

Machine learning for every body

[Source](#)

one hot encoding

[USA,India,Canda,France]

Country	Encoding
USA	[1,0,0,0]
Canada	[0,0,1,0]
India	[0,1,0,0]
France	[0,0,0,1]

Qualitative

- categorical data (finite number of categories or groups).
- Ordinal data (Inherent order)

Quantative

- numerical valued data (could be discrete or continuous)

What are the predictions that our model can output ?

Supervised Learning tasks

1. Classification - Predict discrete classes

Binary classification	Multiclass classification
Positive/negative	cat/dog/lizard/dolphin
cat/dog	orange/apple/pear
Spam/notspam	plant species

2. Regression - predict continuous values

- price of etherium
 - temperation
 - price of real estate
- predict future based on historic values of the features (continuous)

Pregnancies Glucose BloodPressure SkinThickness Insulin BMI Age Outcome									
This is what we would call a feature vector!	6	148	72	35	0	33.6	50	1	
	1	85	66	29	0	26.6	31	0	
	8	183	64	0	0	23.3	32	1	
	1	89	66	23	94	28.1	21	0	
	0	137	40	35	168	43.1	33	1	

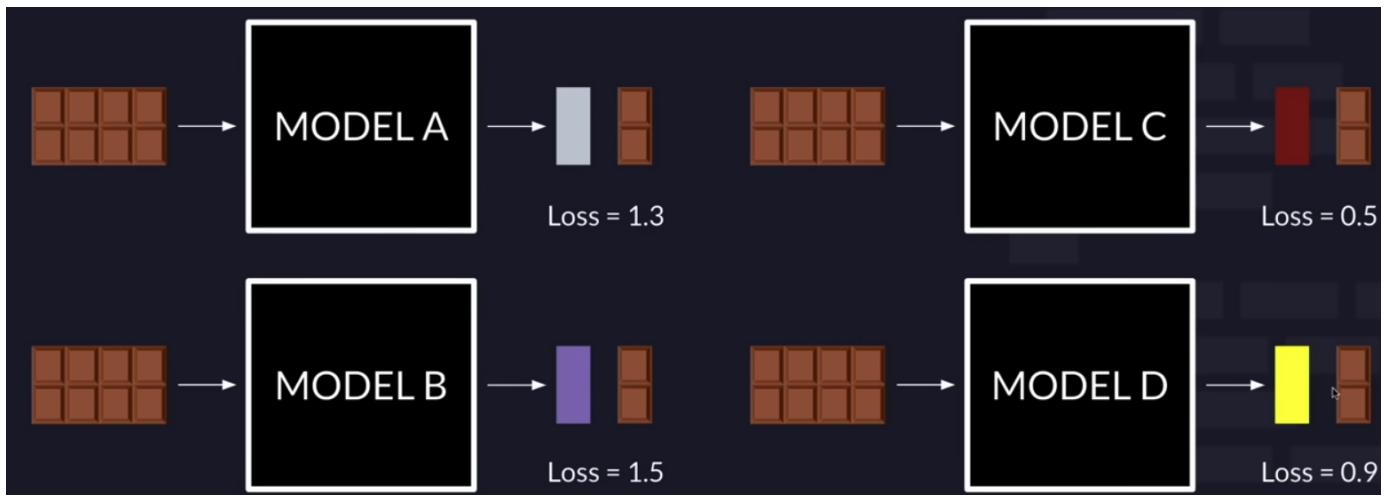
Above image shows, All features of one sample/data set

Divide data into three

- Training dataset
- Validation dataset (used as reality check during/after training to ensure model can handle unseen data)
- Testing dataset

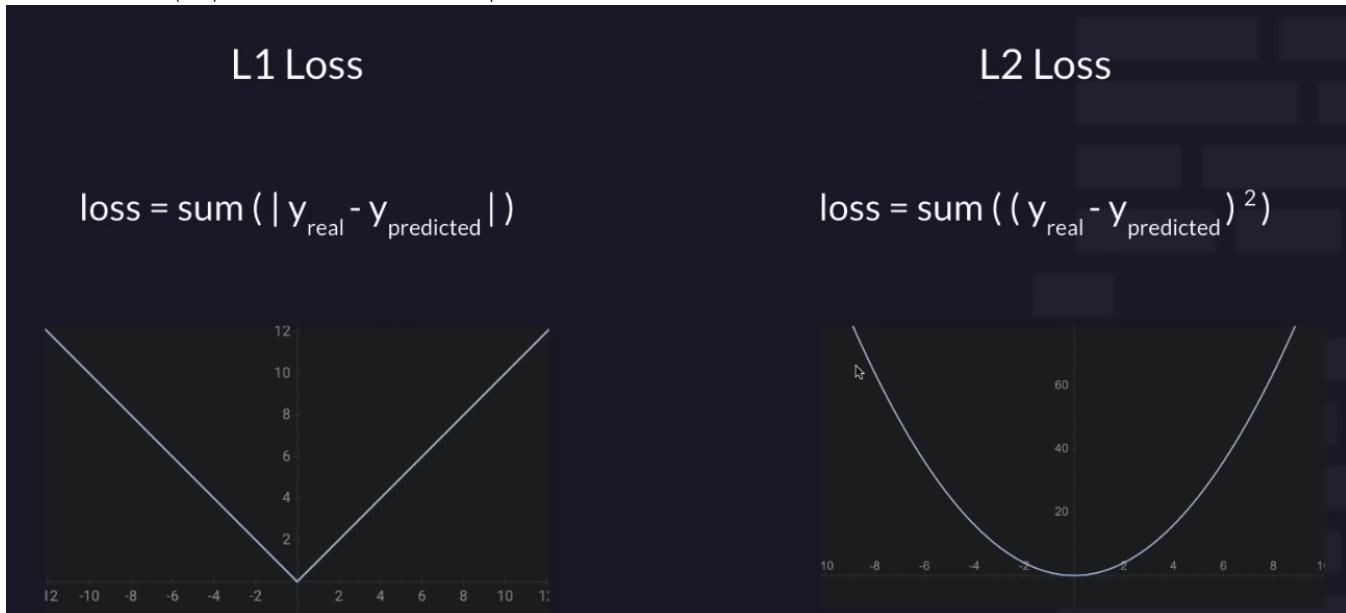
Loss :

Model which has less loss is better performed



Loss functions

Far the value is from 0(zero) over x-axis the more the loss is the poor the model is



Loss decreases as performance increases

Binary Crossentropy is the loss function used when there is a classification problem between 2 categories only.

It is self-explanatory from the name *Binary*, It means 2 quantities, which is why it is constructed in a way that fits the problem of classification of 2 quantities.

