

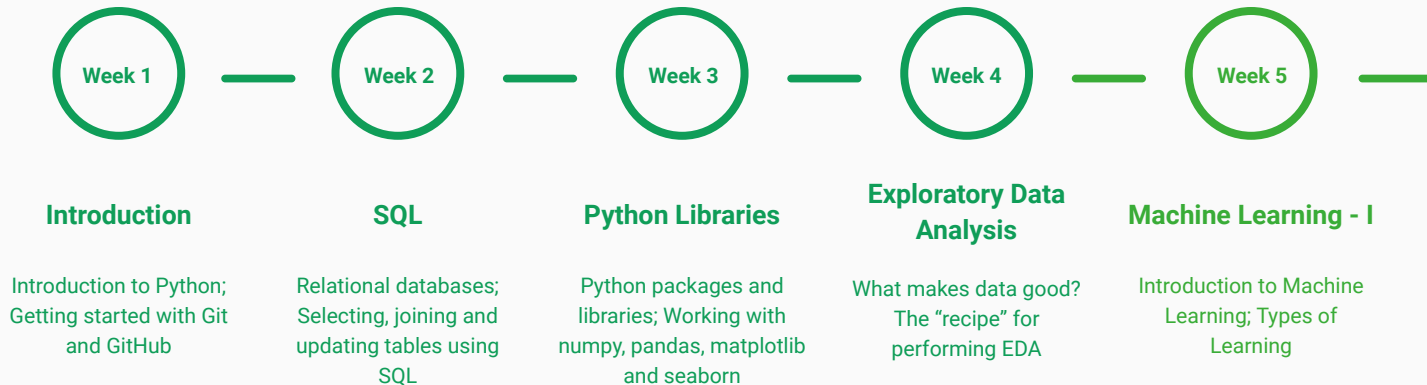
Week 06: Machine Learning - II

Data Science Bootcamp
Fall, 2021

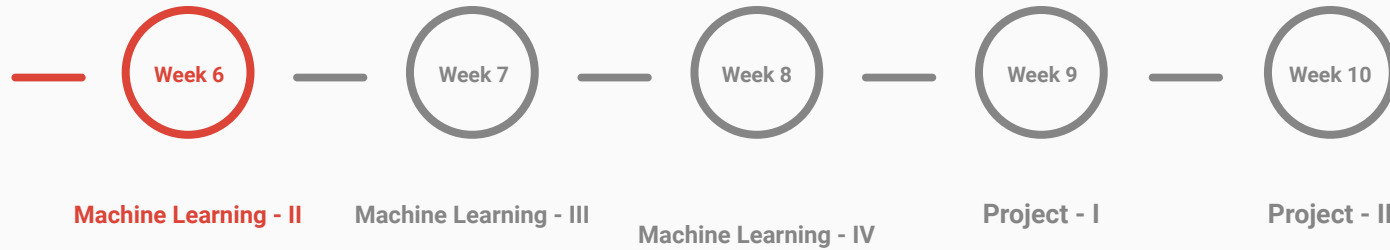
Instructor: Sagar Patel



Halfway there...



Halfway there...



Agenda

- Supervised and Unsupervised Learning
 - Classification; Regression (Linear and Logistic)
 - Clustering
- Model Representation
- Confusion Matrix

