

2-09-2024

Asses2 Syllabus

IR

- TF, DF, IDF, Score of doc for query
- cosine similarity
- Probabilistic Ranking
- Okapi Retrieval model / method
- Language Modelling - unigram - n-gram
- Performance Measure (A, P, R, f1)

CLASSIFICATION

divide into categories.

i/p - document

classification - class A class B

Techniques by classifier

input : $X = x_1, x_2, x_3, \dots, x_n$
class : $C = c_1, c_2, c_3, \dots, c_m$

$m \ll n$

02-09-24

Supervised Classfⁿ:

$(X_1 - C_2)$

$(X_2 - C_2)$

$X_3 - C_1$

$(X_4 - C_2)$

$X_5 - C_1$

$X_6 - C_1$

$(X_7 - C_2)$

$(X_8 - C_2)$

$(X_9 - C_2)$

$X_{10} - C_1$

Training Phase:

- features
 - behaviour
- } will be studied by classifier.

Testing Phase:

New Data

$X_{50} \rightarrow C$

3-09-24

Text Classification problem

A description $d \in X$ is given, X is the document space and a fixed set of class $C = \{c_1, c_2, \dots, c_j\}$.
Class are also called category or level.

$\langle d, c \rangle \Rightarrow$ document & class pair $\langle d, c \rangle \in X \times C$

description \downarrow
a sentence

\langle the main contributors of microprocessor & chips, china \rangle

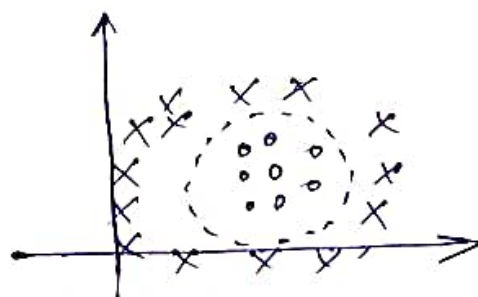
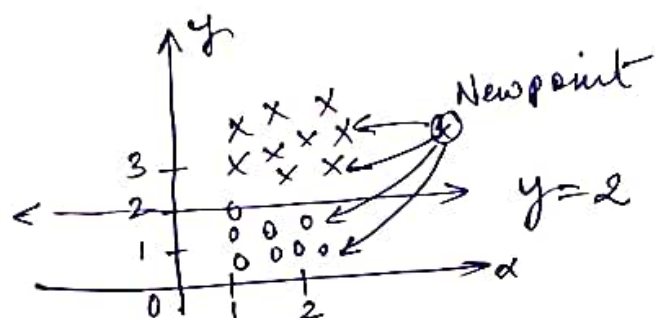
using a learning method, we want to learn a classifier or classification function γ that map the documents to class

$$\gamma: X \rightarrow C.$$

This type of learning is called supervised learning.

Classifiers

$X \rightarrow$ class A
 $O \rightarrow$ class B.



K-Nearest neighbour

- calc distance \Rightarrow then check ^{to} which class does it belong.

Naive Bayes Classification

↳ Probability based algorithm

The probability of a doc d being in class C is:

$$P(c|d) \propto P(c) \prod_{1 \leq k \leq n_d} P(t_k|c)$$

$P(c)$ - prior probability of document occurring & being in class C .

$P(t_k|c)$ - conditional prob of term t_k occurring in class C .

n_d - no. of terms in the documents

↳ Sum of the log of the probability/probabilities

$$P(c|d) = \log_e P(c) + \sum_{1 \leq k \leq n_d} \left[\log_e P(t_k|c) \right]$$

↳ $P(c) = \frac{N_c}{\text{Total of doc.}}$

09.09.24

doc ID	Text	Class
1	Dear friend, lunch money.	N
2	Friend money money	N
3	Dear money money	S
4	Money dear money money	?

$$P(-|N) = \frac{x + x_i}{y + y_i}$$

+1
↓
no. of vocab.