Classification models

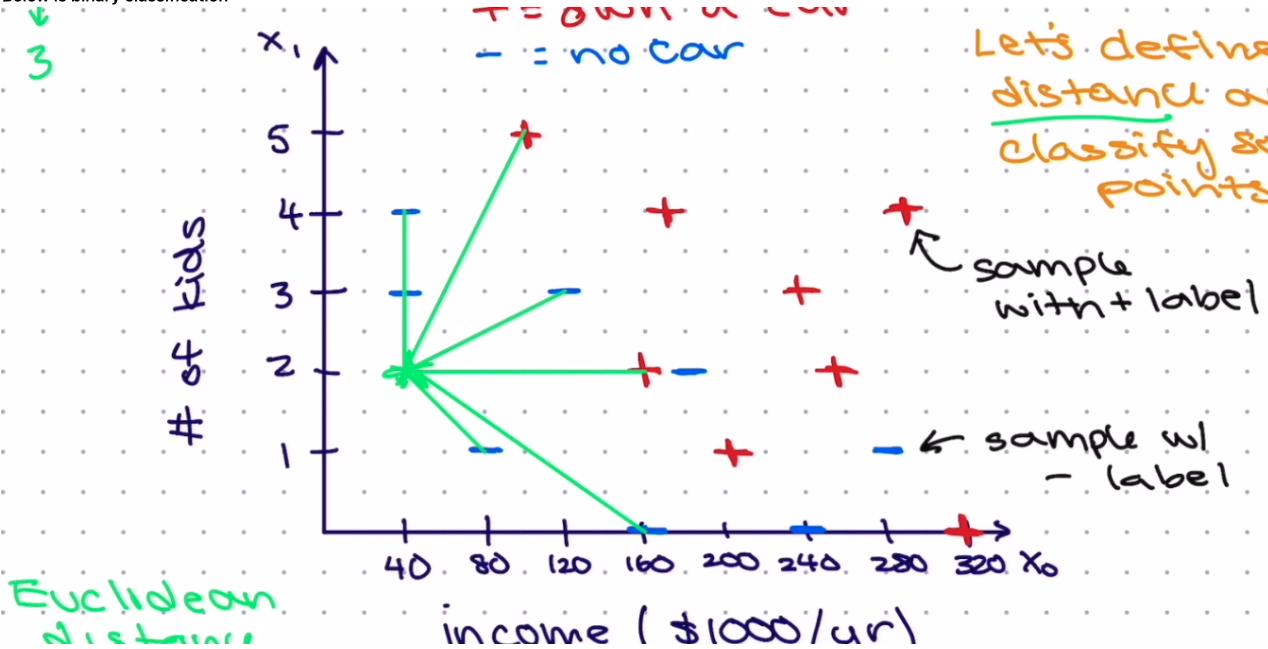
- knn
- naive
- Logistic regression
- svm

K-nearest neighbors Algo

In this example given below we predict if given a person will have car or no. '+' indicates that person has car, '-' indicates that person does not have car.(mapping is from data set). Now we have to predict given a person with income range and children how likely is he can have a car.

Below is binary classification

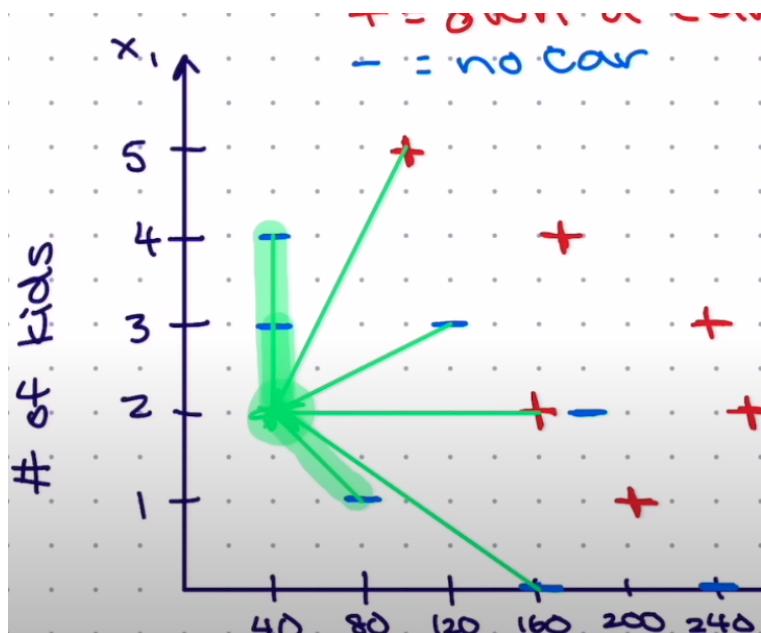
3



K-nearest neighbour algo

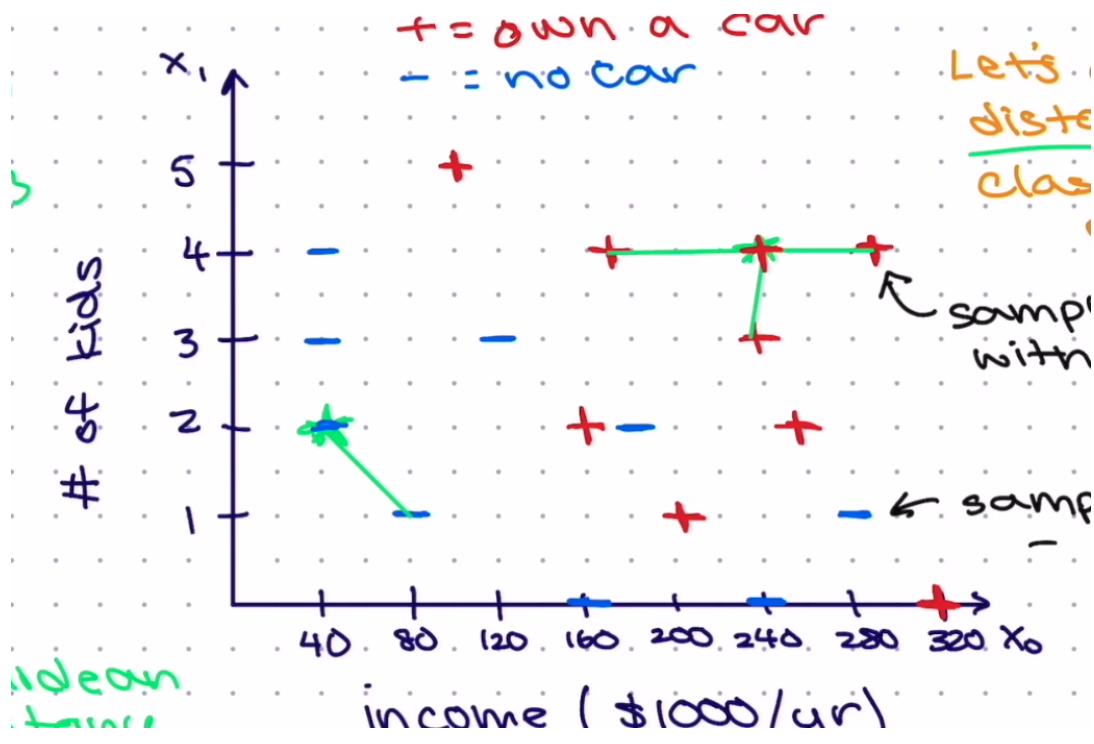
K -> Tell us how many neighbours we are use in order to judge what the label is (usually 3,5)

