

Boolean retrieval is a specific case of DBMS.
go thru algo of skip pointer.

Sample Question

- ① Stemming increases recall in Boolean retrieval model. But if query has not then recall is decreased.
[in other cases stemming never decreases recall]
- ② Relevance feedback improves precision of recall. So same is the case with pseudo. only problem is if top k docs are non relevant then they fail.

③ user is reluctant. query processing takes lot of time. it is hard to understand why new set of docs are retrieved

④ vectors, t-dimensional vector, tf-idf, cosine similarity.

⑤ Local & Global

relevance → Thesaurus.

pseudo relevance → concept cloud

⑥ adv: → more weight to rare terms.
disadv: → short text

Decision tree

| <u>Age</u> | <u>competition</u> | <u>type</u> | <u>profit</u> |
|------------|--------------------|-------------|---------------|
| old | Yes | slow | Down |
| old | No | slow | Down |
| old | No | high | Down |
| mid | Yes | slow | Down |
| mid | Yes | high | Up Down |
| mid | No | high | Up |
| mid | No | slow | Up |
| new | Yes | slow | Up |
| new | No | high | Up |
| new | Yes | slow | Up |

$$P = 5$$

$$N = 5$$

$$\text{Entropy} \approx -P \log_2 \left(\frac{P}{P+N} \right) - N \log_2 \left(\frac{N}{P+N} \right)$$

$$\text{Entropy} = -\frac{P}{P+N} \log_2 \left(\frac{P}{P+N} \right) - \frac{N}{P+N} \log_2 \left(\frac{N}{P+N} \right)$$

$$= -\frac{5}{10} \log_2 \frac{1}{2} - \frac{5}{10} \log_2 \frac{1}{2}$$

$$= 0.5 \cdot 1/2 \log_2 + 1/2 \log_2 = \log_2 \frac{1}{2} = 1$$

$$\text{Entropy}_{(\text{class})} - \text{Entropy}_{(\text{attribute})} = \text{Gain}.$$

Information Gain:-

$$I(P, N) = \sum \frac{P_i + N_i}{P+N} \times I(P_i, N_i)$$

$I(Age)$

Age

| | P_i | N_i | $I(P_i, N_i)$ |
|-----|-------|-------|---------------|
| old | 0 | 3 | 0 |
| mid | 2 | 2 | 1 |
| new | 3 | 0 | 0 |

$$I(Age) = \sum \frac{P_i + N_i}{P+N} I(P_i, N_i)$$

$$= \frac{3}{10} \times 0 + \frac{4}{10} \times 1 + \frac{3}{10} \times 0$$

$$= 0.4$$

$$Gain = 1 - 0.4$$

$$= 0.6$$

Competition

| | P _i | N _i [°] | I(P _i , N _i) |
|-----|----------------|-----------------------------|-------------------------------------|
| Yes | 1 | 3 | 0.81127 |
| No | 4 | 2 | 0.918295 |

$$I(P_i^{\circ}, N_i^{\circ}) = - \frac{1}{4} \log_2 \frac{1}{4} - \frac{3}{4} \log_2 \frac{3}{4}$$

$$= -\frac{1}{4} \times 2 + \frac{3}{4} \times \log_2 \frac{4}{3} = 0.81127$$

$$I(P_i^{\circ}, N_i^{\circ}) = - \frac{4}{6} \log \frac{4}{6} - \frac{2}{6} \log \frac{2}{6}$$

$$= -\frac{2}{3} \log \frac{2}{3} - \frac{1}{3} \log \frac{1}{3}$$

$$= 0.918295$$

$$\begin{aligned} \text{Entropy (comp)} &= \frac{1+3}{10} \times 0.81127 + \frac{4+2}{10} \times 0.918295 \\ &= 0.875484 \end{aligned}$$

$$\text{Gain} = 1 - 0.875484.$$

type

$$P_i \quad N_i \quad I(P_i, N_i)$$

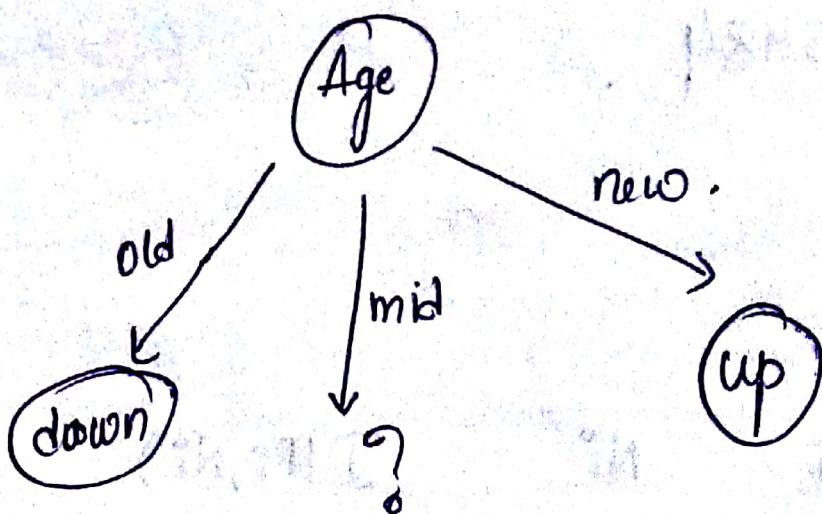
$$S/W \quad 3 \quad 3 \quad 1$$

$$H/W \quad 2 \quad 2 \quad 1$$

$$\text{Entropy (type)} = \frac{6}{10} \times 1 + \frac{4}{10} \times 1$$

$$= 1$$

$$\text{Gain} = 1 - 1 = 0.$$



competition type.

