

Machine Learning

A Case study
Approach

(Supervised) Regression

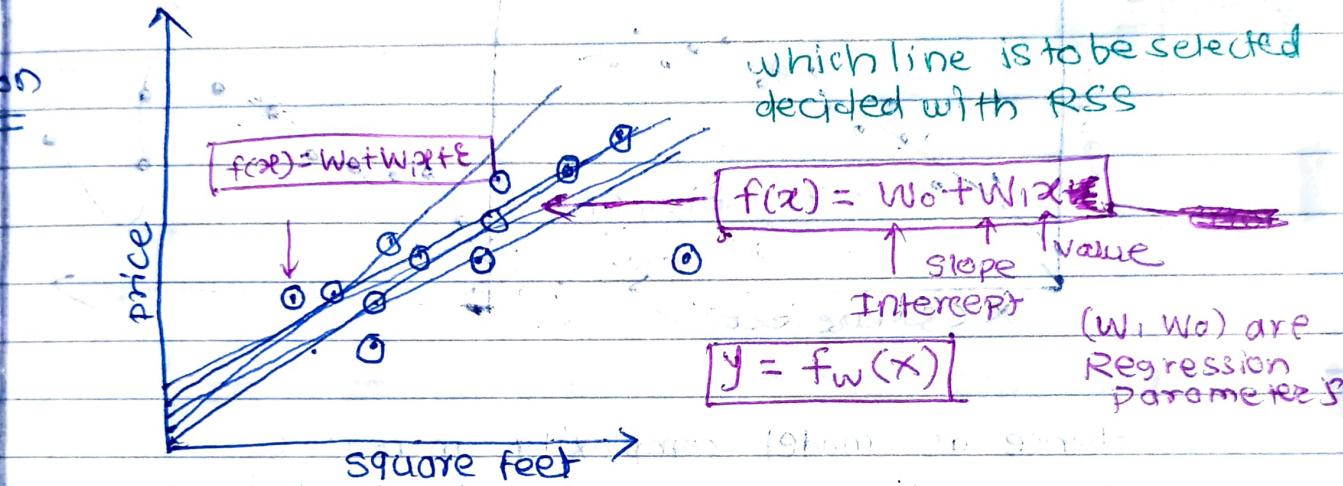
DATA

Date
Price

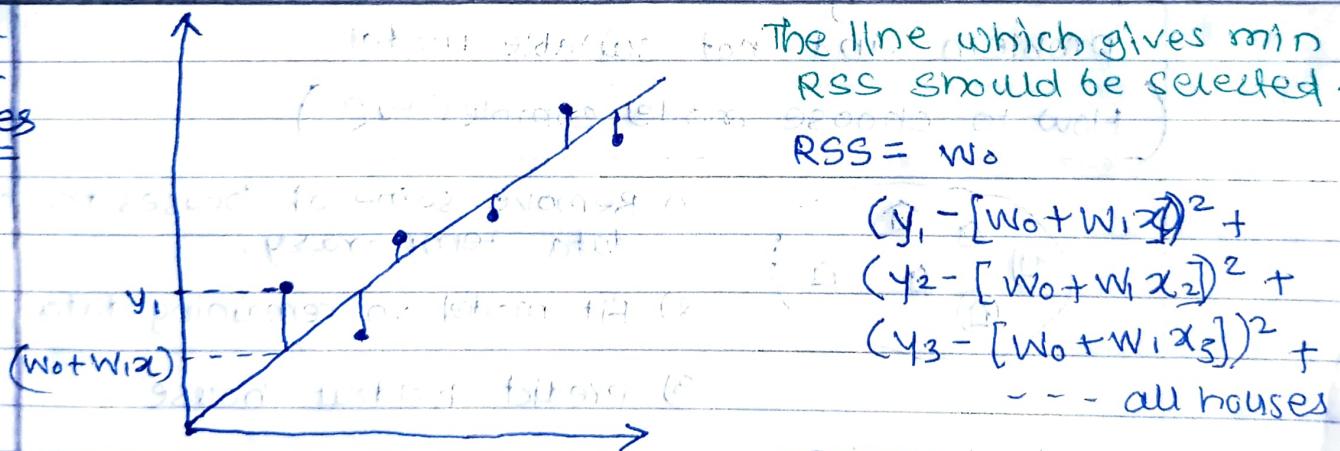
App

Predicting house price.

