

- Popular algorithms: Regression, Classification, and Clustering
- Recommender Systems: Content-Based and Collaborative Filtering
- Popular models: Train/Test Split, Gradient Descent, and Mean Squared Error
- Get ready to do more learning than your machine!

# Syllabus

## Module 1 - Machine Learning

- Python for Machine Learning
- Supervised vs Unsupervised
- Lab & Review

## Module 2 - Regression

- Simple Linear Regression
- Multiple Linear Regression
- Model Evaluation in Regression Models
- Non-Linear Regression
- Lab & Review

## Module 3 - Classification

- K-Nearest Neighbors
- Decision Trees
- Evaluation Metrics in Classification
- Logistic Regression vs Linear Regression
- Support Vector Machine (SVM)
- Lab & Review

## Module 4 - Clustering

- K-Means Clustering
- Hierarchical Clustering
- DBSCAN
- Lab & Review

## Module 5 - Recommender Systems

- Content-Based Recommender Systems
- Collaborative Filtering
- Lab & Review

# Learning Objectives

In this lesson you will learn about:

- Machine Learning applications
- Python libraries for Machine Learning

- Supervised vs Unsupervised Learning

# Module 1

## Introduction to machine learning

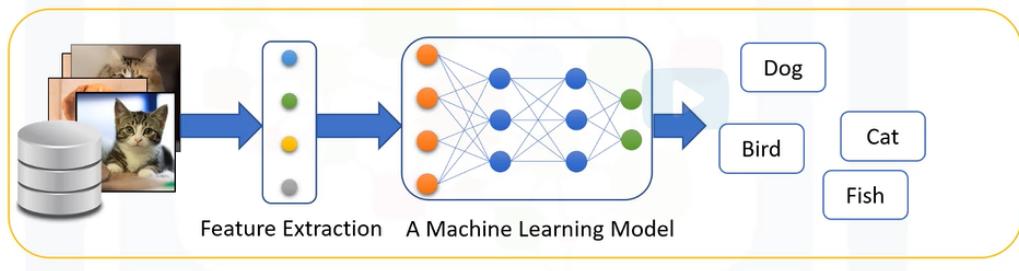
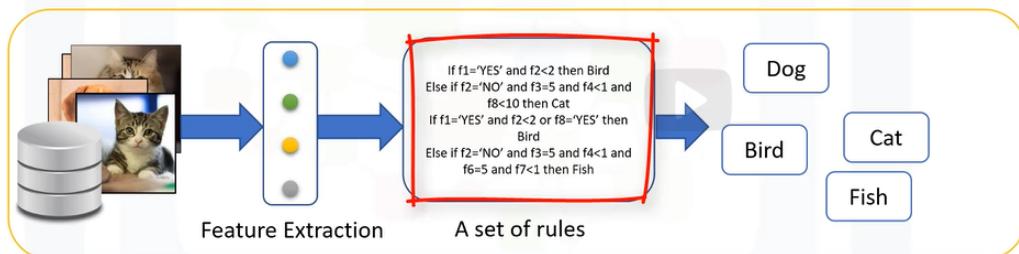
### What is machine learning

**Machine learning is the subfield of computer science that gives “computers the ability to learn without being explicitly programmed.”**

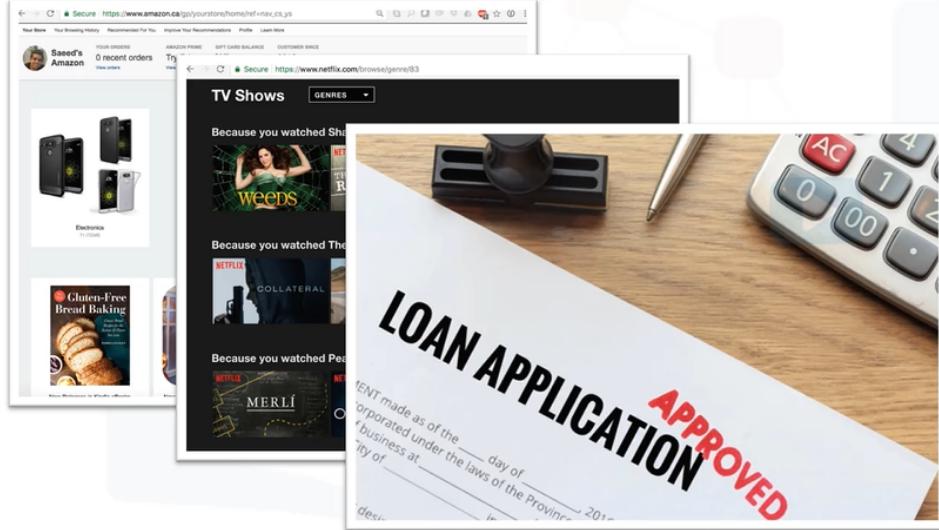
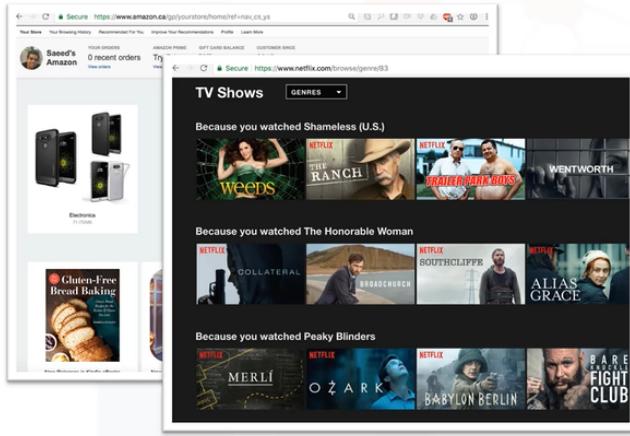
**Arthur Samuel**

American pioneer in the field of computer gaming and artificial intelligence, coined the term "machine learning" in 1959 while at IBM.

### How machine learning works ?



### Examples of machine learning



## Major machine learning techniques

- Regression/Estimation
  - Predicting continuous values
- Classification
  - Predicting the item class / category a case
- Clustering
  - Finding the structure of data; summarization
- Association
  - Association frequent co-occurring items / events
- Anomaly detection
  - Discovering abnormal and unusual cases
- Sequence mining
  - Predicting next events; click-stream (Markov Model, HMM)
- Dimension Reduction
  - Reducing the size of data (PCA)
- Recommendation systems
  - Recommending items

## Difference between artificial intelligence, machine learning and deep learning

- **AI components:**

- Computer Vision
- Language Processing
- Creativity
- Etc.



- **Machine learning:**

- Classification
- Clustering
- Neural Network
- Etc.

- **Revolution in ML:**

- Deep learning

**Let's get started with machine learning!**

- Machine learning applications
- Machine learning algorithms

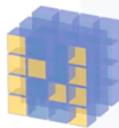


## Python for Machine Learning

### Python libraries for machine learning

- Numpy
- Scipy
- matplotlib
- pandas

- sklearn



NumPy



SciPy



## More about scikit-learn

- Free software machine learning library
- Classification, Regression and Clustering algorithms
- Works with NumPy and SciPy
- Great documentation
- Easy to implement



## Scikit-learn functions

```
from sklearn import preprocessing
X = preprocessing.StandardScaler().fit(X).transform(X)
```

```
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.33)
```

```
from sklearn import svm
clf = svm.SVC(gamma=0.001, C=100.)
```

```
clf.fit(X_train, y_train)
```

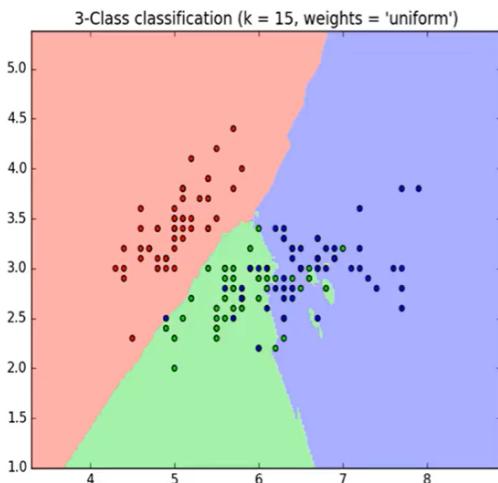
```
clf.predict(X_test)
```

```
from sklearn.metrics import confusion_matrix
print(confusion_matrix(y_test, yhat, labels=[1,0]))
```

```
import pickle
s = pickle.dumps(clf)
```

## Supervised vs Unsupervised

### What is supervised learning

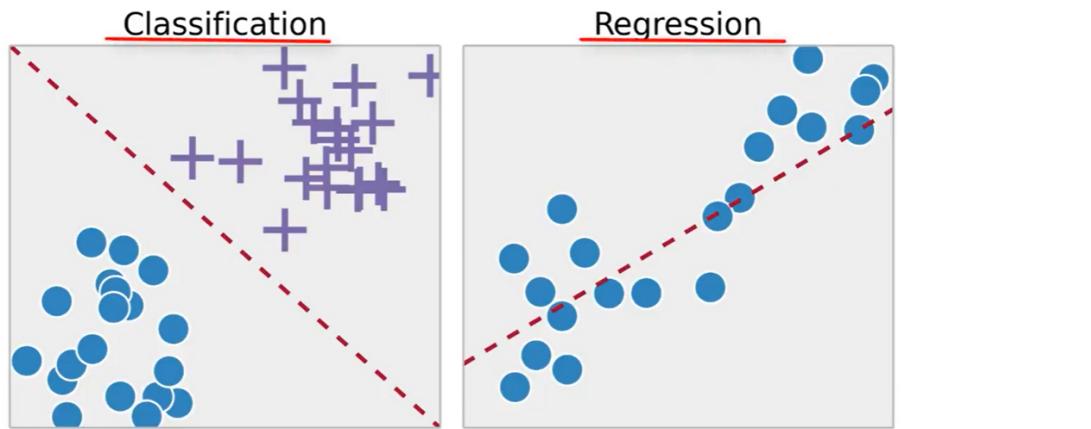


We “teach the model,” then with that knowledge, it can predict unknown or future instances.

## Teaching the model with labeled data

ID	Clump	UnifSize	UnifShape	MargAdh	SingEpiSize	BareNuc	BlandChrom	NormNucl	Mit	Class
1000025	5	1	1	1	2	1	3	1	1	benign
1002945	5	4	4	5	7	10	3	2	1	benign
1015425	3	1	1	1	2	2	3	1	1	malignant
1016277	6	8	8	1	3	4	3	7	1	benign
1017023	4	1	1	3	2	1	3	1	1	benign
1017122	8	10	10	8	7	10		7	1	malignant
1018099	1	1	1	1	2	10	3	1	1	benign
1018561	2	1	2	H	2	1	3	1	1	benign
1033078	2	1	1	1	2	1	1	1	5	benign
1033078	4	2	1	1	2	1	2	1	1	benign

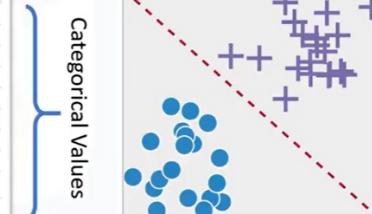
## Type of supervised learning



## What is classification ?

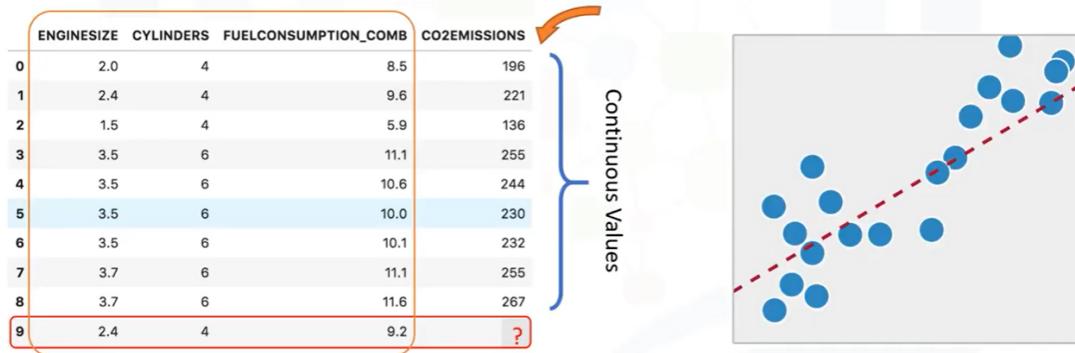
**Classification** is the process of predicting discrete class labels or categories.

ID	Clump	UnifSize	UnifShape	MargAdh	SingEpiSize	BareNuc	BlandChrom	NormNucl	Mit	Class
1000025	5	1	1	1	2	1	3	1	1	benign
1002945	5	4	4	5	7	10	3	2	1	benign
1015425	3	1	1	1	2	2	3	1	1	malignant
1016277	6	8	8	1	3	4	3	7	1	benign
1017023	4	1	1	3	2	1	3	1	1	benign
1017122	8	10	10	8	7	10		7	1	malignant
1018099	1	1	1	1	2	10	3	1	1	benign
1018561	2	1	2	H	2	1	3	1	1	benign
1033078	2	1	1	1	2	1	1	1	5	benign
1033078	4	2	1	1	2	1	2	1	1	benign



## What is regression ?

Regression is the process of predicting continuous values.



## What is unsupervised learning

Customer Id	Age	Edu	Years Employed	Income	Card Debt	Other Debt	Address	DebtIncomeRatio
1	41	2		6	19	0.124	1.073 NBA001	6.3
2	47	1		26	100	4.582	8.218 NBA021	12.8
3	33	2		10	57	6.111	5.802 NBA013	20.9
4	29	2		4	19	0.681	0.516 NBA009	6.3
5	47	1		31	253	9.308	8.908 NBA008	7.2
6	40	1		23	81	0.998	7.831 NBA016	10.9
7	38	2		4	56	0.442	0.454 NBA013	1.6
8	42	3		0	64	0.279	3.945 NBA009	6.6
9	26	1		5	18	0.575	2.215 NBA006	15.5
10	47	3		23	115	0.653	3.947 NBA011	4
11	44	3		8	88	0.285	5.083 NBA010	6.1
12	34	2		9	40	0.374	0.266 NBA003	1.6

Unsupervised learning techniques:

- Dimension reduction
- Density estimation
- Market basket analysis /
- Clustering

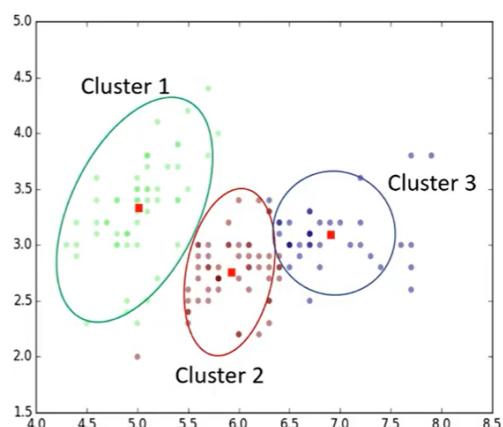
ALL OF THIS DATA IS UNLABELED

The model works on its own to discover information.

## What is clustering ?

Clustering is grouping of data points or objects that are somehow similar by:

- Discovering structure
- Summarization
- Anomaly detection



## Supervised vs unsupervised learning

## Supervised Learning

- **Classification:**  
Classifies labeled data
- **Regression:**  
Predicts trends using previous labeled data
- Has more evaluation methods than unsupervised learning
- Controlled environment

## Unsupervised Learning

- **Clustering:**  
Finds patterns and groupings from unlabeled data
- Has fewer evaluation methods than supervised learning
- Less controlled environment

# Module 2

## Learning Objectives

In this lesson you will learn about:

- Regression Algorithms
- Model Evaluation
- Model Evaluation: Overfitting & Underfitting
- Understanding Different Evaluation Models
- Simple Linear Regression

## Introduction to Regression

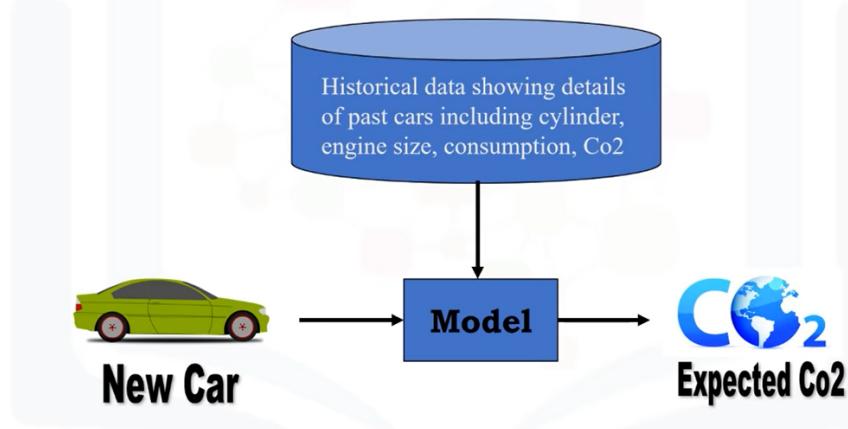
### What is regression ?

	X: Independent variable			Y: Dependent variable
	ENGINESIZE	CYLINDERS	FUELCONSUMPTION_COMB	CO2EMISSIONS
0	2.0	4	8.5	196
1	2.4	4	9.6	221
2	1.5	4	5.9	136
3	3.5	6	11.1	255
4	3.5	6	10.6	244
5	3.5	6	10.0	230
6	3.5	6	10.1	232
7	3.7	6	11.1	255
8	3.7	6	11.6	267
9	2.4	4	9.2	?

A red bracket on the right side of the table is labeled "Continuous Values". A blue bracket on the right side of the table is labeled "Continuous Values". A red arrow points from the "CO2EMISSIONS" column header to the question mark in the last row.

Regression is the process of predicting a continuous value

### What is a regression model ?



## Types of regression

- Simple Regression:
    - Simple Linear Regression
    - Simple Non-linear Regression
  - Multiple Regression:
    - Multiple Linear Regression
    - Multiple Non-linear Regression
- Predict `co2emission` vs `EngineSize` of all cars
- Predict `co2emission` vs `EngineSize` and `Cylinders` of all cars

## Application of regression

- Sales forecasting
- Satisfaction analysis
- Price estimation
- Employment income

## Regression algorithms

- Ordinal regression
- Poisson regression
- Fast forest quantile regression
- Linear, Polynomial, Lasso, Stepwise, Ridge regression
- Bayesian linear regression
- Neural network regression
- Decision forest regression
- Boosted decision tree regression
- KNN (K-nearest neighbors)

# Simple Linear Regression

## Using linear regression to predict continuous values

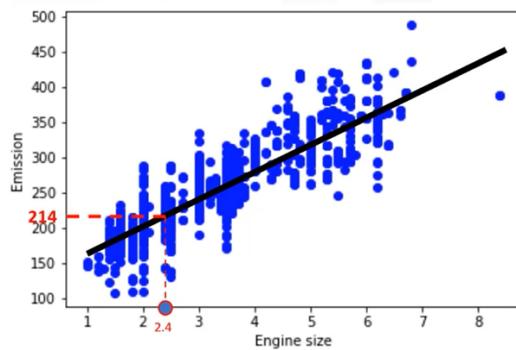
X: Independent variable				Y: Dependent variable
	ENGINESIZE	CYLINDERS	FUELCONSUMPTION_COMB	CO2EMISSIONS
0	2.0	4	8.5	196
1	2.4	4	9.6	221
2	1.5	4	5.9	136
3	3.5	6	11.1	255
4	3.5	6	10.6	244
5	3.5	6	10.0	230
6	3.5	6	10.1	232
7	3.7	6	11.1	255
8	3.7	6	11.6	267
9	2.4	4	9.2	?

## Linear regression topology

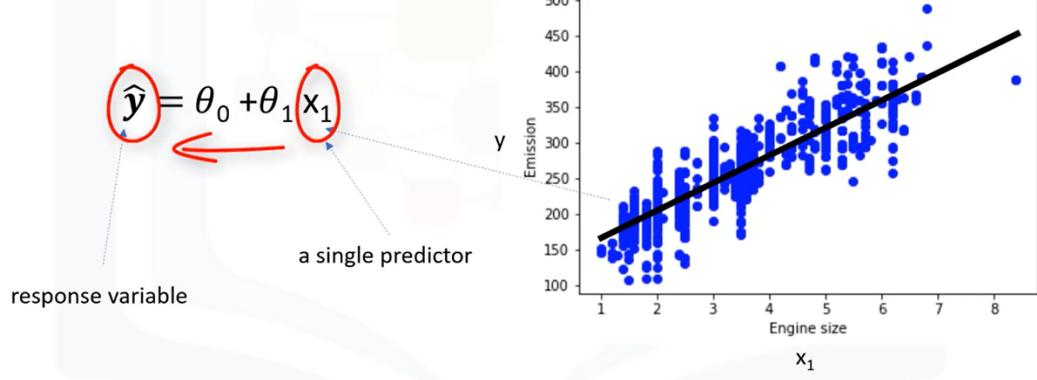
- • Simple Linear Regression:
- Predict co2emission vs EngineSize of all cars
    - Independent variable (x): EngineSize
    - Dependent variable (y): co2emission
- Multiple Linear Regression:
- Predict co2emission vs EngineSize and Cylinders of all cars
    - Independent variable (x): EngineSize, Cylinders, etc
    - Dependent variable (y): co2emission

## How does linear regression works ?

	ENGINESIZE	CYLINDERS	FUELCONSUMPTION_COMB	CO2EMISSIONS
0	2.0	4	8.5	196
1	2.4	4	9.6	221
2	1.5	4	5.9	136
3	3.5	6	11.1	255
4	3.5	6	10.6	244
5	3.5	6	10.0	230
6	3.5	6	10.1	232
7	3.7	6	11.1	255
8	3.7	6	11.6	267
9	2.4	4	9.2	?



## Linear regression model representation



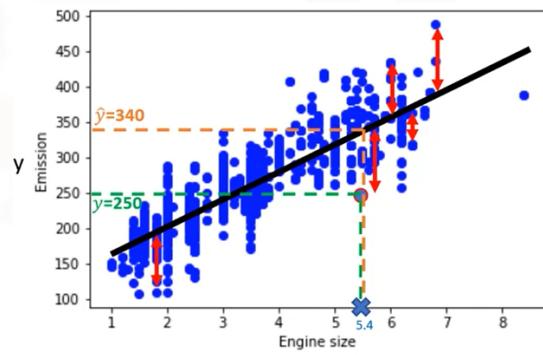
## How to find the best fit ?

