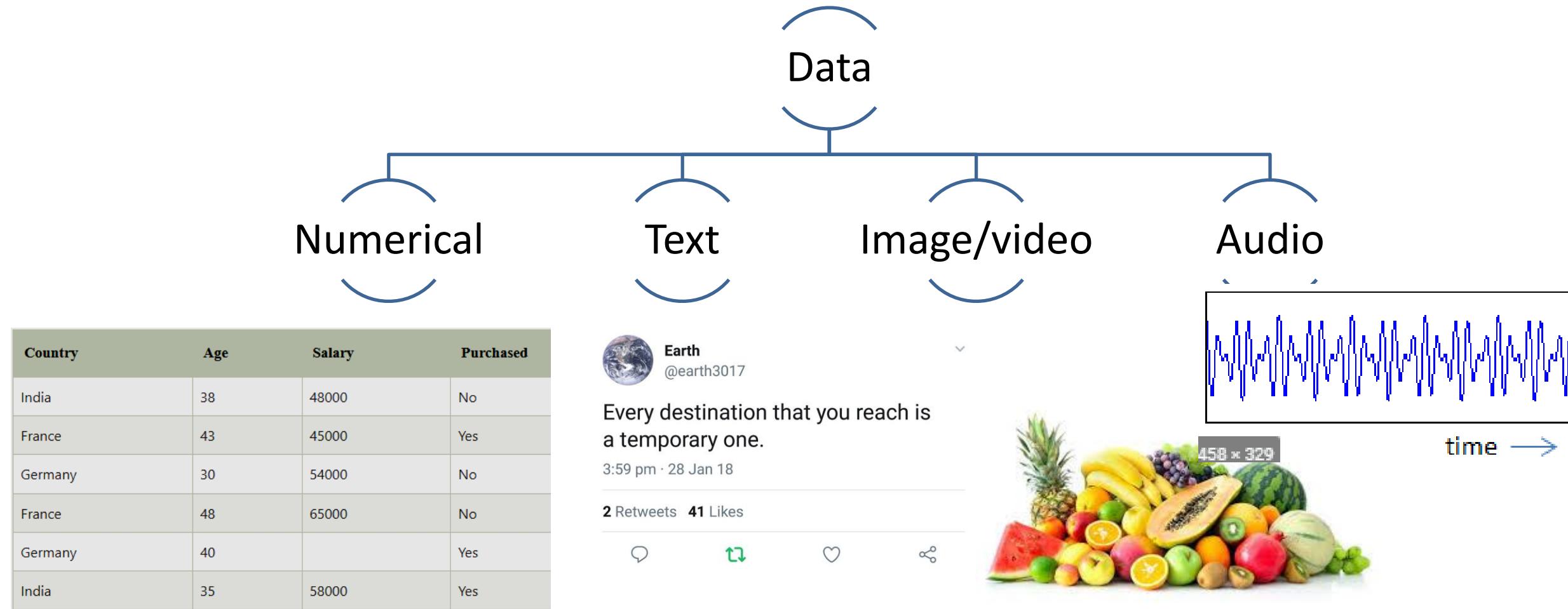
We need to identify the different data sources, as data can be collected from **various sources** such as files, database, internet, or mobile devices.

It is one of **the most important steps** of the life cycle.

→ **The quantity and quality** of the collected data will determine the efficiency of the output. The more will be the data, the more accurate will be the prediction.



# Types of data in datasets



# Popular sources for ML datasets

**Kaggle datasets:** <https://www.kaggle.com/datasets>.

**UCI machine learning repository:** <https://archive.ics.uci.edu/ml/index.php>.

**AWS resources:** <https://registry.opendata.aws/>.

**Google dataset search engine:** <https://toolbox.google.com/datasetsearch>.

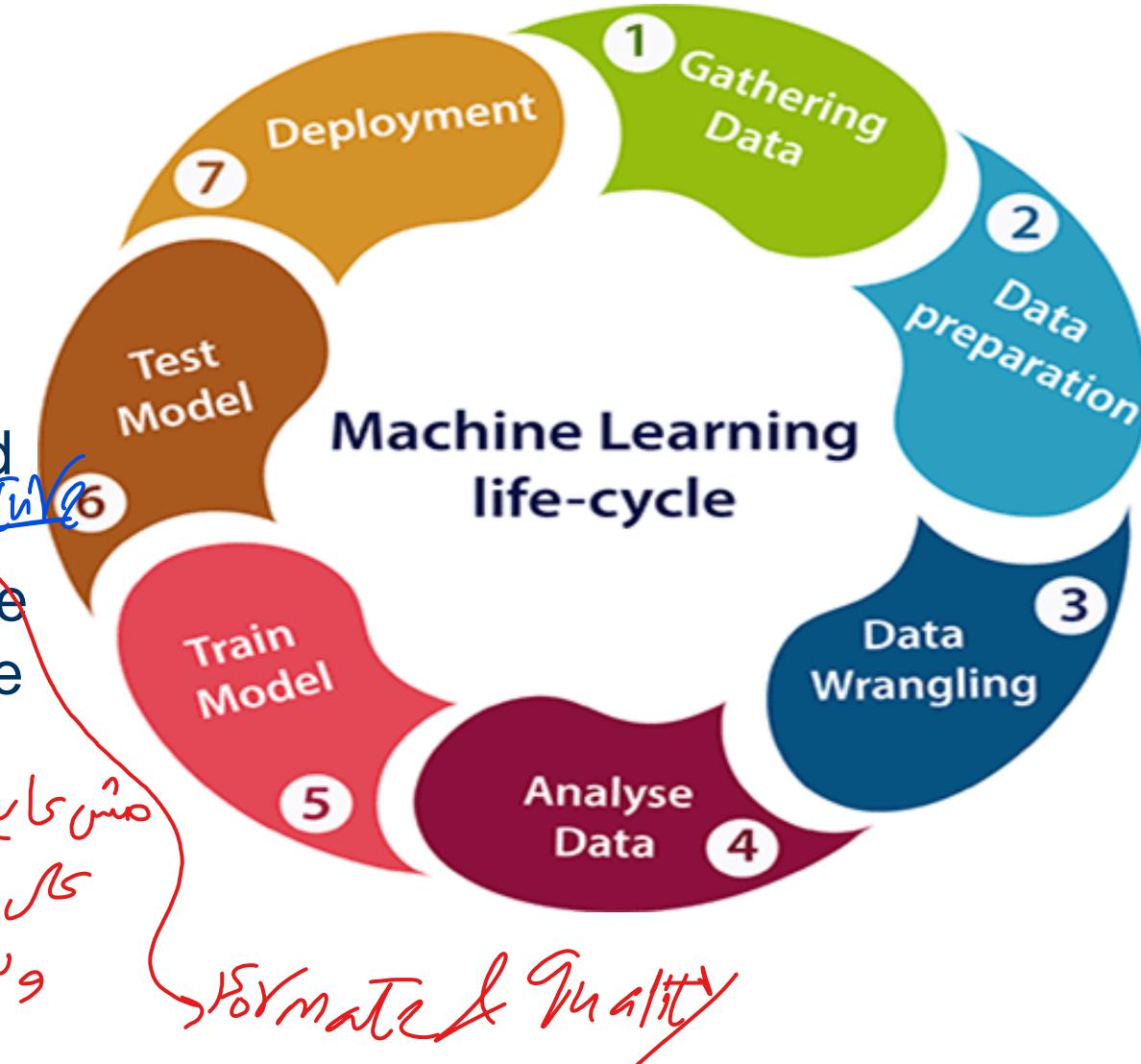
## 2- Data preparation

It is a step where we put our data into a suitable place and prepare it to use in our machine learning training.

In this step, first, we put all data together, and then **randomize the ordering of data**. *STRUCTURE*

Then, **understand the nature of data** that we have to work with. We need to understand the characteristics, **format**, and quality of data.

→ *item file (e.g. sequence) sequence (دروجات)  
لهم اذن لي بالرandonize الـ data لـ*  
*لـ وـ لـ، sequence) modeles (دروجات) لـ*



# 3- Data Wrangling (Cleansing / Preprocessing)

is the process of cleaning and converting raw data into a useable format.