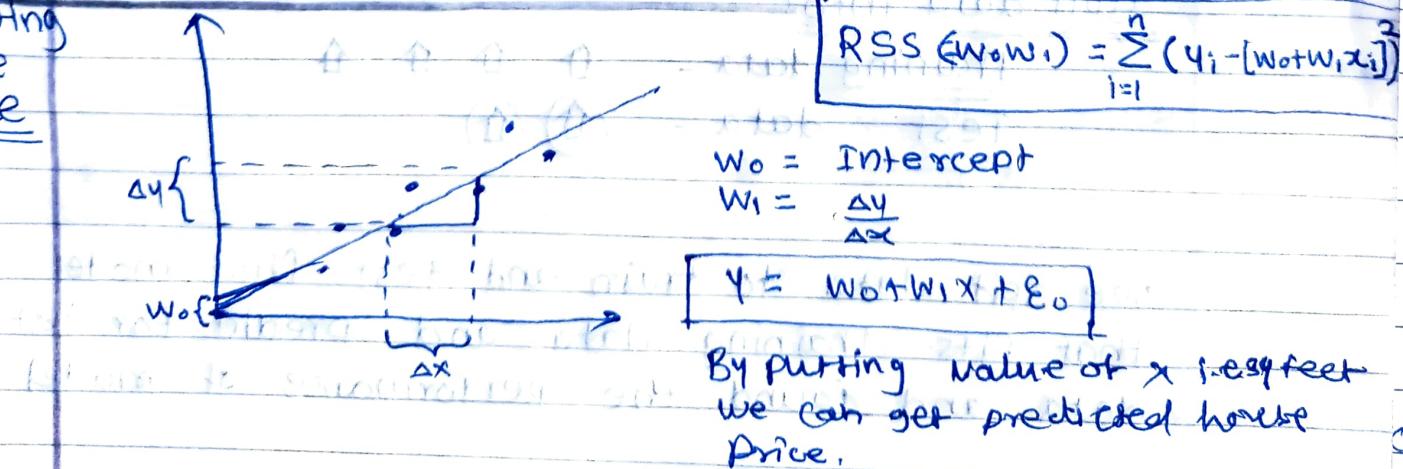
Linear Regression

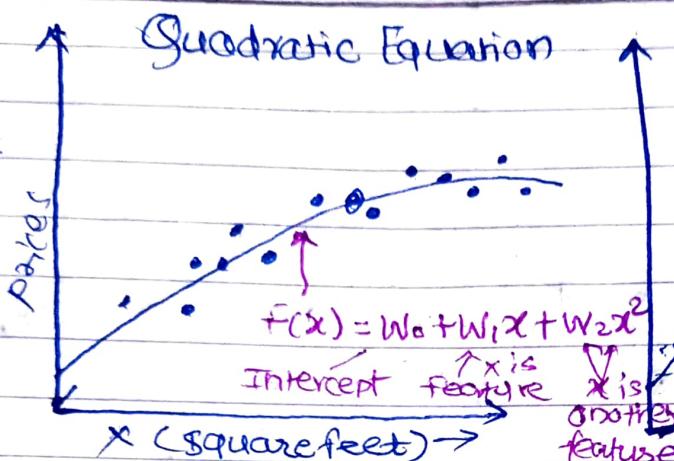


Residual Sum of squares

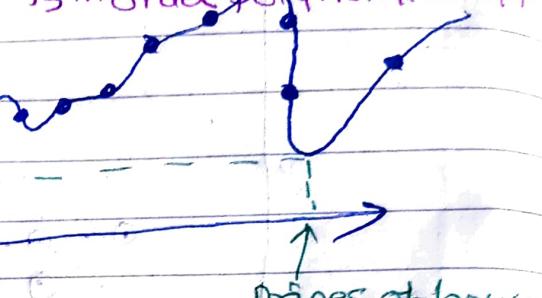


Predicting House price



Higher Order RegressionsOverfitting

overminimized RSS
13th order polynomial fit



Hence as model complexity increases i.e. overfitting occur and which leads to problem and not suitable model

Prices of house cannot be so few hence overfitting give problem.

(How to choose model complexity?)



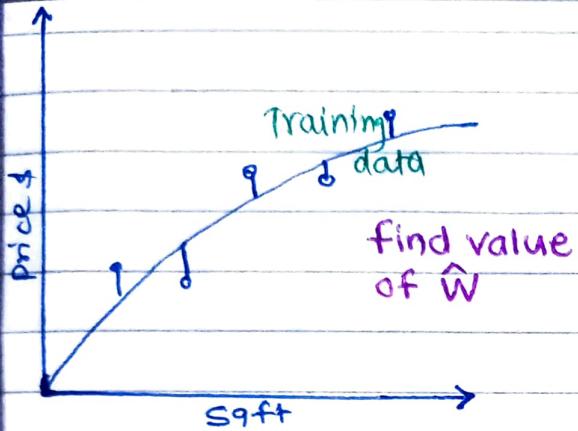
- 1) Remove some of houses from data temporarily,
- 2) Fit model on remaining data
- 3) predict heldout house

Split data into -

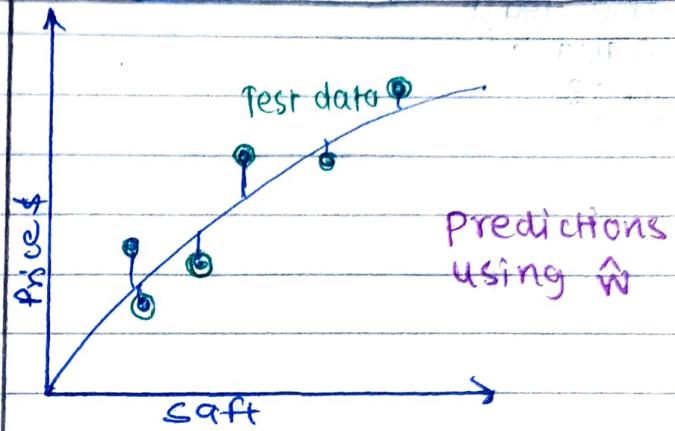
Training data - (Four house icons)

Test data - (Two house icons)

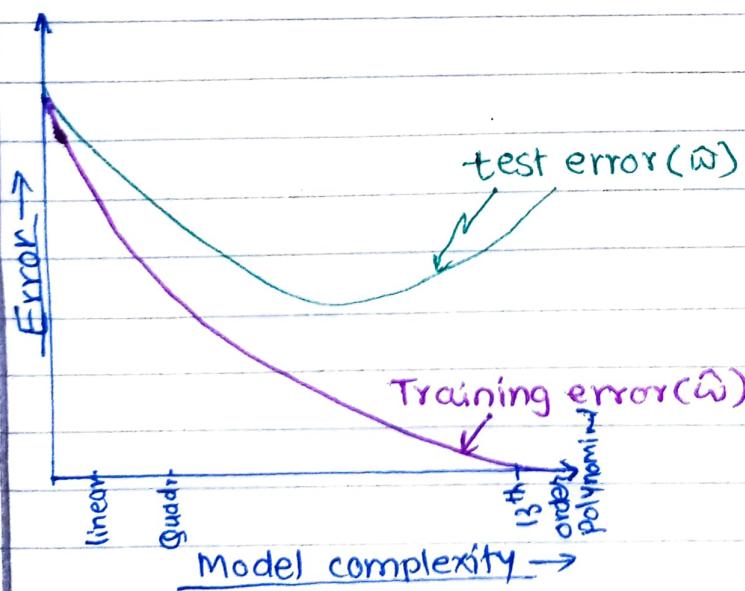
we split data to train and test find model that fits training data and predict for test data and found the performance of model

Training errorTraining error(W) =

$$\begin{aligned} & (\$_{train1} - f_w(\text{sq.ft}_{train1}))^2 \\ & + (\$_{train2} - f_w(\text{sq.ft}_{train2}))^2 \\ & \dots \text{for all data} \end{aligned}$$

Test errorTest error(\hat{W}) =

$$\begin{aligned} & (\$_{test1} - f_{\hat{w}}(\text{sq.ft}_{test1}))^2 \\ & + (\$_{test2} - f_{\hat{w}}(\text{sq.ft}_{test2}))^2 \\ & \dots \text{for all test data} \end{aligned}$$

