

## Naive Bayes Model

relies on the Bayes theorem

$$P(Y|X) = \frac{P(X|Y) \times P(Y)}{P(X)}$$

$$P(Y | w_1, \dots, w_n) = \frac{P(y) \times \prod_{i=0}^n P(w_i | y)}{P(y) \times \prod_{i=0}^n P(w_i | y) + P(\bar{y}) \times \prod_{i=0}^n P(w_i | \bar{y})}$$

Prior probability

Likelihood  
of  
being  
spam

Frequencies

can be done simply by counting  
the words.

But a better way is by doing

TF-IDF → Inverse Doc. Frequency:  
Measures the weight of a word based  
on its corpus

Term Frequency

## Counting Words

$$TF - TDF = TF \times IDF$$

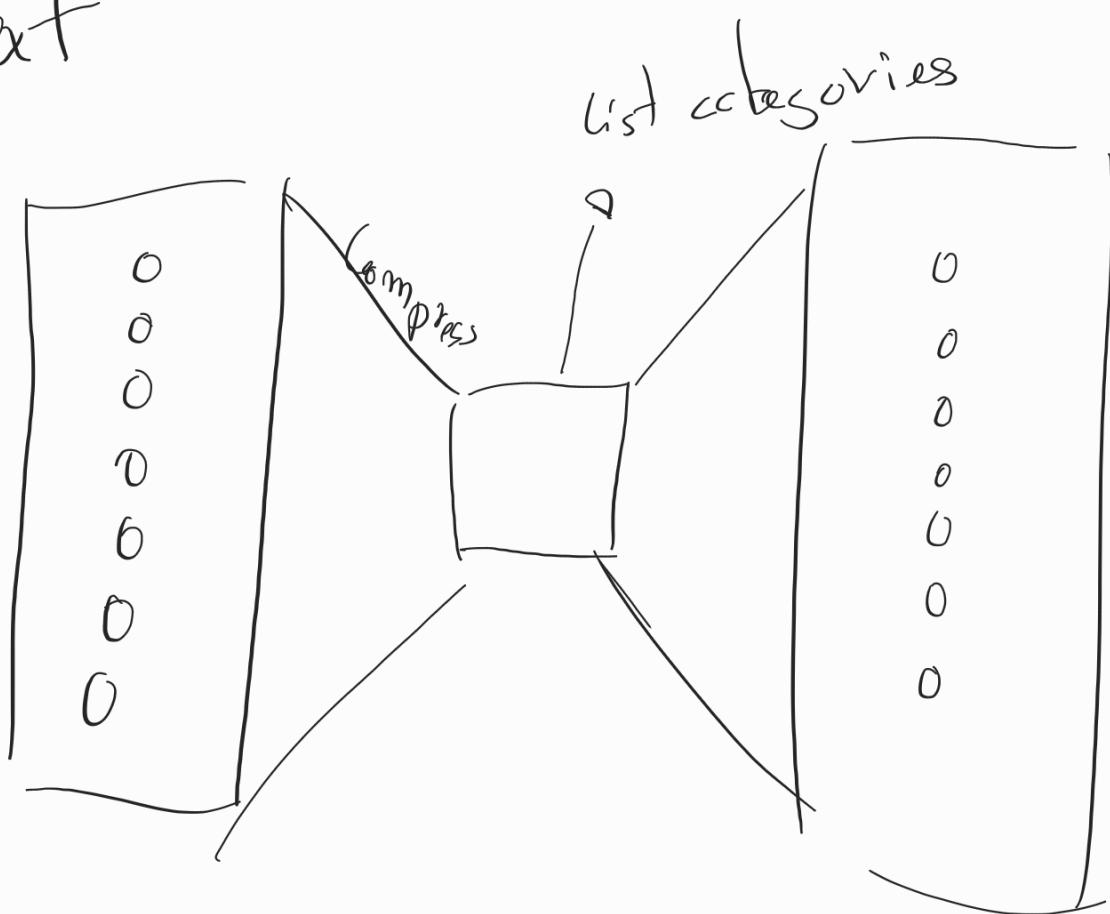
M - Ram J