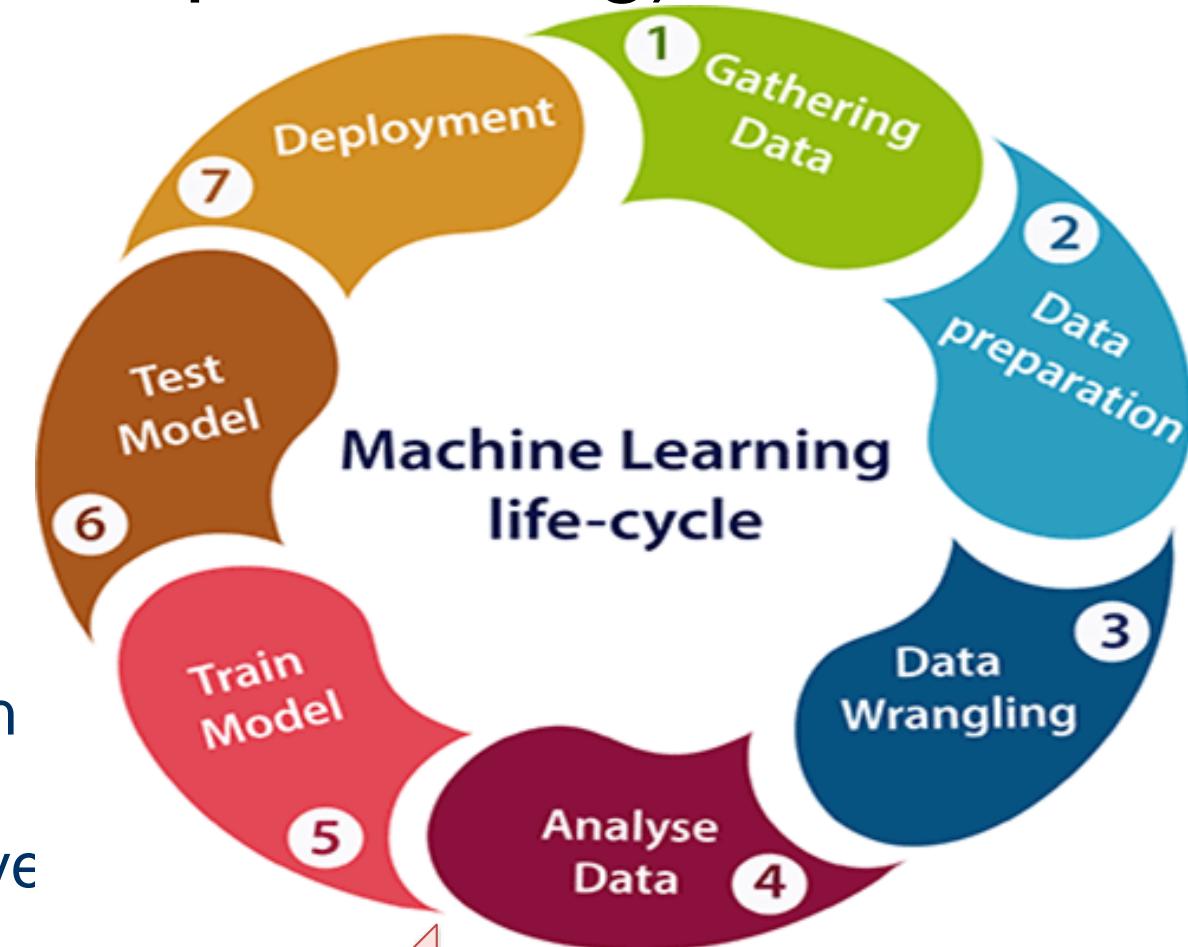
In real-world applications, collected data may have various issues, including:

- Missing Values
- Duplicate data
- Invalid data
- Noise

So, we use various filtering techniques to clean the data.

الزاص

It is **mandatory** to detect and remove the above issues because it can negatively affect the quality of the outcome.



Our Course starting

# Split Dataset

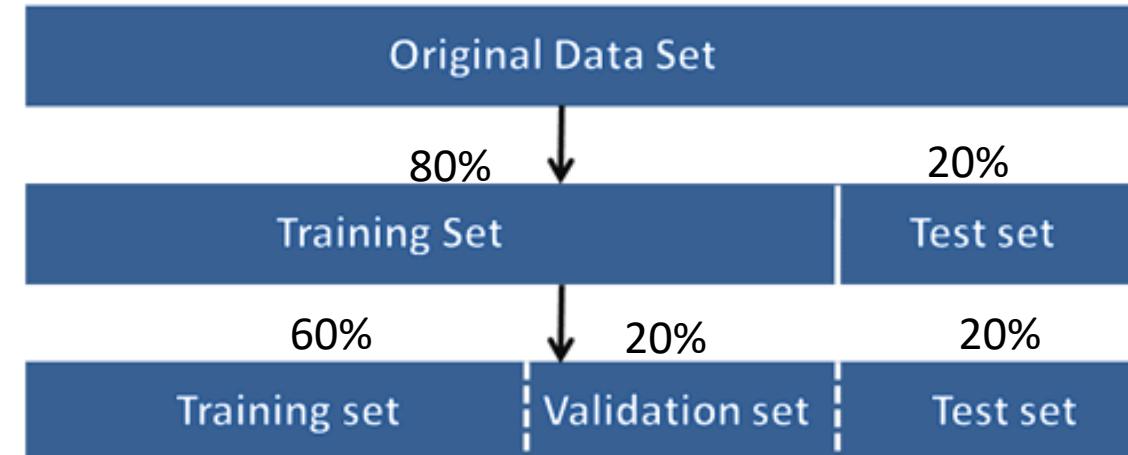
bài

During the development of the ML project, the developers completely rely on the datasets.

In building ML applications, datasets are divided into two parts:

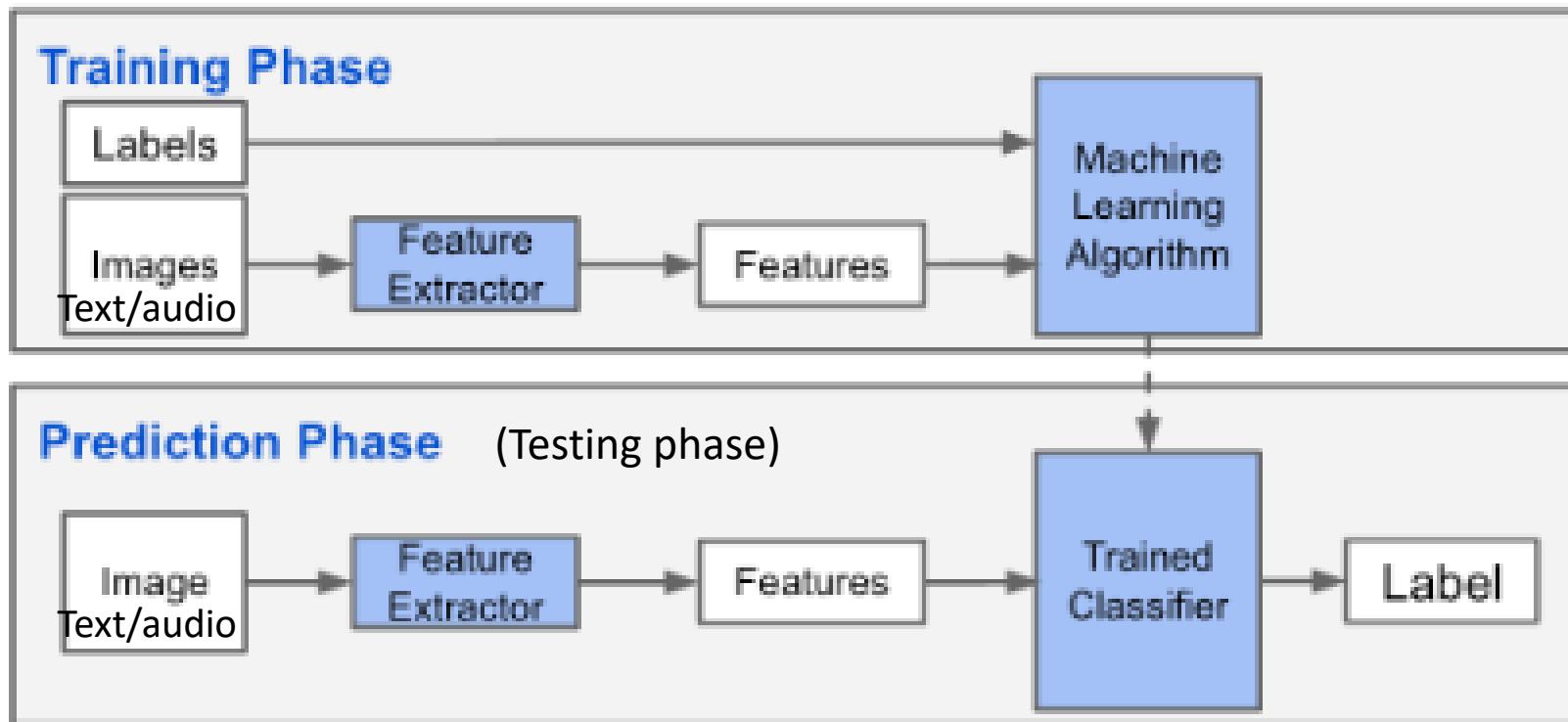
Training dataset:

Test Dataset



# Feature Extraction (Selection)

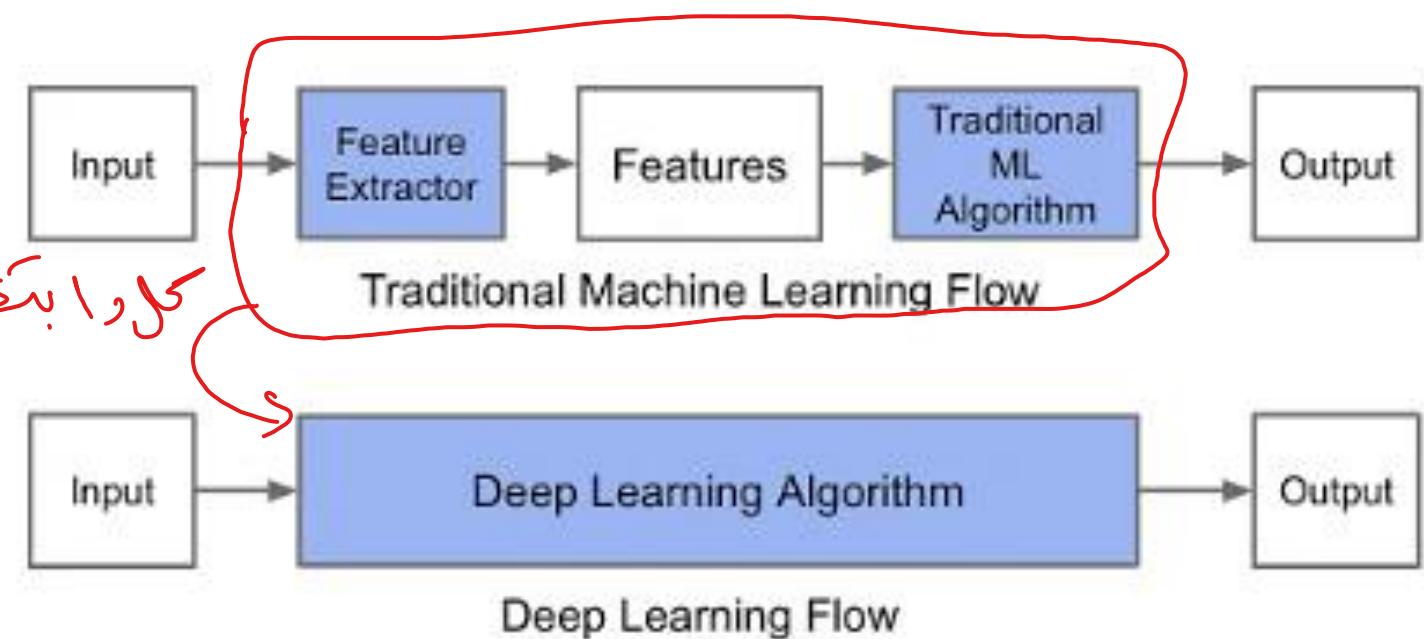
In **traditional machine learning** algorithms, **hand-craft** the features, i.e. select the features, is needed.



# Feature Extraction (Selection) Cont.

*testing / training / isig - mode / isig -> this feature extraction is automatically*

In **traditional machine learning** algorithms, hand-craft the features is needed.  
By contrast, in **deep learning algorithms** feature engineering is done **automatically** by the algorithm.



Deep learning versus Traditional ML algorithms

# 4- Data Analysis

*Select techniques , Build, and evaluate models*

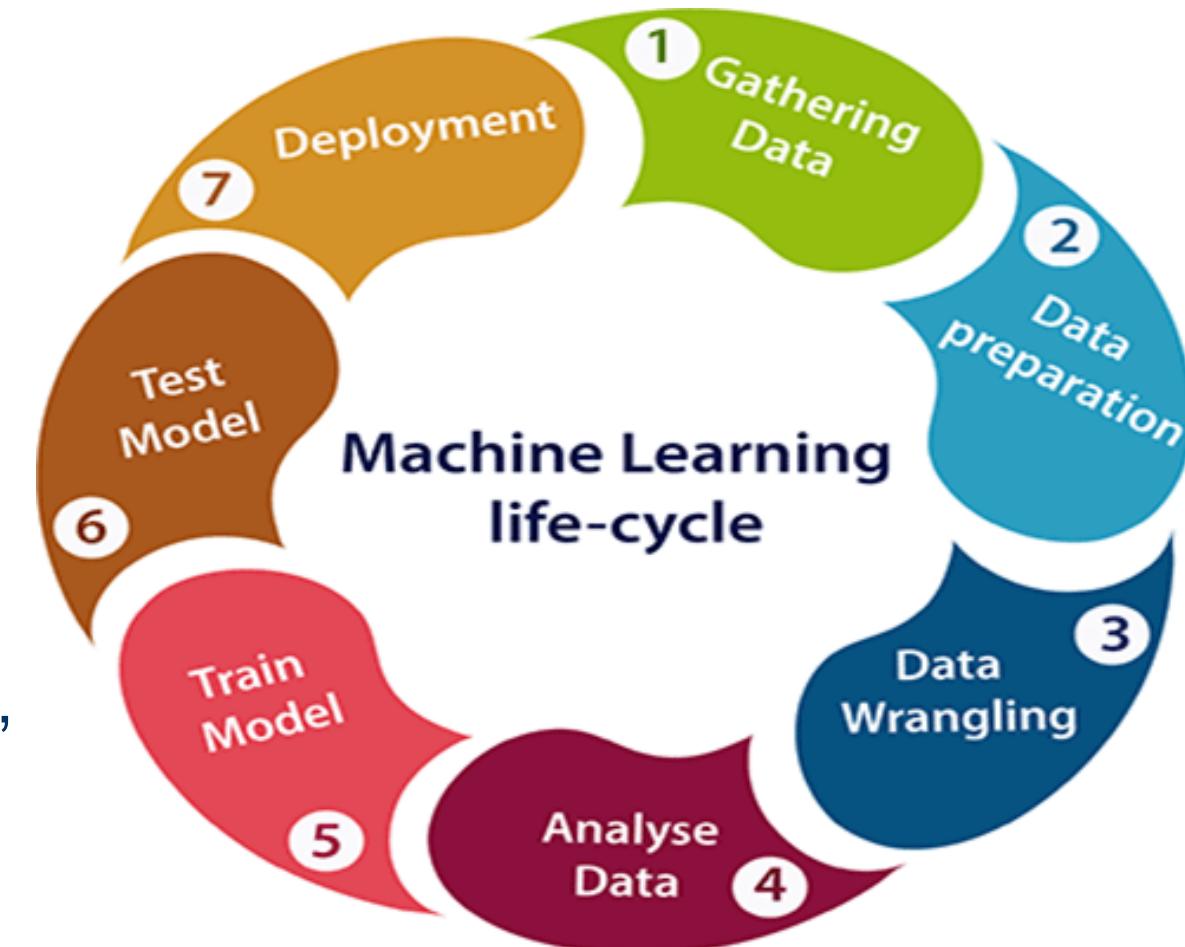
This step involves:

**a) Selection of analytical techniques**

Select ML techniques such as Classification, Regression, Cluster analysis, Association, etc.

**b) Building & evaluate models**

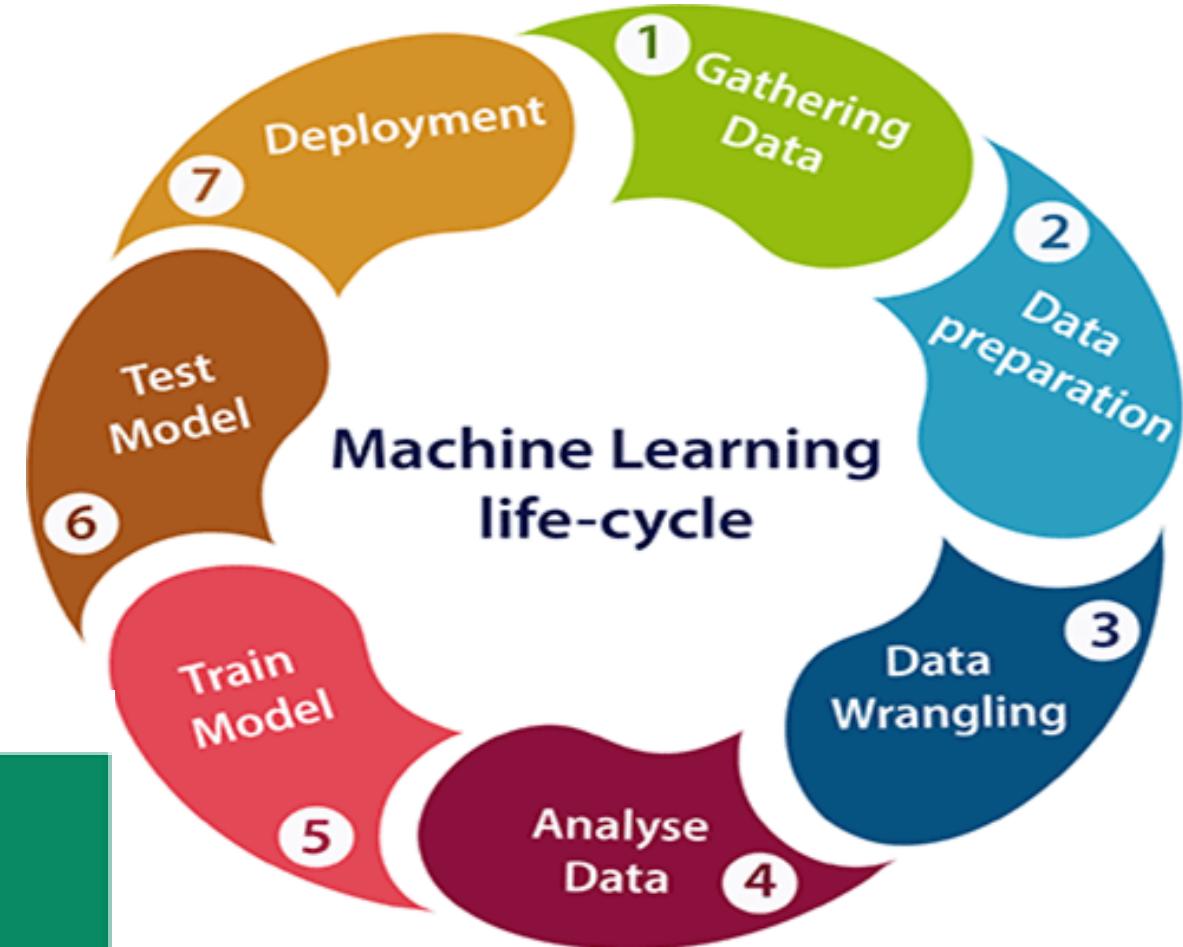
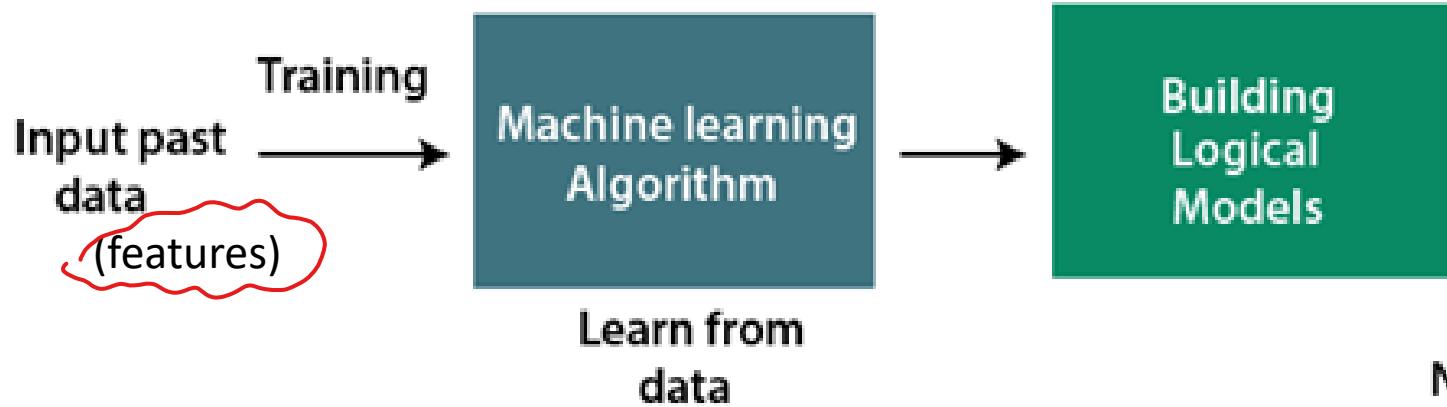
build the selected model using prepared data, and evaluate the model.