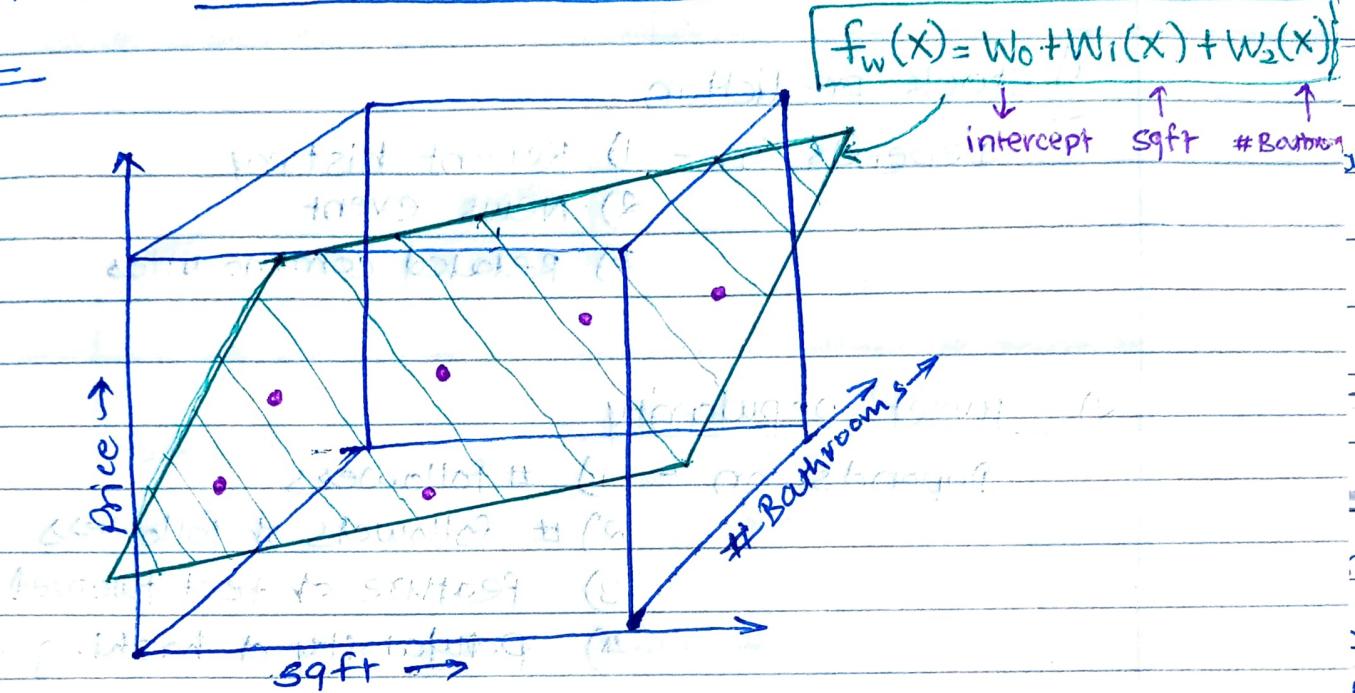
Variation of Test/Training Errors with model Complexity

- * As we increase model complexity we fit the training data more precisely as seen in overfitting for 10th polynomial hence training error decreases.

- * But with increase in complexity test error decrease to optimal and then increases due to overfitting.

Adding other features

we can add more features to model



Here we added no. of Bathrooms which cause model to have two features and which is given by hyperplane

we also can add many more features to model like for house price

- ↳ 1) sqft
- 2) # Bathrooms
- 3) # Bedrooms
- 4) Lot size
- 5) Year of built.