$x_1 = 2.4$  independent variable  
 $y = 250$  actual Co2 emission of  $x_1$

$\hat{y} = \theta_0 + \theta_1 x_1$   
 $\hat{y} = 340$  the predicted emission of  $x_1$

$$\begin{aligned} \text{Error} &= y - \hat{y} \\ &= 250 - 340 \\ &= -90 \end{aligned}$$

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



## Estimating the parameters

$$\hat{y} = \theta_0 + \theta_1 x_1$$

$$\theta_1 = \frac{\sum_{i=1}^s (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^s (x_i - \bar{x})^2}$$

$$\bar{x} = (2.0 + 2.4 + 1.5 + \dots)/9 = 3.34$$

$$\bar{y} = (196 + 221 + 136 + \dots)/9 = 256$$

$$\theta_1 = \frac{(2.0 - 3.34)(196 - 256) + (2.4 - 3.34)(221 - 256) + \dots}{(2.0 - 3.34)^2 + (2.4 - 3.34)^2 + \dots}$$

$$\theta_0 = \bar{y} - \theta_1 \bar{x}$$

$$\hat{y} = \theta_0 + \theta_1 x_1$$

$$\theta_1 = \frac{\sum_{i=1}^s (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^s (x_i - \bar{x})^2}$$

$$\bar{x} = (2.0 + 2.4 + 1.5 + \dots)/9 = 3.34$$

$$\bar{y} = (196 + 221 + 136 + \dots)/9 = 256$$

$$\theta_1 = \frac{(2.0 - 3.34)(196 - 256) + (2.4 - 3.34)(221 - 256) + \dots}{(2.0 - 3.34)^2 + (2.4 - 3.34)^2 + \dots}$$

$$\theta_1 = 39$$

$$\theta_0 = \bar{y} - \theta_1 \bar{x}$$

$$\boxed{\theta_0} = 256 - 39 \cdot 3.34$$

## Predictions with linear regression

	ENGINESIZE	CYLINDERS	FUELCONSUMPTION_COMB	CO2EMISSIONS
0	2.0	4	8.5	196
1	2.4	4	9.6	221
2	1.5	4	5.9	136
3	3.5	6	11.1	255
4	3.5	6	10.6	244
5	3.5	6	10.0	230
6	3.5	6	10.1	232
7	3.7	6	11.1	255
8	3.7	6	11.6	267
9	2.4	4	9.2	?

$$\hat{y} = \theta_0 + \theta_1 x_1$$

$$Co2Emission = \theta_0 + \theta_1 EngineSize$$

$$Co2Emission = 125 + 39 EngineSize$$

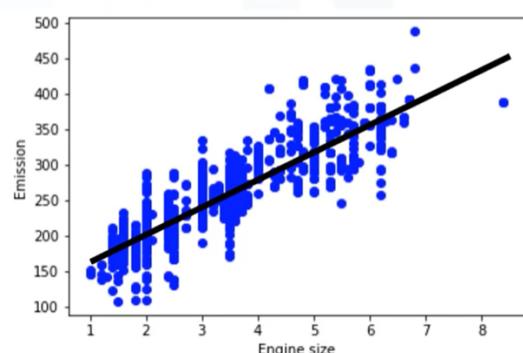
$$Co2Emission = 125 + 39 \times 2.4$$

$$Co2Emission = 218.6$$



## Pros of linear regression

- Very fast
- No parameter tuning
- Easy to understand, and highly interpretable



## Data Source Lab 2

<https://open.canada.ca/data/en/dataset/98f1a129-f628-4ce4-b24d-6f16bf24dd64>

## Multiple Linear Regression

### Examples of multiple linear regression

- Independent variables effectiveness on prediction
  - Does revision time, test anxiety, lecture attendance and gender have any effect on the exam performance of students?

### → • Predicting impacts of changes

- How much does blood pressure go up (or down) for every unit increase (or decrease) in the BMI of a patient?

## Predicting continuous values with multiple linear regression

$$Co2 Em = \theta_0 + \theta_1 Engine size + \theta_2 Cylinders + \dots$$

$$\hat{y} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

	X: Independent variable	Y: Dependent variable
0	2.0	8.5
1	2.4	9.6
2	1.5	5.9
3	3.5	11.1
4	3.5	10.6
5	3.5	10.0
6	3.5	10.1
7	3.7	11.1
8	3.7	11.6

$$Co2 Em = \theta_0 + \theta_1 Engine size + \theta_2 Cylinders + \dots$$

$$\hat{y} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

$$\hat{y} = \theta^T X$$

$$\theta^T = [\theta_0, \theta_1, \theta_2, \dots] \quad X = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \end{bmatrix}$$

	X: Independent variable	Y: Dependent variable
0	2.0	8.5
1	2.4	9.6
2	1.5	5.9
3	3.5	11.1
4	3.5	10.6
5	3.5	10.0
6	3.5	10.1
7	3.7	11.1
8	3.7	11.6

## Using MSE to expose the errors in the model

$$\hat{y} = \theta^T X$$

$\hat{y}_i = 140$  the predicted emission of  $x_i$

$y_i = 196$  actual value of  $x_i$

$y_i - \hat{y}_i = 196 - 140 = 56$  residual error

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

	ENGINESIZE	CYLINDERS	FUELCONSUMPTION_COMB	CO2EMISSIONS
0	2.0	4	8.5	196
1	2.4	4	9.6	221
2	1.5	4	5.9	136
3	3.5	6	11.1	255
4	3.5	6	10.6	244
5	3.5	6	10.0	230
6	3.5	6	10.1	232
7	3.7	6	11.1	255
8	3.7	6	11.6	267

## Estimating multiple linear regression parameters

- How to estimate  $\theta$ ?
  - Ordinary Least Squares
    - Linear algebra operations
    - Takes a long time for large datasets (10K+ rows)
  - An optimization algorithm
    - Gradient Descent
    - Proper approach if you have a very large dataset

## Making prediction with multiple linear regression

	ENGINESIZE	CYLINDERS	FUELCONSUMPTION_COMB	CO2EMISSIONS	
0	2.0	4	8.5	196	$\hat{y} = \theta^T X$
1	2.4	4	9.6	221	$\theta^T = [125, 6.2, 14, \dots]$
2	1.5	4	5.9	136	$\hat{y} = 125 + 6.2x_1 + 14x_2 +$
3	3.5	6	11.1	255	$Co2Em = 125 + 6.2\text{EngSize} + 14 \text{Cylinders} + \dots$
4	3.5	6	10.6	244	$Co2Em = 125 + 6.2 \times 2.4 + 14 \times 4 + \dots$
5	3.5	6	10.0	230	$Co2Em = 214.1$
6	3.5	6	10.1	232	
7	3.7	6	11.1	255	
8	3.7	6	11.6	267	
9	2.4	4	9.2	?	

## A&A - on multiple linear regression

- How to determine whether to use simple or multiple linear regression?
- How many independent variables should you use?
- Should the independent variable be continuous?
- What are the linear relationships between the dependent variable and the independent variables?

# Model Evaluation in Regression Model

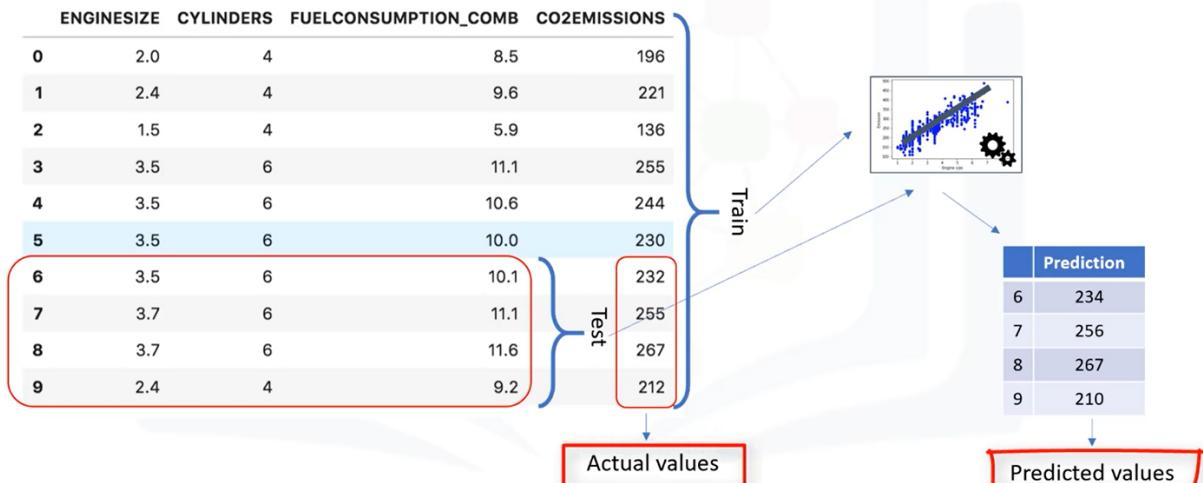
## Model evaluation approaches

- Train and Test on the Same Dataset
- Train/Test Split

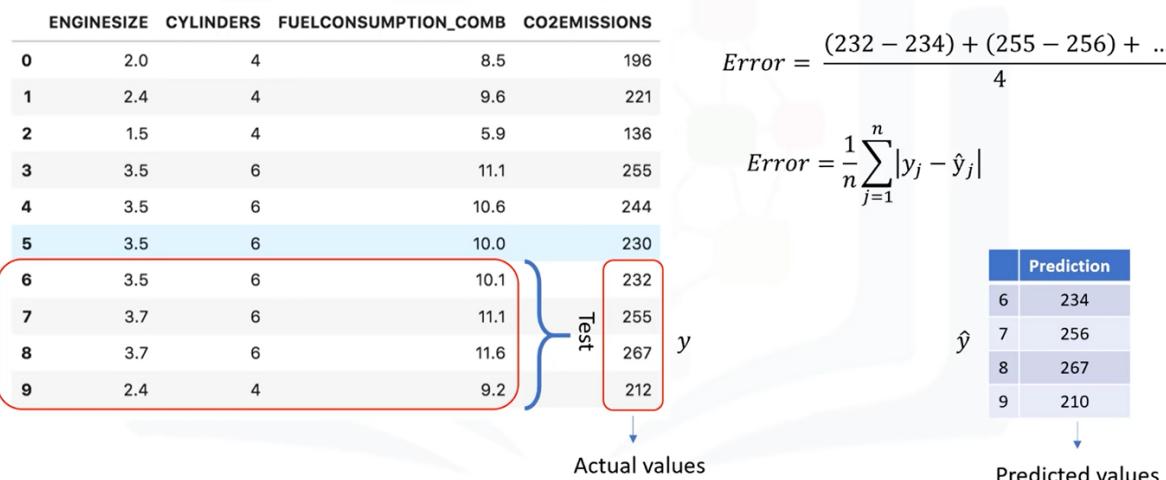
## Regression Evaluation Metrics



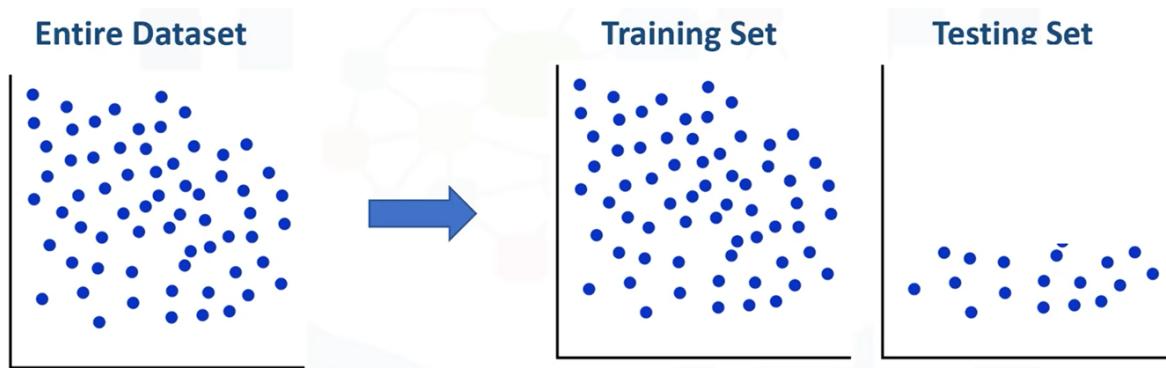
Best approach for most accurate result ?



## Calculating the accuracy of a model



## Train and test on the same dataset



## What is training & out-of-sample accuracy?

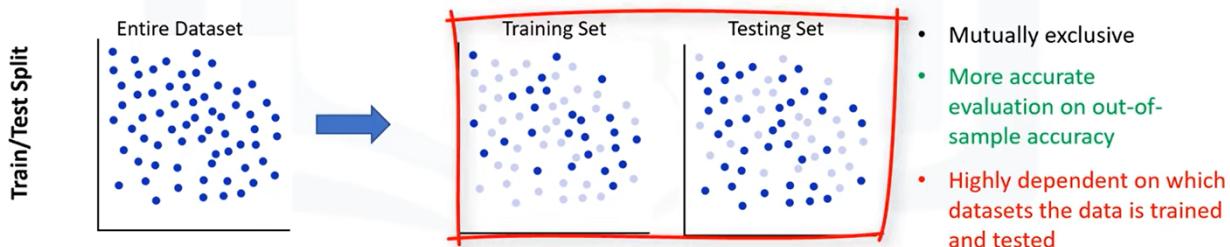
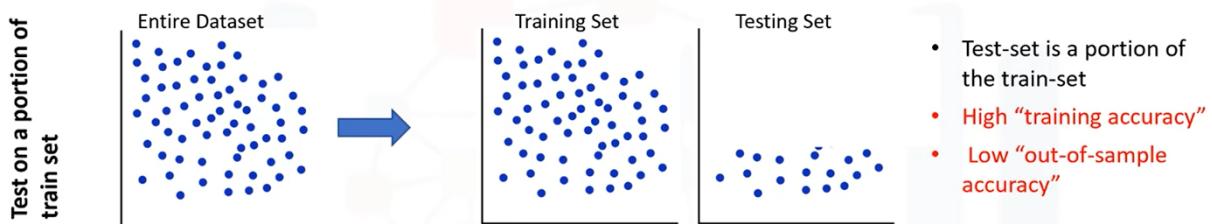
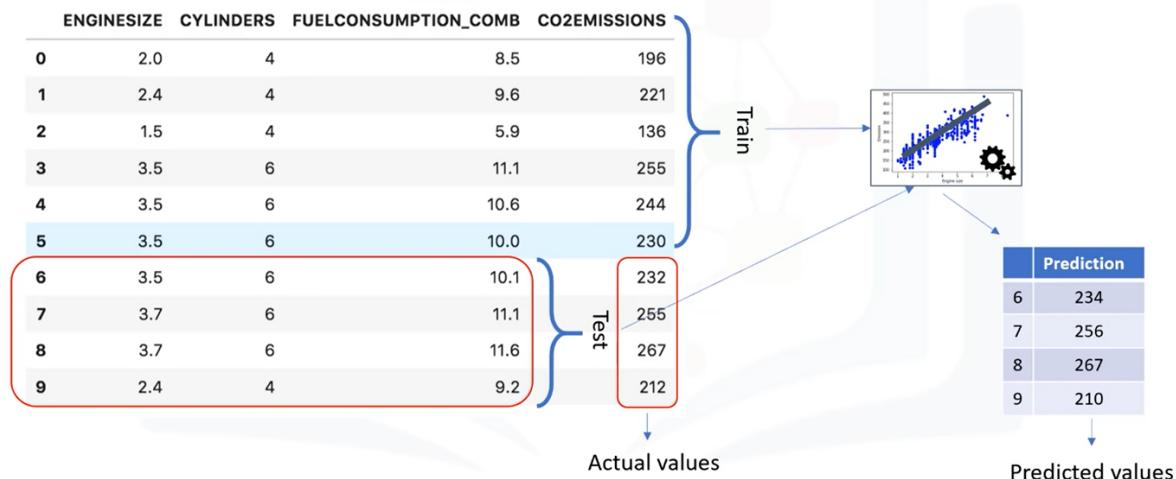
### • Training Accuracy

- High training accuracy isn't necessarily a good thing
- Result of over-fitting
  - Over-fit: the model is overly trained to the dataset, which may capture noise and produce a non-generalized model

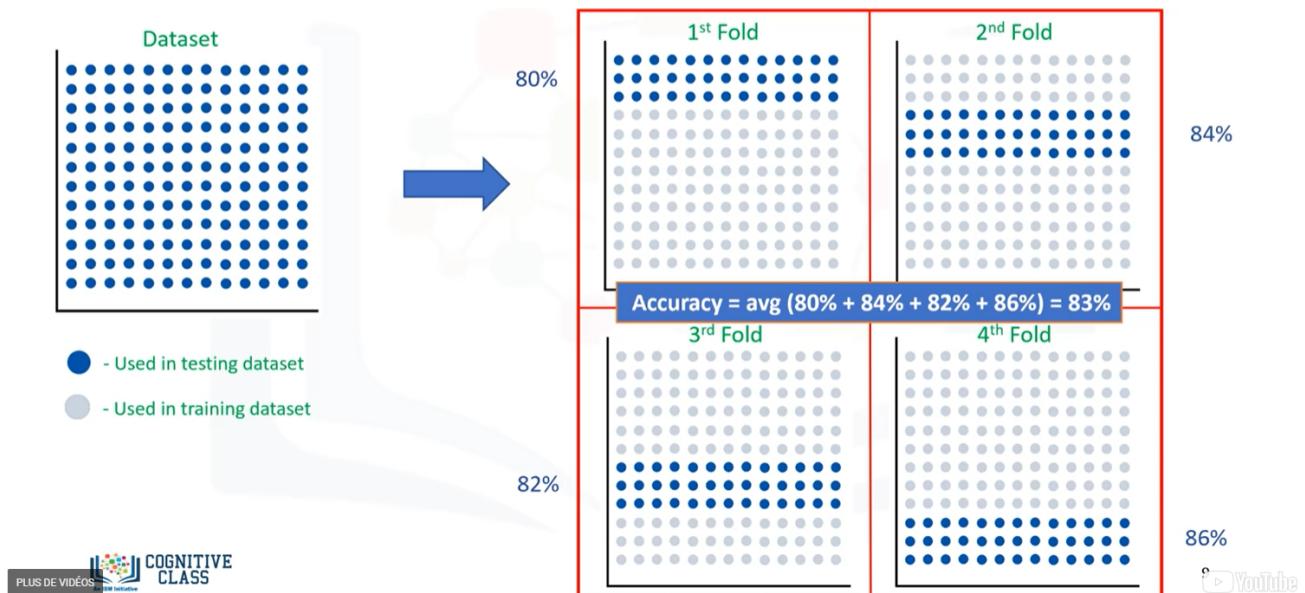
### • Out-of-Sample Accuracy

- It's important that our models have a high, out-of-sample accuracy
- How can we improve out-of-sample accuracy?

## Train/Test split evaluation approach



## How to use K-fold cross-validation ?



## Evaluation Metrics in Regression

### Regression accuracy

	ENGINESIZE	CYLINDERS	FUELCONSUMPTION_COMB	CO2EMISSIONS
0	2.0	4	8.5	196
1	2.4	4	9.6	221
2	1.5	4	5.9	136
3	3.5	6	11.1	255
4	3.5	6	10.6	244
5	3.5	6	10.0	230
6	3.5	6	10.1	232
7	3.7	6	11.1	255
8	3.7	6	11.6	267
9	2.4	4	9.2	212

$$Error = \frac{(232 - 234) + (255 - 256) + \dots}{4}$$

$$Error = \frac{1}{n} \sum_{j=1}^n |y_j - \hat{y}_j|$$

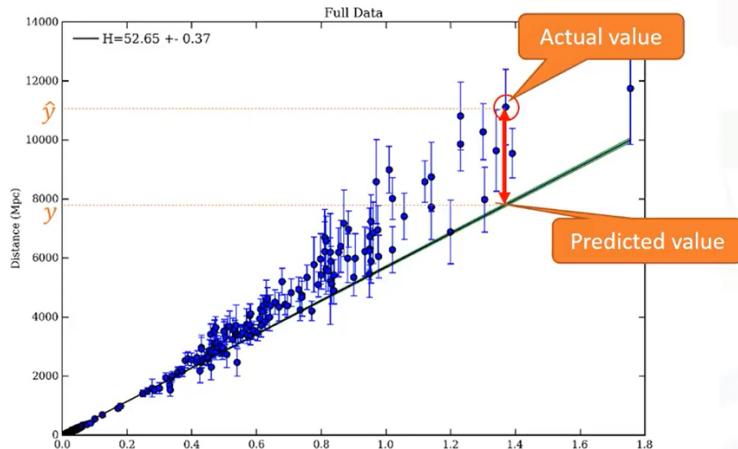
- MAE
- MSE
- RMSE
- ...