11

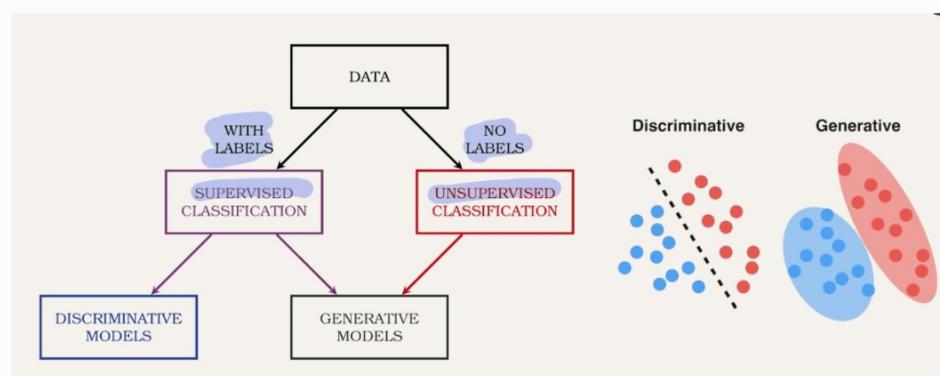
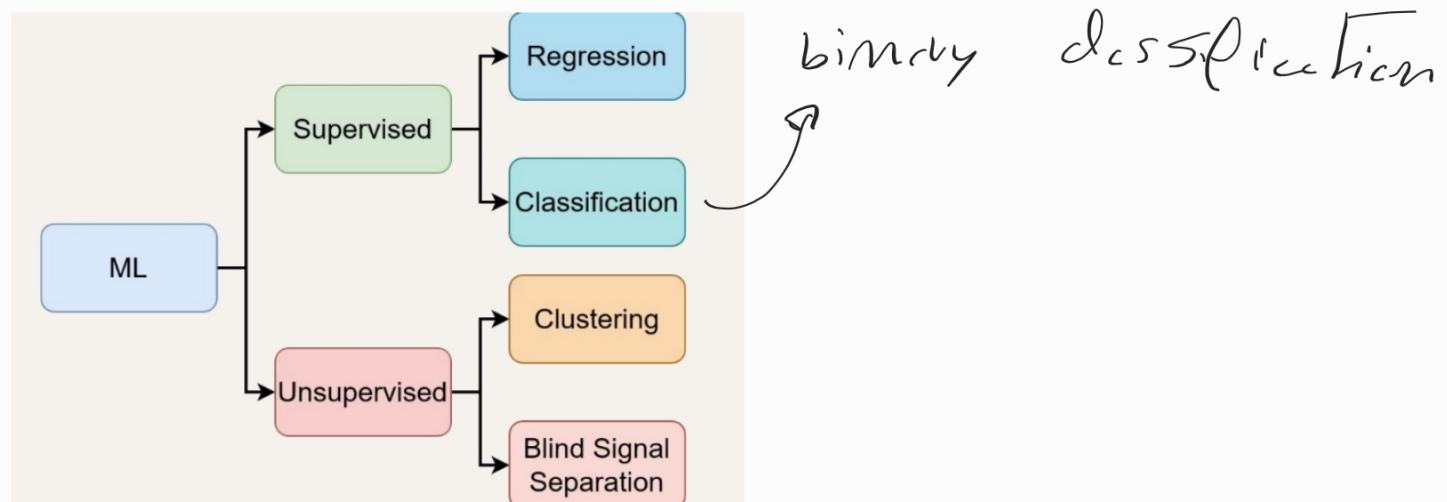
/1

The diagram illustrates the concept of 'Units of Words'. The words 'I am a human' are written in black cursive. Blue curved lines above the words group them into four distinct units: 'I', 'am', 'a', and 'human'. Red curved lines below each of these four units further divide them into smaller units: 'I' into 'I', 'am' into 'am', 'a' into 'a', and 'human' into 'human'. A single yellow bracket is positioned below the entire phrase, grouping all four blue-underlined units together as a single, larger semantic unit.