# 5- Train Model

We use datasets to train the model using various machine learning algorithms.

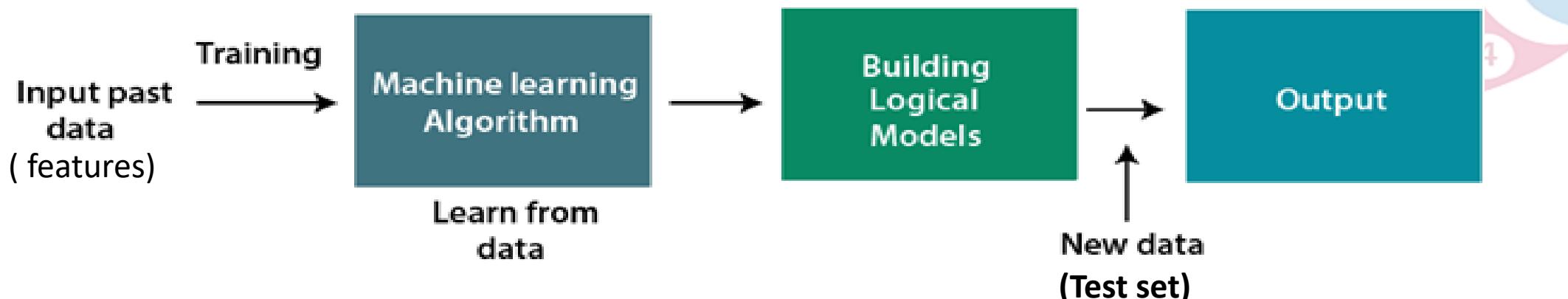
Training a model is required so that it can understand the various patterns, and rules.



# 6- Test Model

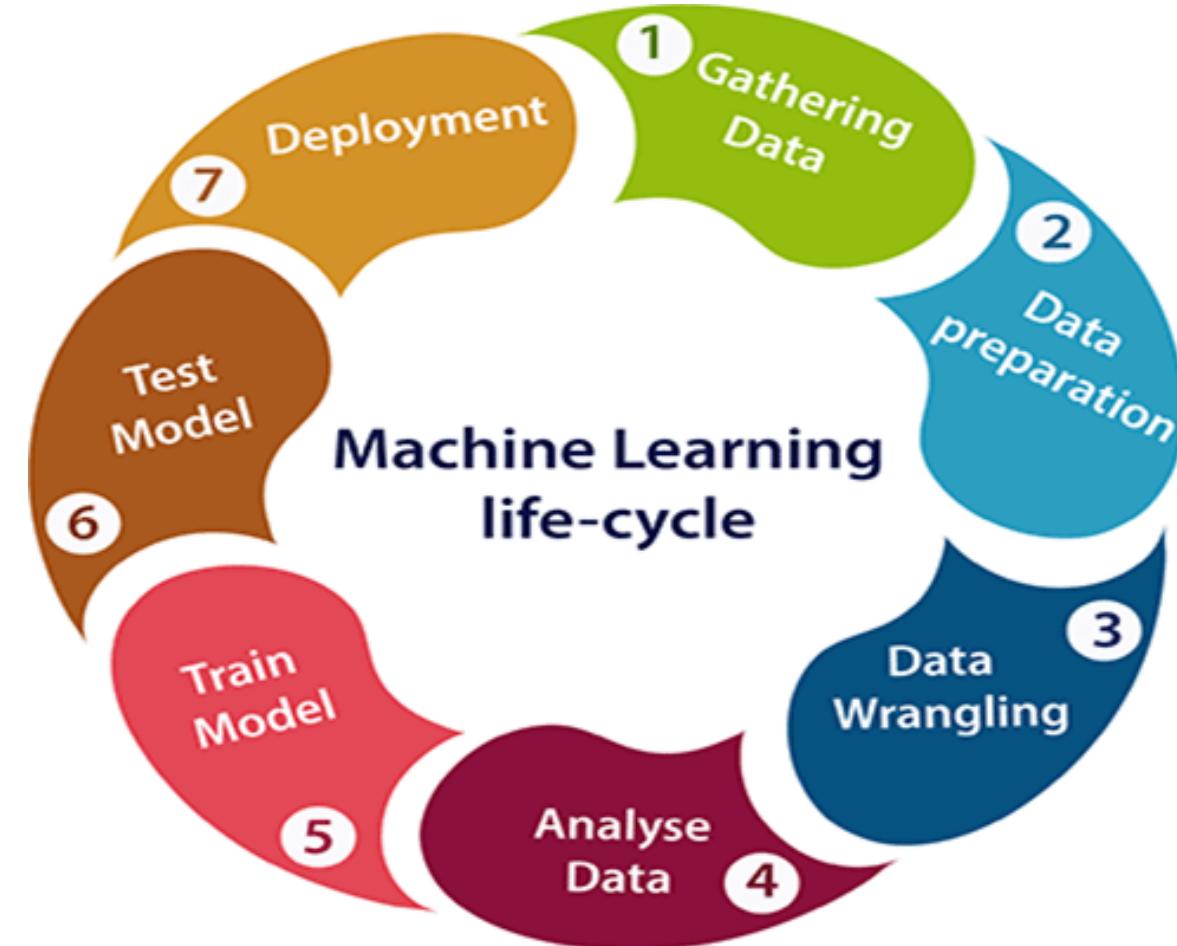
In this step, we check for the accuracy of our model by providing a test dataset to it.

Testing the model determines the **percentage accuracy** of the model as per the requirement of project or problem.



# 7- Deployment

If the above-prepared model is producing an accurate result as per our requirement with acceptable speed, then we deploy the model in the real system.



# **Types of Machine Learning**

# Types of Machine Learning

make prediction  
أُولَئِكَ الْأَوْفَى

→ labeled data

**Supervised**

- Classification
- Regression
- Ranking

Training data includes desired outputs

this technique  
is used in the  
Pre-processing  
data stage

→ unlabeled data  
→ don't make prediction

## Unsupervised

- Clustering
- Dim. Reduction
- Association

Training data does not include desired outputs

Big data suggest  
involves a lot

reduce the number of variables in  
data

the algorithms look for relationships between variables  
in the data

→ the algorithm knows what the correct is but unsupervised doesn't

is the environment) more often  
and update → feedback  
will be given for bad (not seen  
actions (seen)

Reinforcement

Learn from mistakes & rewards

Train + test → is useful  
Testing New training's model  
ans Train Just its own reward  
move data ?