Prediction
6 234
7 256
8 267
9 210

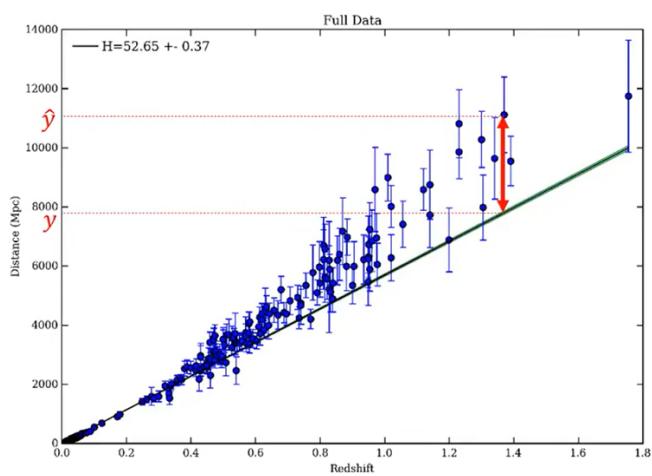
Predicted values

But, before we get into defining these, we

### What is an error of the model ?



Error: measure of how far the data is from the fitted regression line.



$$MAE = \frac{1}{n} \sum_{j=1}^n |y_j - \hat{y}_j|$$

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

$$RMSE = \sqrt{\frac{1}{n} \sum_{j=1}^n (y_j - \hat{y}_j)^2}$$

$$RAE = \frac{\sum_{j=1}^n |y_j - \hat{y}_j|}{\sum_{j=1}^n |y_j - \bar{y}|}$$

$$RSE = \frac{\sum_{j=1}^n (y_j - \hat{y}_j)^2}{\sum_{j=1}^n (y_j - \bar{y})^2}$$

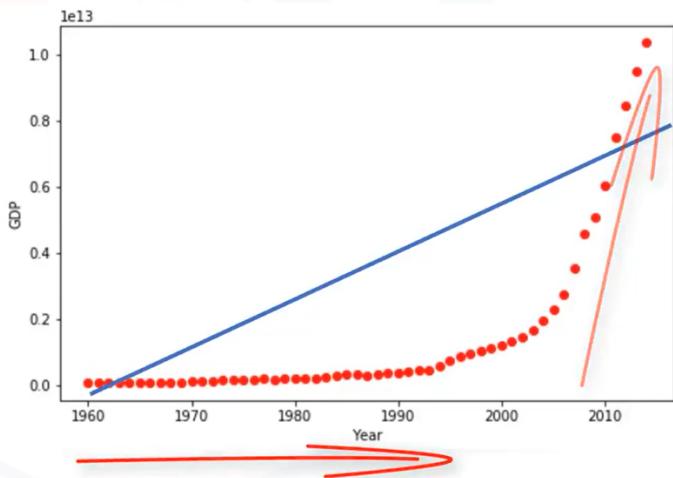
$$R^2 = 1 - RSE$$



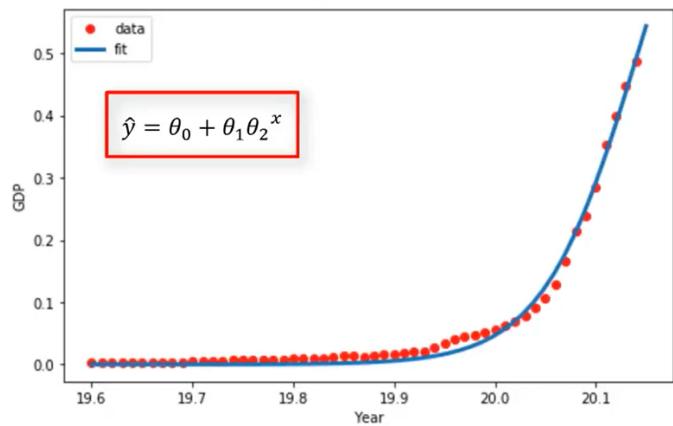
## Non-Linear Regression

Should we use linear regression ?

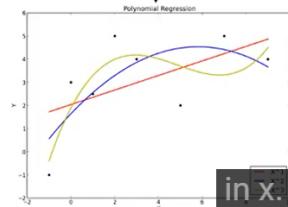
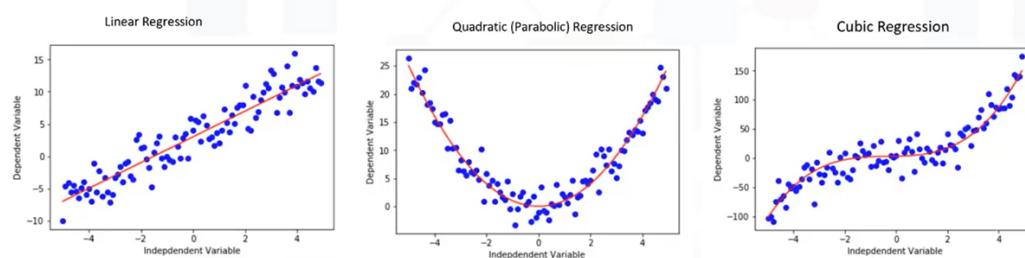
	Year	Value
0	1960	5.918412e+10
1	1961	4.955705e+10
2	1962	4.668518e+10
3	1963	5.009730e+10
4	1964	5.906225e+10
5	1965	6.970915e+10
6	1966	7.587943e+10
7	1967	7.205703e+10
8	1968	6.999350e+10
9	1969	7.871882e+10
...	....	.....



	Year	Value
0	1960	5.918412e+10
1	1961	4.955705e+10
2	1962	4.668518e+10
3	1963	5.009730e+10
4	1964	5.906225e+10
5	1965	6.970915e+10
6	1966	7.587943e+10
7	1967	7.205703e+10
8	1968	6.999350e+10
9	1969	7.871882e+10
...	....	.....



## Different types of regression



## What is polynomial regression ?

- Some curvy data can be modeled by a **polynomial regression**
- For example:

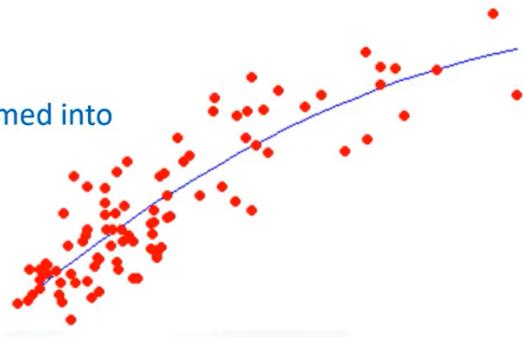
$$\hat{y} = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3$$

- A polynomial regression model can be transformed into linear regression model.

$$\begin{aligned}x_1 &= x \\x_2 &= x^2 \\x_3 &= x^3\end{aligned}$$

$$\hat{y} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$$

Multiple linear regression → Least Squares



## What is non-polynomial regression ?

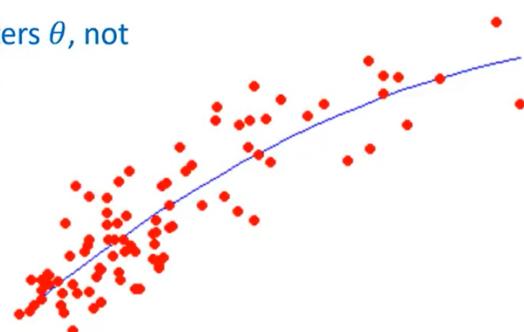
- To model non-linear relationship between the dependent variable and a set of independent variables
- $\hat{y}$  must be a non-linear function of the parameters  $\theta$ , not necessarily the features  $x$

$$\hat{y} = \theta_0 + \theta_2 x^2$$

$$\hat{y} = \theta_0 + \theta_1 \theta_2 x$$

$$\hat{y} = \log(\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3)$$

$$\hat{y} = \frac{\theta_0}{1 + \theta_1^{(x-\theta_2)}}$$



## Linear vs non-linear regression

- How can I know if a problem is linear or non-linear in an easy way?
  - Inspect visually
  - Based on accuracy
- How should I model my data, if it displays non-linear on a scatter plot?
  - Polynomial regression
  - Non-linear regression model
  - Transform your data

# Module 3

## Learning Objectives

In this lesson you will learn about:

- K-Nearest Neighbors
- Decision Trees
- Support Vector Machines
- Logistic Regression

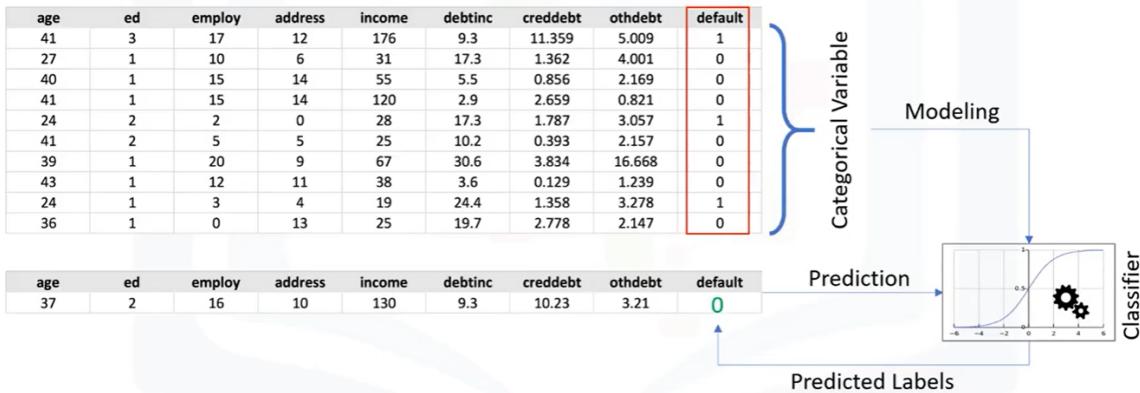
# Intro to Classification

## What is classification ?

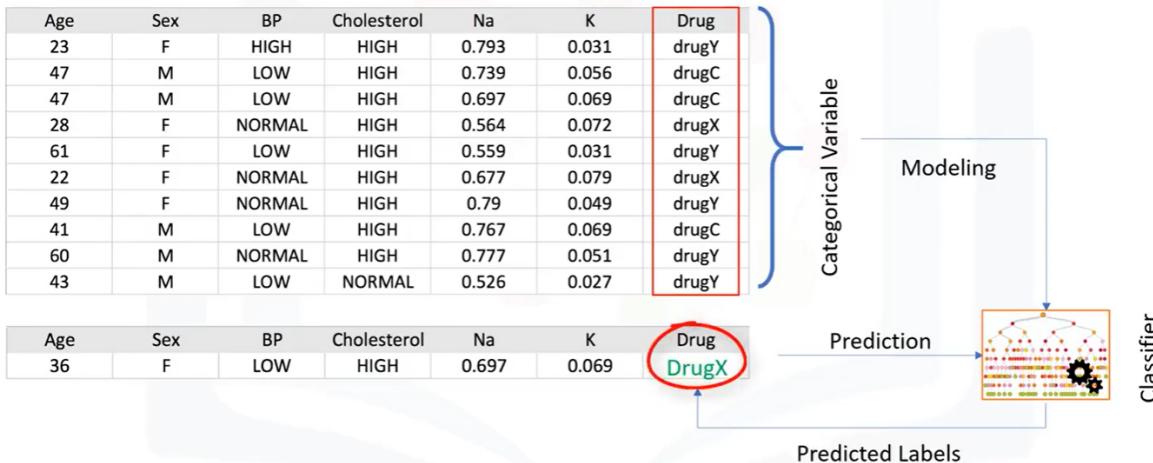
- A supervised learning approach
- Categorizing some unknown items into a discrete set of categories or “classes”
- The target attribute is a categorical variable

## How does classification work ?

Classification determines the class label for an unlabeled test case.



## Example of multi-class classification

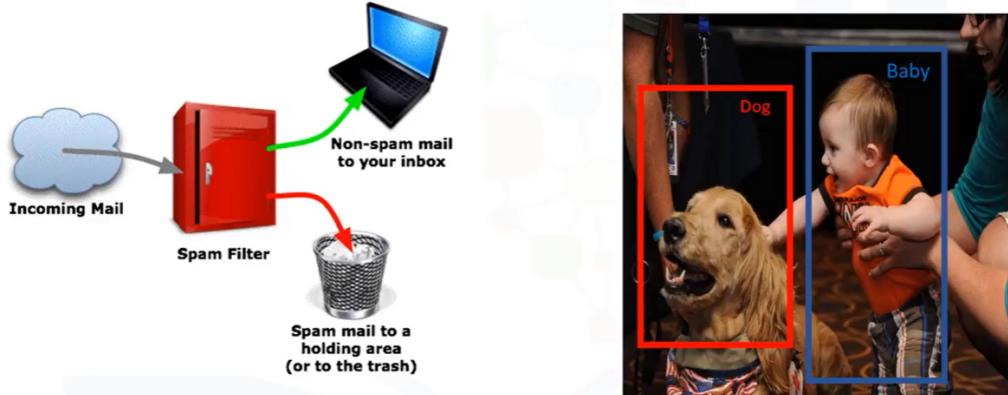


## Classification use cases

	tenure	age	address	income	ed	employ	equip	callcard	wireless	churn
0	11.0	33.0	7.0	136.0	5.0	5.0	0.0	1.0	1.0	Yes
1	33.0	33.0	12.0	33.0	2.0	0.0	0.0	0.0	0.0	Yes
2	23.0	30.0	9.0	30.0	1.0	2.0	0.0	0.0	0.0	No
3	38.0	35.0	5.0	76.0	2.0	10.0	1.0	1.0	1.0	No
4	7.0	35.0	14.0	80.0	2.0	15.0	0.0	1.0	0.0	?

- Which category a customer belongs to?
- Whether a customer switches to another provider/brand?
- Whether a customer responds to a particular advertising campaign?

## Classification applications

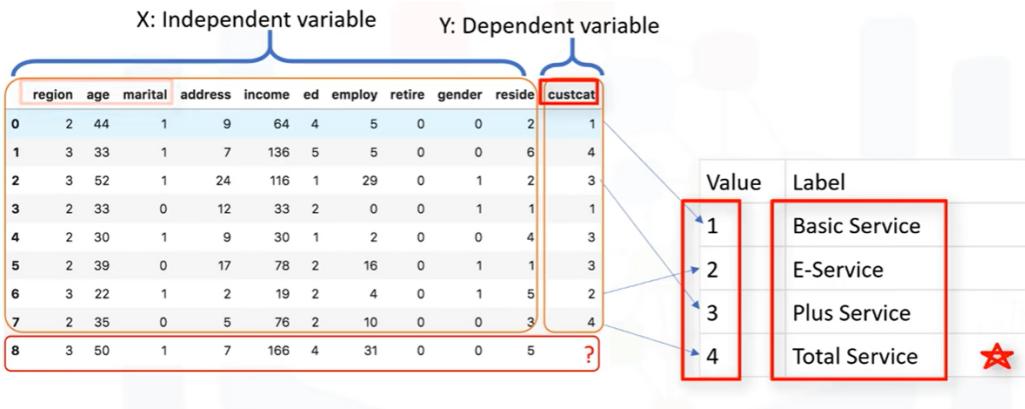


## Classification algorithms in machine learning

- Decision Trees (ID3, C4.5, C5.0)
- Naïve Bayes
- Linear Discriminant Analysis
- $k$ -Nearest Neighbor
- Logistic Regression
- Neural Networks
- Support Vector Machines (SVM)

## K-Nearest Neighbors

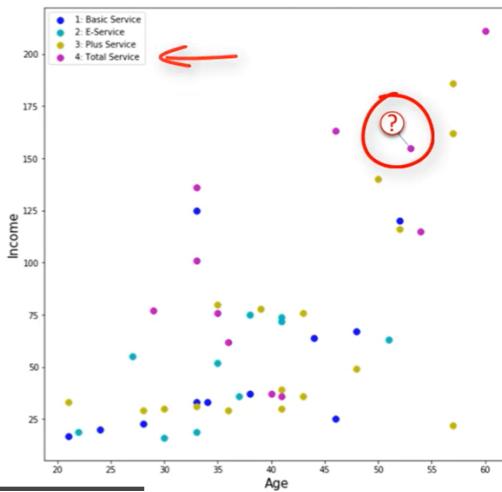
### Intro to KNN



## Determining the class using 1<sup>st</sup> KNN

	region	age	marital	address	income	ed	employ	retire	gender	reside	custcat
0	2	44	1	9	64	4	5	0	0	2	1
1	3	33	1	7	136	5	5	0	0	6	4
2	3	52	1	24	116	1	29	0	1	2	3
3	2	33	0	12	33	2	0	0	1	1	1
4	2	30	1	9	30	1	2	0	0	4	3
5	2	39	0	17	78	2	16	0	1	1	3
6	3	22	1	2	19	2	4	0	1	5	2
7	2	35	0	5	76	2	10	0	0	3	4
8	3	50	1	7	166	4	31	0	0	5	?

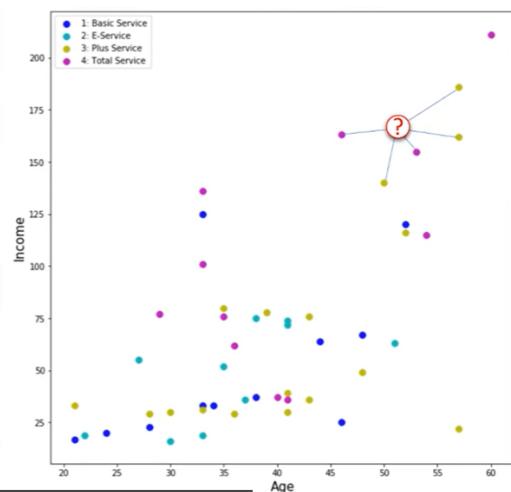
1-NN → 4: Total Service



## Determining the class using 5 KNNs

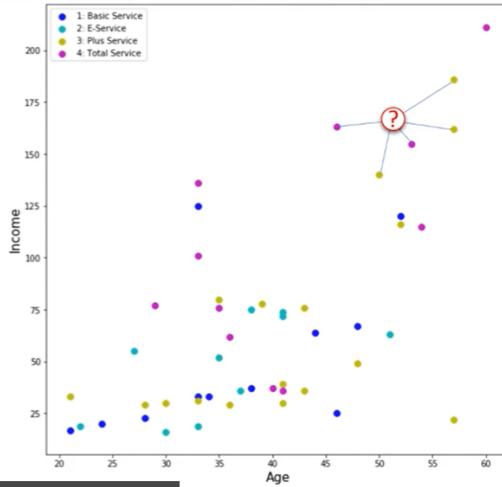
	region	age	marital	address	income	ed	employ	retire	gender	reside	custcat
0	2	44	1	9	64	4	5	0	0	2	1
1	3	33	1	7	136	5	5	0	0	6	4
2	3	52	1	24	116	1	29	0	1	2	3
3	2	33	0	12	33	2	0	0	1	1	1
4	2	30	1	9	30	1	2	0	0	4	3
5	2	39	0	17	78	2	16	0	1	1	3
6	3	22	1	2	19	2	4	0	1	5	2
7	2	35	0	5	76	2	10	0	0	3	4
8	3	50	1	7	166	4	31	0	0	5	?

5-NN → 3: Plus Service



## What is K-Nearest Neighbor ( or KNN ) ?

- A method for classifying cases based on their similarity to other cases
- Cases that are near each other are said to be “neighbors”
- Based on similar cases with same class labels are near each other



#### #### The K-Nearest Neighbors algorithm

1. Pick a value for K.
2. Calculate the distance of unknown case from all cases.
3. Select the K-observations in the training data that are “nearest” to the unknown data point.
4. Predict the response of the unknown data point using the most popular response value from the K-nearest neighbors.

### Calculating the similarity / distance in

#### 1-dimensional space



Customer 1
Age
54



Customer 2
Age
50

$$\text{Dis } (x_1, x_2) = \sqrt{\sum_{i=0}^n (x_{1i} - x_{2i})^2}$$

$$\text{Dis } (x_1, x_2) = \sqrt{(34 - 30)^2} = 4$$

#### multi-dimensional space



Customer 1		
Age	Income	Education
54	190	3

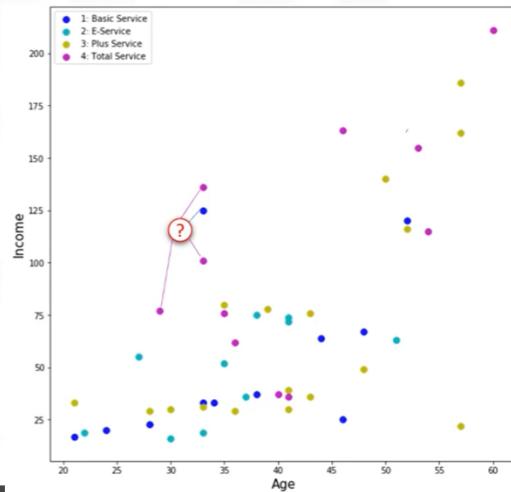
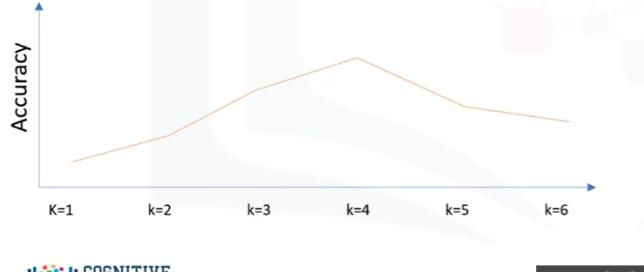
Customer 2		
Age	Income	Education
50	200	8

$$\text{Dis } (x_1, x_2) = \sqrt{\sum_{i=0}^n (x_{1i} - x_{2i})^2}$$

$$= \sqrt{(54 - 50)^2 + (190 - 200)^2 + (3 - 8)^2} = 11.87$$

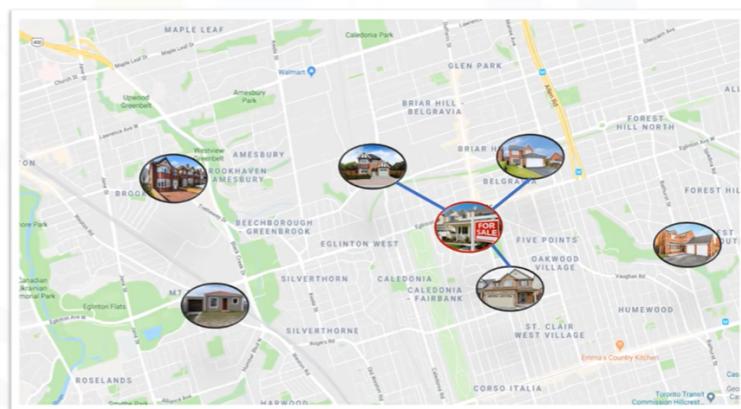
## What is the best value of K for KNN ?

- K =1      class 1
- K =20      ?