| | P_i^o | Ni^o | $I(P_i^o, Ni^o)$ |
|-----|---------|--------|------------------|
| yes | 0 | 2 | 0 |
| No | 2 | 0 | 0 |

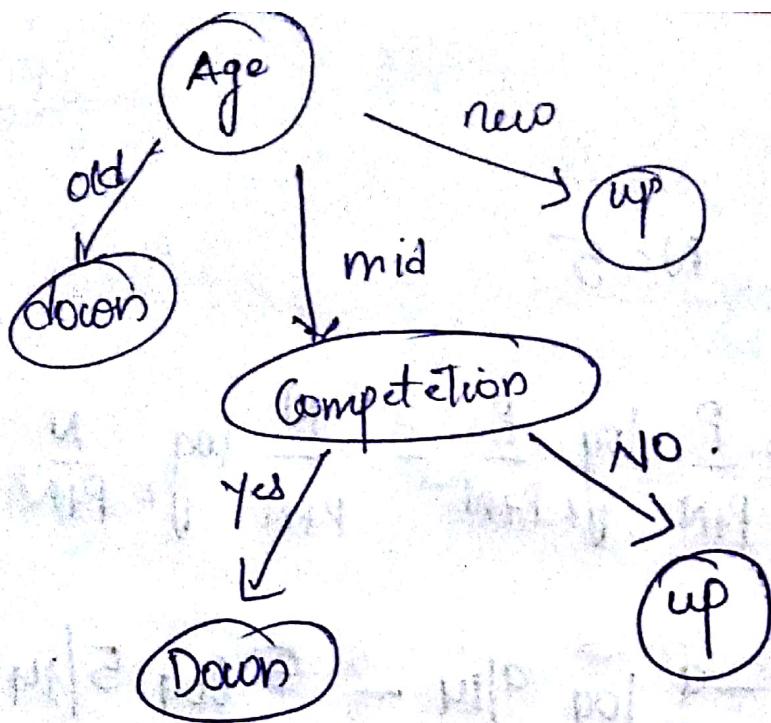
$$\text{Gain} = 1 - 0 = 1$$

type

| | P_i^o | Ni^o | $I(P_i^o, Ni^o)$ |
|-----|---------|--------|------------------|
| S/W | 1 | 1 | 1 |
| H/W | 1 | 1 | 1 |

$$\text{Entropy} = 1$$

$$\text{Gain} = 0.$$



Example - 2

$$P = 9, \quad N = 5.$$

$$\begin{aligned} \text{Entropy}_{\text{class}} &= -\frac{P}{P+N} \log_2 \frac{P}{P+N} - \frac{N}{P+N} \log_2 \frac{N}{P+N} \\ &= -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} \end{aligned}$$

$$\text{Entropy}_{\text{class.}} = 0.4 + 0.53 = 0.93.$$

I(outlook)

| | P | N. | $I(P_i N_i)$ |
|----------|---|----|--------------|
| Sunny | 2 | 3 | 0.96 |
| overcast | 4 | 0 | 0 |
| rain | 3 | 2 | 0.96 |

$$= -\frac{2}{5} \log \frac{2}{5} - \frac{3}{5} \log \frac{3}{5}$$

$$\therefore 0.52 + 0.44 = 0.96.$$

$$\begin{aligned} I(\text{outlook}) &= \sum \frac{P_i + N_i}{P+N} I(P_i, N_i) \\ &= \frac{5}{14} \times 0.96 + \frac{5}{14} \times 0.96 \\ &= \frac{5}{7} \times 0.96 = 0.68. \end{aligned}$$

$$\text{Gain} = 0.93 - 0.68 = 0.25$$

$I(\text{Temp})$

Hot

2

N

$I(P_h, N_i)$

1

Cool

3

2

0.81

Mild

4

2

0.9

$$-\frac{3}{4} \log \frac{3}{4} - \frac{1}{4} \log \frac{1}{4}$$

$$0.31 + 0.5 = 0.81$$

$$-4/6 \log 4/6 - 2/6 \log 2/6$$

$$0.38 + 0.52 = 0.9$$

$$I(\text{temp}) = \frac{4}{14} \times 1 + \frac{4}{14} \times 0.81 + \frac{6}{14} \times 0.9$$

$$= \frac{4 + 3.24 + 5.4}{14} = 0.902$$

$$\text{Gain} = 0.93 - 0.902 = 0.028.$$

I (Humidity)

High

Normal

P

Q

6

N.

Q

1

I (P, N)

0.9

$$-4/7 \log 4/7 - 3/7 \log 3/7 =$$

$$0.46 + 0.52 = 0.98.$$

$$-6/7 \log 6/7 - 1/7 \log 1/7 =$$

$$0.19 + 0.40 = 0.59$$

$$\mathbb{P}(\text{Humidity}) = \frac{7}{14} \times 0.98 + \frac{7}{14} \times 0.59$$