$$P(N|4) = P(N) \prod_{1 \leq k \leq 3} P(t_k | N)$$

$$= P(N) \left[\frac{P(\text{Money}|N) + P(\text{dear}|N)}{P(\text{money}|N) + P(\text{money}|N)} \right]^*$$

$$P(N) = \frac{2}{3}$$

$$P(S) = \frac{1}{3}$$

$$= \frac{2}{3} \times \frac{3}{7} \times \frac{1}{7} \times \frac{2 \times 3}{7 \times 7} = \frac{6 \times 3}{343 \times 7} = 0.0175 \times \frac{3}{7}$$

$$= \underline{0.00749}$$

$$P(\text{money}|N)$$

$$= \frac{3}{7}$$

$$P(\text{dear}|N)$$

$$= \frac{1}{7}$$

$$P(\text{money}|S)$$

$$= \frac{2}{3}$$

$$P(\text{dear}|S)$$

$$= \frac{1}{3}$$

$$P(S|4) = \frac{P(S) + P(\text{money}|S) + P(\text{dear}|S)}{P(\text{money}|S) + P(\text{money}|S)}$$

$$= \frac{1}{3} \times \frac{2}{3} \times \frac{1}{3} \times \frac{2}{3} \times \frac{2}{3}$$

$$= \frac{4 \times 2}{81 \times 3} = 0.0494 \times \frac{2}{3}$$

$$= \underline{0.0329}$$

\therefore doc 4 will be in class S (spam).

8.

	doc ID	words in doc	class
TRAINING	1	Chinese Beijing, chinese	C
	2	Chinese chinese shanghai	C
	3	chinese Macau	C
	4	Tokyo japan chinese	NC
	5	chinese chinese chinese Tokyo japan	?

$$P(\text{Tokyo} | C) = \frac{1}{8+6} = \frac{1}{14}$$

$$\therefore t_{\text{Tokyo} | C} = 0$$

$$P(\text{japan} | C) = \frac{1}{8+6} = \frac{1}{14}$$

$$P(\text{chinese} | C) = \frac{6}{14} = \frac{3}{7}$$

$$L \quad P(C) = 3/4$$

$$\begin{aligned}
 P(C | 5) &= P(C) \times P(\text{chinese} | C)^3 \times P(\text{Tokyo} | C) \times P(\text{Jap} | C) \\
 &= \frac{3}{4} \times \frac{3}{7} \times \frac{3}{7} \times \frac{3}{7} \times \frac{1}{14} \times \frac{1}{14} \\
 &= 0.0003012
 \end{aligned}$$

$$P(NC) = \frac{1}{4}$$

$$P(\text{Tokyo} | NC) = \frac{1+1}{3+6} = \frac{2}{9}$$

$$P(\text{Japan} | NC) = \frac{2}{9}$$

$$P(\text{chin} | NC) = \frac{2}{9}$$

$$P(NC|5) = \frac{1}{4} \times \frac{2}{9} \times \frac{2}{7} \times \frac{2}{9} \times \frac{2}{9} \times \frac{2}{9}$$

$$= 0.000135$$

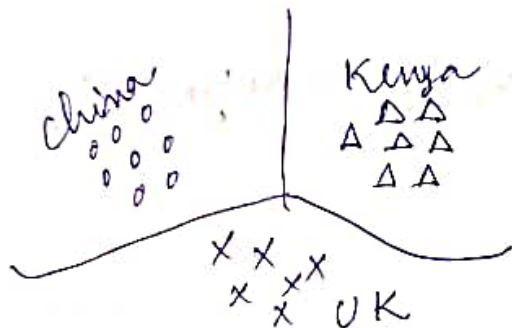
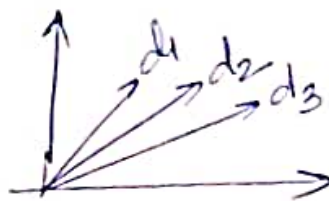
Doc 5 will be in C