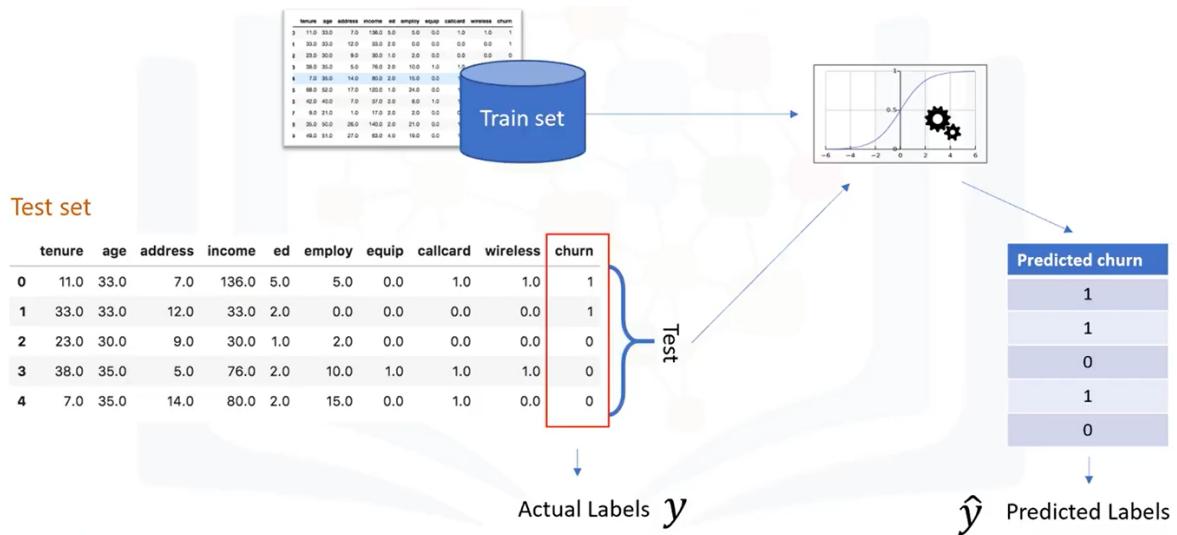
## Computing continuous target using KNN

- KNN can also be used for regression



## Evaluation Metrics in Classification

### Classification accuracy

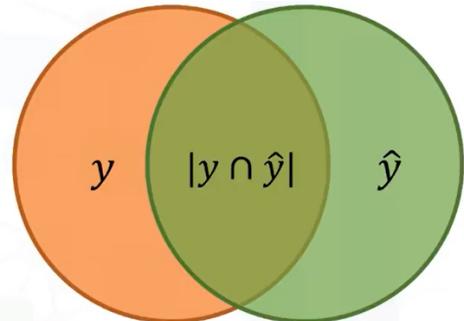


## Jaccard index

$y$ : Actual labels

$\hat{y}$ : Predicted labels

$$J(y, \hat{y}) = \frac{|y \cap \hat{y}|}{|y \cup \hat{y}|} = \frac{|y \cap \hat{y}|}{|y| + |\hat{y}| - |y \cap \hat{y}|}$$



$y: [0, 0, 0, 0, 1, 1, 1, 1]$

$\hat{y}: [1, 1, 0, 0, 0, 1, 1, ,1, 1]$

$$J(y, \hat{y}) = \frac{8}{10+10-8} = 0.66$$

of labels, then the subset accuracy is 1.0;

otherwise it is 0.0.

$$J(y, \hat{y}) = 1.0$$

Higher Accuracy →

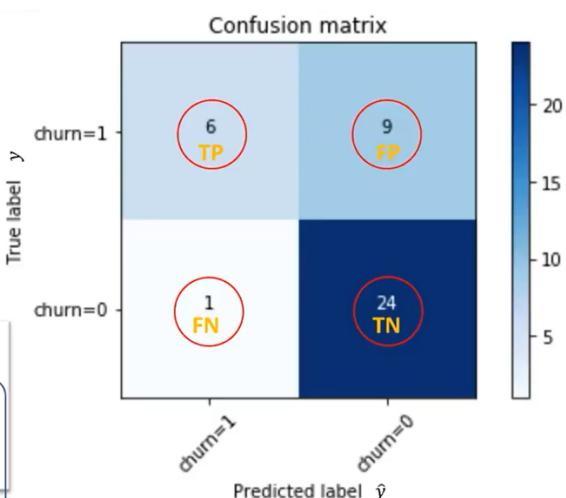


## F1-score

- Precision =  $TP / (TP + FP)$
- Recall =  $TP / (TP + FN)$
- F1-score =  $2x (prc \times rec) / (prc+rec)$

F1-score: 0.00 ... 0.20 .... 0.55 .... 0.83 ... 1.00  
Higher Accuracy →

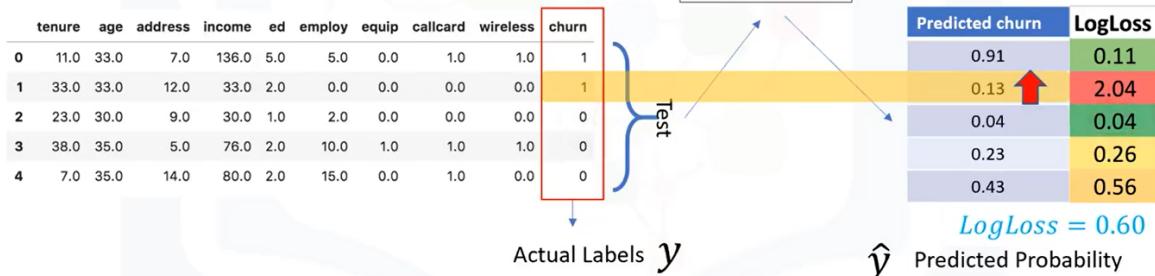
	precision	recall	f1-score
Churn = 0	0.73	0.96	0.83
Churn = 1	0.86	0.40	0.55
			Avg Accuracy = 0.72



## Log Loss

Performance of a classifier where the predicted output is a probability value between 0 and 1.

Test set



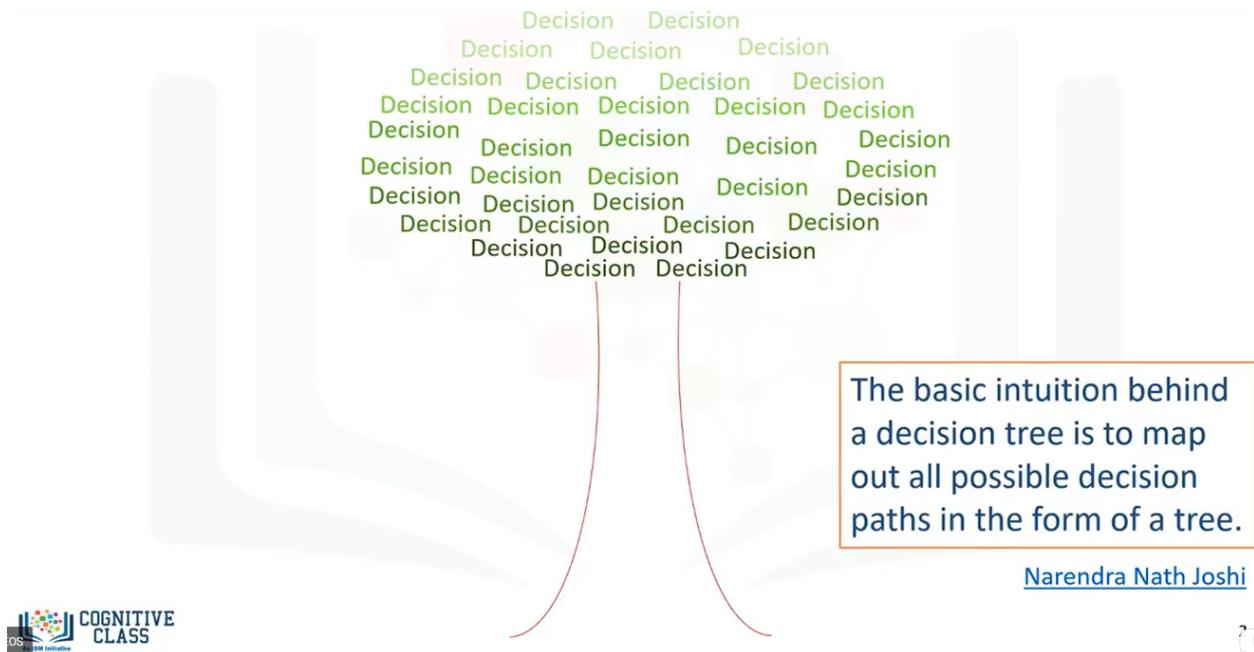
$$\text{LogLoss} = -\frac{1}{n} \sum (y \times \log(\hat{y}) + (1 - y) \times \log(1 - \hat{y}))$$

LogLoss: 0.00 ... 0.35 .... 0.60 ... 1.00

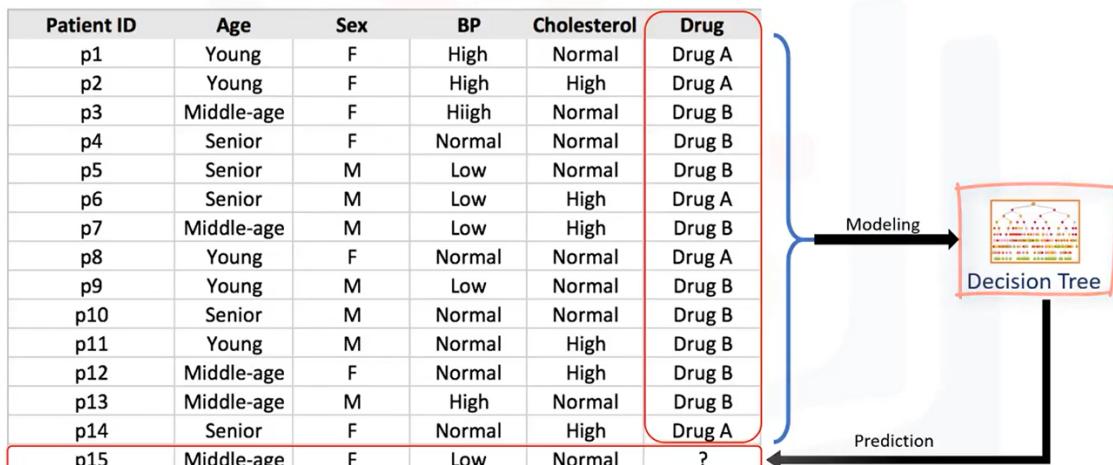
← Higher Accuracy

## Intro to Decision Trees

What is decision tree ?

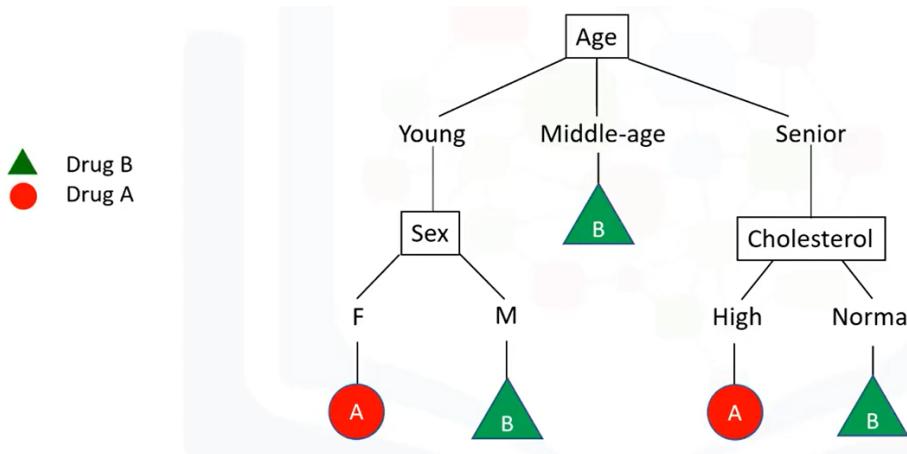


How to build a decision tree ?



How to build a decision tree ?

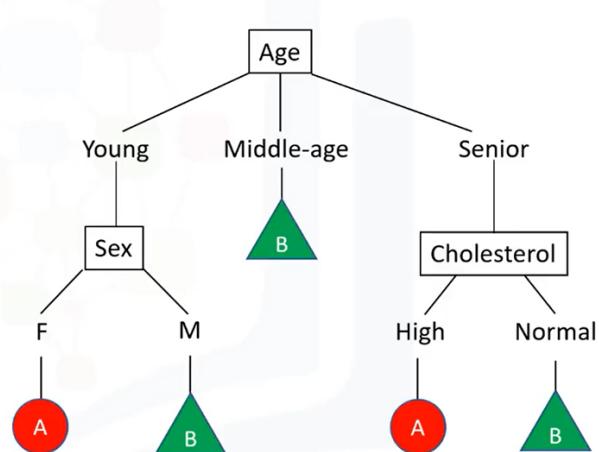
## Building a decision tree with the training set



- Each **internal node** corresponds to a test
- Each **branch corresponds** to a result of the test
- Each **leaf node** assigns a classification

## Decision tree learning algorithm

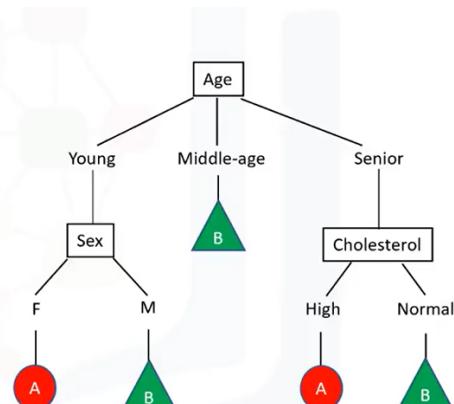
1. Choose an attribute from your dataset.
2. Calculate the significance of attribute in splitting of data.
3. Split data based on the value of the best attribute.
4. Go to step 1.



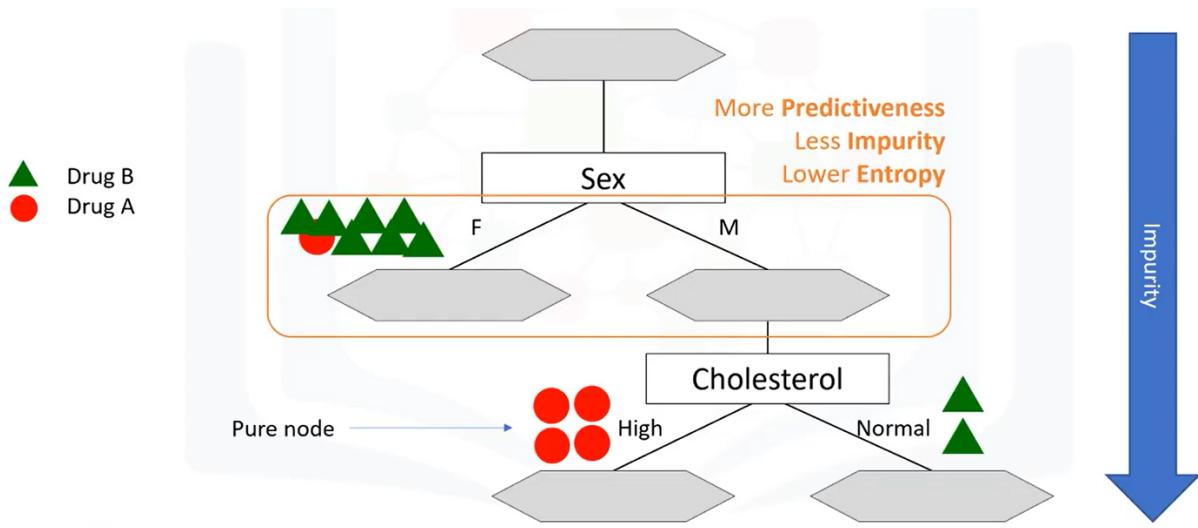
## Building Decision Trees

### How to build decision tree

Patient ID	Age	Sex	BP	Cholesterol	Drug
p1	Young	F	High	Normal	Drug A
p2	Young	F	High	High	Drug A
p3	Middle-age	F	Hiigh	Normal	Drug B
p4	Senior	F	Normal	Normal	Drug B
p5	Senior	M	Low	Normal	Drug B
p6	Senior	M	Low	High	Drug A
p7	Middle-age	M	Low	High	Drug B
p8	Young	F	Normal	Normal	Drug A
p9	Young	M	Low	Normal	Drug B
p10	Senior	M	Normal	Normal	Drug B
p11	Young	M	Normal	High	Drug B
p12	Middle-age	F	Normal	High	Drug B
p13	Middle-age	M	High	Normal	Drug B
p14	Senior	F	Normal	High	Drug A
p15	Middle-age	F	Low	Normal	?



## Which attribute is the best attribute ?

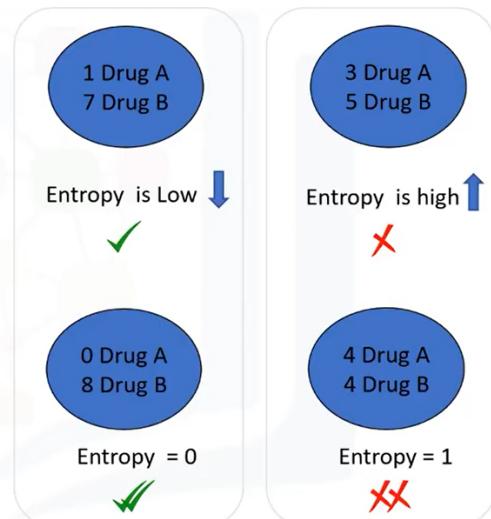


## Entropy

- Measure of randomness or uncertainty

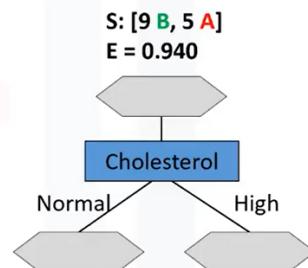
$$\text{Entropy} = - p(A)\log(p(A)) - p(B)\log(p(B))$$

The lower the Entropy, the less uniform the distribution, the purer the node.



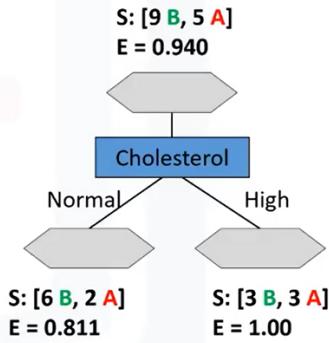
## With attribute is the best one to use ?

Patient ID	Age	Sex	BP	Cholesterol	Drug
p1	Young	F	High	Normal	Drug A
p2	Young	F	High	High	Drug A
p3	Middle-age	F	Hiigh	Normal	Drug B
p4	Senior	F	Normal	Normal	Drug B
p5	Senior	M	Low	Normal	Drug B
p6	Senior	M	Low	High	Drug A
p7	Middle-age	M	Low	High	Drug B
p8	Young	F	Normal	Normal	Drug A
p9	Young	M	Low	Normal	Drug B
p10	Senior	M	Normal	Normal	Drug B
p11	Young	M	Normal	High	Drug B
p12	Middle-age	F	Normal	High	Drug B
p13	Middle-age	M	High	Normal	Drug B
p14	Senior	F	Normal	High	Drug A



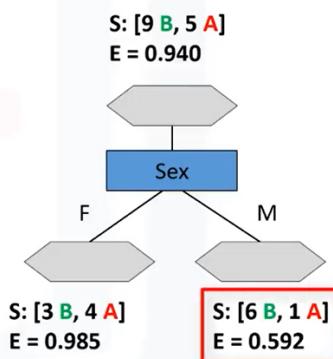
## Is 'Cholesterol' the best attribute ?

Patient ID	Age	Sex	BP	Cholesterol	Drug
p1	Young	F	High	Normal	Drug A
p2	Young	F	High	High	Drug A
p3	Middle-age	F	Hiigh	Normal	Drug B
p4	Senior	F	Normal	Normal	Drug B
p5	Senior	M	Low	Normal	Drug B
p6	Senior	M	Low	High	Drug A
p7	Middle-age	M	Low	High	Drug B
p8	Young	F	Normal	Normal	Drug A
p9	Young	M	Low	Normal	Drug B
p10	Senior	M	Normal	Normal	Drug B
p11	Young	M	Normal	High	Drug B
p12	Middle-age	F	Normal	High	Drug B
p13	Middle-age	M	High	Normal	Drug B
p14	Senior	F	Normal	High	Drug A

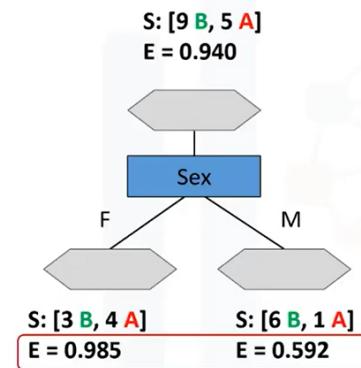


## What about 'Sexe' ?

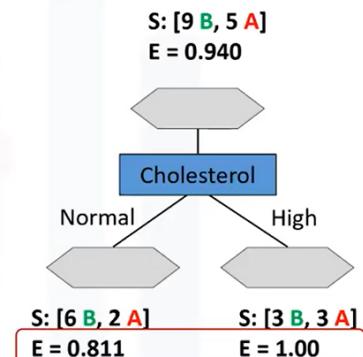
Patient ID	Age	Sex	BP	Cholesterol	Drug
p1	Young	F	High	Normal	Drug A
p2	Young	F	High	High	Drug A
p3	Middle-age	F	Hiigh	Normal	Drug B
p4	Senior	F	Normal	Normal	Drug B
p5	Senior	M	Low	Normal	Drug B
p6	Senior	M	Low	High	Drug A
p7	Middle-age	M	Low	High	Drug B
p8	Young	F	Normal	Normal	Drug A
p9	Young	M	Low	Normal	Drug B
p10	Senior	M	Normal	Normal	Drug B
p11	Young	M	Normal	High	Drug B
p12	Middle-age	F	Normal	High	Drug B
p13	Middle-age	M	High	Normal	Drug B
p14	Senior	F	Normal	High	Drug A



## Which attribute is the best ?



Patient ID	Age	Sex	BP	Cholesterol	Drug
p1	Young	F	High	Normal	Drug A
p2	Young	F	High	High	Drug A
p3	Middle-age	F	Hiigh	Normal	Drug B
p4	Senior	F	Normal	Normal	Drug B
p5	Senior	M	Low	Normal	Drug B
p6	Senior	M	Low	High	Drug A
p7	Middle-age	M	Low	High	Drug B
p8	Young	F	Normal	Normal	Drug A
p9	Young	M	Low	Normal	Drug B
p10	Senior	M	Normal	Normal	Drug B
p11	Young	M	Normal	High	Drug B
p12	Middle-age	F	Normal	High	Drug B
p13	Middle-age	M	High	Normal	Drug B
p14	Senior	F	Normal	High	Drug A



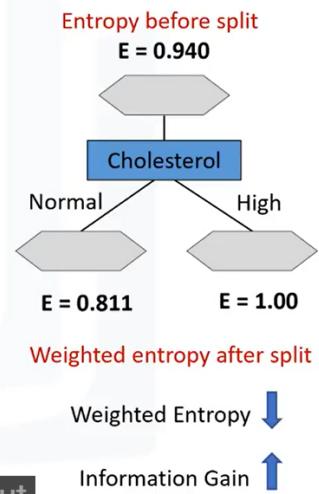
?