$$= 0.785$$

$$\text{Grain} = 0.93 - 0.785 = 0.145$$

$\mathbb{I}(\text{Wind})$

$$P \quad N \quad \mathbb{P}(P_i, N_i)$$

Strong

6

2

0.81

Weak

3

3

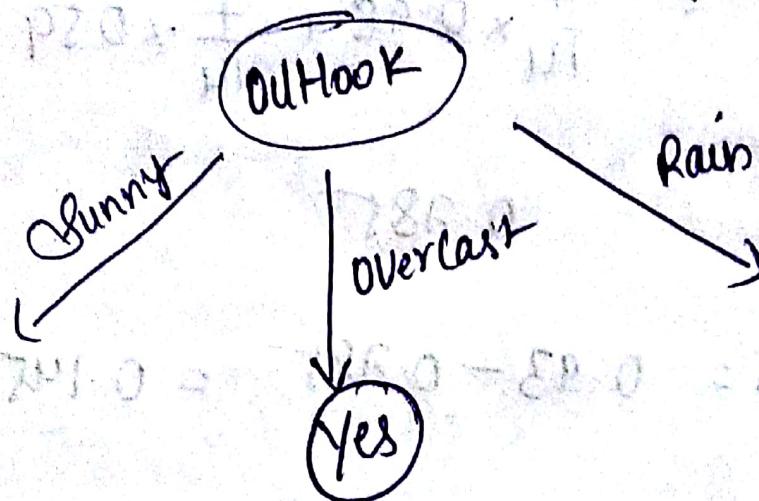
1

$$-6/8 \log 6/8 - 2/8 \log 2/8$$

$$P(\text{wind}) = \frac{8}{14} \times 0.81 + \frac{6}{14} \times 1$$

$$= 0.89$$

$$\text{Gain} = 0.93 - 0.89 = 0.04$$



Sunny

$I(\text{temp})$

Hot

P
2

N
2

$I(P_i, N_i)$

0

Cold

1

0

0

Mild

1

1

1

$$I(\text{temp}) = \frac{2}{5} \times 1 = 0.4.$$

$$\text{Gain} = 0.93 - 0.4 = 0.53.$$

$I(\text{Humidity})$

| | P _i | N | $I(P_i^o, N_i)$ |
|---------|----------------|---|-----------------|
| High | 0 | 3 | 0 |
| Normal. | 2 | 0 | 0 |

$$\text{Gain} = 0.93.$$

$I(\text{wind})$

| | P | N | $I(P_i^o, N_i)$ |
|--------|---|---|-----------------|
| Weak | 2 | 1 | 0.9 |
| Strong | 1 | 1 | 1 |

$$F(P_{NN}^o) = 3 \times 0.9 + 2 \times 5 \times 1 = 0.94$$

Naïve Bayes classifier :-

$$P(h|D) = \frac{P(D|h) P(h)}{P(D)}$$

There is no dependency between attributes i.e.
independent attributes.

Q) Does patient cancer or not? A patient takes a lab test & the result comes bad positive. The result comes correct positive in only 98% cases in which the disease is actually present and a correct negative result in 97%. Only in which the disease is not present. Furthermore 0.008% of the entire population have this cancer. By Bayes using Bayes theorem, say whether patient has cancer or not.

| | |
|---|----|
| 8 | 2 |
| 3 | 97 |

②

$$P(\text{bad test} + \text{Correct test}) = P(\text{cancer}) = 0.008$$

$$P(\text{Correct test}) = P_{\text{pop}} = 0.992$$

$$P(+ | \text{cancer}) = 0.98$$

$$P(- | \text{cancer}) = 0.02$$

$$P(+ | \neg \text{cancer}) = 0.03$$

$$P(- | \neg \text{cancer}) = 0.97$$

$$P(+ | \text{cancer}) \cdot P(\text{cancer}) = 0.98 \times 0.008 = 0.0078$$

$$P(+ | \neg \text{cancer}) \cdot P(\neg \text{cancer}) = 0.03 \times 0.992 = 0.0298$$

Patient does not have cancer.

$$\textcircled{2} \quad P(\text{Play} | \text{Yes}) = 9/14, \quad P(\text{Play} | \text{No}) = 5/14.$$

outlook

$$P(\text{Sunny} | P) = 3/9$$

$$P(\text{Overcast} | P) = 4/9$$

$$P(\text{Rain} | P) = 3/9$$

$$P(\text{Sunny} | N) = 3/5$$

$$P(\text{Overcast} | N) = 0$$

$$P(\text{Rain} | N) = 2/5$$

temp

$$P(\text{Hot} | P) = 2/9$$

$$P(\text{Hot} | N) = 2/5$$

$$P(\text{Mild} | P) = 4/9$$

$$P(\text{Mild} | N) = 2/5$$

$$P(\text{Cool} | P) = 3/9$$

$$P(\text{Cool} | N) = 1/5$$

humidity

$$P(\text{High} | P) = 3/9$$

$$P(\text{High} | N) = 4/5$$

$$P(\text{Normal} | P) = 6/9$$

$$P(\text{Normal} | N) = 1/5$$

wind

$$P(\text{Strong} | P) = 6/9$$

$$P(\text{Strong} | N) = 2/5$$

$$P(\text{Weak} | P) = 3/9$$

$$P(\text{Weak} | N) = 3/5$$

$x: \langle \text{rain}, \text{hot}, \text{high}, \text{false} \rangle$

$$P(\text{play}|x) = P$$

$$P(\text{play } x | P) = P(\text{rain}|P) \times P(\text{hot}|P) \times P(\text{high}|P)$$
$$\times P(\text{Strong}|P)$$

$$= \frac{3}{9} \times \frac{2}{9} \times \frac{3}{9} \times \frac{6}{9}$$

$$= \frac{162}{9^4}$$

$$P(P|x) = \underline{P(x|P)P(P)}$$

$$= \frac{162}{9^4} \times \frac{9}{14} = \frac{1}{9} \times \frac{2}{7}$$

$$= \frac{1}{168} = 0.006$$

$$P(x|N) = P(\text{rain}|N) \times P(\text{hot}|N) \times P(\text{high}|N)$$

$$\times P(\text{Strong}|N)$$

$$= \frac{2}{5} \times \frac{2}{5} \times \frac{4}{5} \times \frac{8}{5}$$

$$P(N|X) = \frac{3^2}{5^4 5^{3/4}} \times \frac{8}{5} = 0.018.$$

Result → NO.