# Supervised vs. Unsupervised Learning

Problem type	Description	
Ranking	Helping users find the most relevant <i>items</i>	
Recommendation	Giving users the <i>items</i> they may be most interested in	
Classification	Figuring out what category does an <i>item</i> belongs to	
Regression	Predicting a numerical value of an <i>item</i>	
Clustering	Putting similar <i>items</i> together	
Anomaly Detection	Finding uncommon <i>items</i>	

**Supervised Learning**

Data is provided **with** the correct labels

**Unsupervised Learning**

Data is provided **without** labels

means that the output is already known to you

means there is no fixed output variable

# Supervised vs. Unsupervised Learning

Problem type	Description	
Ranking	Helping users find the most relevant <i>items</i>	
Recommendation	Giving users the <i>items</i> they may be most interested in	
Classification	Figuring out what category does an <i>item</i> belongs to	
Regression	Predicting a numerical value of an <i>item</i>	
Clustering	Putting similar <i>items</i> together	
Anomaly Detection	Finding uncommon <i>items</i>	

Supervised Learning  
↳ inputs with outputs

Unsupervised Learning  
↳ unlabeled inputs

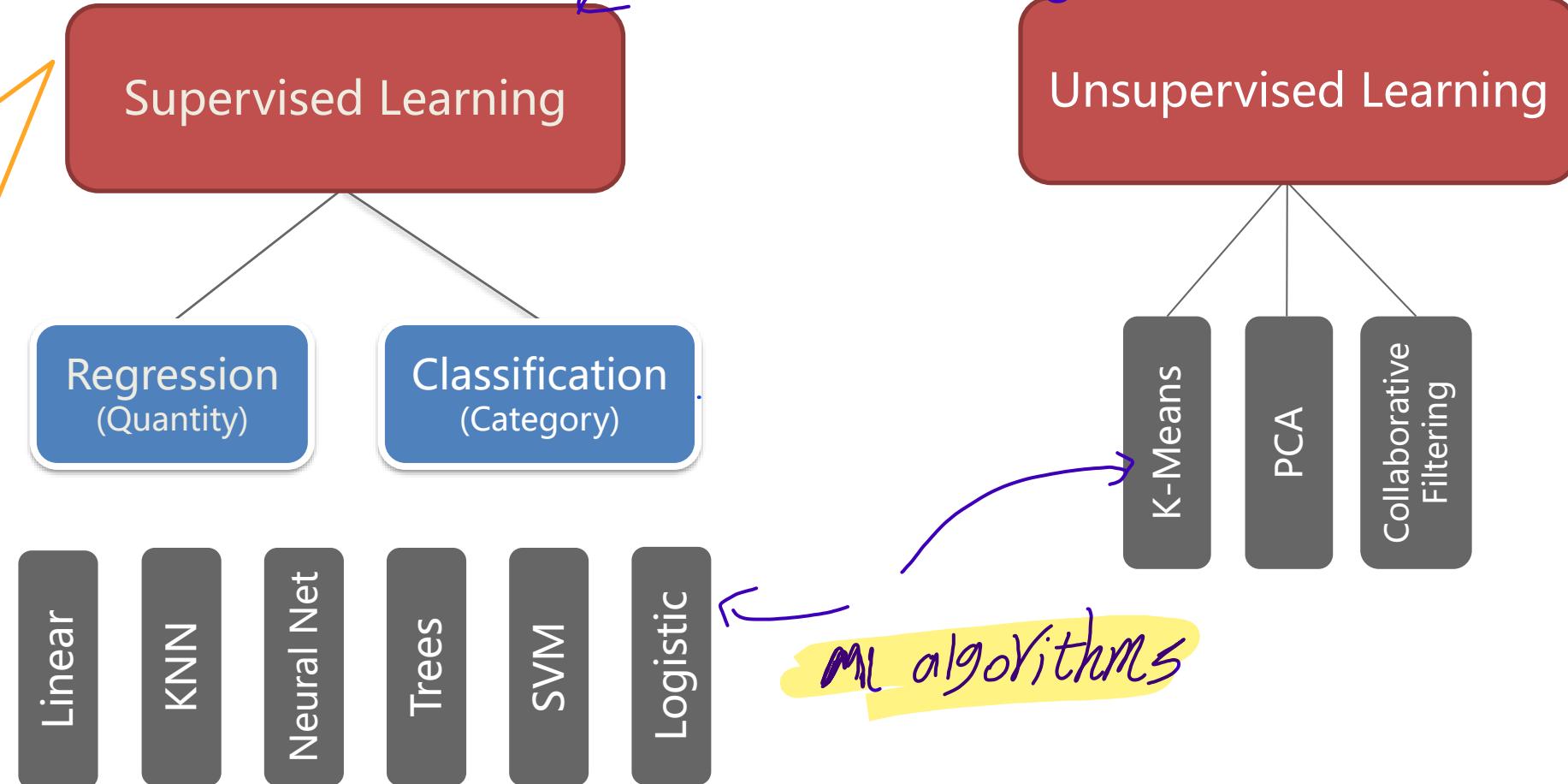
Data is provided **with** the correct labels

Data is provided **without** labels

# Supervised vs. Unsupervised Learning

ML Techniques

Data is provided with the correct labels  
Model learns by looking at these examples

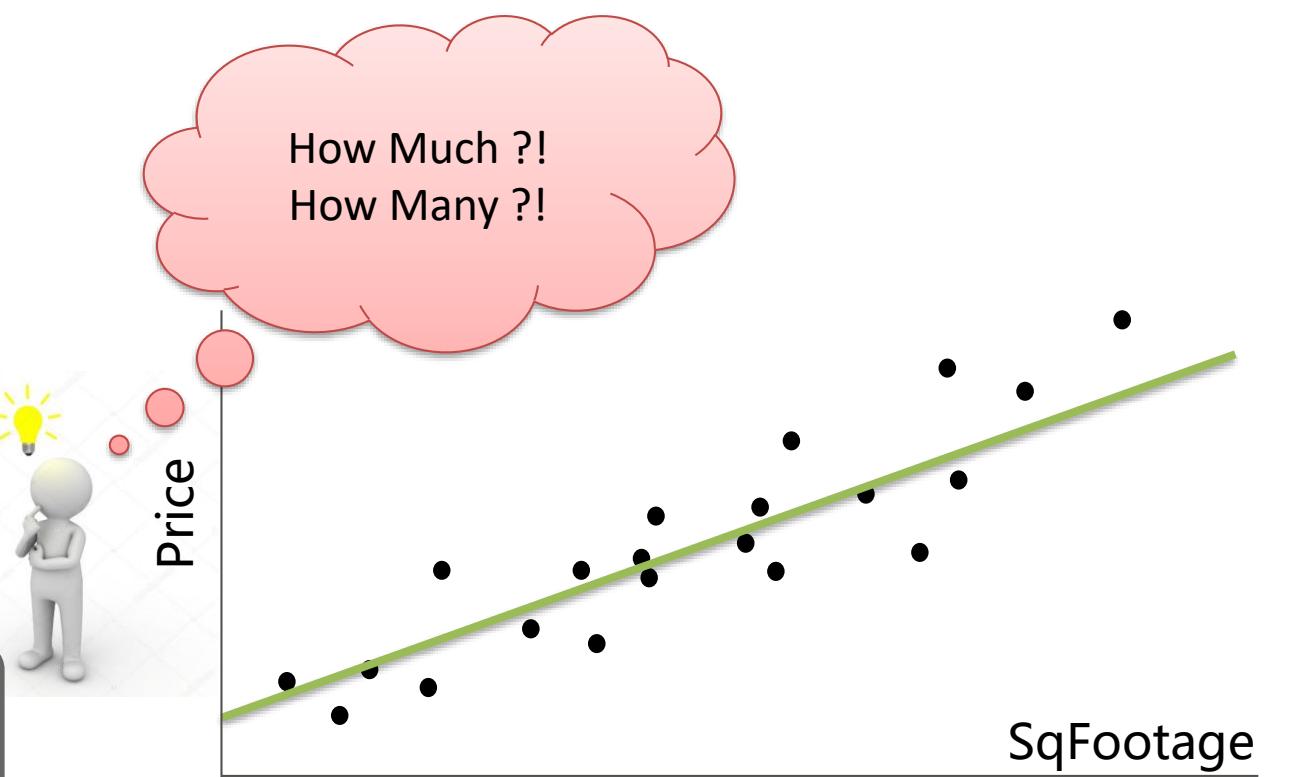
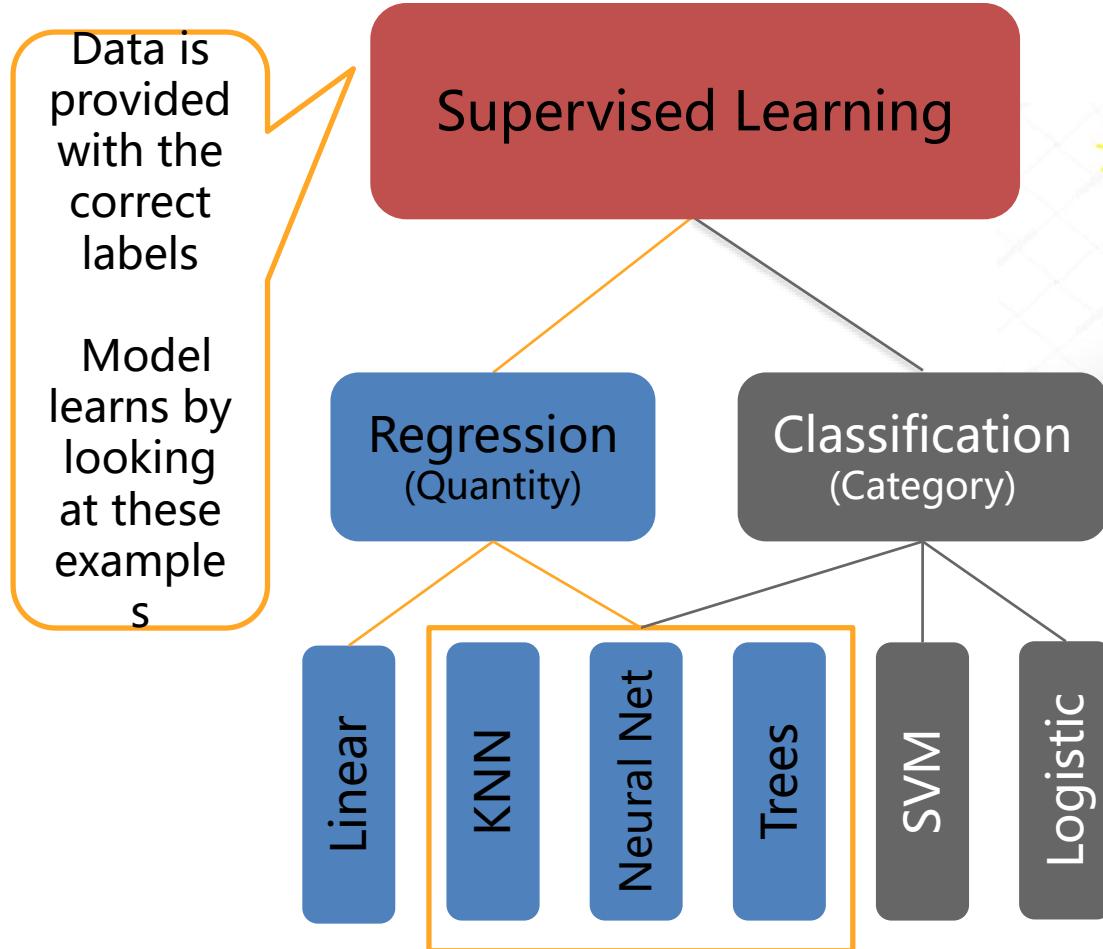