The tree with the higher Information Gain after splitting.

## What is information gain ?

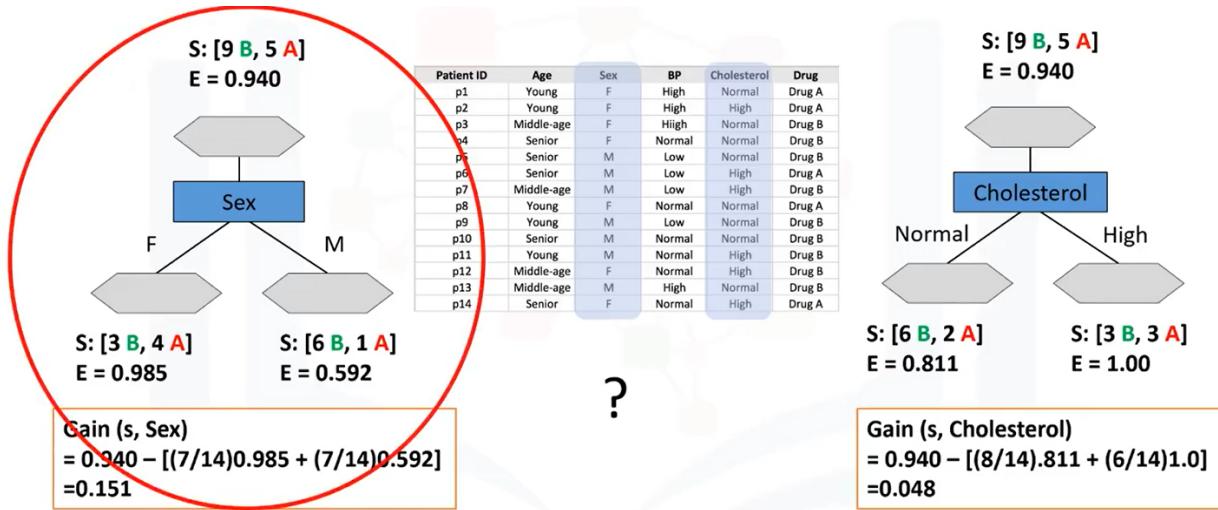
**Information gain** is the information that can increase the level of certainty after splitting.

$$\text{Information Gain} = (\text{Entropy before split}) - (\text{weighted entropy after split})$$

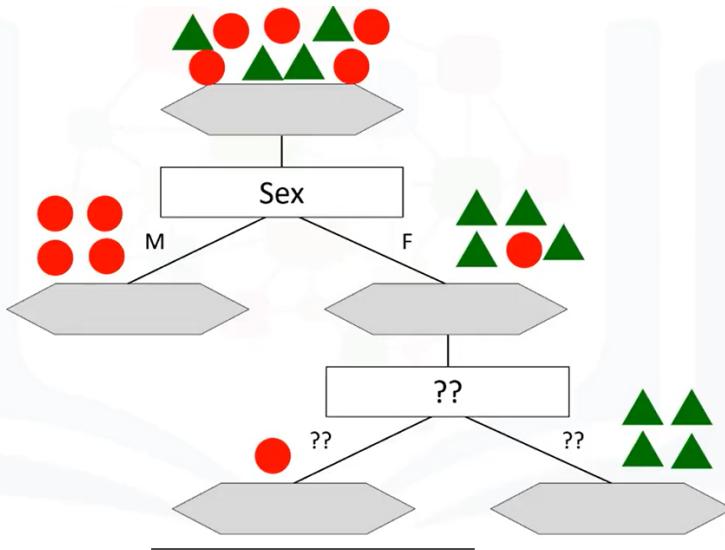


So, constructing a decision tree is all about:

## Calculating information



## Correct way to build a decision tree



## Intro to Logistic Regression

## What is logistic regression ?

Logistic regression is a classification algorithm for categorical variables.

	Independent variables										Dependent variable
	tenure	age	address	income	ed	employ	equip	callcard	wireless	churn	
0	11.0	33.0	7.0	136.0	5.0	5.0	0.0	1.0	1.0	Yes	
1	33.0	33.0	12.0	33.0	2.0	0.0	0.0	0.0	0.0	Yes	
2	23.0	30.0	9.0	30.0	1.0	2.0	0.0	0.0	0.0	No	
3	38.0	35.0	5.0	76.0	2.0	10.0	1.0	1.0	1.0	No	
4	7.0	35.0	14.0	80.0	2.0	15.0	0.0	1.0	0.0	?	

classification, but for simplicity, in this example, we will focus on continuous/categorical variables

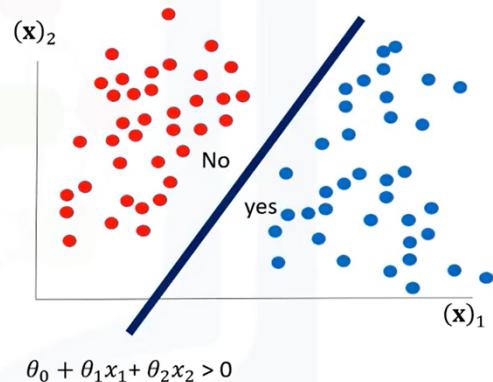
COGNITIVE

## Logistic regression applications

- Predicting the probability of a person having a heart attack
- Predicting the mortality in injured patients
- Predicting a customer's propensity to purchase a product or halt a subscription
- Predicting the probability of failure of a given process or product
- Predicting the likelihood of a homeowner defaulting on a mortgage

## When is logistic regression suitable?

- If your data is binary
  - 0/1, YES/NO, True/False
- If you need probabilistic results
- When you need a linear decision boundary
- If you need to understand the impact of a feature



## Building a model for customer churn

**X**

	tenure	age	address	income	ed	employ	equip	callcard	wireless	churn
0	11.0	33.0	7.0	136.0	5.0	5.0	0.0	1.0	1.0	1.0
1	33.0	33.0	12.0	33.0	2.0	0.0	0.0	0.0	0.0	1.0
2	23.0	30.0	9.0	30.0	1.0	2.0	0.0	0.0	0.0	0.0
3	38.0	35.0	5.0	76.0	2.0	10.0	1.0	1.0	1.0	0.0

$$X \in \mathbb{R}^{m \times n}$$

$$y \in \{0,1\}$$

$$\hat{y} = P(y=1|x)$$

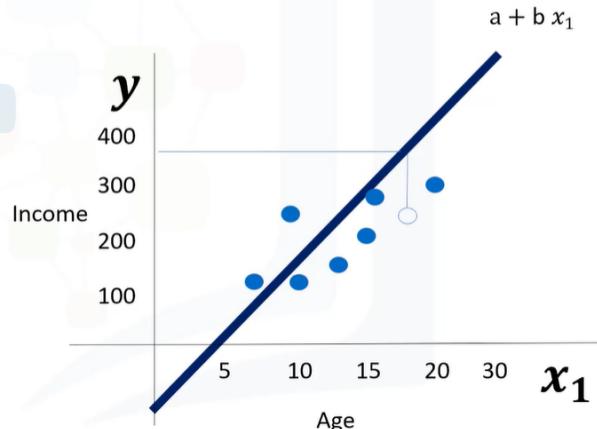
$$P(y=0|x) = 1 - P(y=1|x)$$

DATA SCIENCE FOR BUSINESS

## Logistic Regression vs Linear Regression

### Predicting customer income

	tenure	age	address	income	ed	employ	equip	callcard	wireless	churn
0	11.0	33.0	7.0	136.0	5.0	5.0	0.0	1.0	1.0	1
1	33.0	33.0	12.0	33.0	2.0	0.0	0.0	0.0	0.0	1
2	23.0	30.0	9.0	30.0	1.0	2.0	0.0	0.0	0.0	0
3	38.0	35.0	5.0	76.0	2.0	10.0	1.0	1.0	1.0	0
4	7.0	35.0	14.0	80.0	2.0	15.0	0.0	1.0	0.0	0

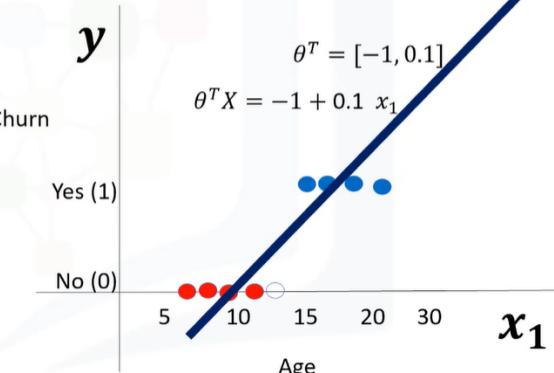


### Predicting churn using linear regression

$$\theta^T X = \theta_0 + \theta_1 x_1$$

$$\theta^T X = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots$$

$$\theta^T = [\theta_0, \theta_1, \theta_2, \dots] \quad X = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \dots \end{bmatrix}$$



## Linear regression classification problems ?

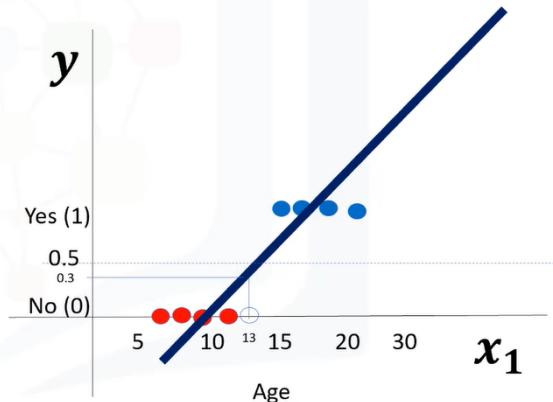
$$\theta^T X = \theta_0 + \theta_1 x_1$$

$$p_1 = [13] \rightarrow \theta^T X = -1 + 0.1 \cdot x_1 \\ = -1 + 0.1 \times 13 \\ = 0.3$$

$$\hat{y} = \begin{cases} 0 & \text{if } \theta^T X < 0.5 \\ 1 & \text{if } \theta^T X \geq 0.5 \end{cases}$$

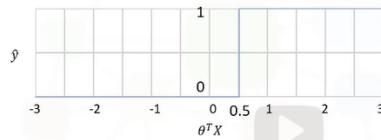
$$\theta^T X = 0.3 \rightarrow \text{Class 0}$$

$$\theta^T X = -1 + 0.1 \cdot x_1$$



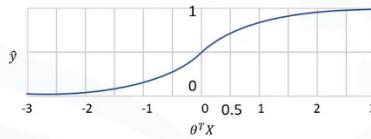
## The problem with using linear regression

$$\theta^T X = \theta_0 + \theta_1 x_1 + \dots$$



$$\hat{y} = \begin{cases} 0 & \text{if } \theta^T X < 0.5 \\ 1 & \text{if } \theta^T X \geq 0.5 \end{cases}$$

$$\sigma(\theta^T X) = \sigma(\theta_0 + \theta_1 x_1 + \dots)$$



$$\hat{y} = \sigma(\theta^T X)$$

$$P(y=1|x)$$

## Sigmoid function in logistic regression

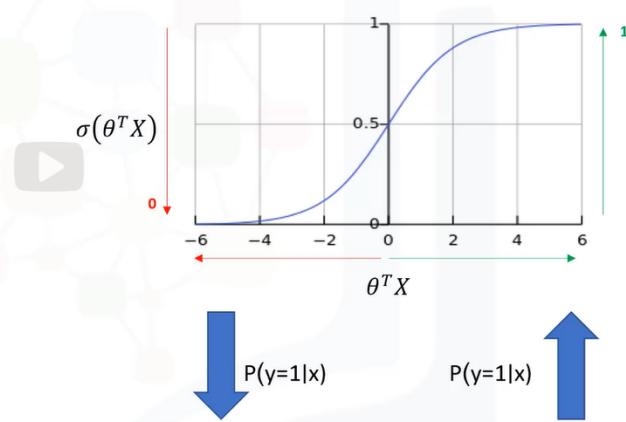
### • Logistic Function

$$\sigma(\theta^T X) = \frac{1}{1 + e^{-\theta^T X}}$$

$$\sigma(\theta^T X) = 1$$

$$\sigma(\theta^T X) = 0$$

[0, 1]



## Clarification customer churn model

## What is the output of our model?

- $P(Y=1|X)$
- $P(y=0|X) = 1 - P(y=1|X)$

- $P(\text{Churn}=1|\text{income},\text{age}) = 0.8$
- $P(\text{Churn}=0|\text{income},\text{age}) = 1 - 0.8 = 0.2$

$$\sigma(\theta^T X) \longrightarrow P(y=1|x)$$

$$1 - \sigma(\theta^T X) \longrightarrow P(y=0|x)$$

## The training process

$$\sigma(\theta^T X) \longrightarrow P(y=1|x)$$

1. Initialize  $\theta$ .
2. Calculate  $\hat{y} = \sigma(\theta^T X)$  for a customer.
3. Compare the output of  $\hat{y}$  with actual output of customer,  $y$ , and record it as error.
4. Calculate the error for all customers.
5. Change the  $\theta$  to reduce the cost.
6. Go back to step 2.

$$\theta = [-1, 2]$$

$$\hat{y} = \sigma([-1, 2] \times [2, 5]) = 0.7$$

$$\text{Error} = 1 - 0.7 = 0.3$$

$$Cost = J(\theta)$$

$$\theta_{new}$$

## Logistic Regression - Training

### General cost function

$$\sigma(\theta^T X) \longrightarrow P(y=1|x)$$

- Change the weight -> Reduce the cost

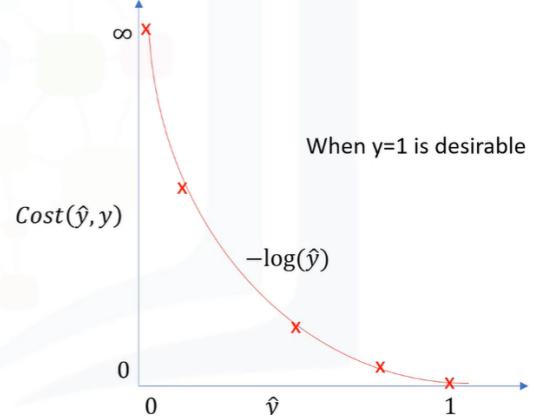
- Cost function

$$Cost(\hat{y}, y) = \frac{1}{2} (\sigma(\theta^T X) - y)^2$$

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m Cost(\hat{y}_i, y_i)$$

### Plotting cost function of the model

- Model  $\hat{y}$
- Actual Value  $y=1$  or  $0$
- If  $Y=1$ , and  $\hat{y}=1 \rightarrow \text{cost} = 0$
- If  $Y=1$ , and  $\hat{y}=0 \rightarrow \text{cost} = \text{large}$



## Logistic regression cost function

- So, we will replace cost function with:

$$Cost(\hat{y}, y) = \frac{1}{2} (\sigma(\theta^T X) - y)^2$$

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m Cost(\hat{y}, y)$$

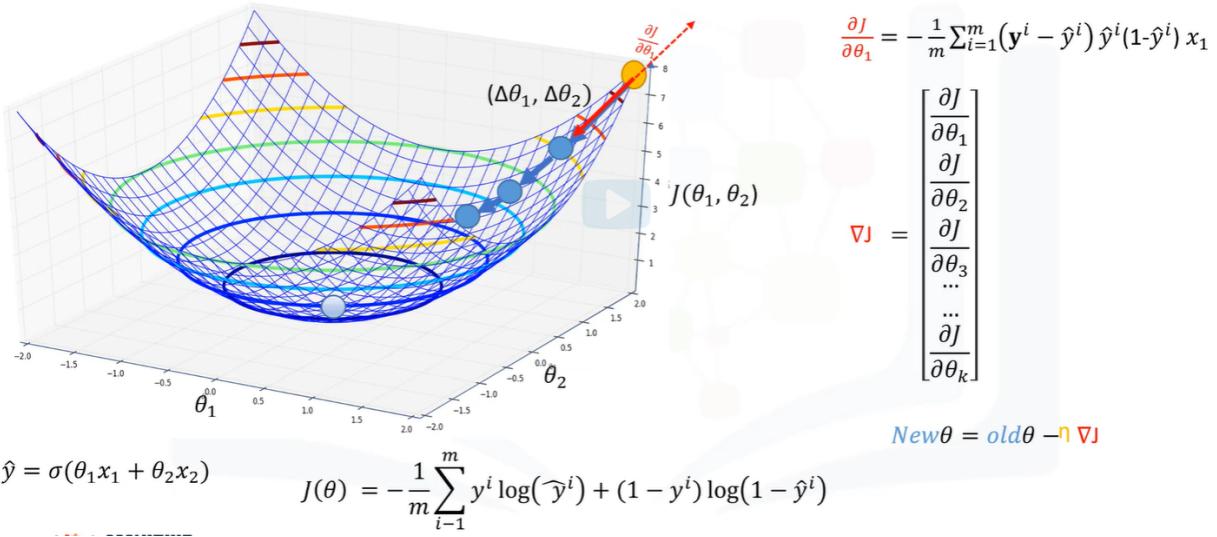
$$Cost(\hat{y}, y) = \begin{cases} -\log(\hat{y}) & \text{if } y = 1 \\ -\log(1 - \hat{y}) & \text{if } y = 0 \end{cases}$$

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m y^i \log(\hat{y}^i) + (1 - y^i) \log(1 - \hat{y}^i)$$

## Minimizing the cost function of the model

- How to find the best parameters for our model?
  - Minimize the cost function
- How to minimize the cost function?
  - Using Gradient Descent
- What is gradient descent?
  - A technique to use the derivative of a cost function to change the parameter values, in order to minimize the cost

## Using gradient descent to minimizing the cost



## Training algorithm recap

1. initialize the parameters randomly.
2. Feed the cost function with training set, and calculate the error.
3. Calculate the gradient of cost function.
4. Update weights with new values.
5. Go to step 2 until cost is small enough.
6. Predict the new customer X.

$$\theta^T = [\theta_0, \theta_1, \theta_2, \dots]$$

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m y^i \log(\hat{y}^i) + (1 - y^i) \log(1 - \hat{y}^i)$$

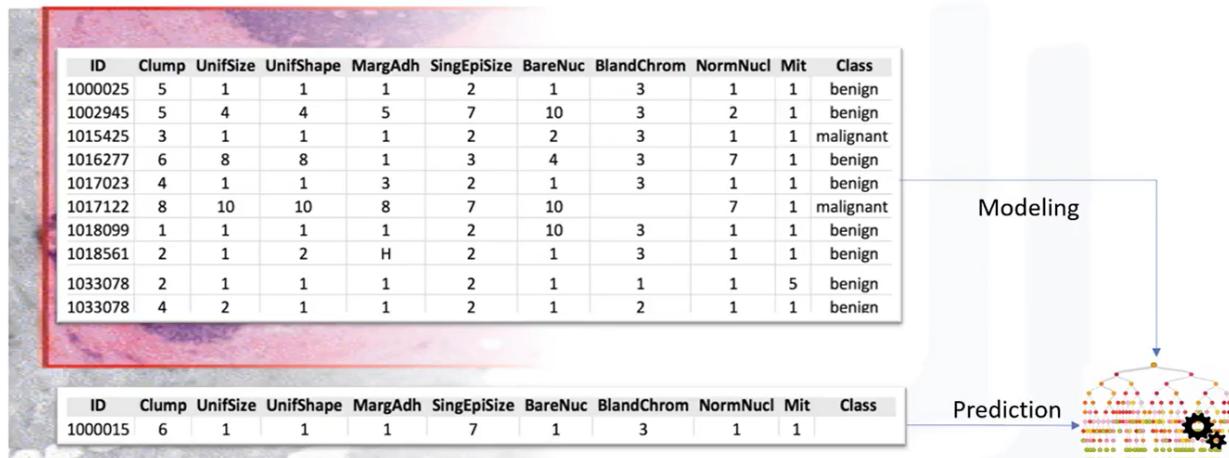
$$\nabla J = \left[ \frac{\partial J}{\partial \theta_1}, \frac{\partial J}{\partial \theta_2}, \frac{\partial J}{\partial \theta_3}, \dots, \frac{\partial J}{\partial \theta_k} \right]$$

$$\theta_{new} = \theta_{prev} - \eta \nabla J$$

$$P(y=1|x) = \sigma(\theta^T x)$$

## Support Vector Machines

### Classification with SVM



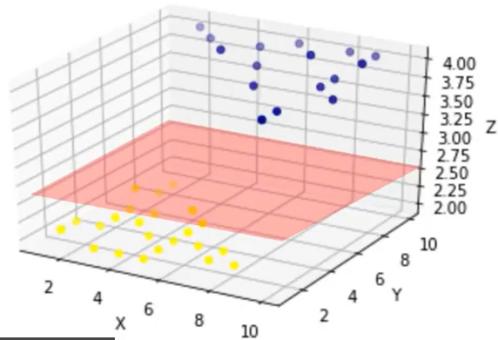
## What is SVM ?

SVM is a supervised algorithm that classifies cases by finding a separator.