Example + Training Set

doc id

1. chinese beijing chinese Yes

2. chinese Chinese shanghai Yes

3. chinese macau Yes

4. Tokyo Japan Chinese No.

test 5. Chinese Chinese Chinese Tokyo Japan.

P(chinese | p(Yes)) = 3/4 , P(No) = 1/4

P(chinese | yes) = 5/8 , P(chinese | No) = 1/3.

P(beijing | yes) = 1/8 , P(beijing | No) = 0

P(Tokyo | yes) = 0 , P(Tokyo | No) = 1/3.

$$P(\text{to} \cancel{\text{to}} \text{Japan} | \text{yes}) = 0$$

$$P(\text{Japan} | \text{No}) = 1/3.$$

Ans

Laplacean correction:-

$$P(t|c) = \frac{T_{ct} + 1}{\sum_{t \in V}^i(t) + B}$$

$\sum_{t \in V}^i(t) + B \rightarrow$ distinguished words.

$$P(\text{yes}) = 3/4, P(\text{No}) = 1/4$$

$$P(\text{Chinese} | \text{yes}) = \frac{5+1}{8+6}$$

$$P(\text{Chinese} | \text{No}) = \frac{1+1}{3+6} = 2/9$$

$$= 6/14$$

$$P(\text{Tokyo} | \text{yes}) = \frac{0+1}{8+6} = \frac{1}{14}, P(\text{Tokyo} | \text{No}) = \frac{1+1}{3+6} = 2/9$$

$$P(\text{Japan} | \text{yes}) = \frac{0+1}{8+6} = 1/14, P(\text{Tokyo} | \text{No}) = \frac{1+1}{3+6} = 2/9$$

$$P(X|yes) = P(\text{Chinese}|yes)^3 + P(\text{Other}|yes) \times P(\text{Japan}|yes)$$
$$= \frac{(6)^3}{145} =$$

$$P(yes|X) = \frac{6^3}{145} \times \frac{3}{4} = \frac{6 \times 8}{4 \times 145} = 0.0003.$$