If k=3, Below image shows what three neighbours we can take. All those points are blue. Chances are prediction is no car

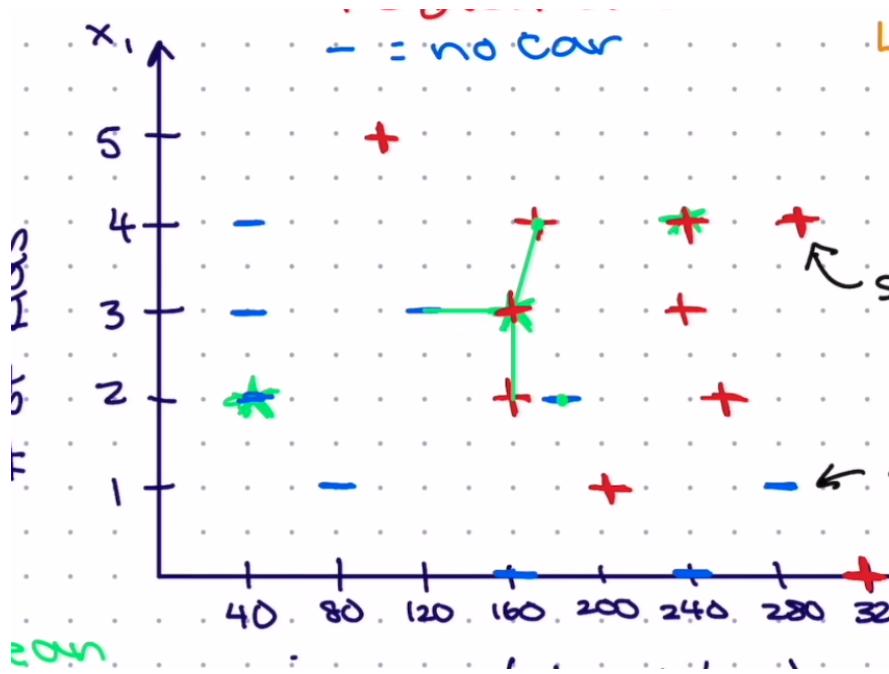


Prediction for -> What if a person has 4 kids and earns around 240k per year.

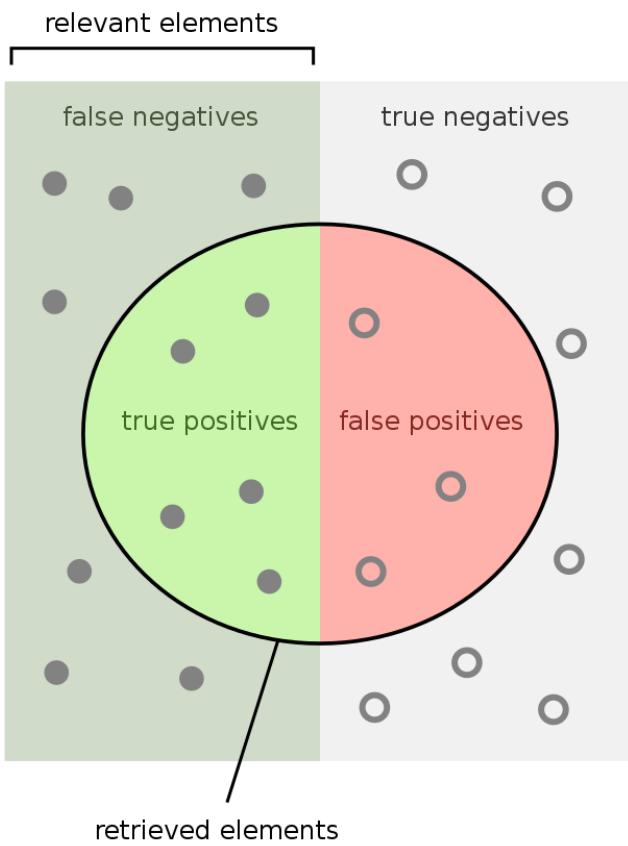
From below image, its more likely that a person will have car



Below image is surrounded by "+" and "-" (red and blue), But two "+" are close to the point, hence conclusion is that the person might have car



Precision and recall



How many retrieved items are relevant?

$$\text{Precision} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}}$$

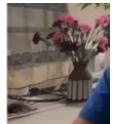
How many relevant items are retrieved?

$$\text{Recall} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$

Precision -> Out of all the ones we labeled positives how many are true positives

Recall -> Out of all the ones we know truly positive how many do we actually get right

Bayes formula



Covid test result

		+	-	
Has covid?	Y	531	6	537
	N	20	9443	9463
Total	551		9449	

What is the probability of having covid given a positive test?

$$P(\text{Covid} | + \text{test}) = 531/551 \Rightarrow 96.4\%$$

| -> given that

Bayes Formula

$$P(A|B) = P(AnB)/P(B) = P(A) \cdot P(B|A)/P(B)$$

where:

P(A)= The probability of A occurring

P(B)= The probability of B occurring

P(A|B)=The probability of A given B

P(B|A)= The probability of B given A

P(AnB)= The probability of both A and B occurring

From above example : $P(\text{Covid} | + \text{test}) = P(A|B)$

Probability will sum up to 1

$$P(+ | \text{no disease})$$

$P(\text{false positive}) = 0.05$

$$P(- | \text{disease})$$

$P(\text{false negative}) = 0.01$



$$P(\text{disease}) = 0.1$$

disease	0.99	0.01
no disease	0.05	0.95

$$P(\text{disease} | (+) \text{test}) = ?$$

$$= \frac{P(+ | \text{disease}) \cdot P(\text{disease})}{P(+)}$$

$$\begin{aligned}
 P(\text{User} | +) &= \frac{P(+) | \text{User})P(\text{User})}{P(+)} \\
 &= \frac{P(+) | \text{User})P(\text{User})}{P(+) | \text{User})P(\text{User}) + P(+) | \text{Non-user})P(\text{Non-user})} \\
 &= \frac{0.99 \times 0.005}{0.99 \times 0.005 + 0.01 \times 0.995} \\
 &\approx 33.2\%
 \end{aligned}$$



$$P(\text{disease} | (+) \text{ test}) = ?$$

$$= \frac{P(+) | \text{disease}) \cdot P(\text{disease})}{P(+) | \text{disease}) \cdot P(\text{disease}) + P(+) | \text{no disease}) \cdot P(\text{no disease})}$$

$$= \frac{0.99 \cdot 0.1}{0.99 \cdot 0.1 + 0.05 \cdot 0.9}$$

$$= 0.6875 \quad \text{or} \quad 68.75\%$$

Naive Bayes

$$P(C_k | x_1, x_2, \dots, x_n) \propto p(C_k) \prod_{i=1}^n p(x_i | C_k)$$

↑ ↑ ↑ proportional to



Derivation:

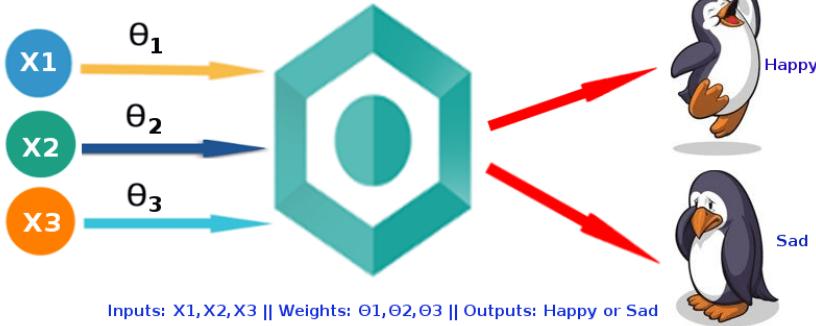
$$\begin{aligned}
 P(C_k | x_1, x_2, \dots, x_n) &= \frac{P(x_1, x_2, \dots, x_n | C_k) \cdot P(C_k)}{P(x_1, x_2, \dots, x_n)} \\
 P(C_k | x_1, x_2, \dots, x_n) &\propto P(x_1, x_2, \dots, x_n | C_k) P(C_k) \\
 &\propto [P(x_1 | C_k) \cdot P(x_2 | C_k) \cdot \dots \cdot P(x_n | C_k)] P(C_k) \\
 &\propto P(C_k) \prod_{i=1}^n P(x_i | C_k)
 \end{aligned}$$

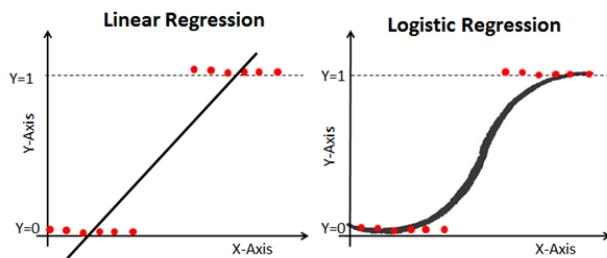
\hat{y} (predicted y):

Logistic Regression:

- Classify using regression
- we will try to fit data in to sigmoid function so the probability is between 0 and 1