Other eg. where regression used

1) Stock prediction

- Depends on - 1) Recent history
2) News event
3) Related commodities

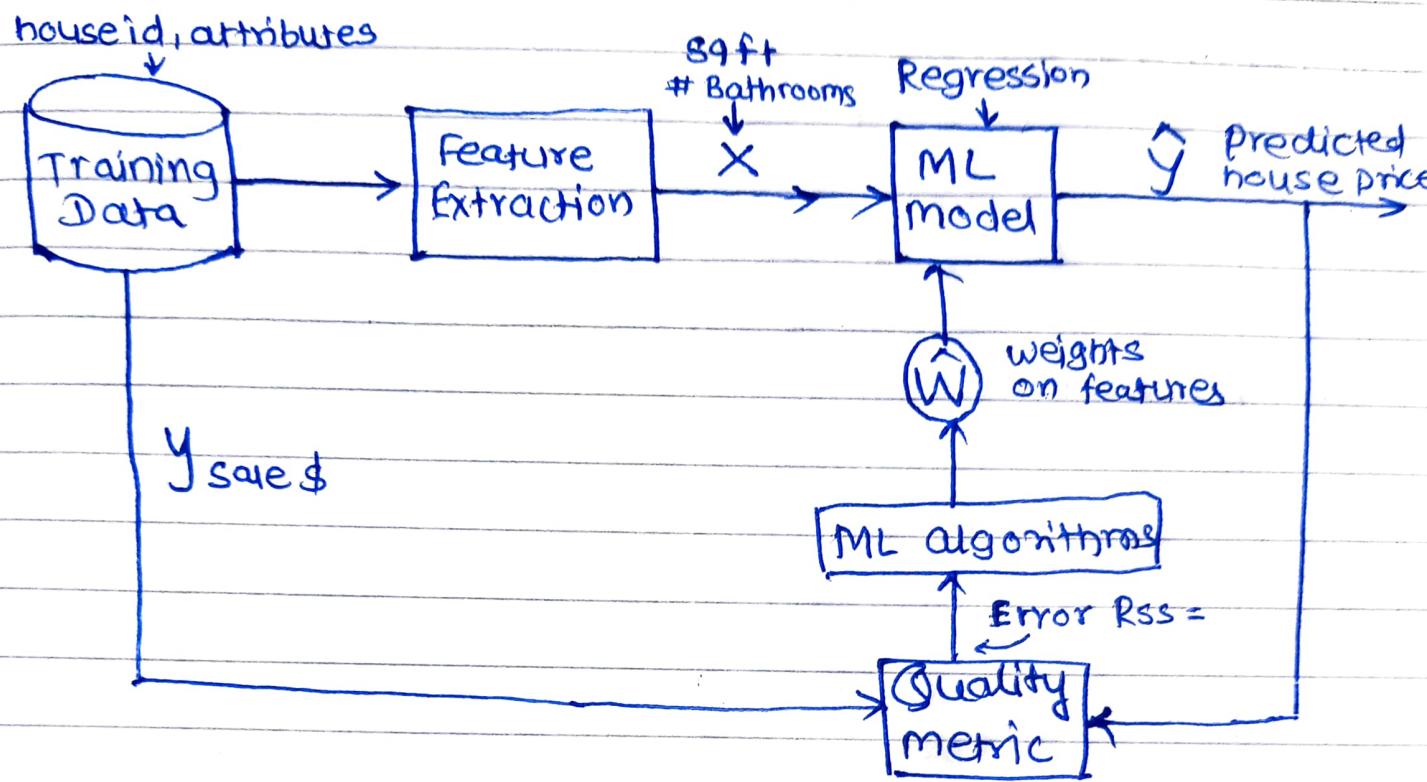
2) Tweet popularity

- Depends on - 1) # followers
2) # followers or follower
3) feature or text tweeted
4) popularity or hashing.

3) Temp of particular space in smart house

- Depends on - 1) thermostat settings
2) open/closed windows
3) vents
4) Temp outside
5) Time of day

Regression ML Block Diagram



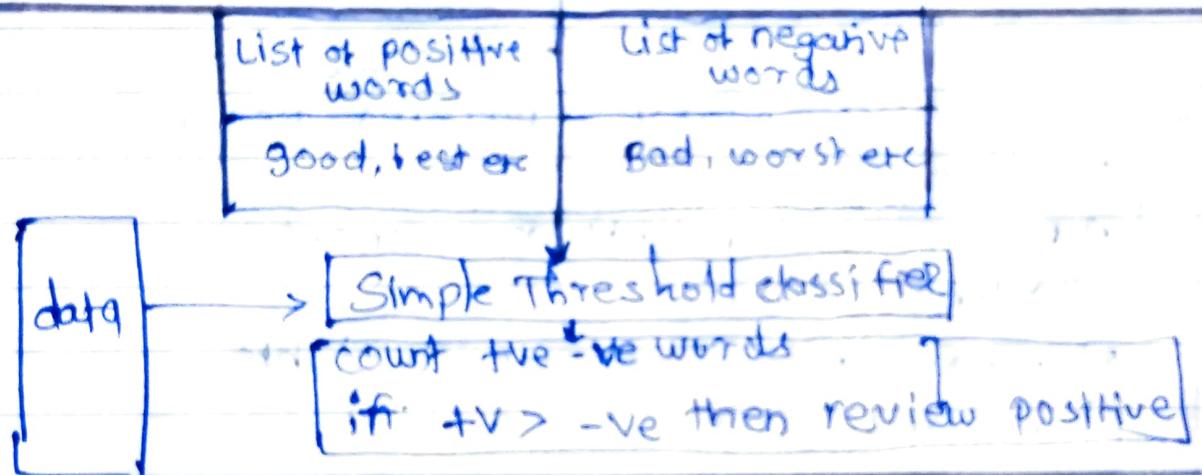
(Supervised)

Classification

IANI

- Application : spam filtering , suggesting the categories of article,
video etc , Image classification , disease classification
Reading our mind.

near classifier To classify review as sentence is positive or negative
we count all positive words negative words
if #positive > #negative review is positive



Simple threshold have problems like -

- 1) How to get list of +ve or -ve
- 2) words have diff. degree of sentiment
e.g. good < great
- 3) single words aren't enough we need to address it by more features.
e.g. good
not good

Hence we will use training data with adding weight to it as follow and examine the sentence as follow

$$\text{good} = 1.0$$

$$\text{great} = 1.5$$

$$\text{awesome} = 2.7$$

$$\text{terrible} = -2.1$$

$$\text{bad} = -2$$

food was good; taste was awesome
but service was terrible

$$\begin{aligned}\text{score}(x) &= 1.2 + 1.7 - 2.1 \\ &= 0.8\end{aligned}$$

As score > 0

Review is positive.

Decision Boundary

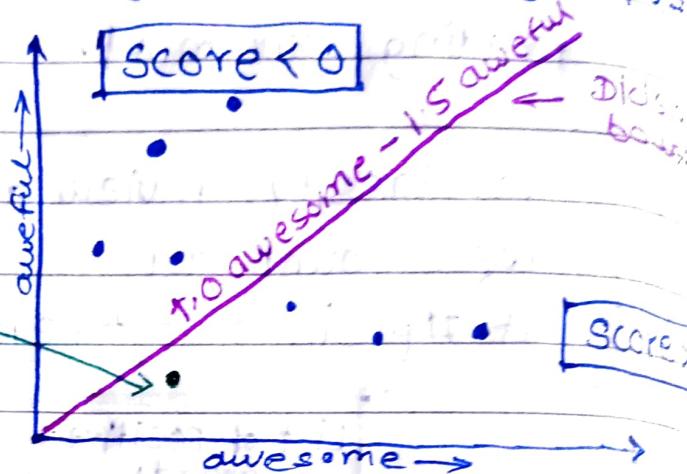
If only two words have non zero weight

words	weight
-------	--------

Awesome 1.0

Awful -1.5

$$\text{score} = 1.0 \# \text{Awesome} - 1.5 \# \text{Awful}$$



* Sushi was awesome
Food was awesome
But service was
awful

* Decision Boundary separates positive or negative predictions.

for linear

when two weights are non zero

↳ line (2D)

a) when three non zero weights

↳ plane (3D)

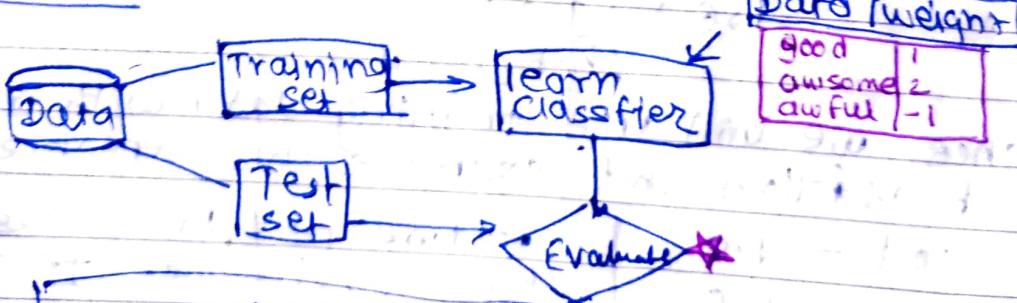
b) when many non zero weights

↳ Hyperplane

for general

complicated shapes.