Data is provided without labels  
Model finds patterns in data and features

ML algorithms

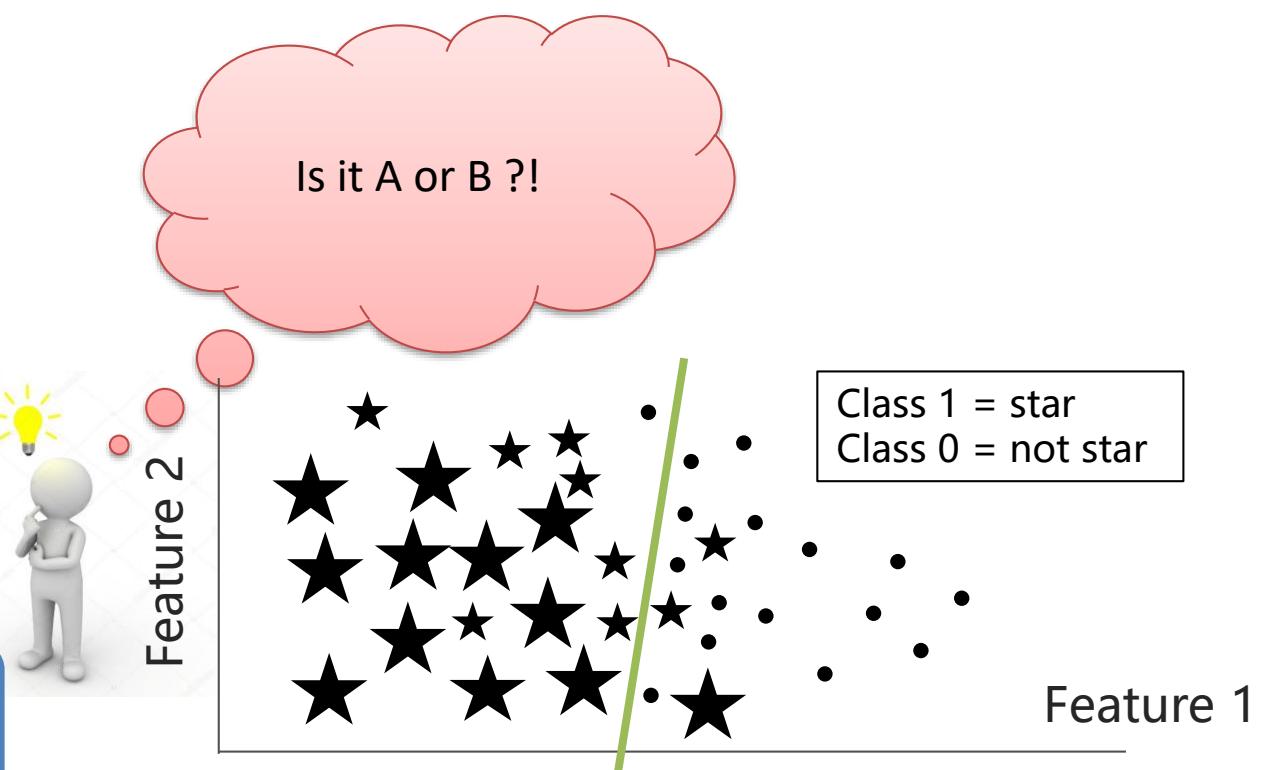
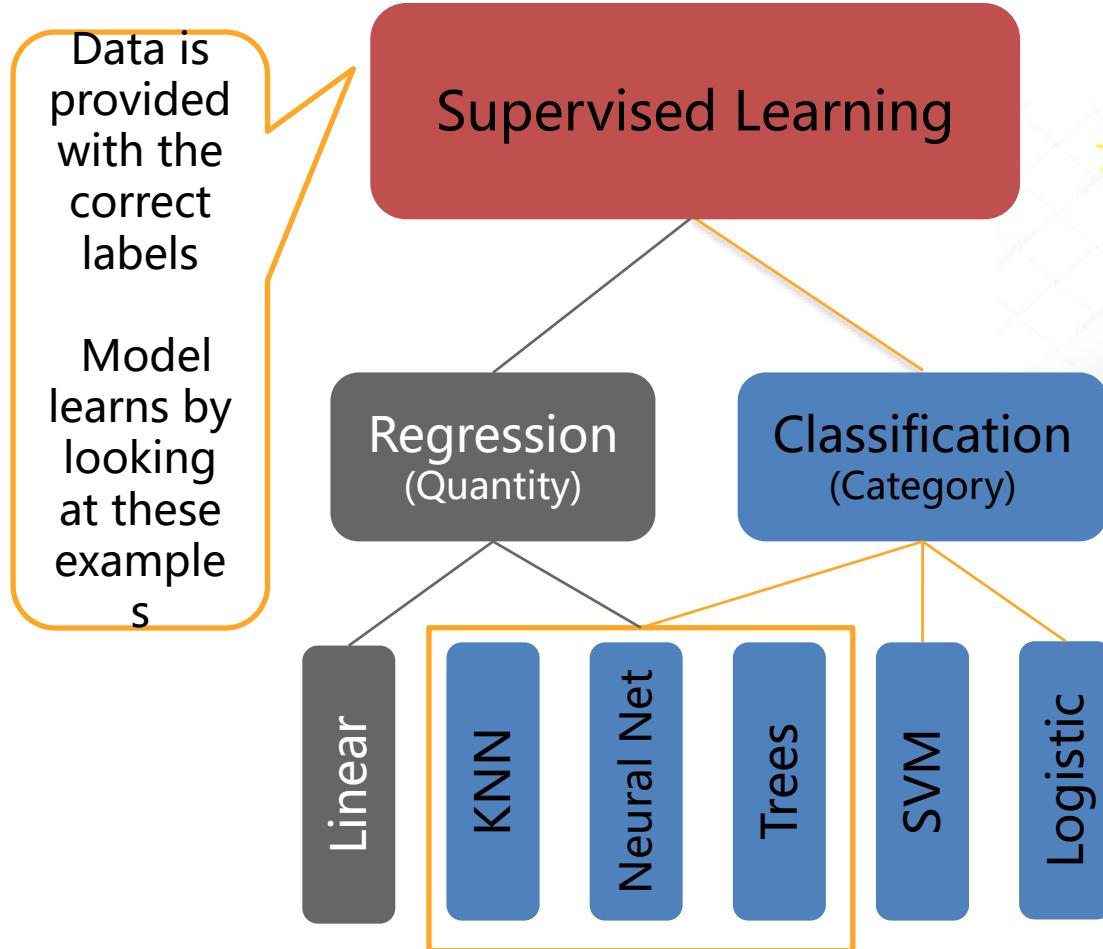
and features