Logistic Regression Model

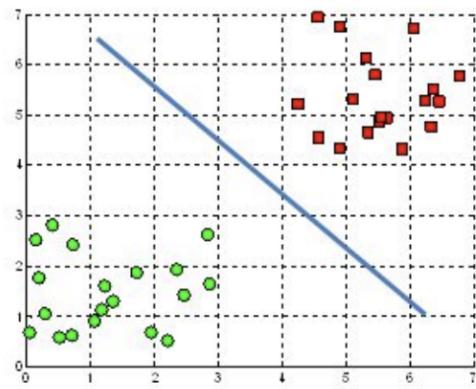




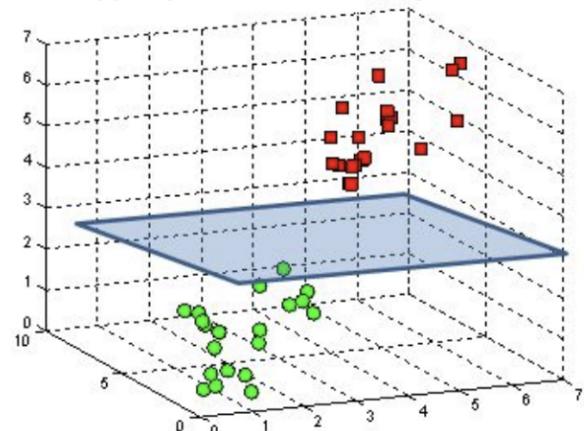
Support Vector Machine - SVM

- In 2d it looks like a line dividing but in 3d its a plane that divides data

A hyperplane in \mathbb{R}^2 is a line

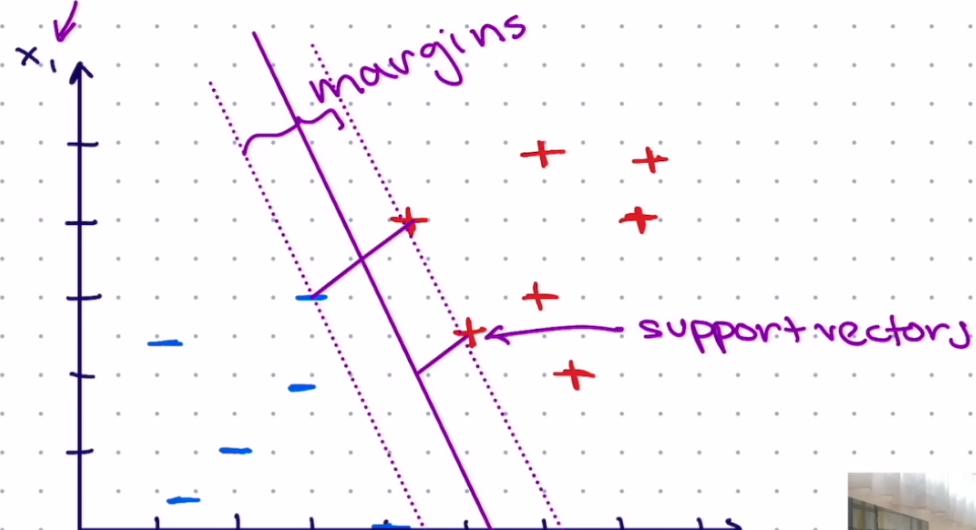


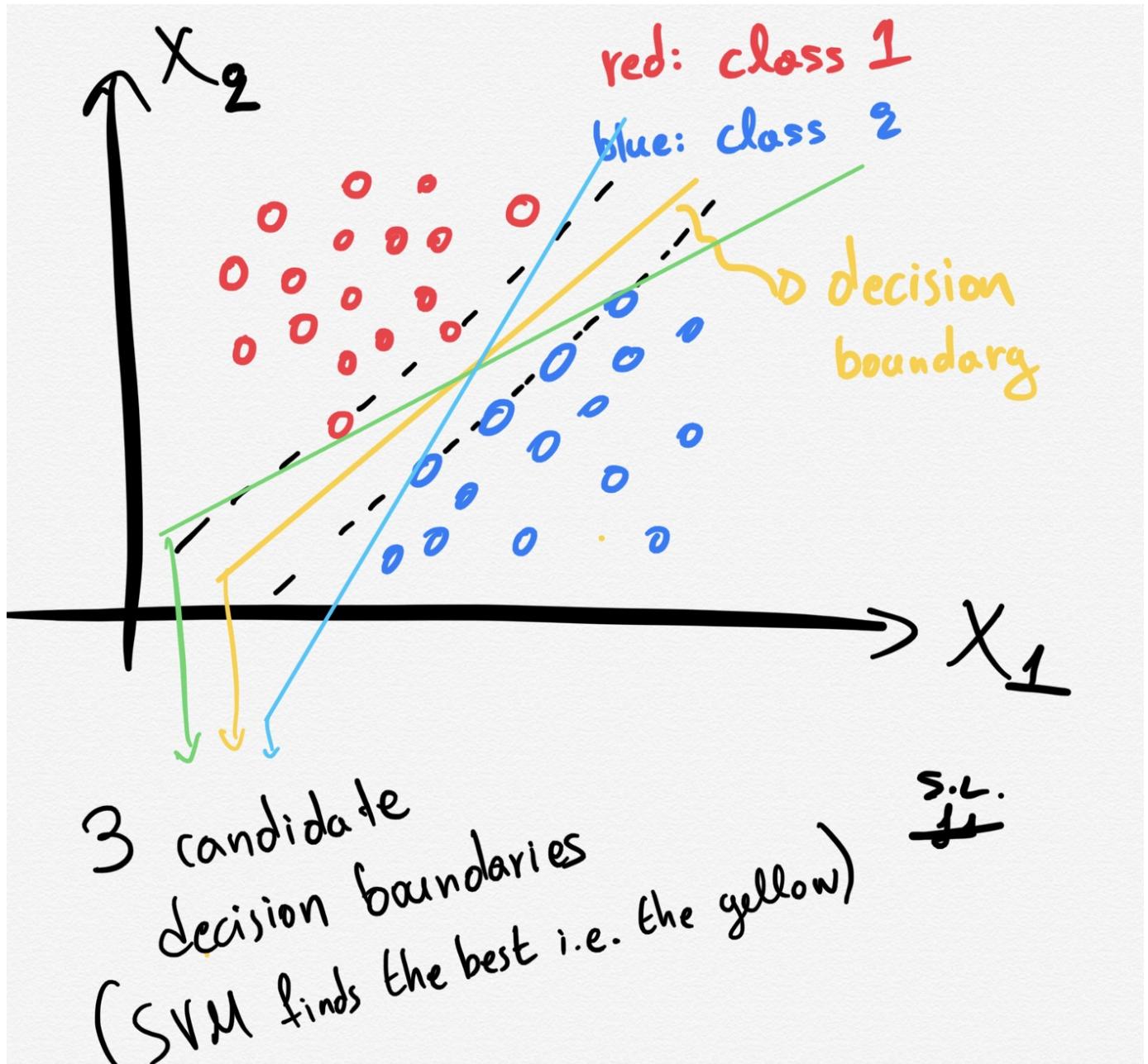
A hyperplane in \mathbb{R}^3 is a plane



In SVM its not just classifying data, we also look at the boundary(margin) of the classifying line. (It should not be close to one type of data)

Support Vector Machines (SVM)



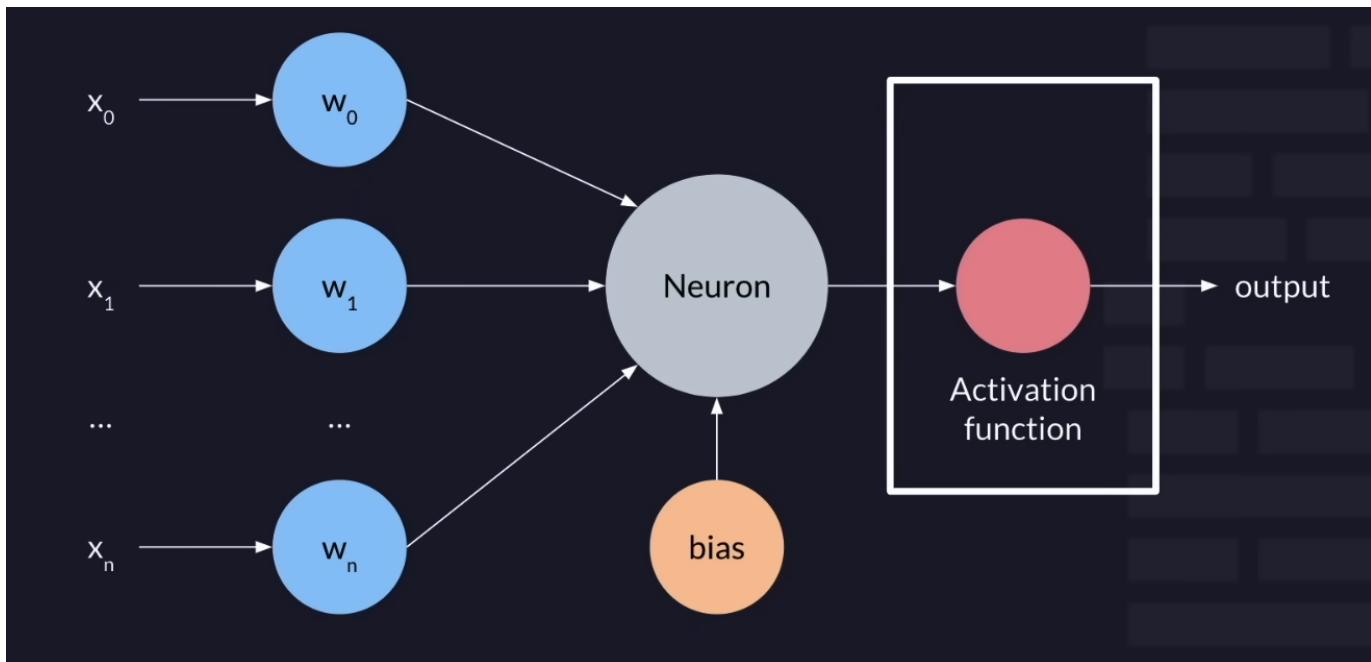


Kernel trick $\rightarrow x$ to (x, x^2)