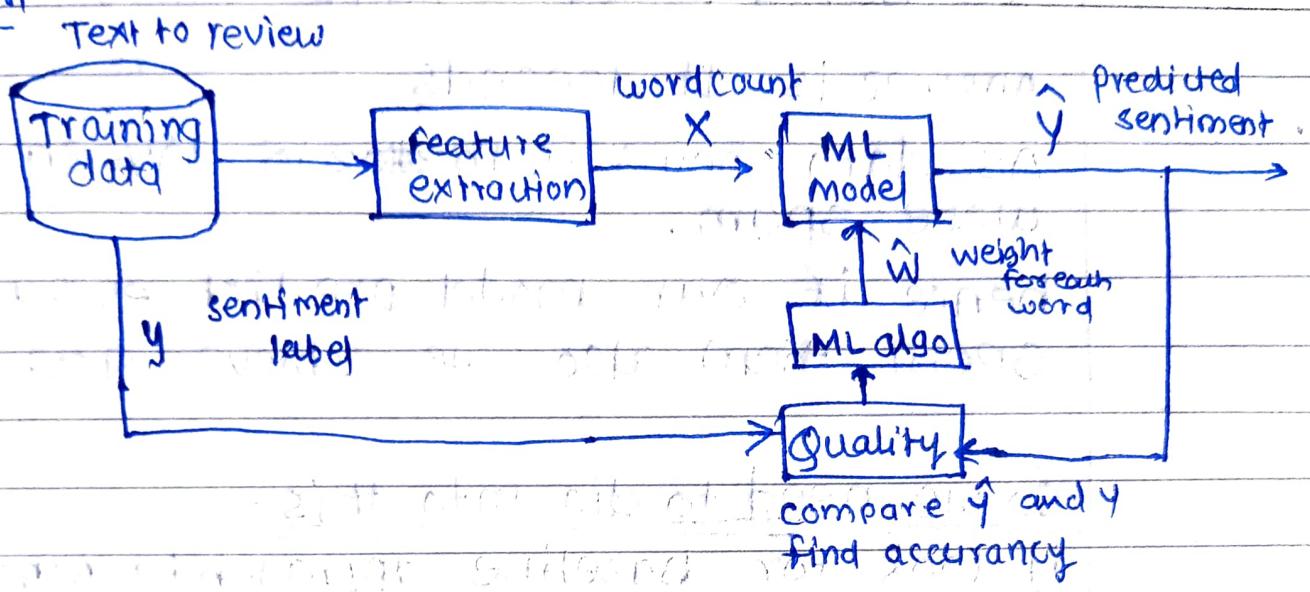
Training classifier



$$\text{Error} = \frac{\#\text{ wrong prediction}}{\#\text{ total sentence}}$$

$$\text{Accuracy} = \frac{\#\text{ Right Prediction}}{\#\text{ total sentence}}$$

$$\text{Accuracy} = 1 - \text{Error}$$

Block diagram

* We evaluate the model by finding error and accuracy.

What is good accuracy?

When we have Binary classification randomly selecting gives 0.5 accuracy

For k classes accuracy is

$$\text{accuracy} = \frac{1}{k}$$

So as with random guessing we get such accuracy we least, least, least beat the accuracy getting by random guessing.

e.g. If we have 3 classes we have

$$\text{accuracy} = \frac{1}{3} = 0.333 \text{ for random guessing.}$$

\therefore we must be having at least

$$\text{Accuracy} > 0.333$$

Is 90% accuracy is good?

Answer is It depends

As data shows from ~~2010~~ 2010 90% emails sent were spam.

hence if our model predict every email as spam, then also we achieve 90% accuracy.

so we need to dig into this -

- i) look for baseline approach i.e. random guess majority class i.e. in email case.
- ii) what is the impact of mistakes/errors if important enough email sent to spam if disease get undetected,

Types of mistakes

		Predicted label	
		+	-
True label	+	True positive	False negative
	-	False positive	True negative

* Different mistakes have different impact

- 1) True Positive } Accuracy
- 2) True negative }
- 3) False positive } Error
- 4) False negative } mistake

For email filtering

False negative \rightarrow spam email in inbox \rightarrow Annoying

False positive \rightarrow important emails lost \rightarrow Harmful

for disease detection

False negative \rightarrow Disease untreated \rightarrow Harmful
 False positive \rightarrow wasteful treatment \rightarrow less harm