A quick Recap

slido



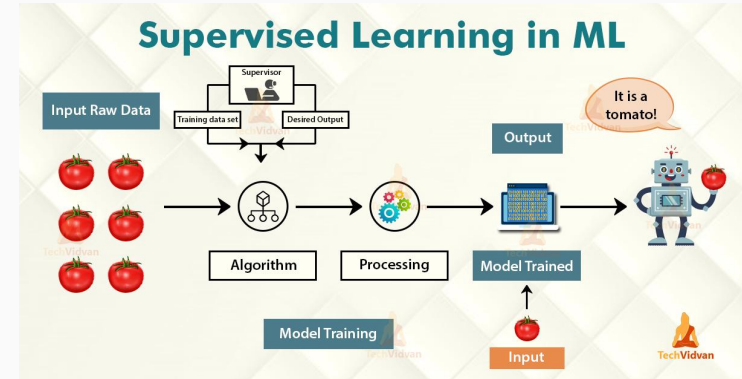
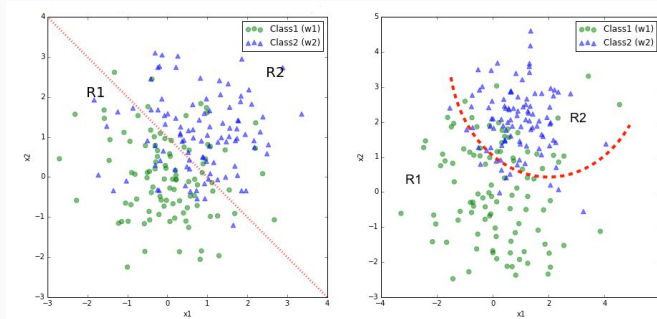
Tomato Detector

① Start presenting to display the poll results on this slide.

Types of “Learning”

- **Supervised Learning**

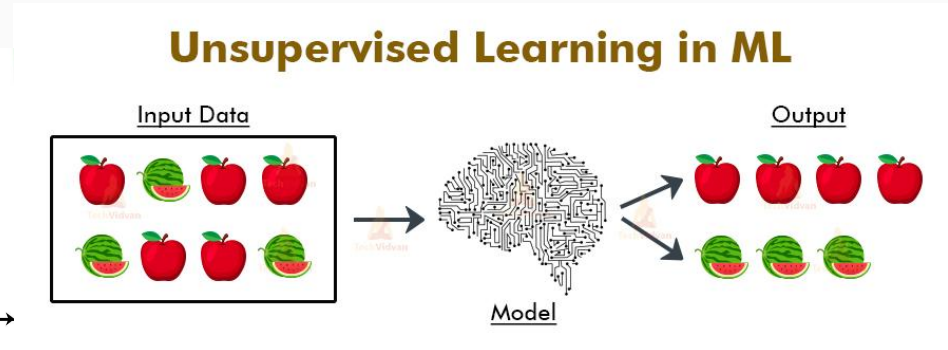
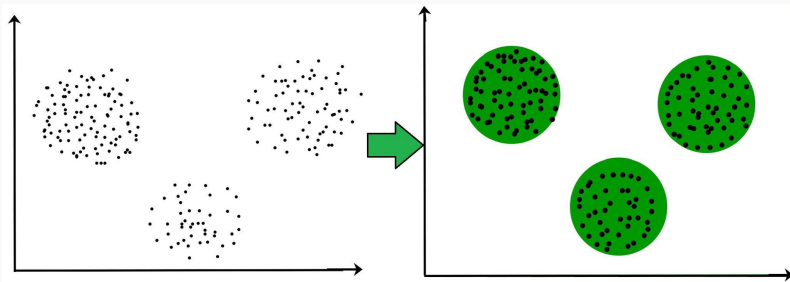
- Algorithms trained using **labeled** data
- The model takes direct **feedback** to check if it's predicting the correct output or not
- Can be categorized in **Classification** and **Regression**
- Example: Tomato Detector



Types of “Learning”

- **Unsupervised Learning**

- Algorithms trained using **unlabeled** (/unknown) data
- There is no **feedback**
- Can be categorized in **Clustering** and **Association**
- Example: Fruit classifier