Neural Network

Classification using Neural network

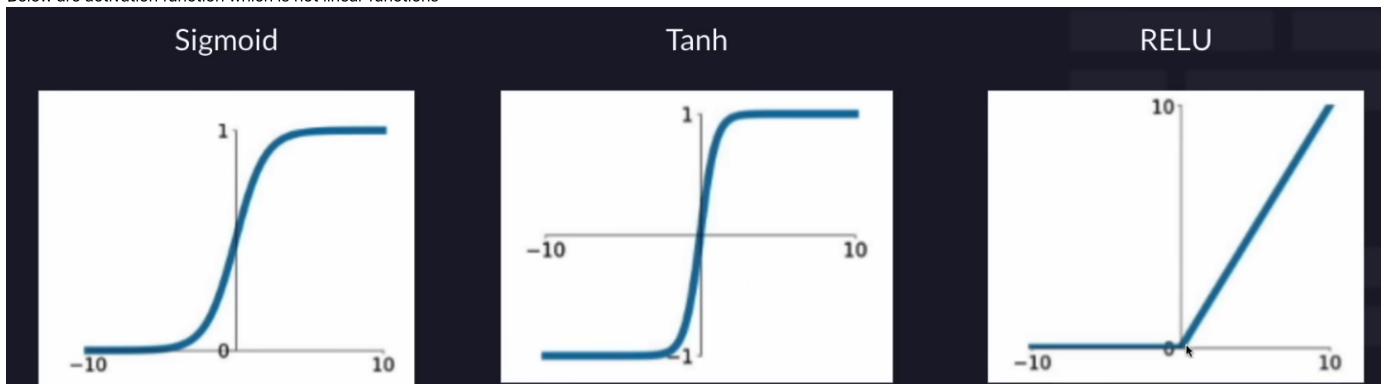
A neural network is a series of algorithms that endeavors to recognize underlying relationships in a set of data through a process that mimics the way the human brain operates



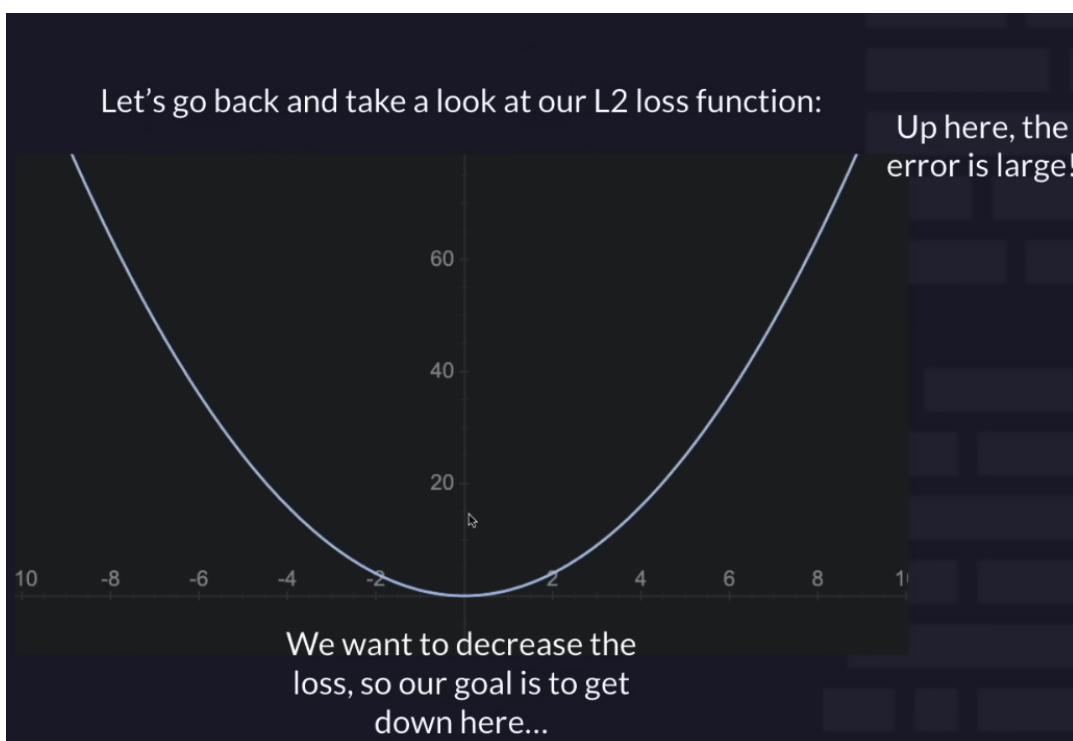
| With out activation funtion neuralnetwork is like linear combination of models

Activation Function

Below are activation function which is not linear functions

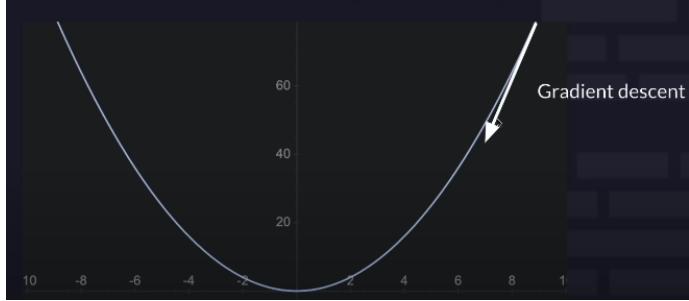


Lossfunction

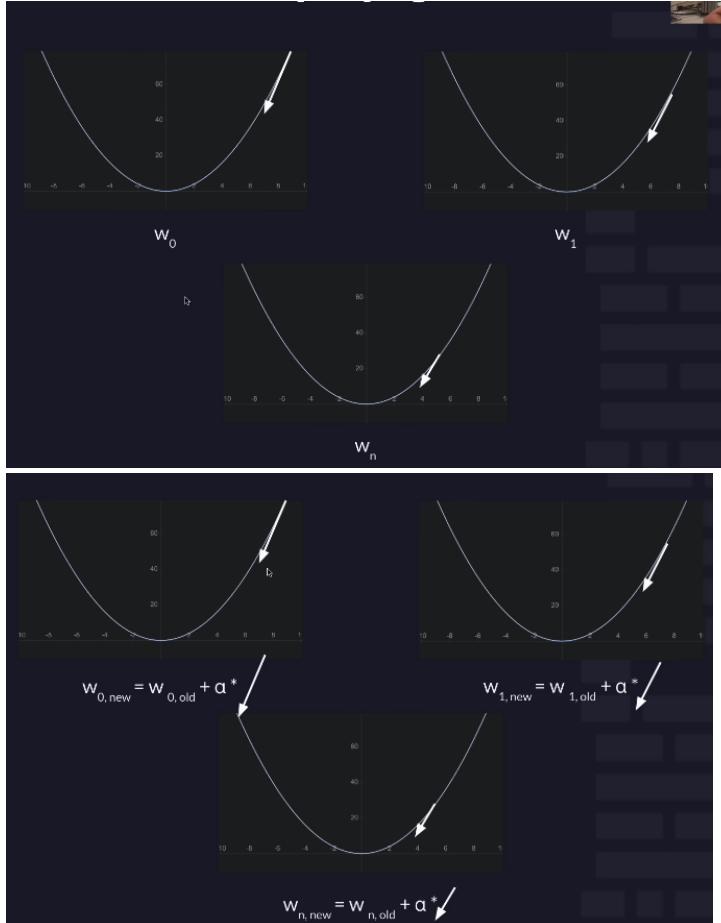


Since top right mean error is more, we use loss function to bring it down. Property to use it is call gradient decent