- * Hence we need to look for the cost of wrong prediction

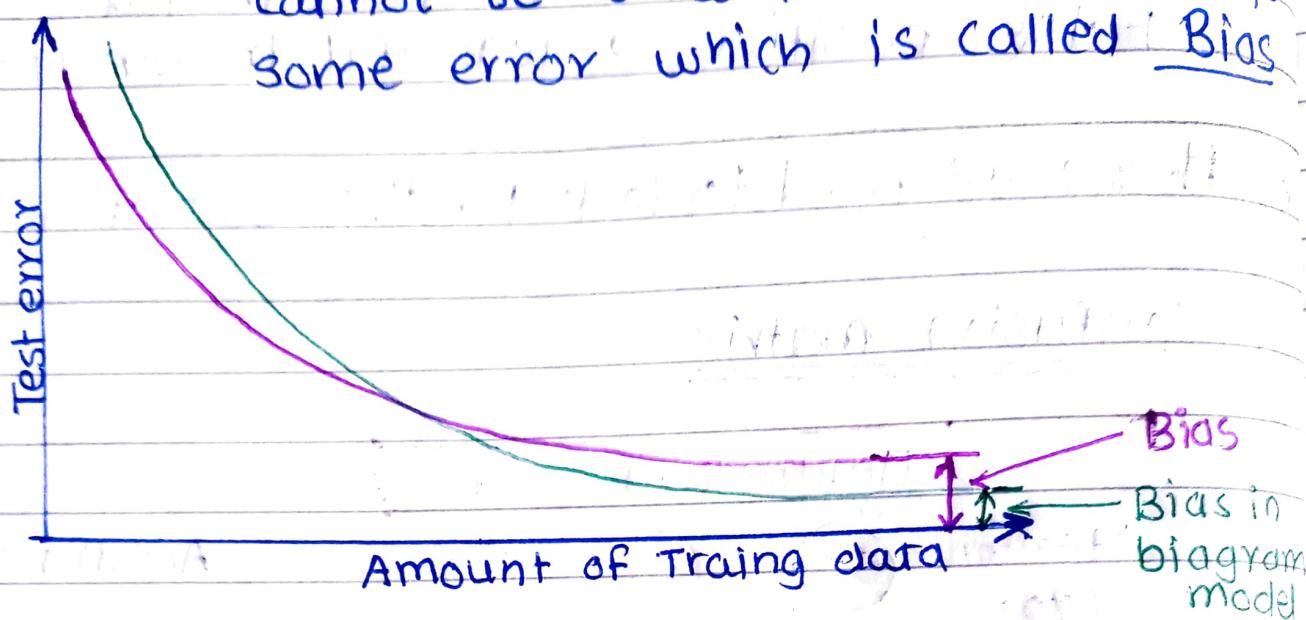
Confusion matrix

	Healthy	Cold	Flue
Healthy (70)	60	8	2
Cold (20)	4	12	4
Flue (10)	0	2	8

$$\text{Accuracy} = \frac{60+12+2}{100} = 0.80$$

Here we need to look for the wrong prediction making what impact and what is consequences which will cost most and prepare model accordingly.

Bias - Even if we provide infinite data error cannot be 0 and there is always some error which is called Bias



In our model of sentiment classifiers single words are OK but when it comes to combination of two words model makes mistake.
e.g. Sushi was not good.

Hence we need more complex model as

word	weight
good	+1.5
Not good	-2.1

← Biagram model

↓
less errors
more accurate
less Bias

↓
But Not completely
accurate

↓
Still have some
Bias.

(unsupervised)

Clustering and Retrieval

IANT

Date:

Page:

pp

To find similar articles like documents

- retrieval
search
- 1) we find distance between Query document and all other documents.
 - 2) Nearest neighbours returned.

Algorithms 1) 1 NN algorithm - most neighbour nearest only one

2) KNN algorithm - group/set of nearest neighbours

Nearest Neighbor Search

- 1) We get one article out of corpus
- 2) Specify distance of each article with every article and form distance matrix.
- 3) articles having lesser distances retrieved for query article.

Word count document Representation.

3	2	1	5	
---	---	---	---	--

the soccer messi or

We count # instances for each word and form vector for each documents and look for similarity of documents based on word counts vector products

Similarity for Doc 1 & 2 & 3

$$\text{Doc 1} \Rightarrow [1|0|0|0|5|3|0|0|1|0|0|0] \quad \text{similarity}_{\text{Doc1} \times 0} \\ \text{Doc 2} \Rightarrow [3|0|0|0|2|0|0|1|0|1|0|0] \quad = 13$$

$$\text{DOC 1} \Rightarrow [1|0|0|0|5|3|0|0|1|0|0|0] \quad \text{similarity} = \text{Doc1} \times 0 \\ \text{DOC 3} \Rightarrow [3|0|0|0|6|4|0|0|2|0|0|0] \quad = 45$$

hence similar document for 1 is 3. Hence

Issues with word count

- 1) As there are common words like 'the', 'or', and which are present in all documents
- 2) As short document have less count of words as compare to large documents hence it hampers the results.

2) To handle short and long documents we use method
Normalization

DOC 1 \Rightarrow

1	0	0	0	5	3	0	0	1	0	0
---	---	---	---	---	---	---	---	---	---	---

$$\text{normalization} \Rightarrow \sqrt{1^2 + 5^2 + 3^2 + 1^2} = 6$$

Normalised DOC 1 \Rightarrow

$\frac{1}{6}$	0	0	0	$\frac{5}{6}$	$\frac{3}{6}$	0	0	$\frac{1}{6}$	0	0
---------------	---	---	---	---------------	---------------	---	---	---------------	---	---

We use this normalised vector for document retrieval as regardless of document size it treat document counts equivalently.

1) To handle frequency of rare words & common words use method called -
TF-IDF - Term Frequency - Inverse Doc Frequency

The most common words which comes in every document must be downgraded like 'the', 'and', 'also' etc.

Common locally
Rare globally

Rare words \Rightarrow words appear in few doc

common words \Rightarrow words appear in every doc

Common
globally

calculating TF-IDF

TF - Term frequency

100	5	15	10	1	1	1	1	1
the	messi							

IDF - Inverse document frequency

10	4	1	1	1	1	1	1
the	messi						

$$\log \frac{\# \text{DOC}}{1 + \# \text{Doc having word 'the'}}$$

$$\log = \frac{\text{large}}{1 + \# \text{large}}$$

* hence globally common words are downgraded and Rare words upgraded by TF-IDF

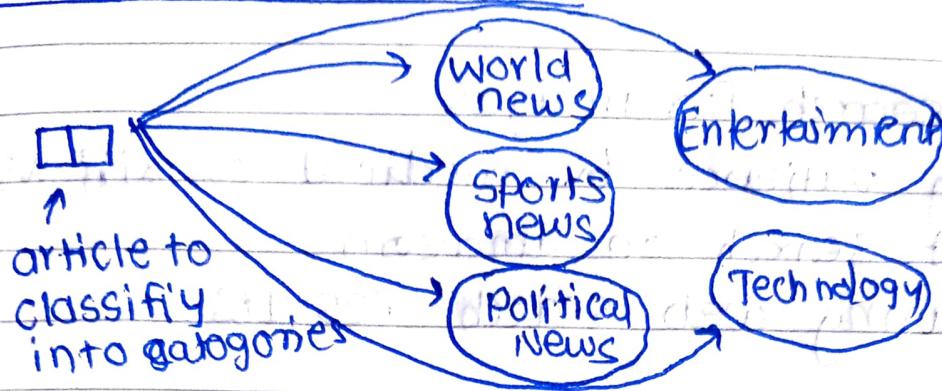
$$\text{Now } \underline{\text{TF-IDF}} = \text{TF} * \text{IDF}$$

$$\text{TF} * \text{IDF} = \boxed{0} \quad \boxed{20} \quad \boxed{10} \quad \boxed{10} \quad \boxed{1} \quad \boxed{1} \quad \boxed{1} \quad \boxed{1} \quad \boxed{1}$$

tf * IDF

use for Document retrieval.