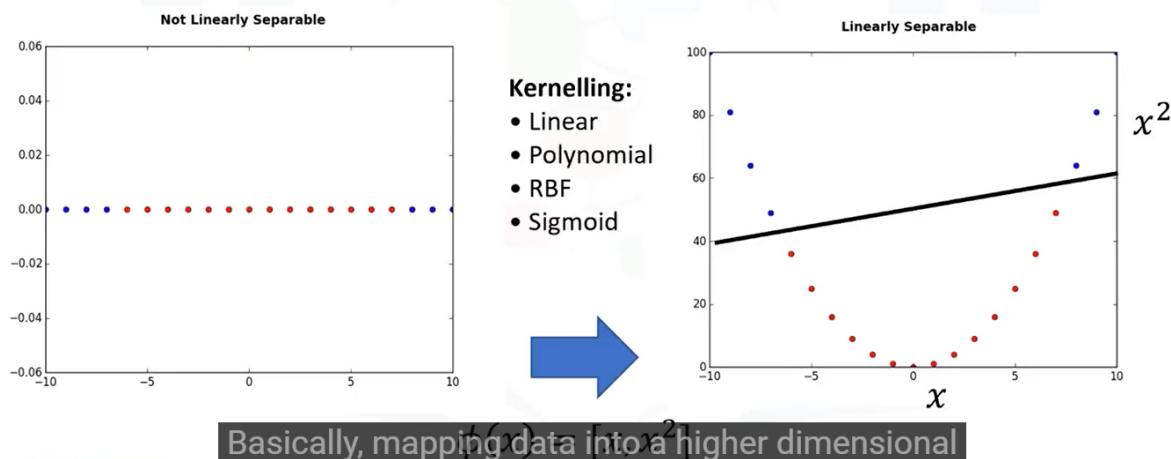
## 1. Mapping data to a **high-dimensional** feature space

## 2. Finding a separator

Clump	UnifSize	UnifShape	MargAdh	SingEpiSize	BareNuc	BlandChrom	NormNucl	Mit	Class
5	1	1	1	2	1	3	1	1	benign
5	4	4	5	7	10	3	2	1	benign
3	1	1	1	2	2	3	1	1	malignant
6	8	8	1	3	4	3	7	1	benign
4	1	1	3	2	1	3	1	1	benign
8	10	10	8	7	10		7	1	malignant
1	1	1	1	2	10	3	1	1	benign
2	1	2	H	2	1	3	1	1	benign
2	1	1	1	2	1	1	1	5	benign
4	2	1	1	2	1	2	1	1	benign

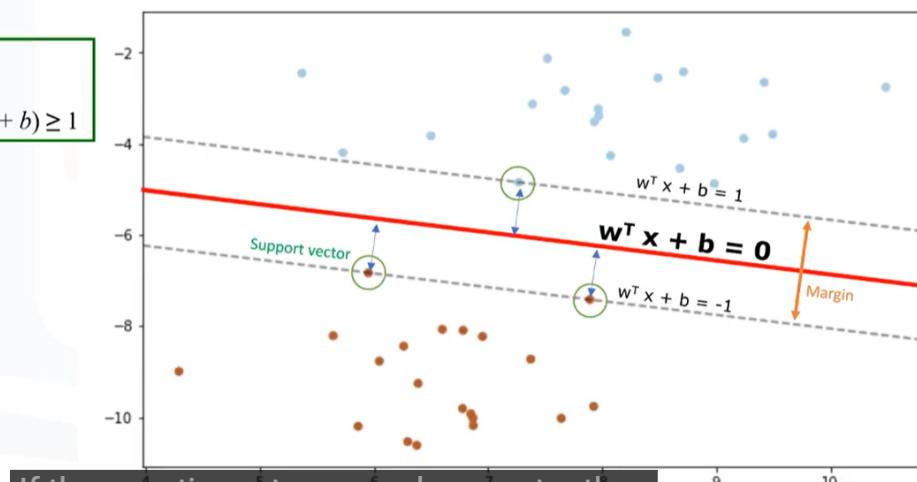


## Data transformation



## Using SVM to find the hyperplane

Find  $\mathbf{w}$  and  $b$  such that  
 $\Phi(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w}$  is minimized;  
and for all  $\{(\mathbf{x}_i, y_i)\}$ :  $y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1$



## Pros and cons of SVM

- Advantages:
  - Accurate in high-dimensional spaces
  - Memory efficient
- Disadvantages:
  - Prone to over-fitting
  - No probability estimation
  - Small datasets

## SVM applications

- Image recognition
- Text category assignment
- Detecting spam
- Sentiment analysis
- Gene Expression Classification
- Regression, outlier detection and clustering

# Module 4

## Learning Objectives

In this lesson you will learn about:

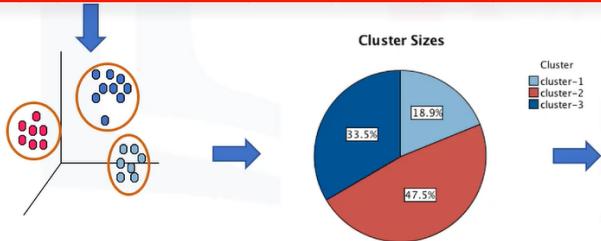
- K-Means Clustering plus Advantages & Disadvantages
- Hierarchical Clustering plus Advantages & Disadvantages
- Measuring the Distances Between Clusters - Single Linkage Clustering
- Measuring the Distances Between Clusters - Algorithms for Hierarchy Clustering
- Density-Based Clustering

## Intro to Clustering

### Clustering for segmentation

Customer ID	Age	Edu	Years Employed	Income	Card Debt	Other Debt	Address	DebtIncomeRatio	Defaulted
1	41	2	6	19	0.124	1.073	NBA001	6.3	0
2	47	1	26	100	4.582	8.218	NBA021	12.8	0
3	33	2	10	57	6.111	5.802	NBA013	20.9	1
4	29	2	4	19	0.681	0.516	NBA009	6.3	0
5	47	1	31	253	9.308	8.908	NBA008	7.2	0
6	40	1	23	81	0.998	7.831	NBA016	10.9	1
7	38	2	4	56	0.442	0.454	NBA013	1.6	0
8	42	3	0	64	0.279	3.945	NBA009	6.6	0
9	26	1	5	18	0.575	2.215	NBA006	15.5	1

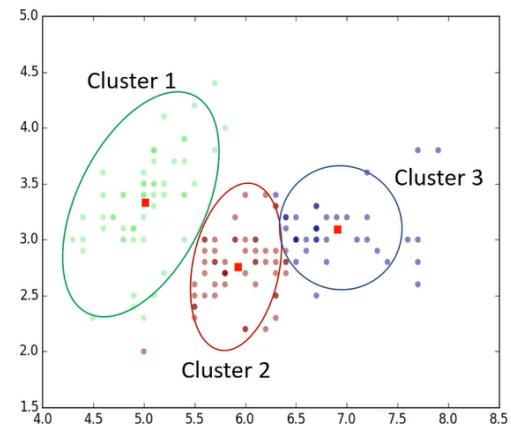
Customer ID	Segment
1	YOUNG AND LOW INCOME
2	AFFLUENT AND MIDDLE AGED
3	AFFLUENT AND MIDDLE AGED
4	YOUNG AND LOW INCOME
5	AFFLUENT AND MIDDLE AGED
6	AFFLUENT AND MIDDLE AGED
7	YOUNG AND LOW INCOME
8	YOUNG AND LOW INCOME
9	AFFLUENT AND MIDDLE AGED



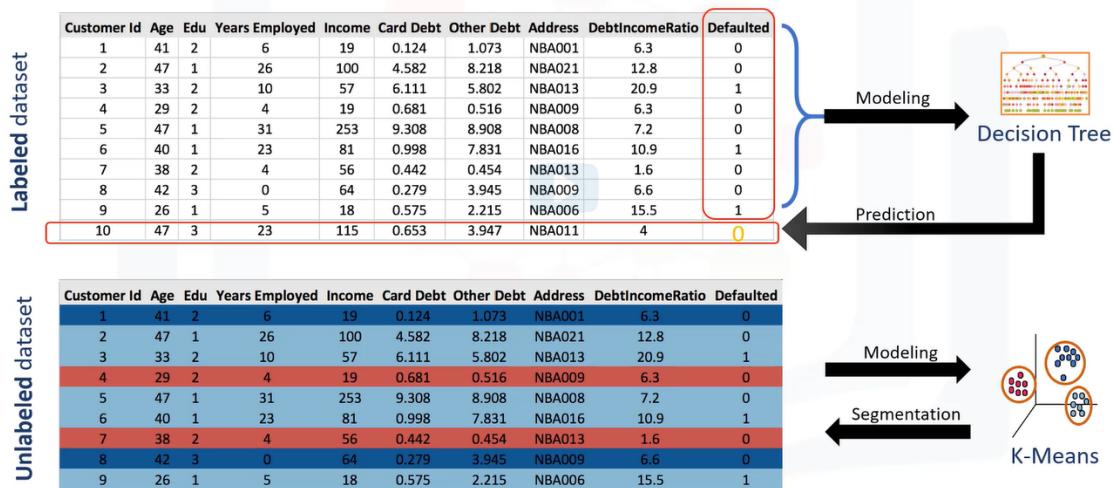
## What is clustering ?

### What is a cluster?

A group of objects that are similar to other objects in the cluster, and dissimilar to data points in other clusters.



## Clustering vs classification



## Clustering application

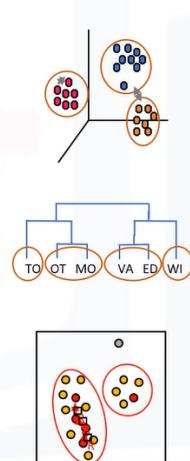
- **RETAIL/MARKETING:**
  - Identifying buying patterns of customers
  - Recommending new books or movies to new customers
- **BANKING:**
  - Fraud detection in credit card use
  - Identifying clusters of customers (e.g., loyal)
- **INSURANCE:**
  - Fraud detection in claims analysis
  - Insurance risk of customers
- **PUBLICATION:**
  - Auto-categorizing news based on their content
  - Recommending similar news articles
- **MEDICINE:**
  - Characterizing patient behavior
- **BIOLOGY:**
  - Clustering genetic markers to identify family ties

## Why clustering ?

- Exploratory data analysis
- Summary generation
- Outlier detection
- Finding duplicates
- Pre-processing step

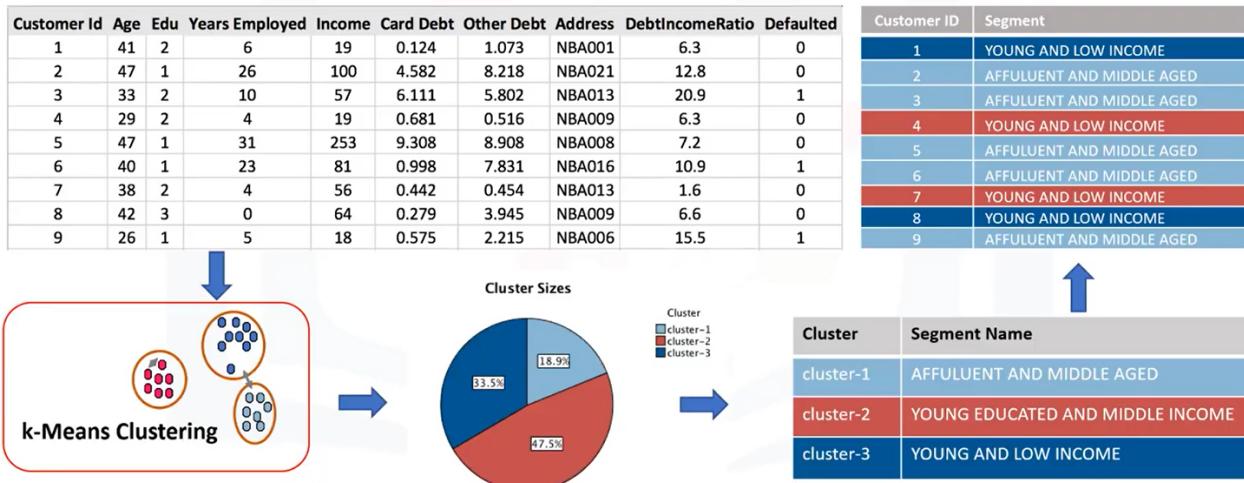
## Clustering algorithms

- **Partitioned-based Clustering**
  - Relatively efficient
  - E.g. k-Means, k-Median, Fuzzy c-Means
- **Hierarchical Clustering**
  - Produces trees of clusters
  - E.g. Agglomerative, Divisive
- **Density-based Clustering**
  - Produces arbitrary shaped clusters
  - E.g. DBSCAN



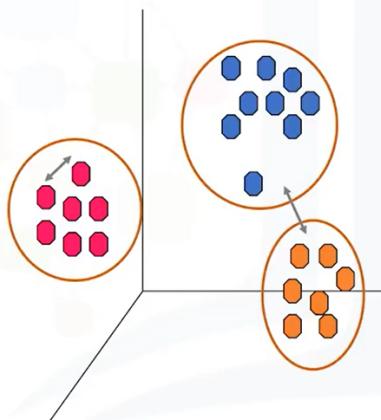
# K-Means Clustering

## What is K-means clustering ?

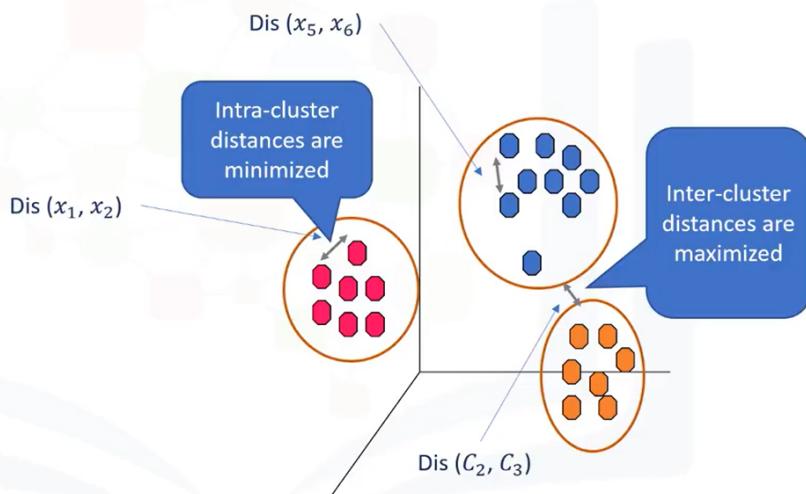


## K-means algorithms

- Partitioning Clustering
- K-means divides the data into non-overlapping subsets (clusters) without any cluster-internal structure
- Examples within a cluster are very similar
- Examples across different clusters are very different



## Determine the similarity or dissimilarity



## 1-dimensional similarity / distance



Customer 1	
Age	
54	



Customer 2	
Age	
50	

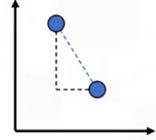
$$\text{Dis } (x_1, x_2) = \sqrt{\sum_{i=0}^n (x_{1i} - x_{2i})^2}$$

$$\text{Dis } (x_1, x_2) = \sqrt{(34 - 30)^2} = 4$$

## 2-dimensional similarity / distance



Customer 1	
Age	Income
54	190



Customer 2	
Age	Income
50	200

$$\text{Dis } (x_1, x_2) = \sqrt{\sum_{i=0}^n (x_{1i} - x_{2i})^2}$$

$$= \sqrt{(54 - 50)^2 + (190 - 200)^2} = 10.77$$

## Multi-dimensional similarity / distance



Customer 1		
Age	Income	education
54	190	3



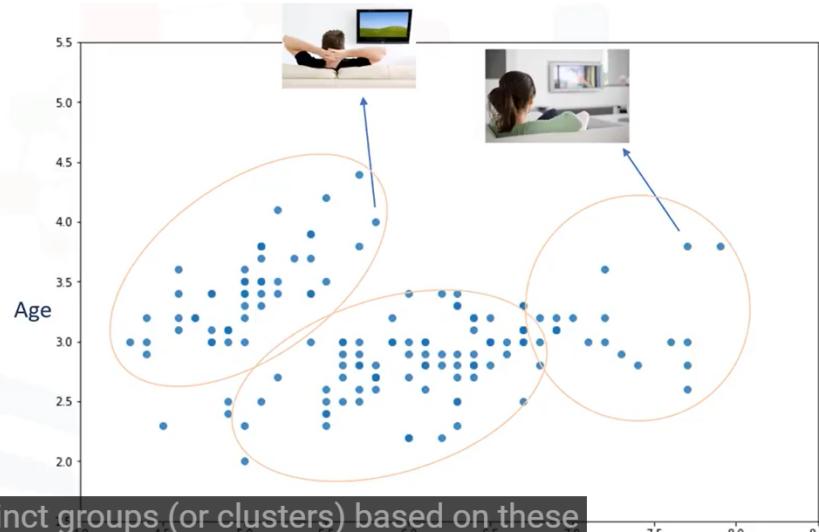
Customer 2		
Age	Income	education
50	200	8

$$\text{Dis } (x_1, x_2) = \sqrt{\sum_{i=0}^n (x_{1i} - x_{2i})^2}$$

$$= \sqrt{(54 - 50)^2 + (190 - 200)^2 + (3 - 8)^2} = 11.87$$

## How does K-means clustering work ?

Customer ID	Age	Income
1	3	4
2	2	6
3	3.5	2
...	...	..

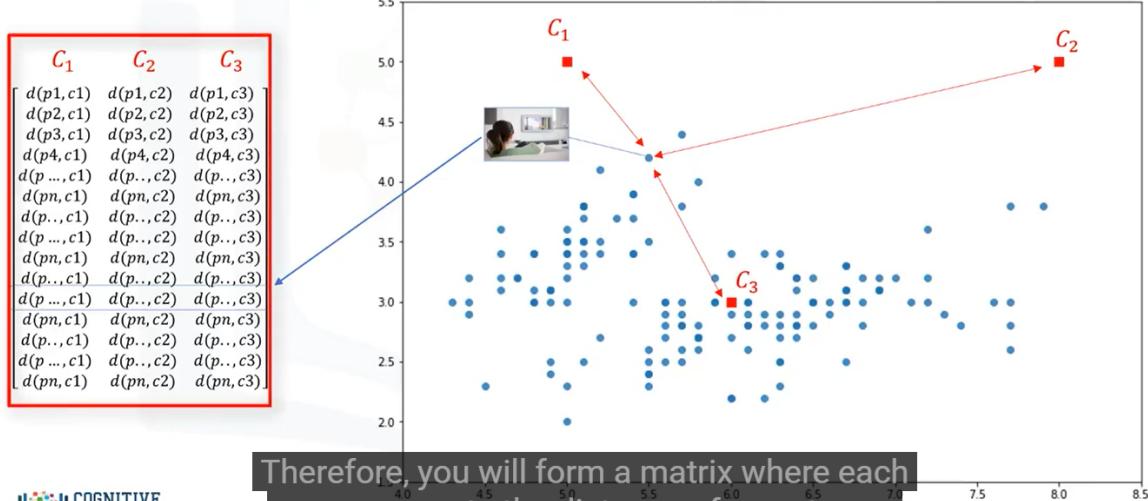


## K-Means clustering - initialize K

!/[k-means initialize](images/k-means-initialize.png)

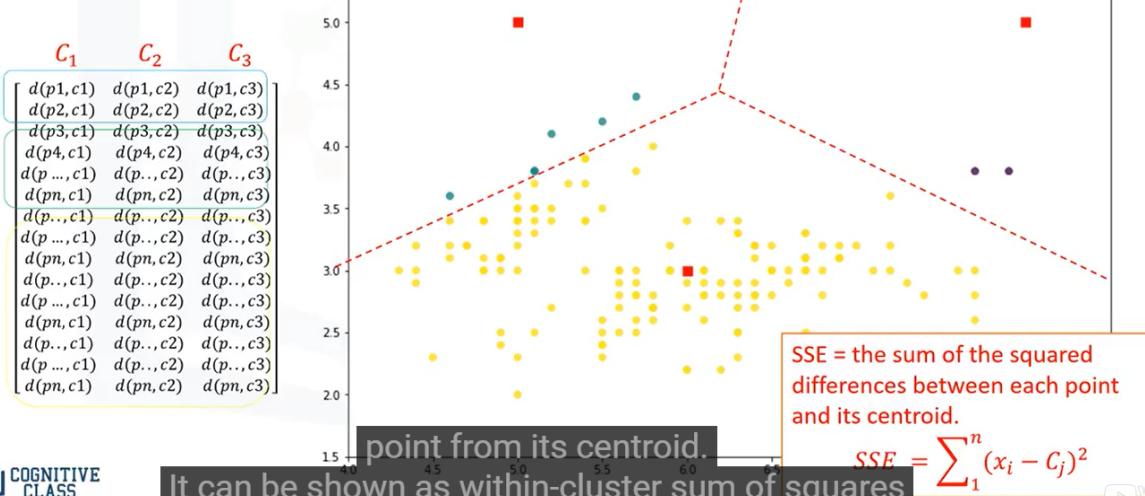
## K-Means clustering - calculate the distance

### 2) Distance calculation



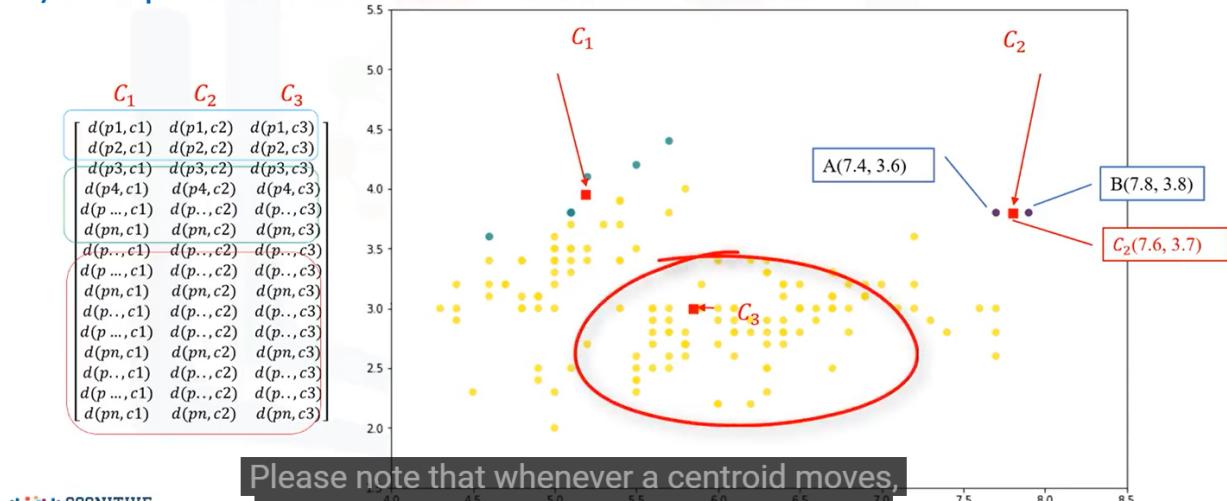
## K-Means clustering - assign to centroid

### 3) Assign each point to the closest centroid



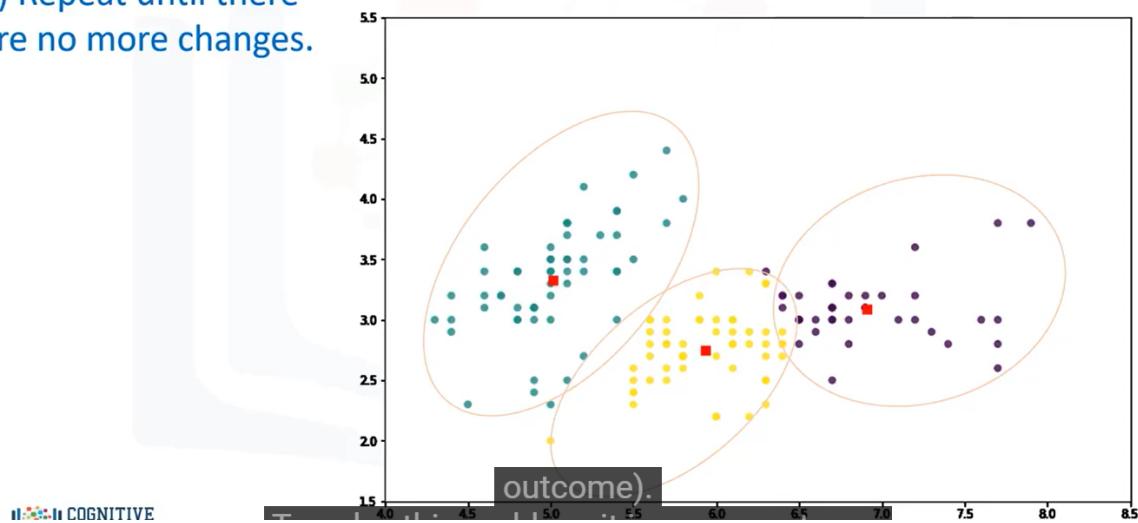
## K-Means clustering - compute new centroids

4) Compute the new centroids for each cluster.