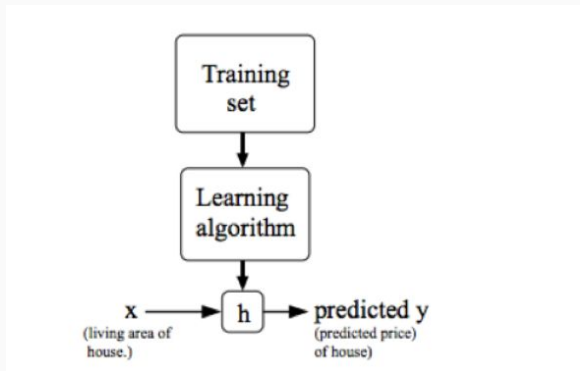


Supervised Learning

- Given a training set, learn a function h (**Hypothesis**) : $X \rightarrow Y$ so that $h(x)$ is a “good” predictor for the corresponding value of y

Regression: When the target variable is continuous - Housing Price Prediction

Classification: When the target variable can only take a small number of discrete values such as predicting if it's a house or an apartment

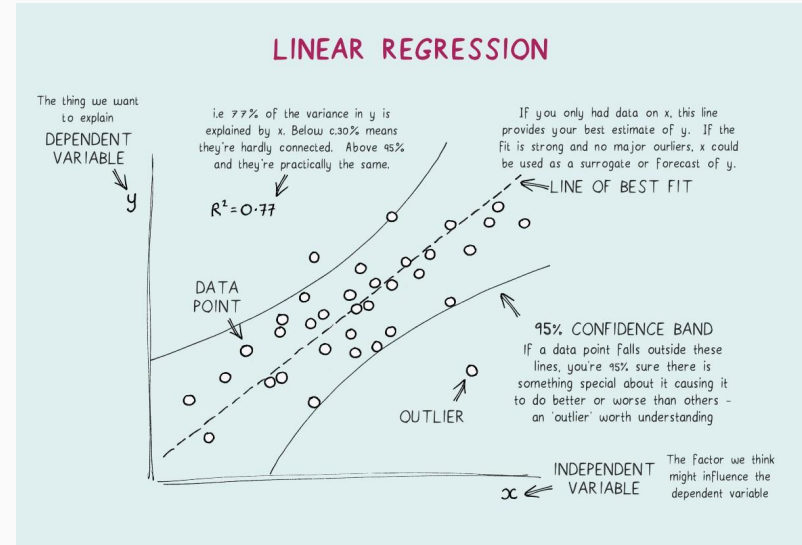


Linear Regression

- Model the relationship between two variables by fitting a **linear equation** to **observed data**
- A **continuous variable** being available
- **Estimating** the dependant

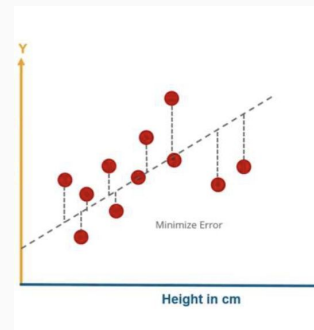
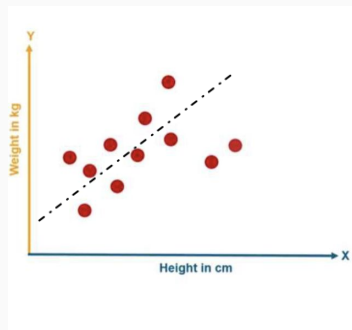
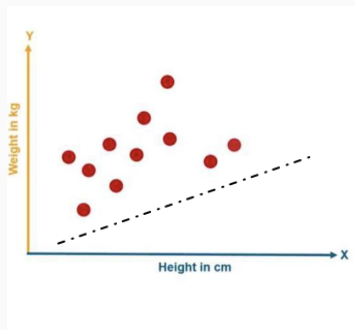
The general equation of a linear regression would be