## Supervised Learning: Regression



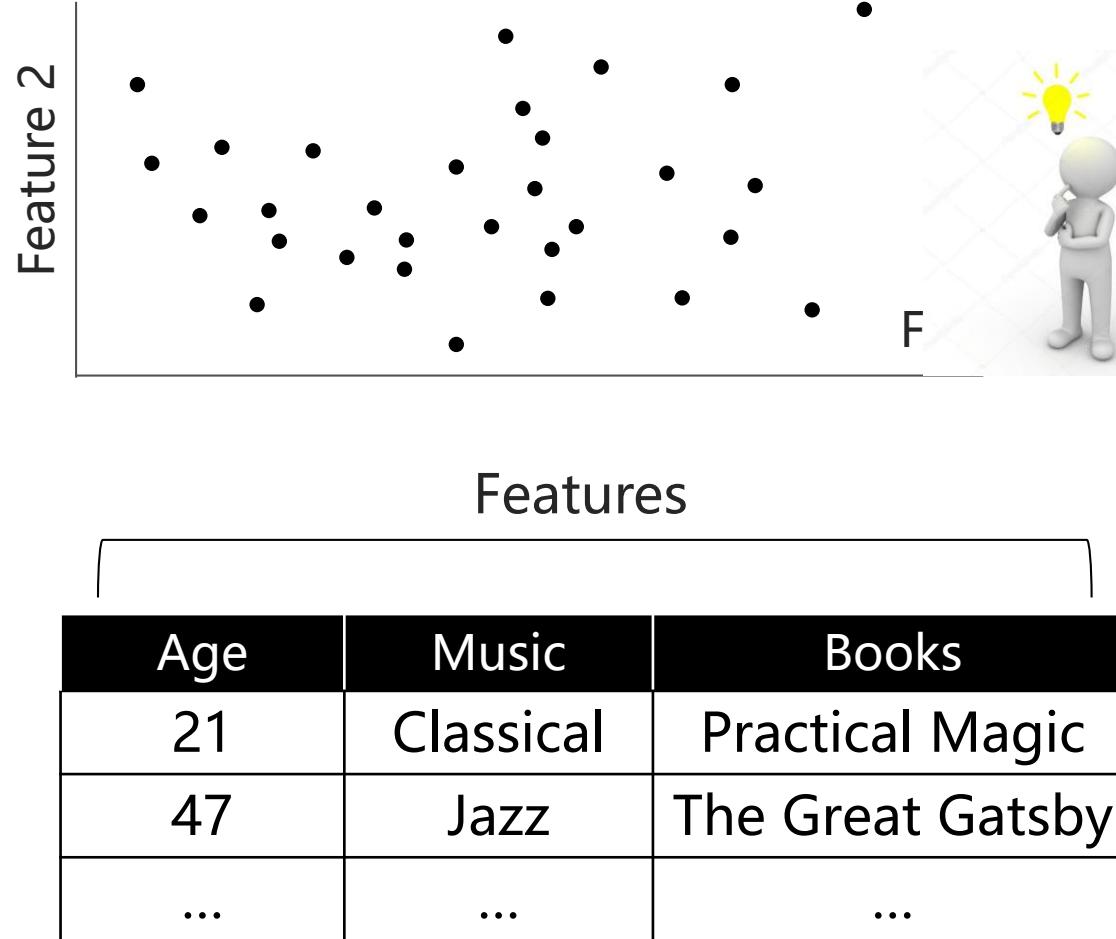
Label		Features		
	Price	Bedrooms	SqFootage	Age
	280.000	3	3292	14
	210.030	2	2465	6
	...	...	...	...

# Supervised Learning: Classification

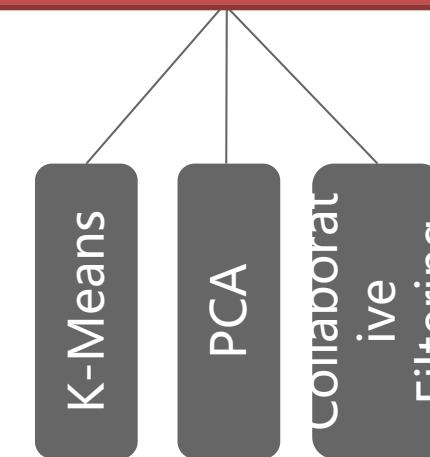


Label		Features		
	Star	Points	Edges	Size
1	1	5	10 <	750
0	0	0	> 9	150
...	...	...	...	...

## Unsupervised Learning: Clustering



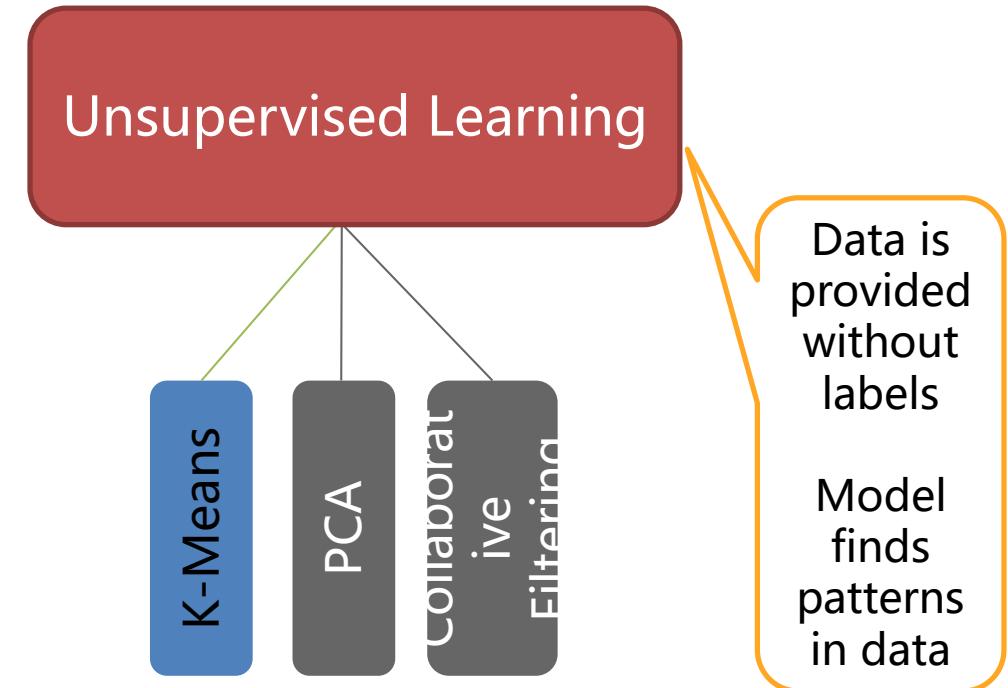
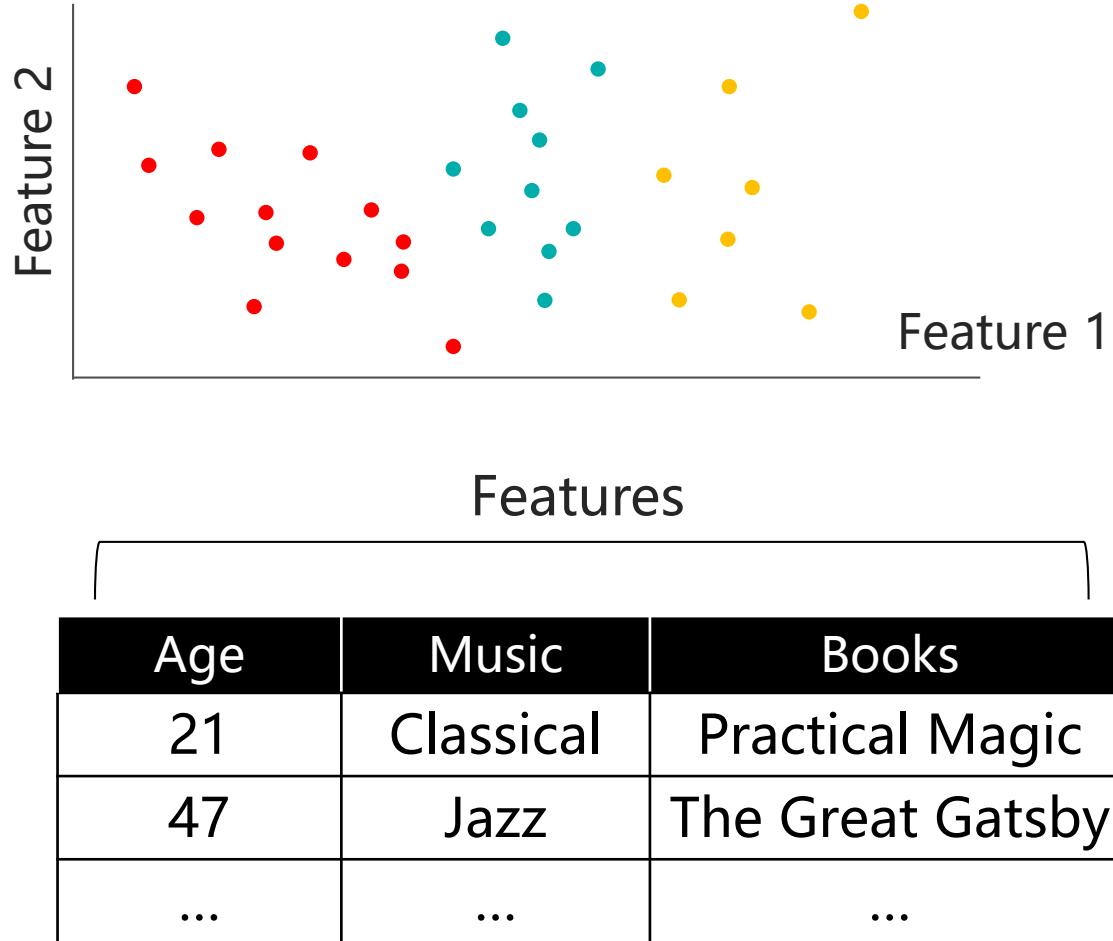
Unsupervised Learning



Data is provided without labels

Model finds patterns in data

## Unsupervised Learning: Clustering



What should I do now ?!