## K-Means clustering - repeat

5) Repeat until there are no more changes.



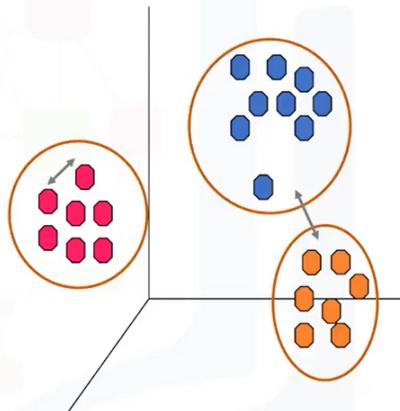
# More on K-Means

# K-Means clustering algorithm

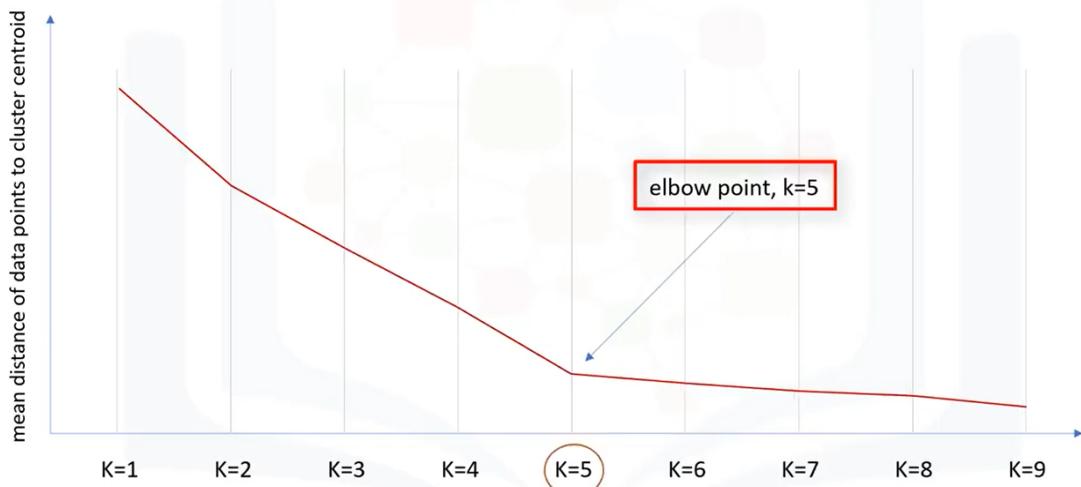
1. Randomly placing  $k$  centroids, one for each cluster.
  2. Calculate the distance of each point from each centroid.
  3. Assign each data point (object) to its closest centroid, creating a cluster.
  4. Recalculate the position of the  $k$  centroids.
  5. Repeat the steps 2-4, until the centroids no longer move.

## K-Means accuracy

- External approach
  - Compare the clusters with the ground truth, if it is available.
- Internal approach
  - Average the distance between data points within a cluster.



## Choosing k



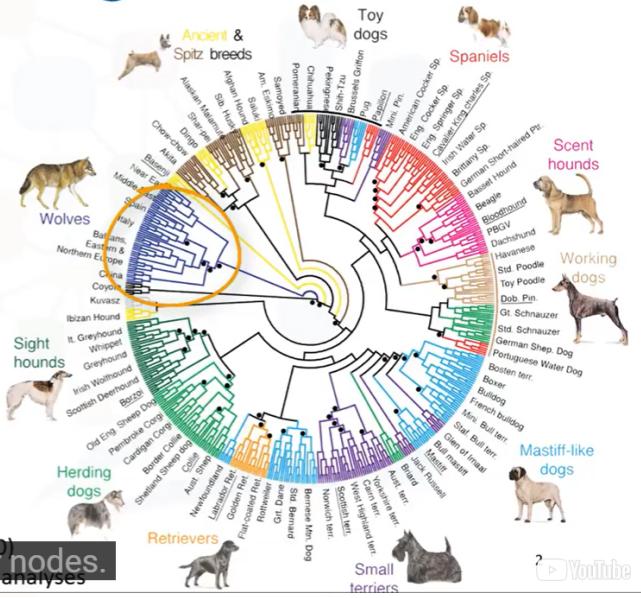
## K-Means recap

- Med and Large sized databases (*Relatively efficient*)
- Produces sphere-like clusters
- Needs number of clusters (k)

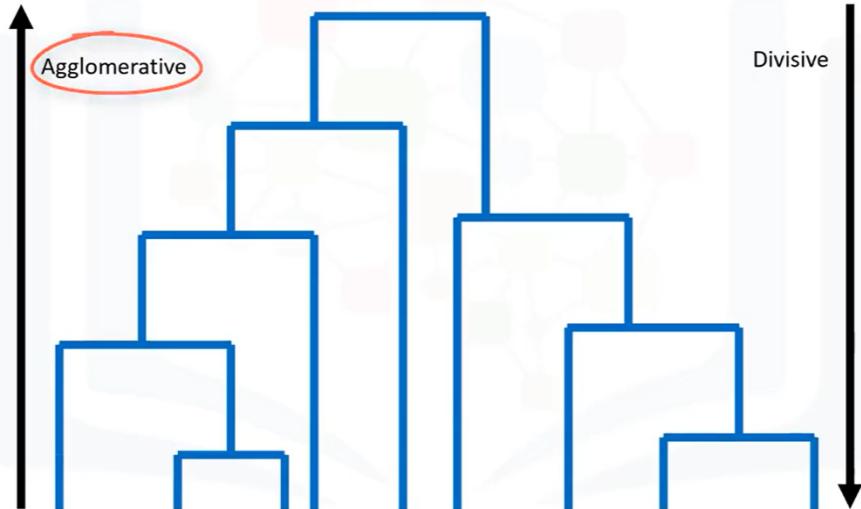
## Hierarchical Clustering

### Hierarchical clustering

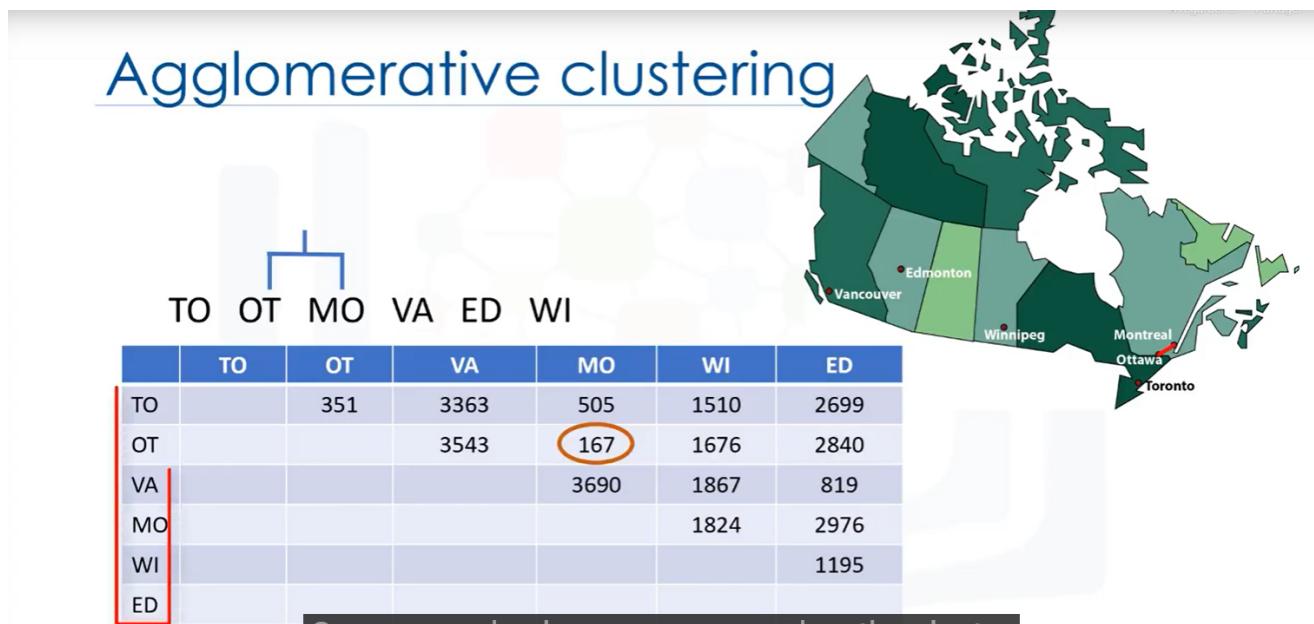
Hierarchical clustering algorithms build a hierarchy of clusters where each node is a cluster consists of the clusters of its daughter nodes.



Source: von Holdt B.M. et al. (2010)  
daughter nodes.  
Genome-wide SNP and haplotype analyses.

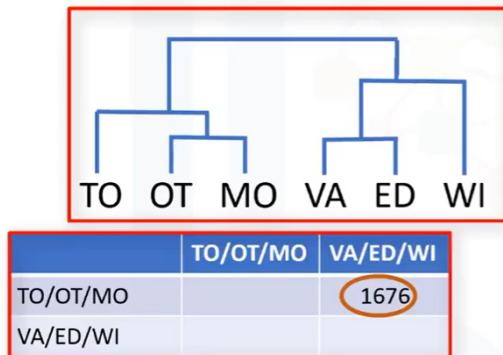


### Agglomerative clustering



## Hierarchical clustering

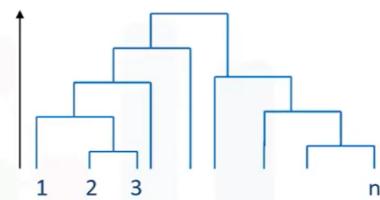
# Hierarchical clustering



## More on Hierarchical Clustering

### Agglomerative algorithm

1. Create  $n$  clusters, one for each data point
2. Compute the Proximity Matrix
3. Repeat
  - i. Merge the two closest clusters
  - ii. Update the proximity matrix
4. Until only a single cluster remains



$$\begin{bmatrix} 0 & & & & \\ d(2,1) & 0 & & & \\ d(3,1) & d(3,2) & 0 & & \\ \vdots & \vdots & \vdots & \ddots & \\ d(n,1) & d(n,2) & \dots & \dots & 0 \end{bmatrix}$$

## Similarity / Distance



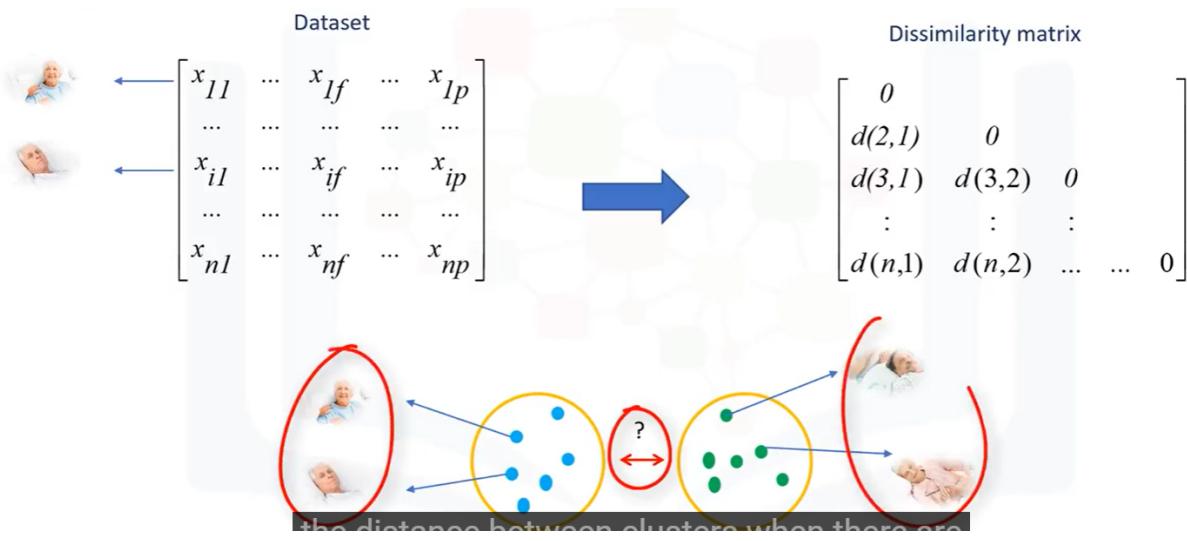
Patient 1	
Age	BMI
54	190

Patient 2	
Age	BP
50	125

Dis (p1,p2)

$$\begin{aligned}
 &= \sqrt{\sum_{i=0}^n (x_i - y_i)^2} \\
 &= \sqrt{(54 - 50)^2 + (190 - 200)^2 + (120 - 125)^2} \\
 &= \text{and Blood Pressure.}
 \end{aligned}$$

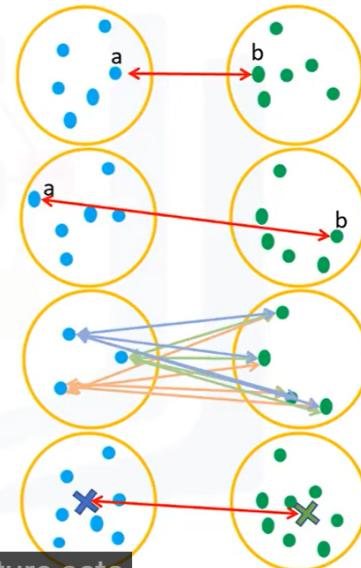
## How can calculate Distance



## Distance between clusters

- Single-Linkage Clustering
  - Minimum distance between clusters
- Complete-Linkage Clustering
  - Maximum distance between clusters
- Average Linkage Clustering
  - Average distance between clusters
- ★ • Centroid Linkage Clustering
  - Distance between cluster centroids

**Linkage Clustering.**  
Centroid is the average of the feature vector.



## Advantages vs. disadvantages

Advantages	Disadvantages
Doesn't required number of clusters to be specified.	Can never undo any previous steps throughout the algorithm.
Easy to implement.	Generally has long runtimes.
Produces a dendrogram, which helps with understanding the data.	Sometimes difficult to identify the number of clusters by the dendrogram.

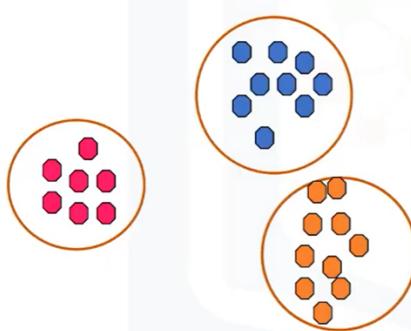
## Hierarchical clustering vs. K-means

K-means	Hierarchical Clustering
1. Much more efficient	1. Can be slow for large datasets
2. Requires the number of clusters to be specified	2. Does not require the number of clusters to run
3. Gives only one partitioning of the data based on the predefined number of clusters	3. Gives more than one partitioning depending on the resolution
4. Potentially returns different clusters each time it is run due to random initialization of centroids	4. Always generates the same clusters

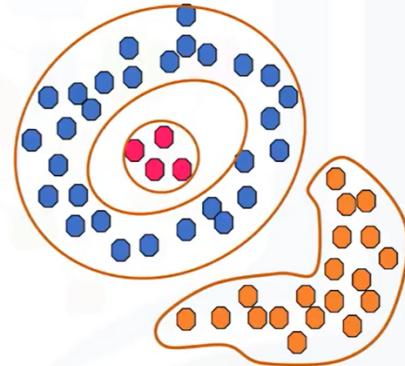
## DBSCAN Clustering

### Density-based clustering

- Spherical-shape clusters



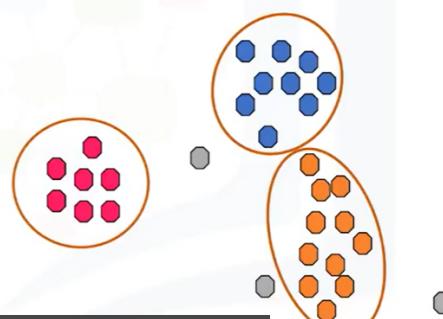
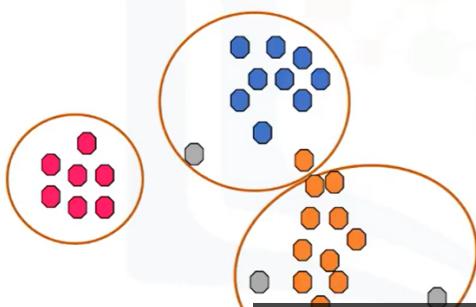
- Arbitrary-shape clusters



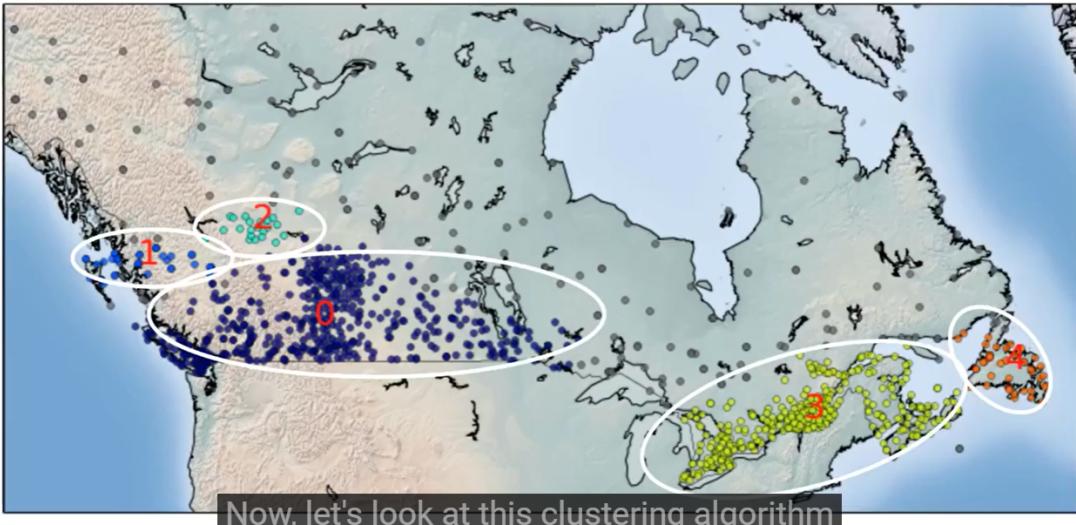
### K-means vs. density-based clustering

- k-Means assigns all points to a cluster even if they do not belong in any

- Density-based Clustering locates regions of **high density**, and separates outliers

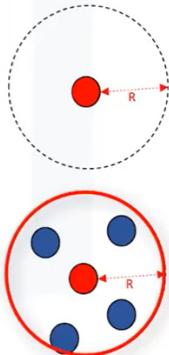


## DBSCAN for class identification

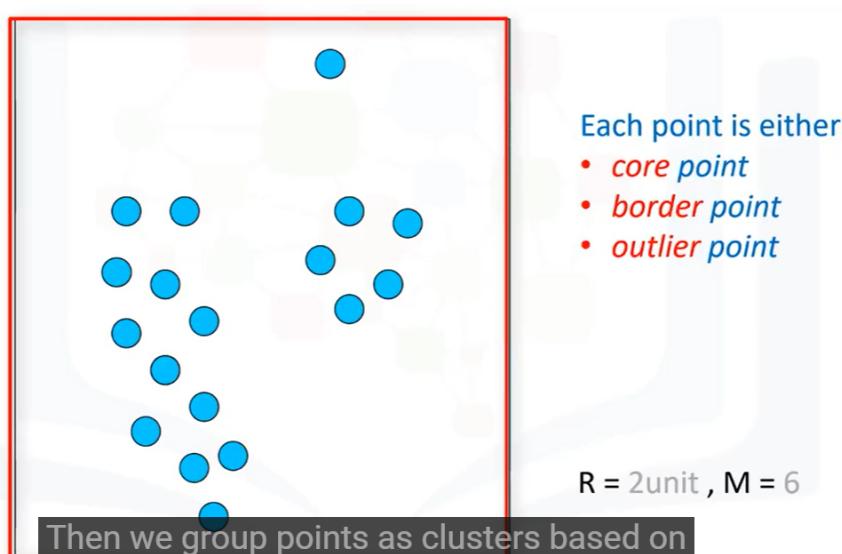


## What is DBSCAN ?

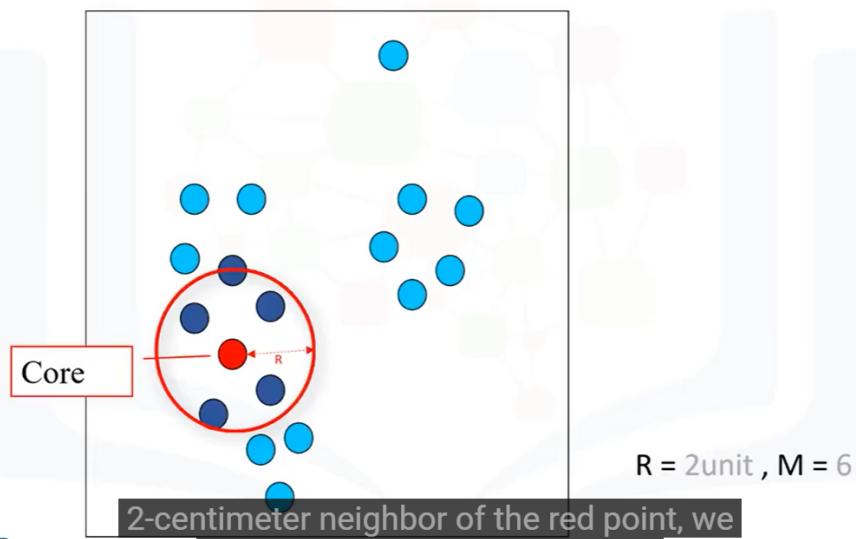
- DBSCAN (Density-Based Spatial Clustering of Applications with Noise)
  - Is one of the most common clustering algorithms
  - Works based on density of objects
- R (Radius of neighborhood)
  - Radius (R) that if includes enough number of points within, we call it a dense area
- M (Min number of neighbors)
  - The minimum number of data points we want in a neighborhood to define a cluster