3-09-2024

Vector Space Collection

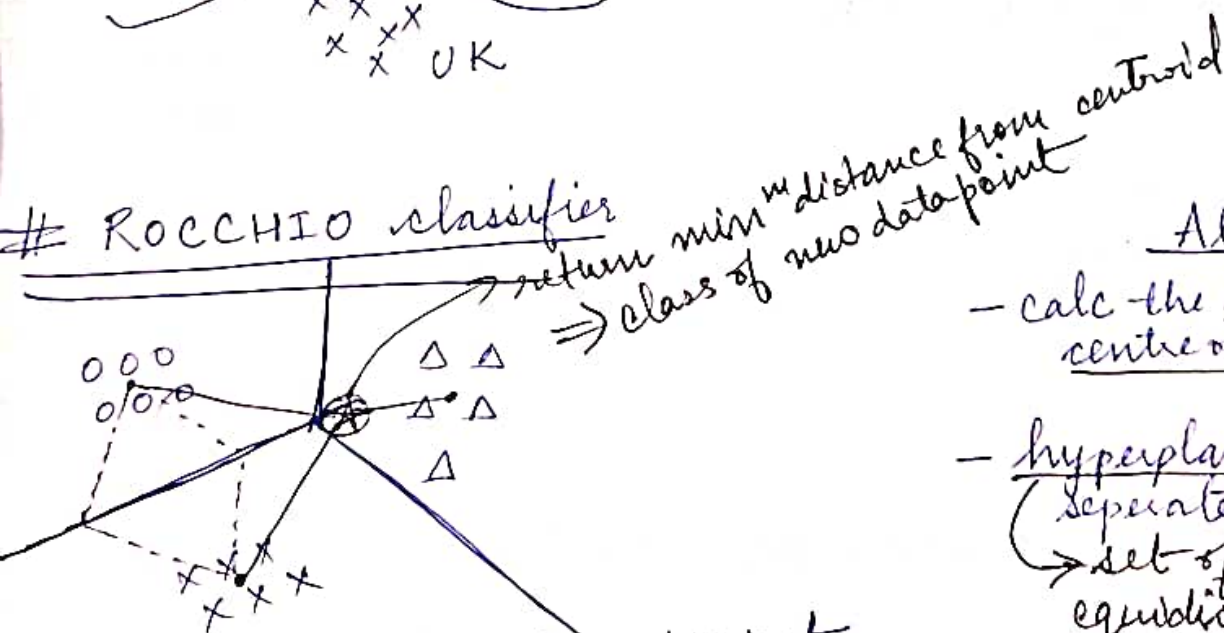
Assumption: contiguity hypothesis

Doc of same class form a contiguous region & region of diff. classes do not overlap.



Based on this
i) Rocchio
ii) KNN.

ROCCHIO classifier



Algorithm:

- calc the centroid / centre of mass.

- hyperplane to separate the classes.
(set of points i.e., equidistant from the centroid).

Advantage: easy to calc; efficient

Disadvantage: if pts belonging to same class are scattered, then centroid may not be proper.

09. 13.09.2024

KNN classifier

1. Algorithm:

- design the value of k
- Calc the dist. from new data pt. to all the existing data pt.
- arrange them in ascending order & consider the top k values
- Perform majority voting & assign the class accordingly.

Eg: let $k=3$ Top 3.

$$\boxed{d_1 < d_2 < d_3 < d_4 \dots < d_n}$$

$\underbrace{\quad \quad \quad}_{C_1} \quad \underbrace{\quad \quad}_{C_1} \quad \underbrace{\quad}_{C_2}$

majority = C_1
 \therefore return C_1

Data pts.

$$\begin{array}{cc} x_1 & y_1 \\ x_2 & y_2 \\ \vdots & \vdots \\ x_n & y_n \end{array} \quad \begin{array}{c} \frac{2d_1}{3} \\ \frac{2d_2}{3} \\ \vdots \\ \frac{2d_n}{3} \end{array}$$

New Data pt.