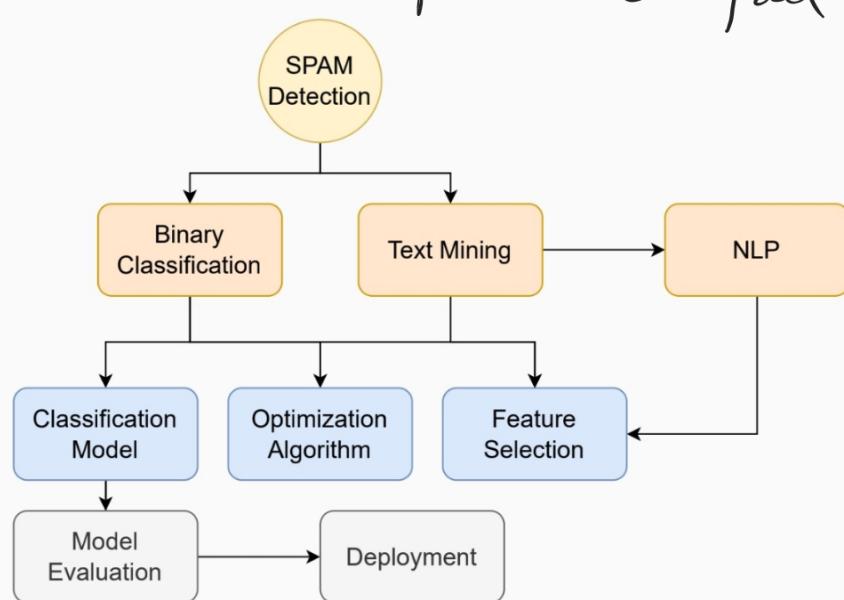
Fast fast



ML → improve its own performance  
→ programmed to learn through trial and error



Limitations → biased assumptions exist  
↳ simplification of reality  
↳ Assumptions can fail



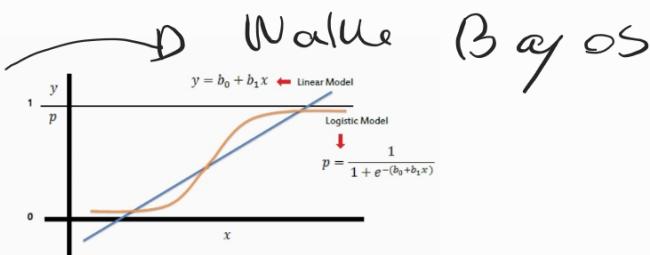
Binary classification → 1 or 0 → spam or not  
→ task of classifying the elements of a set into 2 groups

Text mining → make high quality info from text  
 ↳ uses bag model (bag of words)

NLP → method to see computers ability to understand text

$$P(H|E) = \frac{P(E|H) * P(H)}{P(E)}$$

Likelihood of the Evidence given that the Hypothesis is True  
 Prior Probability of the Hypothesis  
 Posterior Probability of the Hypothesis given that the Evidence is True  
 Prior Probability that the evidence is True



$$\text{acc} = \frac{\text{TP} + \text{TN}}{(\text{TP} + \text{TN} + \text{FP} + \text{FN})}$$

bow vs word embeddings → ordered

↑ struggle on finding center  
 ↓ b.d semantically  
 ↓ n. relationship b/w words

→ similar clo closer in the embedding space  
 unseen words are connected more easily

NB  
 text-mining  
 Tokenization → word → token of → sentence  
 sentence → token of → paragraph

lemmatization → rocks → rock  
 better → good

Simplifying word to is "base"

Term frequency  $\rightarrow$  count word, more effective TF-IDF

TF-IDF  $\rightarrow$  How frequent a word is on a document  
 $\rightarrow$  How relevant a word is on the doc

$\rightarrow$  How many times appears on the doc +  
the inverse doc freq  $\rightarrow$  which is how rare or  
common a word is, closer to 0 the more  
common