Let's go back and take a look at our L2 loss function:



Different weight have different loss value, which is call back propagation

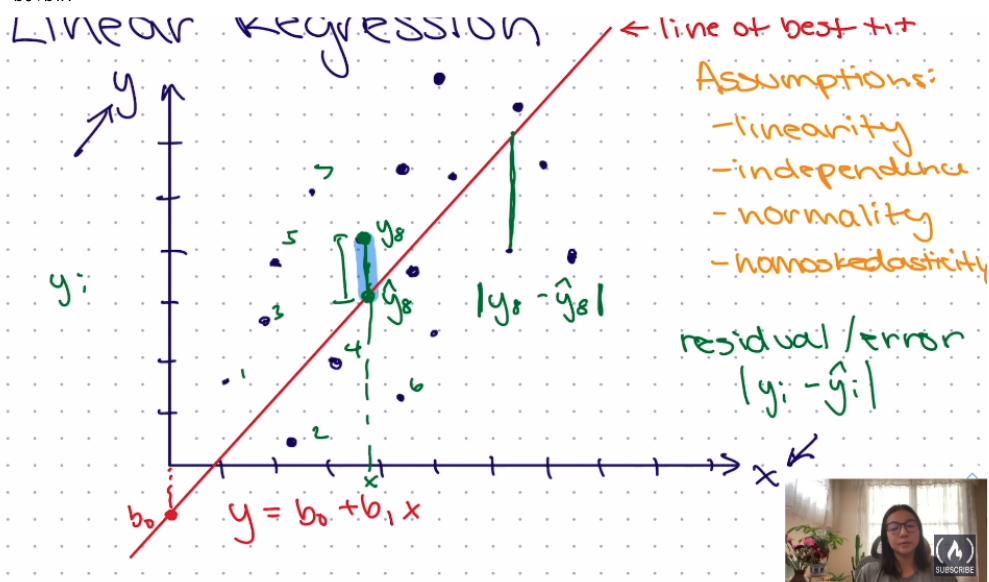


Tensor flow is predefined opensource code for neural network. Its a open source library to build models

Linear Regression

Simple linear regression

$$y = b_0 + b_1 x$$



Minimising error means minimising residual values. So we need to sum up all residual value and get the lowest sum.

Multiple linear regression

$$y = b_0 + b_1 X_1 + b_2 X_2 + b_3 X_3$$

Evaluating linear regression model

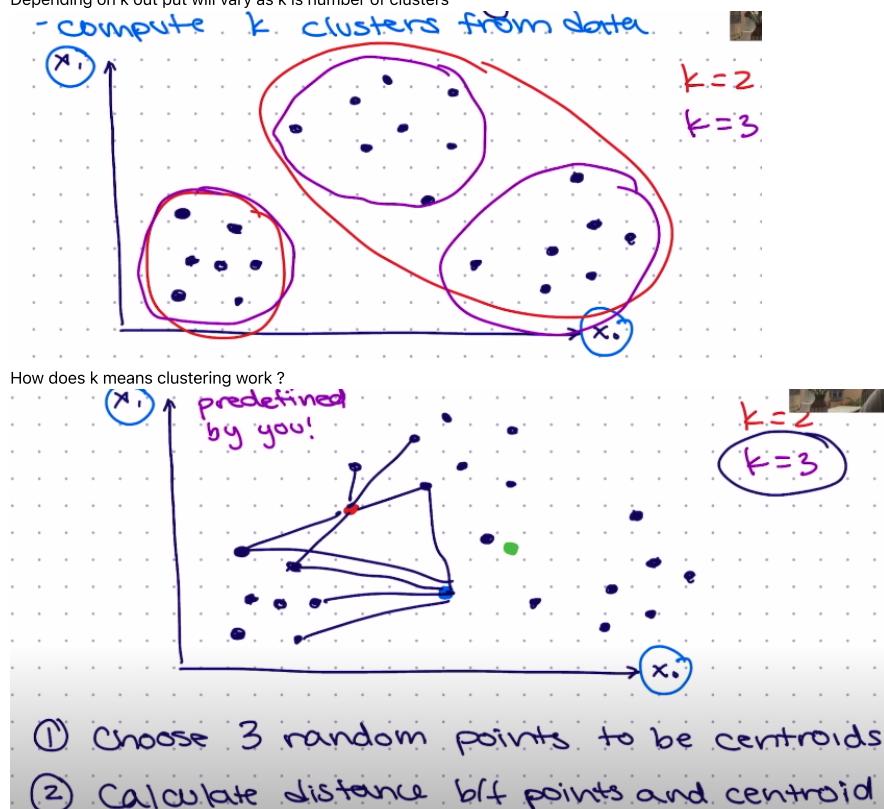
1. Mean Absolute Error (MAE)
2. Mean Squared Error (MSE)
3. Root mean square error (RMSE)
4. Co efficient of Determination (R^2)

We can use neural net (tensor flow keras) and calculate linear regression. Difference is in NN we are using back propagation, in regular linear regression we don't have back propagation

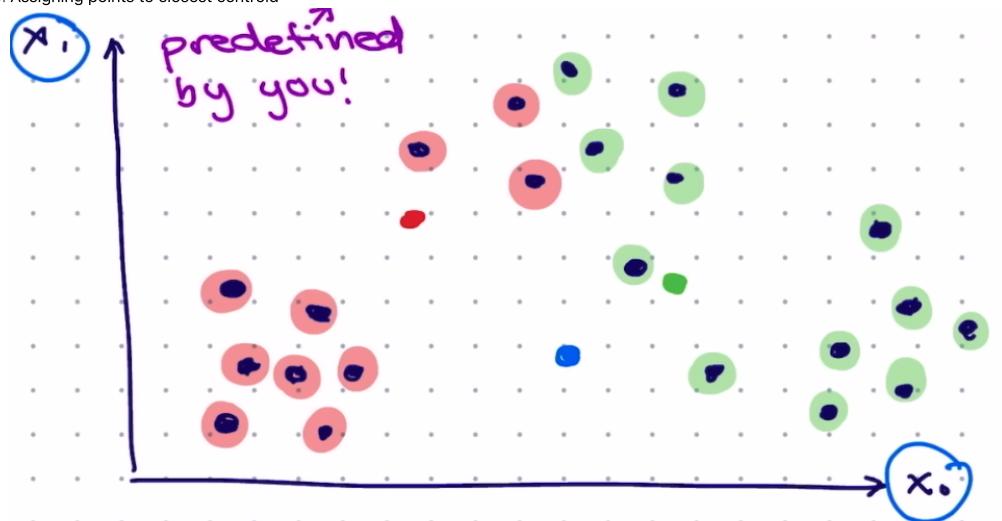
Unsupervised Learning

K-Means clustering

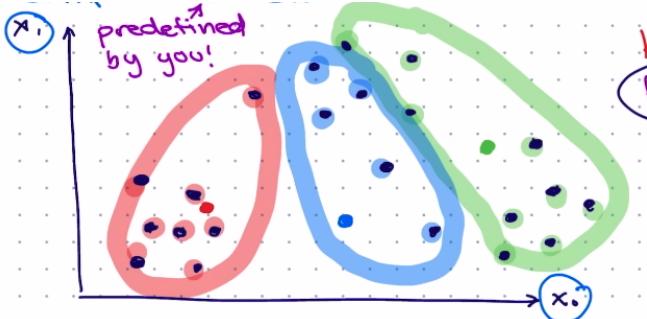
- unsupervised learning
- Compute k clusters from data
- Depending on k output will vary as k is number of clusters