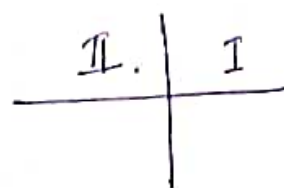
$$x_k, y_k$$

$$\begin{array}{l} (x_k y_k - x_1 y_1) - d_1 \\ (x_k y_k - x_2 y_2) - d_2 \\ \vdots \\ (x_k y_k - x_n y_n) - d_n \end{array}$$

- \downarrow
- i) Euclidean
 - ii) Manhattan
 - iii) others

13.09.2024

x	y	Quadrant
1	2	I
2	3	I
3 4	4	I I
-2	3	II
-3	1	II
-1	+5	?



Calc using
Euclidean dist
 $k=3$

$$x_k = -1 \quad y = 5$$

$$d_1 = \sqrt{4+9} = \sqrt{13}$$

$$d_2 = \sqrt{9+4} = \sqrt{13}$$

$$d_3 = \sqrt{25+1} = \sqrt{26}$$

$$d_4 = \sqrt{1+4} = \sqrt{5}$$

$$d_5 = \sqrt{4+16} = \sqrt{20}$$

$$\left[\begin{array}{cccc} d_4 & d_1 & d_2 & d_5 \\ | & | & | & | \\ \text{II} & \text{I} & \text{I} & \text{II} \end{array} \right] < d_3$$

(Note: In the original image, the 'I' labels under d_1 and d_2 are underlined and connected by a bracket.)

majority = II

$(-1, 5) - \text{II}$



#

Q.

docID

Words in doc

in C = china

1

chinese Beijing chinese

Yes

2

chinese chinese Shanghai

Yes

3

chinese Macau

Yes

4

Tokyo japan chinese

No

5

chinese chinese Tokyo japan
chinese

Using Rocchio algorithm.

Step-1 (Centroid)

$$\text{Term weight table (tf-idf)} = \frac{N=4}{(1 + \log_{10} \frac{tf_{t,d}}{1})} \left(\log_{10} \frac{N}{N_{dt}} \right)$$

Vector | chinese | japan | Tokyo | Macau | Beijing | Shanghai

d₁

0

0

0

0

0.602

0

d₂

0

0

0

0

0

0.602

d₃

0

0

0

0.602

0

0

d₄

0

0.602

0.602

0

0

0

d₅

0

0.602

0.602

0

0

0

M_C
(centroid)
(for C=Yes)

0

0

0

0.2

0.2

0.2

M_{C̄}
(centroid)
(for C=No)

0

0.602

0.602

0

0

0

$$M_C = \frac{1}{3} (d_1 + d_2 + d_3)$$

$$M_{\bar{C}} = \frac{1}{1} (d_4)$$

Basically avg.

$$\left. \begin{array}{l} |M_C - d_5| \\ |M_{\bar{C}} - d_5| \end{array} \right\} \text{min}^m \text{distance}$$

Manhattan distance.