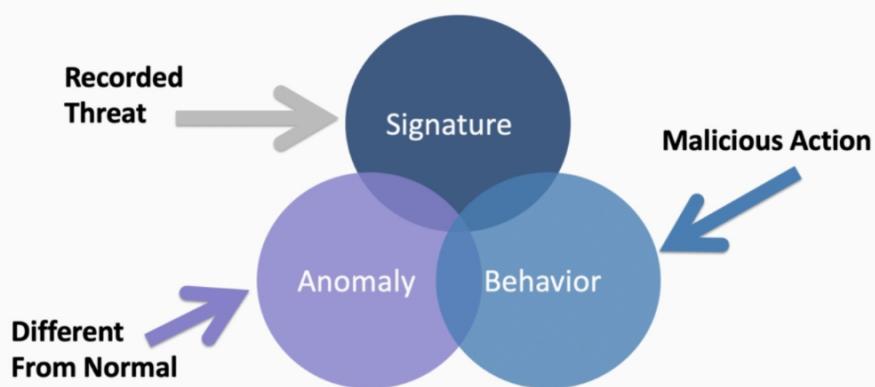
like 0.1 common ->  
never bits of 0.1 I think

Partial swarm  $\rightarrow$  optimization algorithm  
 $\rightarrow$  find min and max

Gradient descent  $\rightarrow$  opt alg  
 $\star$  find min with derived func

Linear Regression  $\rightarrow$

Logistic Regression  $\rightarrow$  Logistic curve is the inverse,



1.0 0.5 0.0

Signature  $\rightarrow$  compare hashes of virus for example or for  
similar ones

Behaviour detection  $\rightarrow$  malicious behaviour  $\rightarrow$  similarities are  
found alarms  $\rightarrow$  on

Anomaly Detection  $\rightarrow$  zero days, threat intelligence  
 $\rightarrow$  outliers, find alarm that differ from  
norm

fraud  $\rightarrow$  outliers  $\rightarrow$  noise  $\rightarrow$  fake arrows

Point anomalies ( $\rightarrow$  from the rest)

contextual  $\rightarrow$  abnormality is context specific

collective anomaly  $\rightarrow$  Group of outliers

clustering  $\rightarrow$  find a structure with meaning

$\hookrightarrow$  Density based - clusters as a dense region

$\hookrightarrow$  hierarchical Based  $\rightarrow$  tree type, new one formed  
using old ones

$\hookrightarrow$  Partitioning

BSS  $\rightarrow$  Blind Signal Separation, separation of signals  
in  $\hookrightarrow$  mixed signals

PCA  $\rightarrow$  principal components, makes indices

ICA  $\rightarrow$  Independent //, , // //

Auto encoders  $\rightarrow$  Giving natural def if can  
reconstruct it if the def is complex & it will

Use a hard line to dropping, making it complex

Missing values

- Replace "" with the mean of the column
  - Can only be done on numerical values
- Computer

Encode

- Hot encoder
- Ordinal encoder

Feature scaling

Obj has 100 0 vs 5 it will  
be more weight on 100 so we scale them  
so it doesn't give importance to one over the other  
feature.

Confusion matrix

	1	0
1	TP	FP
0	FN	TN

Accuracy

Quants uses accuracy

Precision  $\rightarrow$  good when FP high

Acc of positive predictions  $\rightarrow$   $TP/(TP+FP)$

Recall  $\rightarrow$  good when FN high

$$\frac{TP}{TP + FN}$$

$F_1$  - score

mix of recall and precision

MCC from -1 to 1 and takes in consideration all of 4 squares of Conf matrix

Roc graph of TP on top understand the TP on FP with 0.1 threshold.

McNemar

Perz compare the below

		2c	2w
1c	A	B	
1w	C	D	

Pe → Encapsulates all necessary for Windows to run the app, in classmate, dynamic libraries, API.

Feature importance

Logistic Regression (LR)

Random forest Feature (RFFP)

Decision tree (DT) → uses CART

## Permutation Feature

↳ Feature Importance

↳ Scratches a column alone, some to all

↳ Repeats 3, 5, 10 or more times