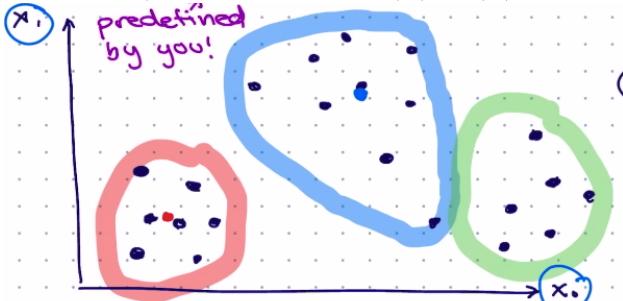
3. Assigning points to closest centroid



4. Compute new centroid and repeat step 2, we get new centroid



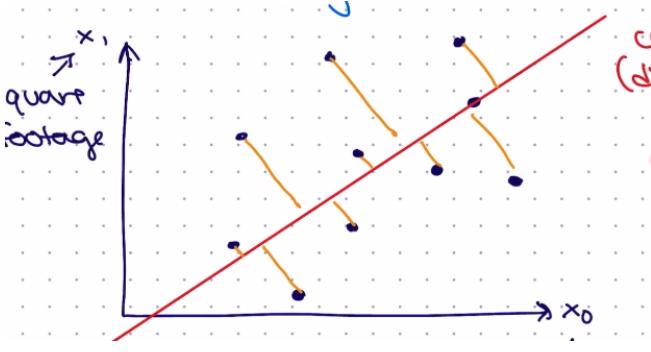
5. Process continues, recalculate centroid and cluster, repeat step 2,3



We stop when the points are stable and nothing is changing (centroid and cluster are same). This process is called expectation maximisation

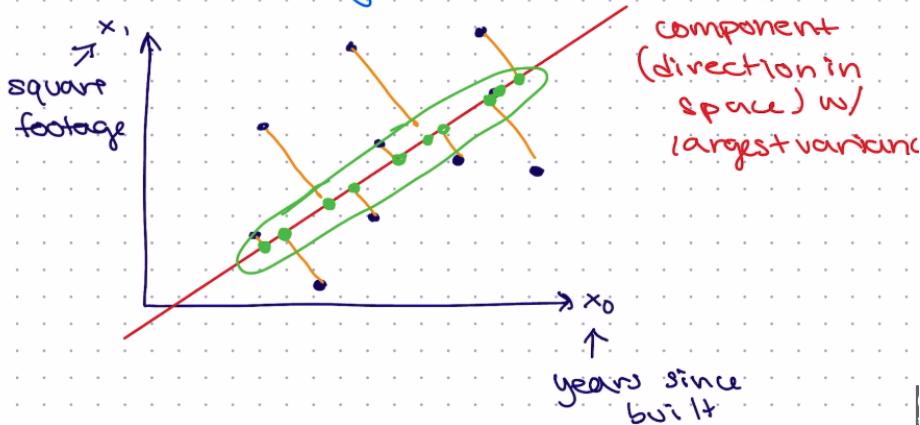
Principal component Analysis(PCA)

- Used for dimensionality reduction
- Technique used in Unsupervised learning
- Definition: Component direction in space with largest variance
- Looks similar to linear regression but not, In this we don't have y component



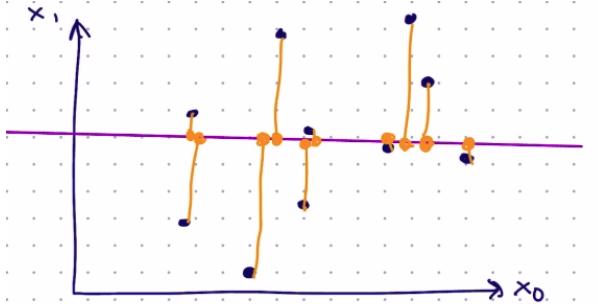
We started with 2d data set and end with 1d

Principal Component Analysis -dimensionality reduction

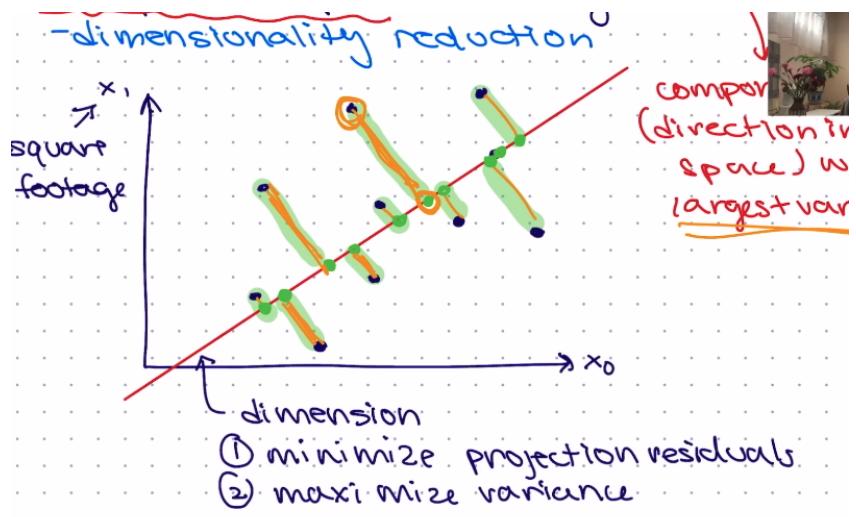


From above picture, we plotted blue ones(data set) in two dimensional graph. We calculated green dots which is one dimensional data set, Red line is the dimensional (directionality). Green dots are the one number to represent two dimensional points (blue)

If directionality/PCA changes (red line) Then it looks like in below image



In Linear regression we find the distance between y and x components, But in PCA we find the distance between x and projected point (dot on red line)



If we have 100 different features and then we need to take 5 features , we can calculate PCA and take top ones. Example if we have sphere and we have 2d area to plot , what is the best way to represent sphere ?? its circle. (Extract max info from many features/columns)