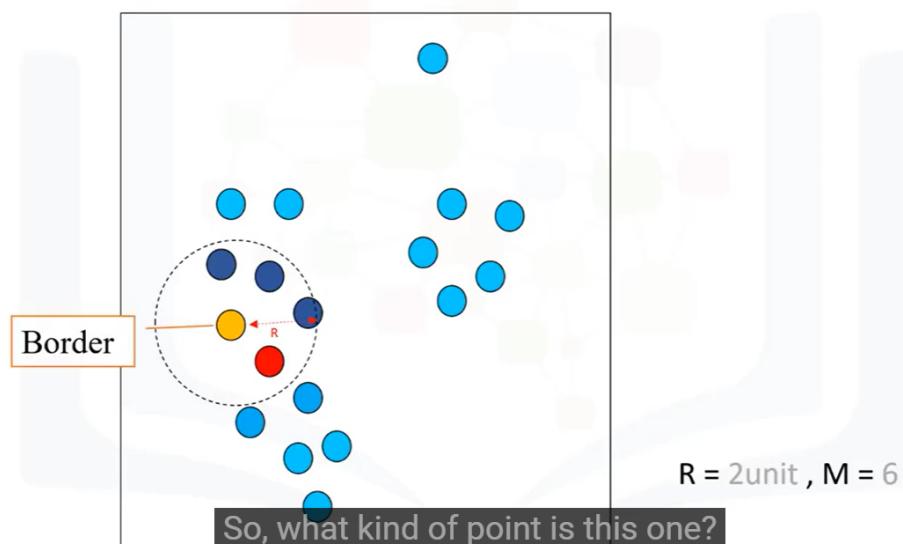
## How DBSCAN works



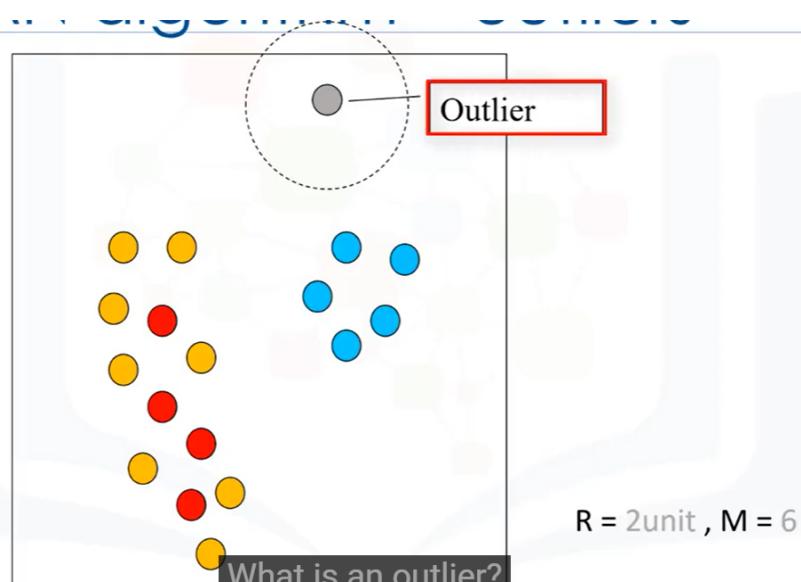
## DBSCAN algorithm - core point ?



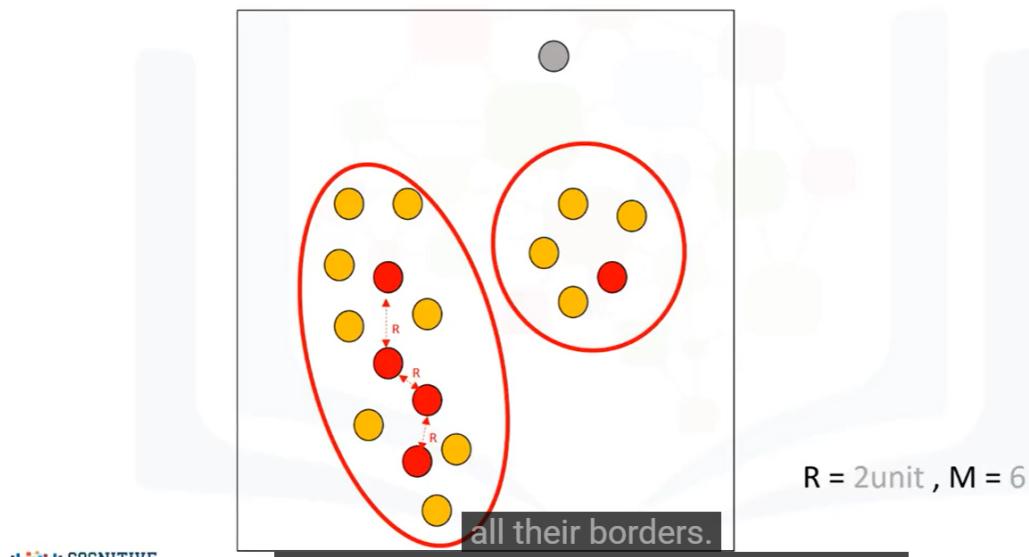
## DBSCAN algorithm - border point ?



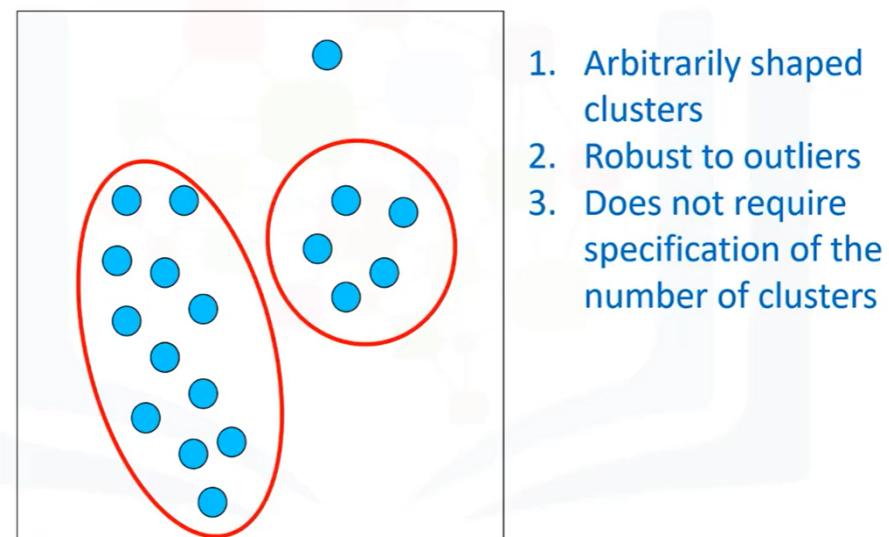
## DBSCAN algorithm - outliers ?



## DBSCAN algorithm - clusters ?



## Advantage of DBSCAN



## Module 5

### Learning Objectives

In this lesson you will learn about:

- To understand the purpose and mechanism of recommendation systems.
- To understand different types of recommender systems.
- To implement recommender system on a real dataset.

### Intro to Recommender Systems

#### What are recommender systems ?

Recommender systems capture the pattern of peoples' behavior and use it to predict what else they might want or like.



g)

## Applications

- What to buy?
  - E-commerce, books, movies, beer, shoes
- Where to eat?
- Which job to apply to?
- Who you should be friends with?
  - LinkedIn, Facebook, ...
- Personalize your experience on the web
  - News platforms, news personalization

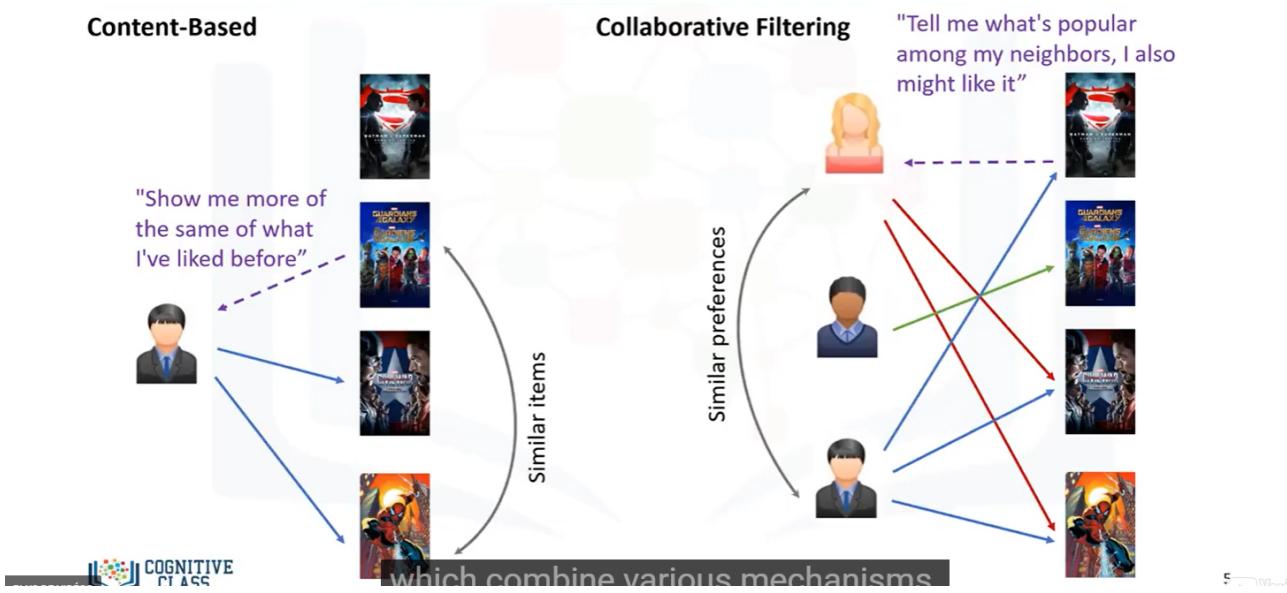


So, let's take a closer look at the main benefits of using a recommendation system.

## Advantage of recommender systems

- Broader exposure
- Possibility of continual usage or purchase of products
- Provides better experience

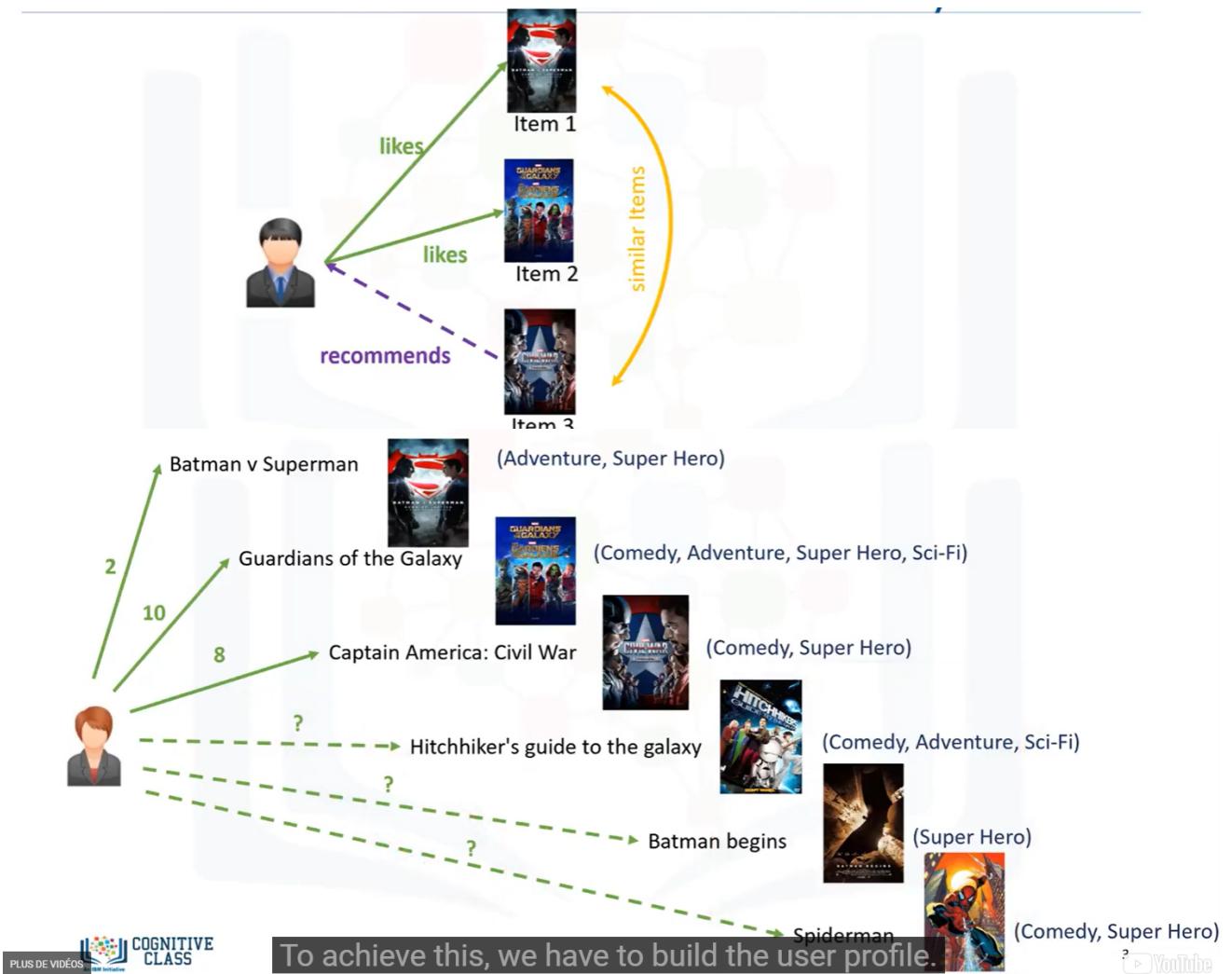
## Two types of recommender systems



## Implementing recommender systems

- **Memory-based**
  - Uses the entire user-item dataset to generate a recommendation
  - Uses statistical techniques to approximate users or items  
e.g., Pearson Correlation, Cosine Similarity, Euclidean Distance, etc.
- **Model-based**
  - Develops a model of users in an attempt to learn their preferences
  - Models can be created using Machine Learning techniques like regression, clustering, classification, etc.

## Content-based recommender systems



## Weighing the genres

### Weighing the genres

Weighted Genre Matrix				
	Comedy	Adventure	Super Hero	Sci-Fi
	0	2	2	0
2	10	10	10	10
10	8	0	8	0
8				

Input User Ratings X Movies Matrix =

User Profile

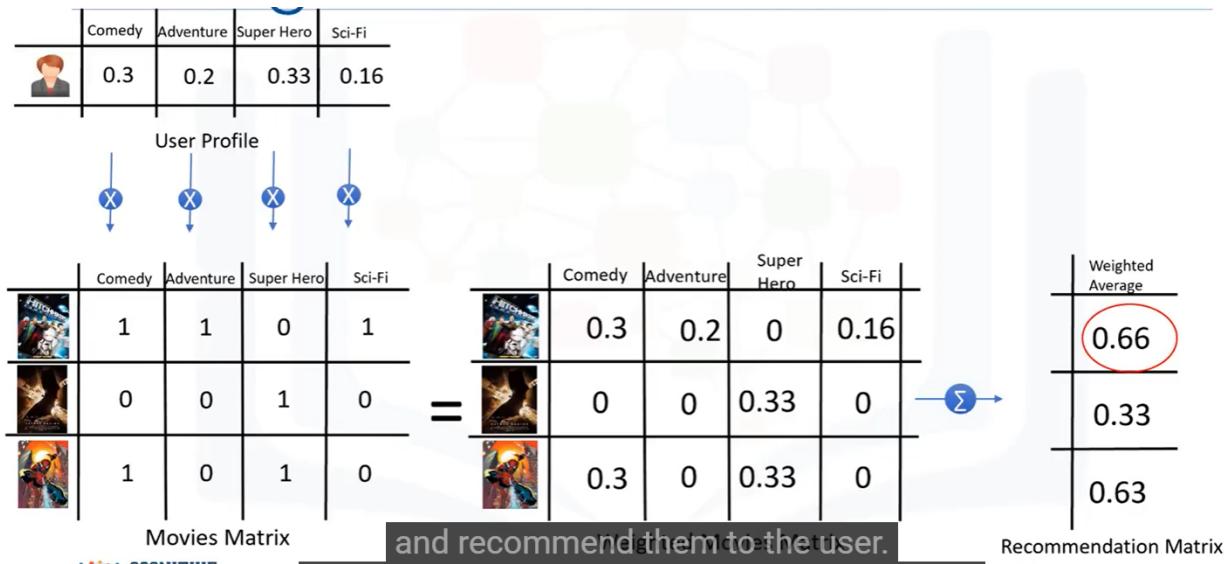
It clearly indicates that she likes "Super hero" movies more than other genres.

PLUS DE VIDÉOS COGNITIVE CLASS

## Candidate movies for recommendation

	Comedy	Adventure	Super Hero	Sci-Fi
	1	1	0	1
	0	0	1	0
	1	0	1	0

## Finding the recommendation



## Come back to recommended



# Collaborative Filtering

- **User-based collaborative filtering**

- Based on users' neighborhood

- **Item-based collaborative filtering**

- Based on items' similarity



Let's first look at the intuition behind.

## User-based collaborative filtering

- **User-based collaborative filtering**



that her neighbor has already seen.

## User ratings matrix

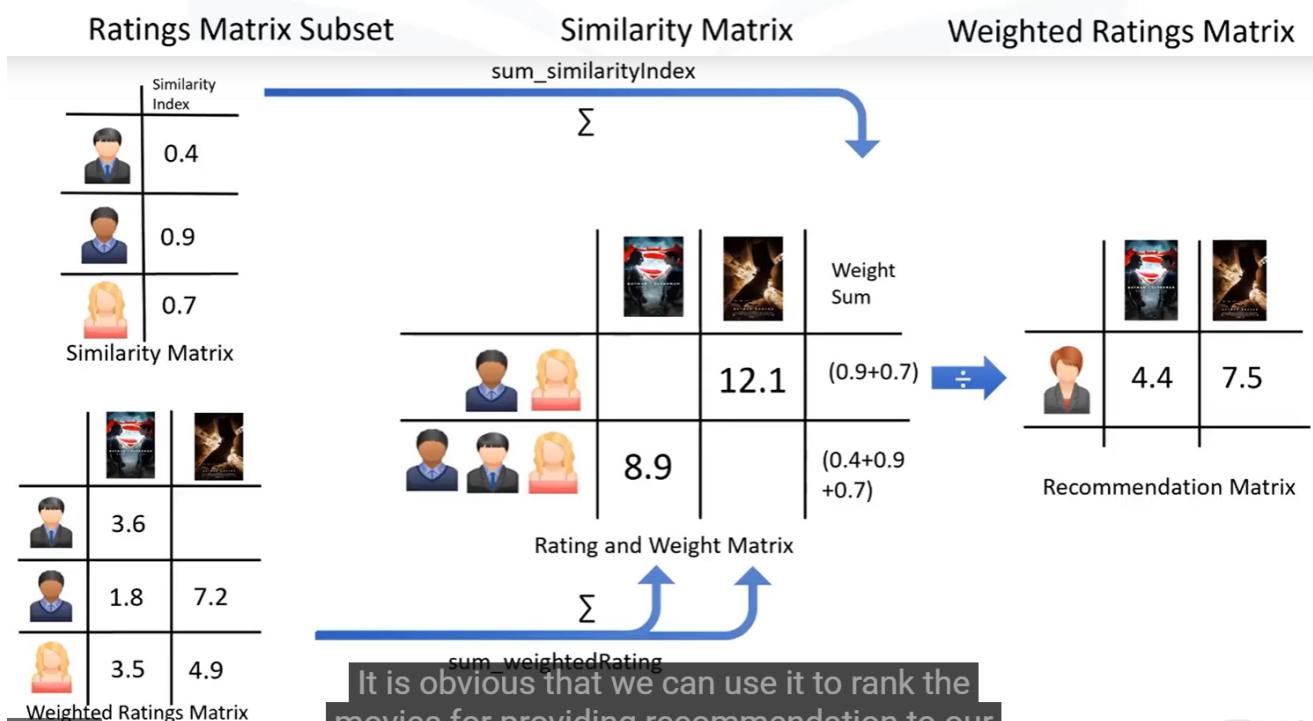
	9	6	8	4	
	2	10	6		8
	5	9		10	7
Active user	(?)	10	7	8	(?)

## Learning the similarity weights



## Creating the weighted ratings matrix

	Movie 1	Movie 2		Similarity Index	
User 1	9		$\times$	0.4	
User 2	2	8	$\times$	0.9	
User 3	5	7	$\times$	0.7	
				=	
				3.6	
				1.8	7.2
				3.5	4.9



# Collaborative filtering



## Challenges of collaborative filtering

- **Data Sparsity**
  - Users in general rate only a limited number of items
- **Cold start**
  - Difficulty in recommendation to new users or new items
- **Scalability**
  - Increase in number of users or items