min distance : (Using Euclidean distance) . 17.09.24

d_5 & μ_c

$$= \sqrt{0^2 + 0^2 + 0^2 + 0^2 + 0^2 + 0^2} = \sqrt{0} = 0.$$

d_5 & μ_c

$$= \sqrt{0^2 + 0.602^2 + 0.602^2 + 0.2^2 + 0.2^2 + 0.2^2}$$

$$= 0.918.$$

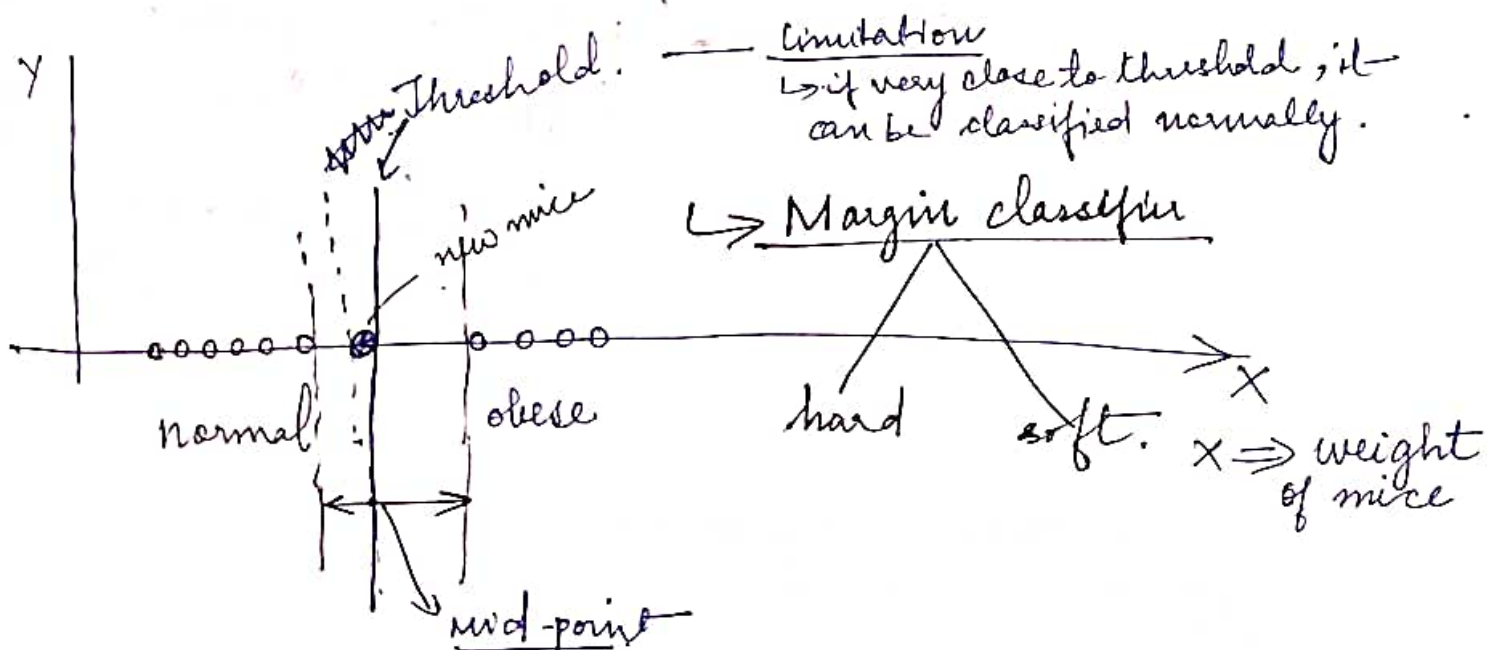
min (0, 0.918)

min is μ_c

d_5 belongs to \bar{c} i.e No.

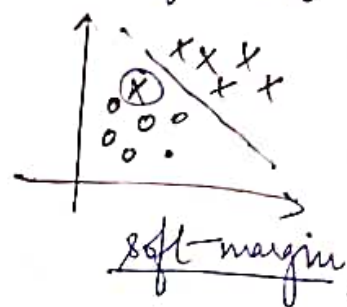
14.09.24

Support Vector Machine (SVM)

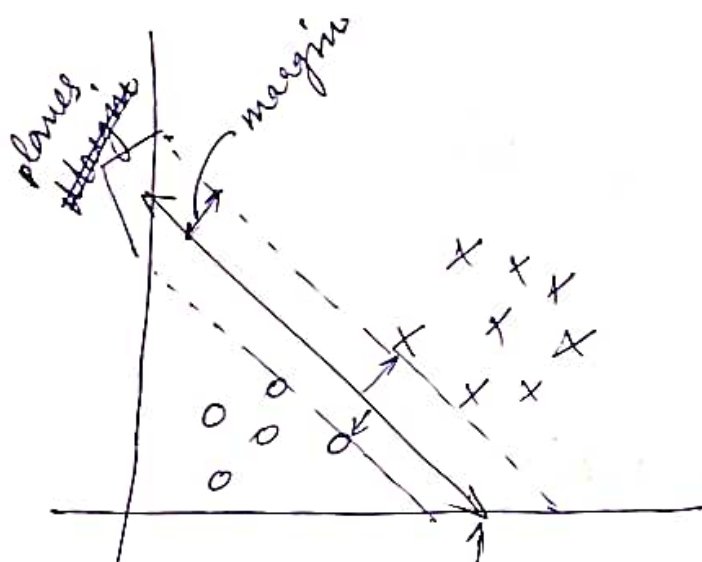


Hard Margin \rightarrow No outlier is allowed.