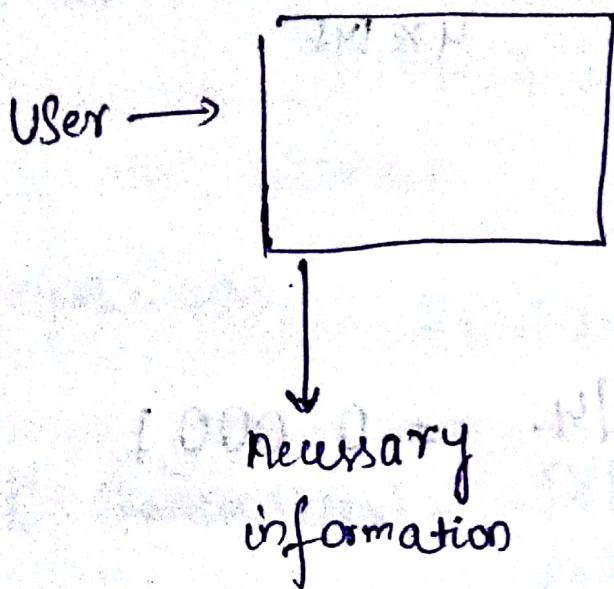
$$P(X|No) = (2/9)^3 \times (2/9)^5$$

$$P(No|X) = (2/9)^5 \times 1/4 = 0.0001$$

Result = yes

Recommender System :-

Information overload. in 1990

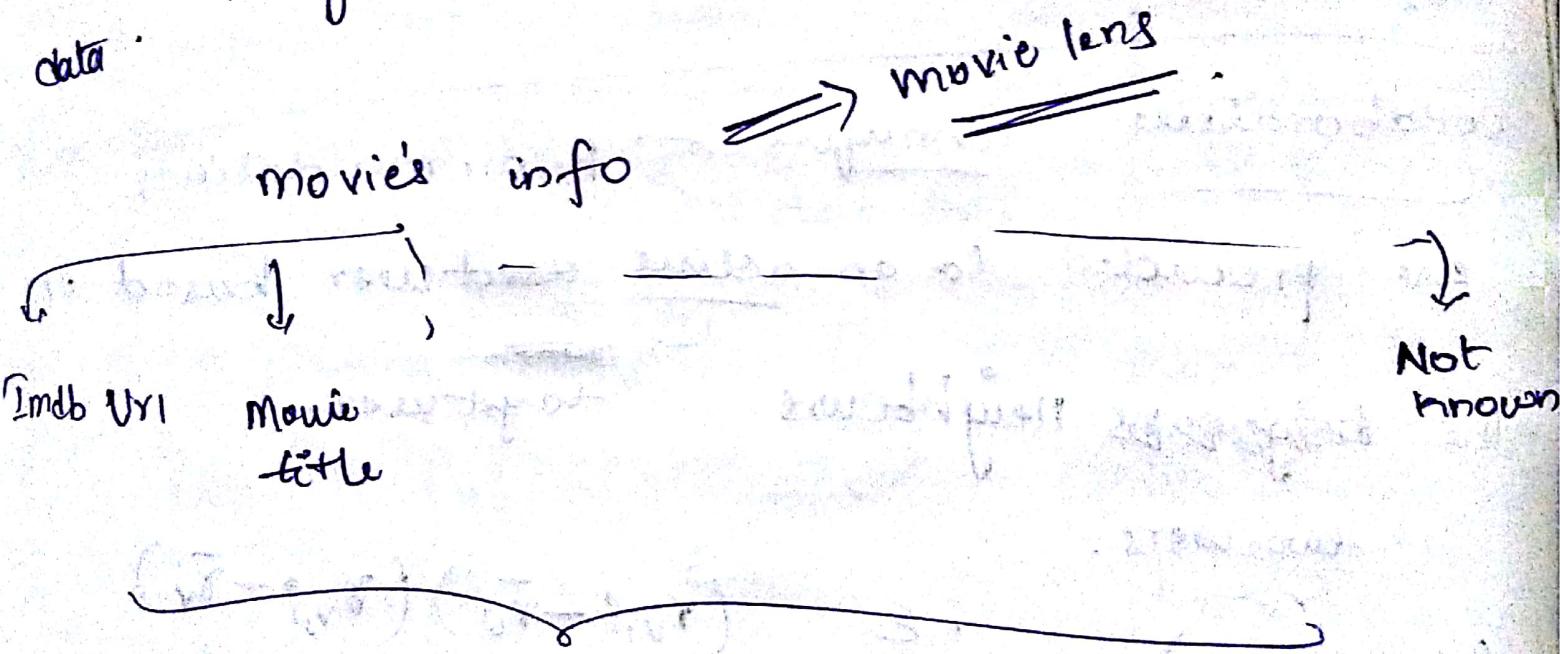


Group lens , movie lens. } → RS are used.

↓
Bench mark data etc

- ① user's age
 - ② user's occupation
 - ③ user's rating on experienced movies.
 - ④ movie's information.
- } by user.
IDS
CWE

for building any ~~not~~ recommender system. we need data.



Recommendation Techniques

- ① collaborative ~~Technique~~ filtering
- ② content based filtering
- ③ hybrid filtering -

Collaborative Technique :- Recommendations are provided to an active user based on the neighbours neighbours.

two ways.

$$\text{sim}(u, v) = \frac{\sum_{i \in \text{Common}} (r_{ui} - \bar{r}_u)(r_{vi} - \bar{r}_v)}{\sqrt{\sum_{i \in \text{Common}} (r_{ui} - \bar{r}_u)^2} \sqrt{\sum_{i \in \text{Common}} (r_{vi} - \bar{r}_v)^2}}$$

$$\text{sim}(u, v) = \frac{\sum_{i \in \text{Common}} (r_{ui} - \bar{r}_u)(r_{vi} - \bar{r}_v)}{\sqrt{\sum_{i \in \text{Common}} (r_{ui} - \bar{r}_u)^2} \sqrt{\sum_{i \in \text{Common}} (r_{vi} - \bar{r}_v)^2}}$$

$$\text{range} = (-1, 1)$$

actions ↑

| | T_1 | T_2 | T_3 | T_4 | T_5 | T_6 | T_7 | T_8 | → movies |
|---------|-------|-------|-------|-------|-------|-------|-------|-------|------------------------|
| U_1 | 3 | 0 | 4 | 0 | 5 | 0 | 2 | 5 | Rate at least 4 movies |
| U_2 | 3 | 4 | 5 | 3 | 4.5 | 0 | 0 | 0 | Seen. |
| U_3 | 1 | 2 | 3 | 4 | 5 | 0 | 0 | 0 | |
| U_4 | 5 | 2 | 3 | 3 | 3 | 3 | 2 | 5 | |
| U_5 | 1 | 2 | 3 | 4 | 3 | 4 | 2 | 4 | |
| | (2) | | | 4 | | 4 | | | |
| Users ↓ | | | | | | | | | |

Data into two parts

→ training data
→ test data

divide in
60 - 40 %.

for $U_1 \rightarrow$ Training movie

(5)

(T_1, T_3, T_5)

Test movie

(2)

(T_7, T_8)

find $\text{Sim}(U_1, U_2)$, $\text{Sim}(U_1, U_3)$, $\text{Sim}(U_1, U_4)$, $\text{Sim}(U_1, U_5)$

$\text{Sim}(U_1, U_2)$

Common movie — T_1, T_3, T_5 .

$$\bar{x}_{U_1} = 4, \bar{x}_{U_2} = \frac{3+4+5+3+4+5}{6} = \frac{24}{6} = 4$$

$$\text{Sim}(U_1, U_2) = \frac{(3-4)(3-4) + (4-4)(5-4) + (5-4)(4-4)}{\sqrt{(3-4)^2 + (4-4)^2 + (5-4)^2} \sqrt{(3-4)^2 + (5-4)^2 + (4-4)^2}}$$

$$= \frac{1}{\sqrt{2} \sqrt{2}} = \frac{1}{2} = 0.5$$

for U_1, U_3

Common modes = I_1, I_3, I_5

$$\bar{x}_{U_1} = 4, \quad \bar{x}_{U_3} = \frac{1+2+3+4+5}{5} = 3$$

$$\text{Sim}(U_1, U_2) = \frac{(3-4)(1-3) + (4-4)(3-3) + (5-4)(5-3)}{\sqrt{(3-4)^2 + (4-4)^2 + (5-4)^2} \sqrt{(1-3)^2 + (3-3)^2 + (5-3)^2}}$$

$$= \frac{2+2}{\sqrt{2} \times 2\sqrt{2}} = \frac{4}{2 \times 2} = 1$$

for U_1, U_4

$$\frac{17}{2}, \frac{18}{4}$$

(5)