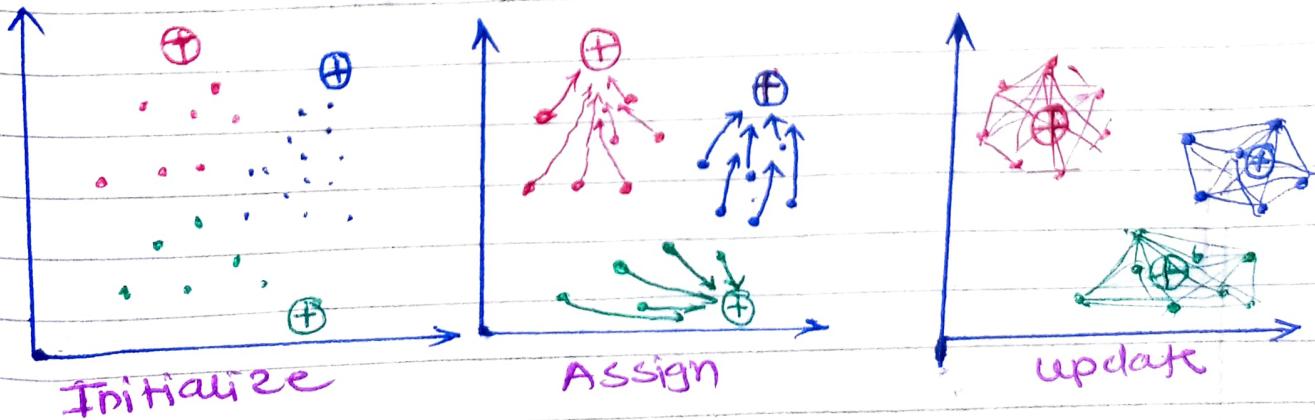
Multiclass classification (Supervised)



Here we need to give labels to clusters which causes it becomes supervised learning.

K-means Algo

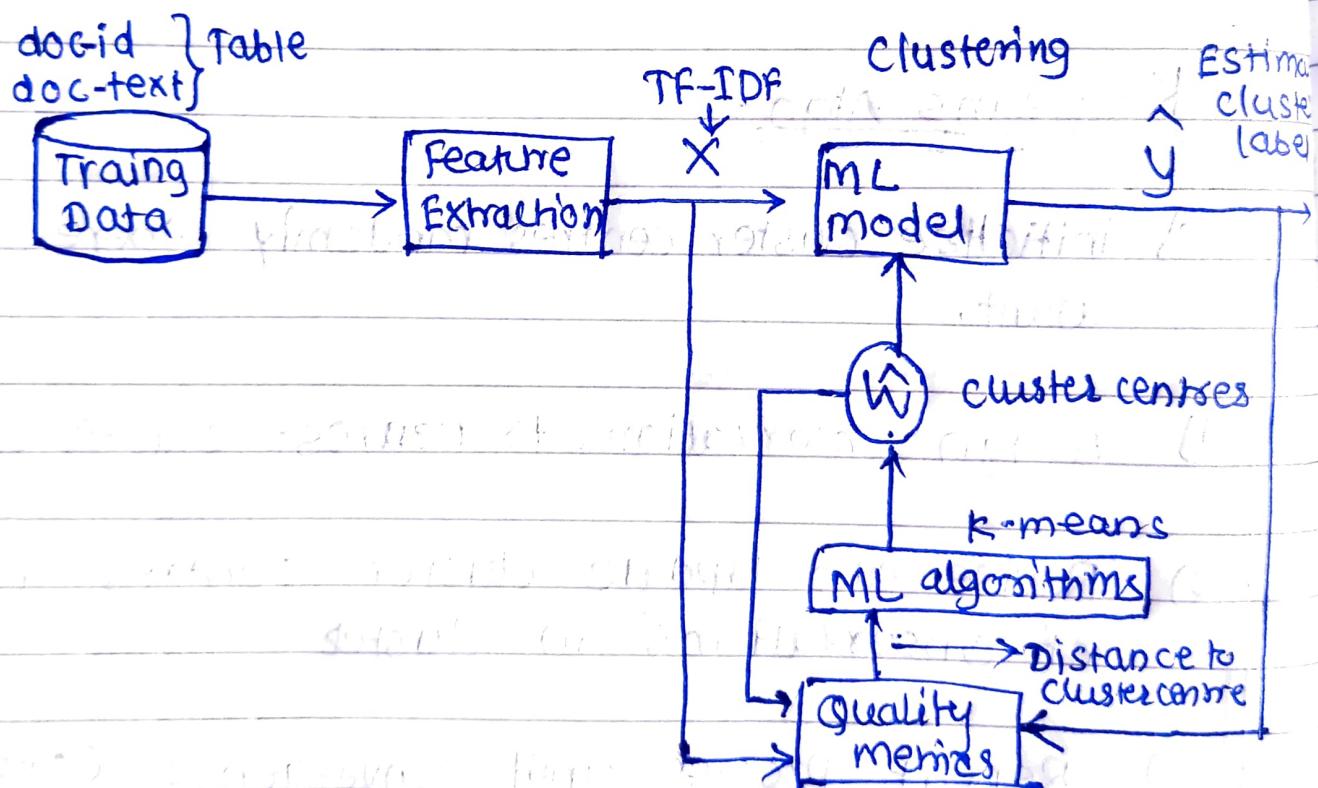
- 1) initialize cluster centres randomly where we want.
- 2) Assign observations to nearest cluster centre
- 3) Revise and update cluster centres as mean of observations in cluster.
- 4) Repeat process until convergence (step 1 & 2)



Examples

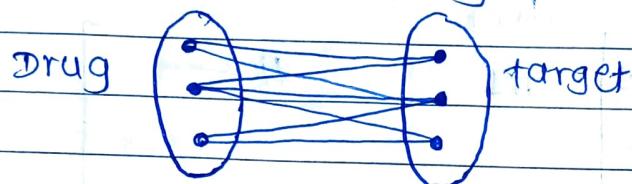
- 1) Google search for images
- 2) Grouping patients by medical conditions
- 3) Product search on Amazon
- 4) Structuring web search results

Clustering Block diagram



Recommender Systems

Application - Youtube, Netflix, Amazon, music, friend suggestions, Drug target interactions i.e. if one medicine is useful in treatment of one disease it can be useful treating some other diseases.