$$y = ax + b$$



Example

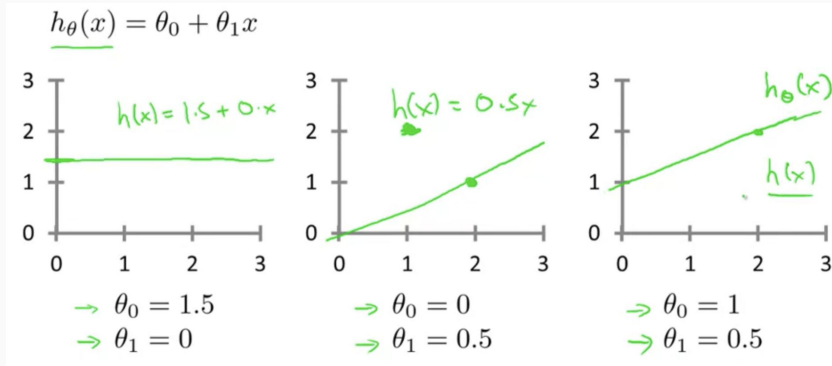
- Start with a random line $w = h$
- Check error in prediction (hypothesis) and update
- Stop when no change or threshold is reached



Hypothesis

Training Set	Size in feet ² (x)	Price (\$) in 1000's (y)
	2104	460
	1416	232
	1534	315
	852	178

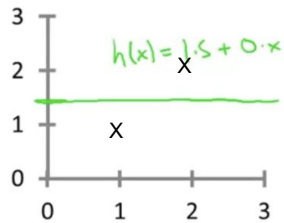
Hypothesis: $h_{\theta}(x) = \theta_0 + \theta_1 x$



Cost Function

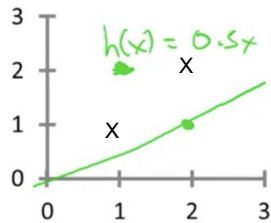
$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (\hat{y}_i - y_i)^2 = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i)^2$$

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$



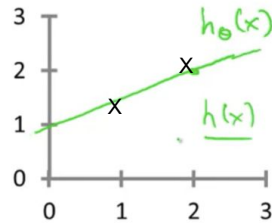
$$\begin{aligned} \rightarrow \theta_0 &= 1.5 \\ \rightarrow \theta_1 &= 0 \end{aligned}$$

$$\text{Cost: } ((1.5-1)^2 + (1.5-2)^2)/2 \cdot 2 = 0.5/4$$



$$\begin{aligned} \rightarrow \theta_0 &= 0 \\ \rightarrow \theta_1 &= 0.5 \end{aligned}$$

Cost: ?



$$\begin{aligned} \rightarrow \theta_0 &= 1 \\ \rightarrow \theta_1 &= 0.5 \end{aligned}$$

$$\text{Cost: } ((1-1)^2 + (2-2)^2)/2 \cdot 2 = 0$$

Cost Function

- Accuracy of hypothesis can be measured by using a **cost function**
- In the previous example, select the values such that the cost is minimum

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (\hat{y}_i - y_i)^2 = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i)^2$$

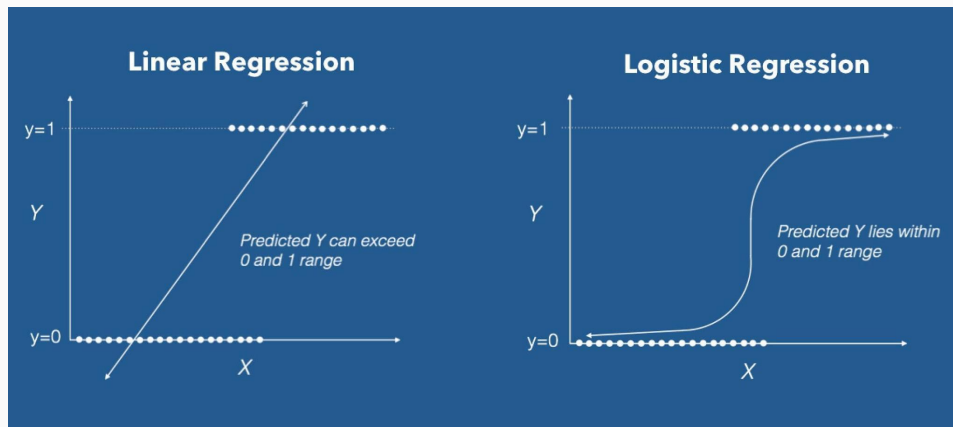
$$\text{Hypothesis: } h_{\theta}(x) = \theta_0 + \theta_1 x$$