Common movies = $\bar{I}_1, \bar{I}_3, \bar{I}_5$

$$\bar{r}_{U_1} = 4 \quad \bar{r}_{U_4} = \underbrace{5+2+3 \times 4+2+5}_{8}$$

$$= \frac{14+12}{8} = 26/8$$

After all similarities. ~~use~~ we use K nearest neighbour with active users.

if similarities are $(U_1, U_2) = 0.5$
 $(U_1, U_3) = 0.7$ ✓
 $(U_1, U_4) = 0.3$
 $(U_1, U_5) = 0.8$ ✗

K=2
generally

Similar user = U_3, U_5

We think test data as unseen movies for performance measure.

Test cases:-

- ① all similar user have seen the movie
- ② some fraction have seen the movie
- ③ No. similar user have seen the movie.

↙
System can't recommend.

Ques - 1 :-

| | I_7 | I_8 |
|-------|-------|-------|
| U_3 | 3 | 4 |
| U_5 | 2 | 4 |

$$\text{Avg of } I_7 = 2.5$$

$$" " I_8 = 4.$$

Predicted rating for $I_7 = 2.5$

$$" " " I_8 = 4.$$

mean absolute error for $I_7 = 0.5$
for $I_8 = 1$

0.25 TRUE
.6
.3
.5
.7 one

We do this by considering other users as active
and find out Mean absolute error. We ~~have~~ all this
error and say ~~the~~ \oplus it as error for collaborative
filtering.

two clauses :- ham / Spam.

ham d1: "good"

ham d2: "Very good"

Spam d3: "bad"

Spam d4: "Very bad"

Spam d5: "Very bad, very bad"

Test: d6: "good? bad! very bad"

$$P(\text{Spam}) = 3/5, P(\text{ham}) = 2/5$$

$$P(\text{good} | \text{Spam}) = 0, P(\text{bad} | \text{Spam}) = 2/3$$

$$P(\text{good} | \text{ham}) = 2/3, P(\text{bad} | \text{Spam}) = 4/7$$

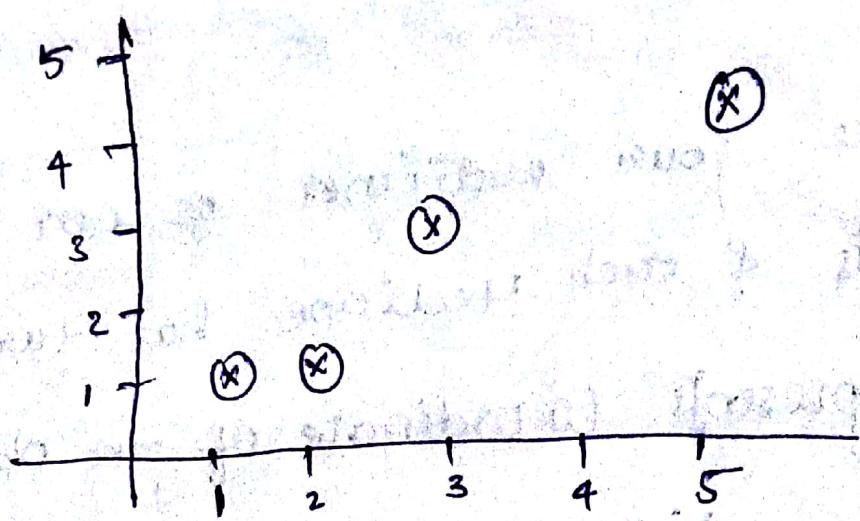
$$P(\text{very} | \text{Spam}) = 8/7, P(\text{bad} | \text{ham}) = 0.$$

$$P(\text{Very} | \text{ham}) = 1/3$$

clustering :-

Suppose we have four medicines as our training data point objects & each Medicine has two attributes, each attribute represents co-ordinate of the object, we have to determine which medicine belongs to medicine 1 & which one belongs medicine 2.

| object | Weighted Index (x) | $p^*(y)$ |
|--------|--------------------|----------|
| Med A | 1 | 1 |
| Med B | 2 | 1 |
| Med C | 3 | 3 |
| Med D | 5 | 4 |



$C_1 (1, 1)$] \rightarrow assumption of
 $C_2 (2, 1)$ centroids as first two data sets.

| A | B | C | D |
|---|---|---|---|
| 1 | 2 | 3 | 5 |
| 1 | 1 | 3 | 4 |

$C_1 (1, 1)$ $C_2 (2, 1)$

| A | B | C | D |
|---------|---|---|------|
| C_1^0 | 0 | 1 | 3.61 |
| C_2^0 | 1 | 0 | 2.83 |

$$G^0 =$$

| | | | |
|---|---|---|---|
| 1 | 0 | 0 | 0 |
| 0 | 1 | 1 | 1 |

Group 1

Group 2

$$C_1 = (1, 1)$$

$$C_2 = \left(\frac{2+3+5}{3}, \frac{1+3+4}{3} \right) = \left(\frac{10}{3}, \frac{8}{3} \right)$$

| | A | B | C | D |
|----|------|------|------|-------|
| C1 | 0 | 1 | 361 | 5 |
| C2 | 3.14 | 2.36 | 0.47 | 1.89. |

| | | |
|---------|-------------------------------------------------------------------------------------------------|---------|
| $G^1 =$ | $\begin{array}{ c c c c c } \hline 1 & 1 & 0 & 0 \\ \hline 0 & 0 & 1 & 1 \\ \hline \end{array}$ | Group 1 |
| | | Group 2 |

$$C_1 = \left(\frac{1+2}{2}, \frac{1+1}{2} \right) = (3/2, 1)$$

$$C_2 = \left(\frac{3+5}{2}, \frac{3+4}{2} \right) = (4, 7/2)$$

D^2_2

C1

C2

| | | | |
|------|------|------|------|
| 0.5 | 0.5 | 3.20 | 4.61 |
| 4.30 | 3.54 | 0.71 | 0.71 |