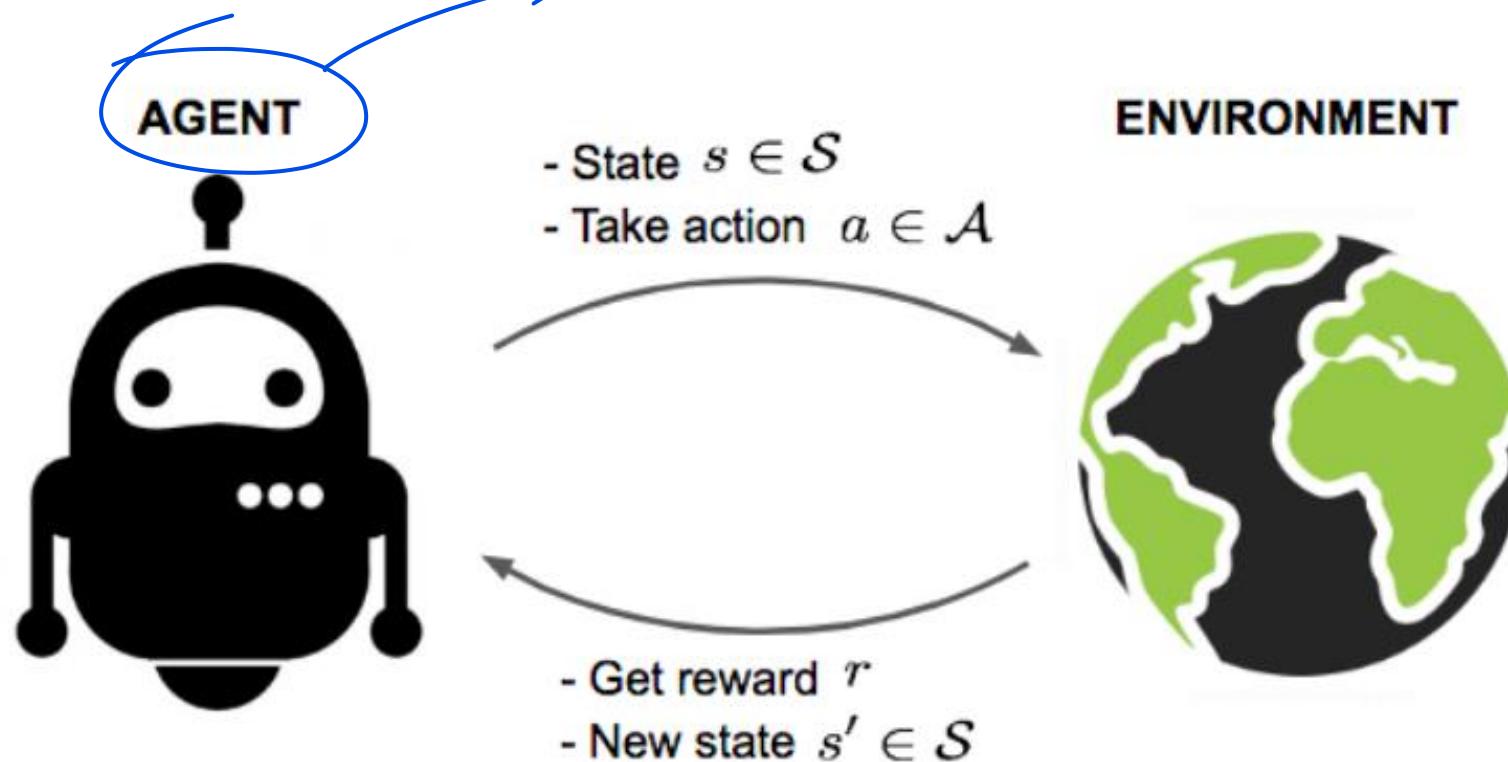


## Reinforcement Learning



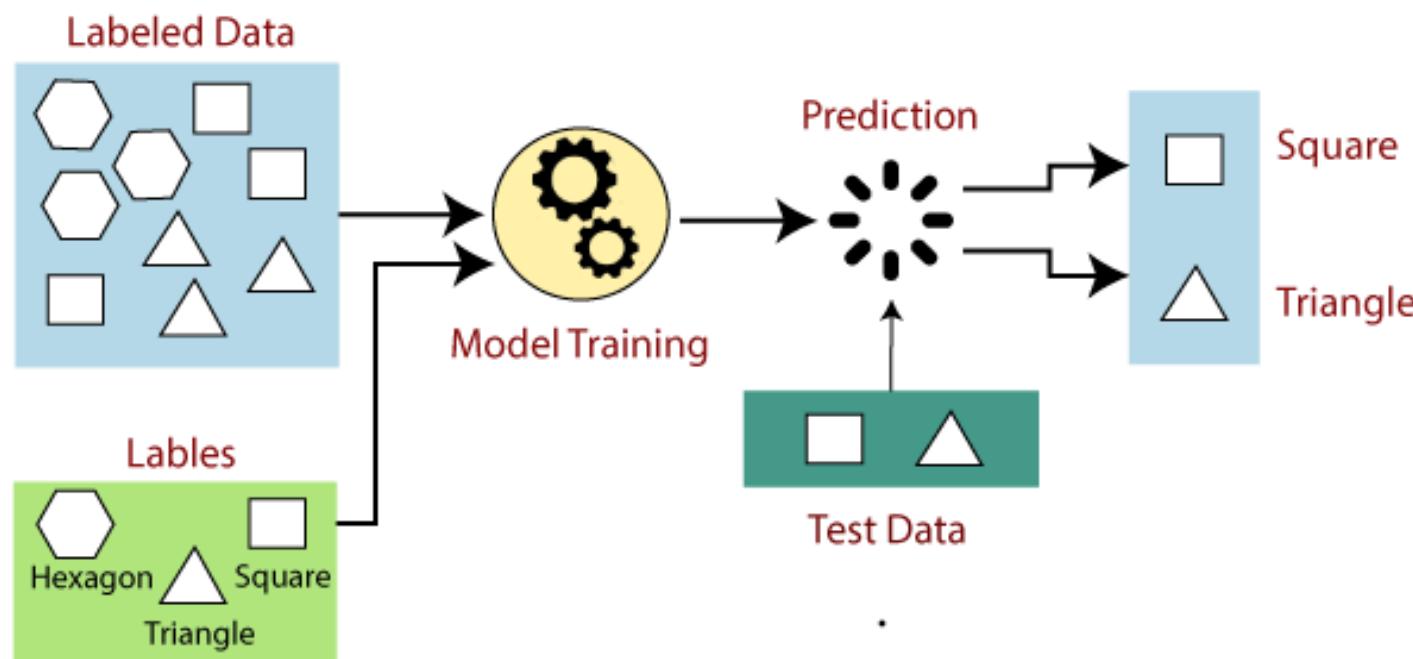
# Reinforcement Learning

model  $\pi$  is learned  
with update,  $\rightarrow$  feedback from environment  $\rightarrow$  model  $\pi$   
so policy  $\pi$  is updated based on  
Training & Testing



# Supervised Learning

In supervised learning, models are trained using **labeled dataset**, where the model learns about each type of data. Once the training process is completed, the model is tested on the basis of test data (a subset of the training set), and then it predicts the output.



# Supervised Learning (Cont.)

## Regression

It is used for the prediction of **continuous variables**

Popular regression algorithms:

**Linear Regression**

Regression Trees

Non-Linear Regression

Bayesian Linear Regression

Polynomial Regression

## Classification

It is used when the output variable is **categorical** (classes).

Popular classification algorithms:

Random Forest

Decision Trees

**Logistic Regression**

**Support vector Machines (SVMs)**

**Neural Networks (NNs)**

# Advantages and Disadvantages of Supervised learning

## Advantages of Supervised Learning

With the help of supervised learning, the model can predict the output on the basis of prior experiences.

Supervised learning model helps us to solve various real-world problems such as **fraud detection, spam filtering**, etc.

## Disadvantages of Supervised Learning

Supervised learning cannot predict the correct output if the test data is different from the training dataset.

Training required lots of computation times.

In supervised learning, we need enough knowledge about the classes of object.

# Unsupervised Learning

- Here, we have taken **unlabeled input data**, which means it is not categorized and corresponding outputs are also not given.

**Clustering:** Clustering is a method **of grouping objects into clusters** such that objects with the most similarities remain in a group and have fewer or no similarities with the objects of another group.

**Association:** An association rule is an unsupervised learning method that is used for **finding the relationships between variables** in a large database. It determines the set of items that occurs together in the dataset. The Association rule makes marketing strategy more effective.

**Example:** People who buy X items (bread) also tend to purchase Y (Butter/Jam) items.

# Advantages and Disadvantages of Unsupervised Learning

## Advantages of Unsupervised Learning

Unsupervised learning is used for **more complex tasks** as compared to supervised learning because, in unsupervised learning, we don't have labeled input data.

Unsupervised learning is preferable as it is **easy to get unlabeled data** in comparison to labeled data.

## Disadvantages of Unsupervised Learning

Unsupervised learning is intrinsically **more difficult** than supervised learning as it **does not have corresponding output**.

The result of the unsupervised learning algorithm might be **less accurate** as input data is not labeled, and algorithms do not know the exact output in advance.

# Thanks