- It is an **average of squared difference** between actual output (y) and the hypothesis results in the input (x)

Logistic Regression

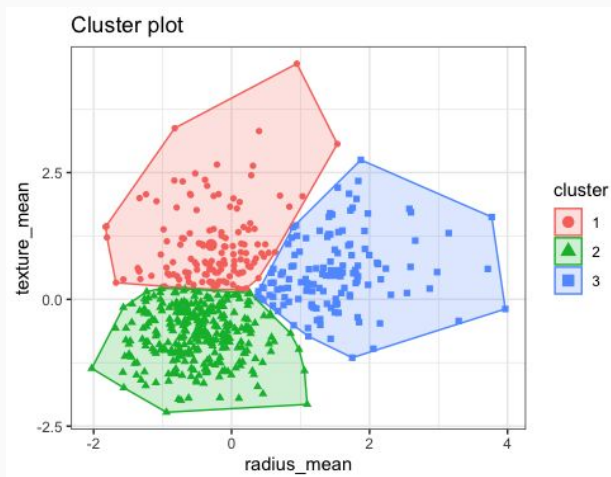
- Modeling the **probability** of a certain class or event
- Can be extended to model **several classes of events**
- Generally used for **classification**
- Has a categorical response

$$P(\text{Event}) = \frac{P(\text{Occurrence})}{P(\text{Not Occurrence})}$$



Clustering

- Automatically discovering **natural grouping** in data
- Unlike **supervised learning**, clustering only interprets the input data and find natural clusters in feature space



Model Representation

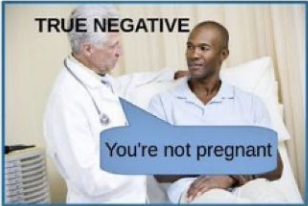



- $\mathbf{x}^{(i)}$ \Leftarrow Input variables / features
- $\mathbf{y}^{(i)}$ \Leftarrow Output variables / features
- $(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})$ \Leftarrow Training Example
- \mathbf{X}/\mathbf{x} \Leftarrow Space of input values
- \mathbf{Y}/\mathbf{y} \Leftarrow Space of output values

Encoding

- Converting categorical data into numerical data to fit and evaluate in the model
- There are multiple methods of encoding

Label Encoding			One Hot Encoding			
Food Name	Categorical #	Calories				
Apple	1	95	Apple	Chicken	Broccoli	Calories
Apple	1	95	1	0	0	95
Chicken	2	231	0	1	0	231
Broccoli	3	50	0	0	1	50

Confusion Matrix

	$\hat{Y} = 0$ NEGATIVE	$\hat{Y} = 1$ POSITIVE
$Y = 0$ NOT PREGNANT	TRUE NEGATIVE 	FALSE POSITIVE  TYPE 1 ERROR
$Y = 1$ PREGNANT	FALSE NEGATIVE  TYPE 2 ERROR	TRUE POSITIVE 

		Predicted class	
		P	N
Actual Class	P	True Positives (TP)	False Negatives (FN)
	N	False Positives (FP)	True Negatives (TN)

Summary

- Supervised learning uses labeled input and output data, while an unsupervised learning algorithm does not
 - Unsupervised learning models work on their own to discover the inherent structure of unlabeled data
- Confusion matrix provides accurate insight into how correctly the model has classified the classes depending upon the data fed or how the classes are misclassified

That's all Folks!

See you in the next session :)

Give us a feedback: <https://bit.ly/3q6ZDID>