G^2_2

| | | | |
|---|---|---|---|
| 1 | 1 | 0 | 0 |
| 0 | 0 | 1 | 1 |

Group 1

Group 2

final
centroids.

$$\left\{ \begin{array}{l} C1 = (3/2, 1) \\ C2 = (4, 9/2) \end{array} \right.$$

K-Means is not suitable for short text analysis,
is not good like tweets.

Hierarchical clustering

| | BA | FI | MI | NA | RM | TO | MI TO |
|----|-----|-----|-----|-----|-----|-----|-------|
| BA | 0 | 662 | 877 | 255 | 412 | 996 | 877 |
| FI | 662 | 0 | 295 | 468 | 268 | 400 | 295 |
| MI | 877 | 295 | 0 | 384 | 564 | 138 | 0 |
| NA | 255 | 468 | 384 | 0 | 219 | 869 | 384 |
| RM | 412 | 295 | 564 | 219 | 0 | 669 | 564 |
| TO | 996 | 400 | 138 | 869 | 669 | 0 | 0 |

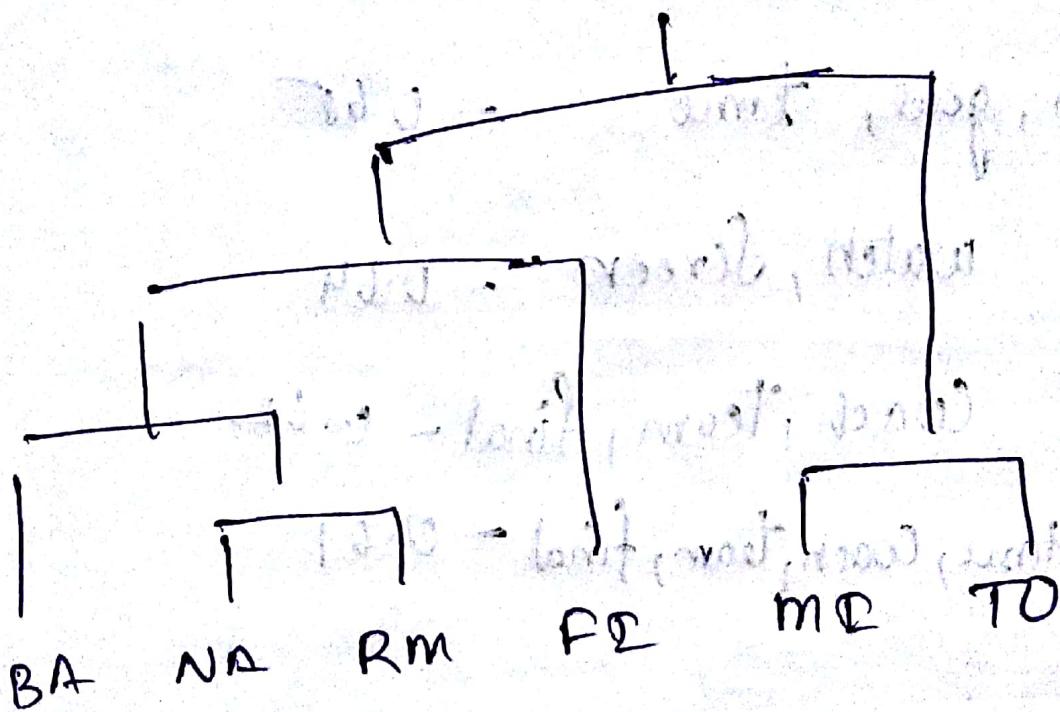
① club two cities, based on 1min distance

| | BA | FI | MI TO | NA | RM |
|-------|-----|-----|-------|-----|-----|
| BA | 0 | 662 | 877 | 255 | 412 |
| FI | 662 | 0 | 295 | 468 | 268 |
| MI TO | 877 | 295 | 0 | 266 | 564 |
| NA | 255 | 468 | 266 | 0 | 219 |
| RM | 412 | 268 | 564 | 219 | 0 |

| | BA | FI | MI/TO | RM/NA |
|-------|-----|-----|------------|----------|
| BA | 0 | 662 | 877 | (255) |
| FI | 0 | 0 | 295 | 268 |
| MI/TO | 877 | 295 | 0 | 268, 564 |
| RM/NA | 255 | 268 | 564 268 | 6 |

| | FI | MI/TO | BA/RM/NA |
|----------|-------------------|------------|-------------------|
| FI | 0 | 295 | 268 268 268 |
| MI/TO | 295 | 0 | 564 |
| BA/RM/NA | 268 268 268 | 564 268 | 0 |

| | MI / TO | BA / RMINA FI |
|------------------|------------|------------------|
| MI / TO | 0 | 895 995 |
| BA / RMINA FI | 895 295 | 0 |



Q2)