Soft Margin \rightarrow Allow some certain misclassification



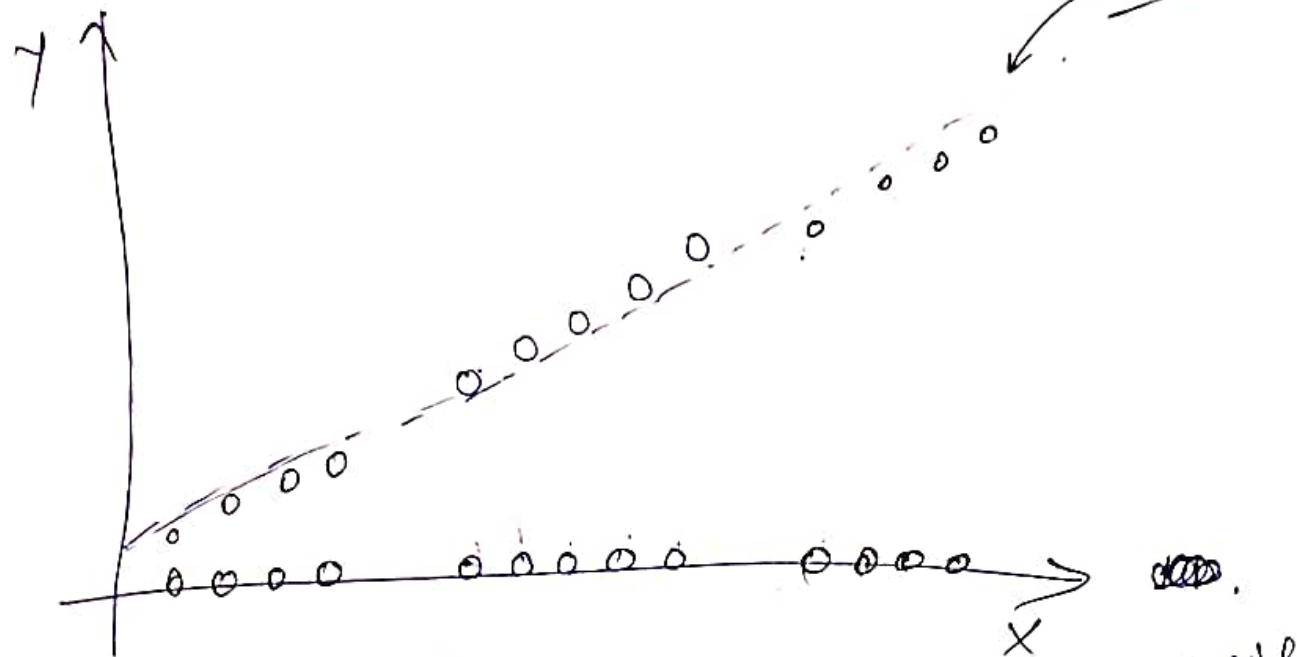
SVM \rightarrow Max^m margin classifier



hyperplane - (max^m dist. from the support vector)

\downarrow
data points which helps to create hyperplane

24.09.24
SVM \rightarrow kernel $\begin{cases} \text{linear} \\ \text{polynomial} \\ \text{RBF} \end{cases}$



- Transform the data points to diff dimension using kernel so as it is easy to classify the scattered datapoints.

$x = \text{weight}$
 $y = (\text{weight})^2$
plot (w_1, w_2) .