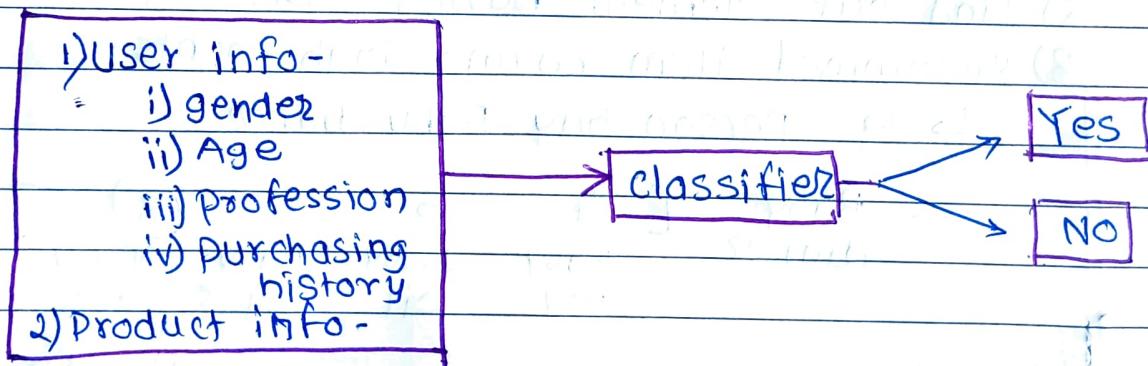


Building Recommender System

Approach 0 - Based on popularity products are ranked and suggested accordingly (global popularity) which has limitations for personalised suggestions.

Approach 1 - classification model.

- To find whether person would like or not purchasing the product.



Limitations - features like user info, product info may not be available.

Approach 2 - Collaborative filtering -

- people brought this also buy -

Co-occurrence Matrix

		product 1	2	3	4	
		product 1	#	#	#	No. of people purchasing some i.e. + products
matrix	C =	2	#	#	#	Product 1 & Product 3
		3	#	#	#	
		4	#	#	#	

given matrix

is a symmetric

matrices

Hence no. people buying product 1 and no. people buying product 3 is same

$$C_{ij} = C_{ji}$$

To recommend the product

- 1) we go to product row
- 2) find out element having most/largest counts
- 3) Recommended item having high counts.

↳ e.g. Person buy Baby diaper

Baby diaper	[4 5 7 8]
Baby toy	
Baby oil	
Baby wipes	
Baby food	

Limitations -

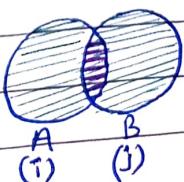
But this causes recommendation

which is based on popularity and drown down other effects hence we need to normalisation of popularity as follows -

Normalisation -

Purchased i and j

Purchased i or j



this item is more likely to be recommended

Weighted average of purchased items.

User bought items { diaper, milk }

$$\text{Weighted score for (Baby Wipes)} = \frac{1}{2} (S_{\text{babywipes, diaper}} + S_{\text{babywipes, milk}})$$

- 1) We look for purchased history i.e. diaper, milk.
- 2) We look for row of this products
- 3) If any item is to be recommended from cofactor matrix we look for it for e.g. baby wipe
- 4) We find # People purchased baby wipe & baby diaper and # people purchased baby wipe & milk and get average for weighted score.

Limitations

It doesn't look for other features like time of day, user's information & product information.

Cold start problem - When there is no history.

Approach 3 - Discovering hidden structure by matrix factorization.

Let some users have watched some movies and rated those movie we have matrix with user and their rating for movie but not every movie rated by every user and this blank space can be filled up and based on that movies recommended.

1	4	5	
2			
3	3	5	
4		2	
5		4	1
6	5		
	movies	1	2

Now white spaces filled up by looking for movies user rated and predicted the rating for unrated movie Completed matrix used further for Recommendation

Let user 1 rated movie 4 as 5 rating.

movies [2 3 4 1]
 Action Romance Comedy drama

vector of feature and their weight is known

User -1 [2 3 4 1]
 Action Romance Comedy drama

vector of feature and their weight of that liked by user

for movie Rating of unknown movie

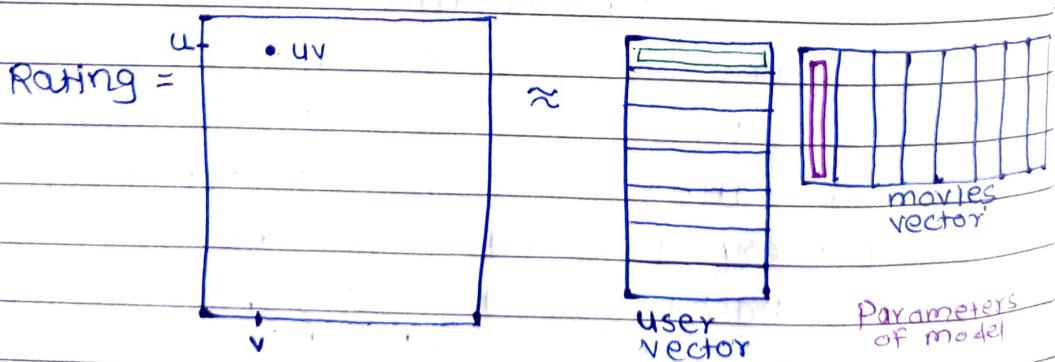
	Action	Romance	comedy	drama
movie 4	1	3	4	5
user 1	2	3	4	1

Rating that user may give is decided by product of vectors.

$$\begin{aligned} \text{Rating for movie 4} &= \frac{\text{movie4} \cdot [1 3 4 5]}{\text{user1} \cdot [2 3 4 1]} \\ &= \frac{1 \times 2 + 3 \times 3 + 4 \times 4}{2 \times 1 + 3 \times 4 + 4 \times 1} \\ &= 3.2 \end{aligned}$$

Based on the this white spaces filled.

Predictions in matrix form.



When we get product of User vector and movies vector we get all the values of matrix.

$$\text{User vector} = Lu = [\text{Action Romance comedy drama}]$$

$$\text{Movie vector element} = Rv = [\text{Action Romance comedy drama}]$$

$$\text{RSS (LR)} = \sum_{\text{all rated movies}} (\text{Rating given} - [\text{L}_u \text{R}_v])^2$$

↑ ↑
actual Rating Predicted Rating

Limitation

cold start problem.

 ↑ ↑
 New movie New user