=

replays, offside - 0.69

goal, time - 0.68

team, final - 0.65

match, goal, time - 0.65

watch, screen - 0.64

coach, team, final - 0.62

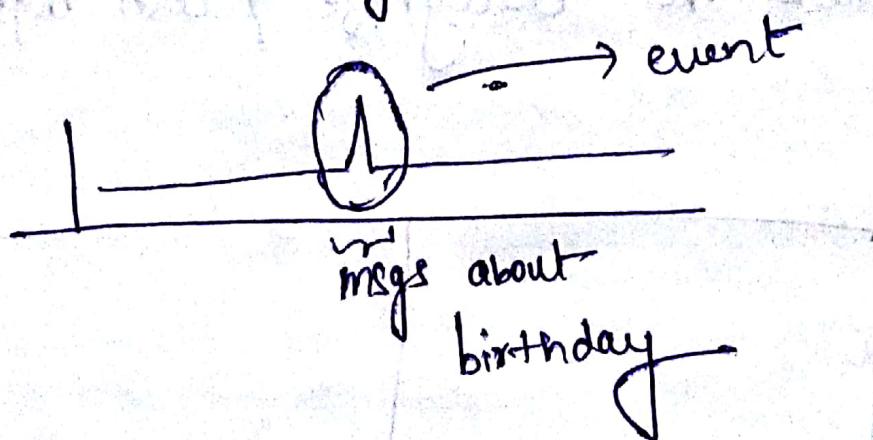
match, goal, time, coach, team, final - 0.61

Social Network Analysis

trends, events, personalization, social connectivity,
digital marketing, knowing people, observe
monitor people,

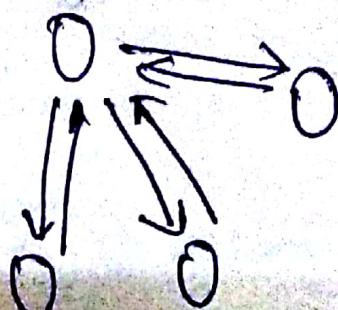
trend → what people are interested in

event → like birthday.

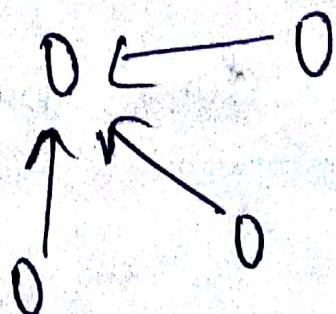


Essentially analysing human behaviour.
(groups of humans)

If there is content analysis, we need NLP.



facebook



-twitter

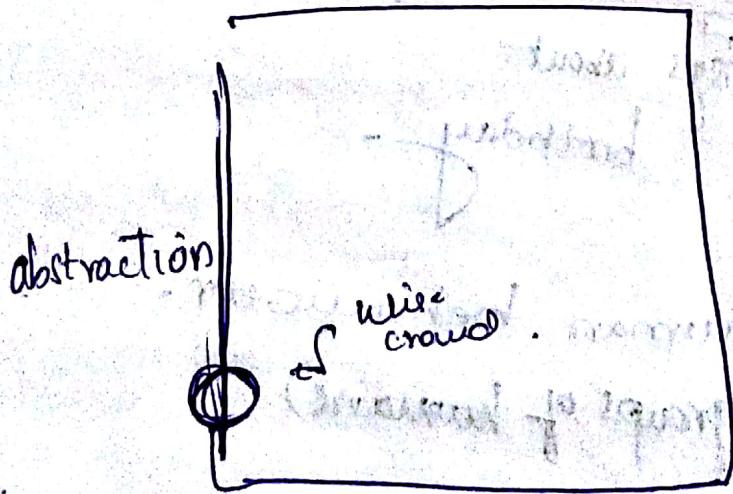
Qualitative analysis of online concepts! -

Herding :- following someone.

Crowd Sourcing :-

By reddit → disastrous

asking crowd to recognize / ask information



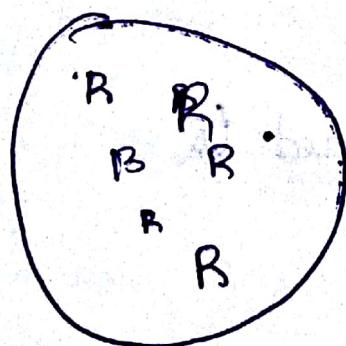
expression

Wise Crowd → have independency, but agree upon
some basic rules.

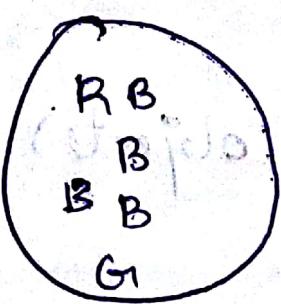
Purity of a cluster :— max count of similar objects in a class of the cluster, divided by total no. of objects.

$$\text{Purity} = \frac{\text{max count in a class}}{\text{total no. of objects}}$$

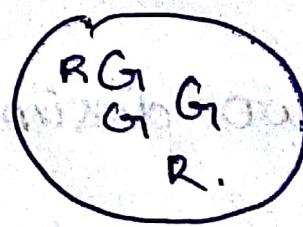
cluster 1



cluster -2



cluster -3



$$\text{cluster 1} = \max(5, 0, 1) / 6. = 5/6$$

$$\text{cluster 2} = \max(4, 1, 1) / 6. = 4/6$$

$$\text{cluster 3} = \max(3, 0, 1) / 5 = 3/5$$

$$\begin{aligned} \text{Purity of clustering algo} &= \frac{5/6 + 4/6 + 3/5}{6} \\ &= \cancel{\frac{5}{6} + \cancel{\frac{4}{6}} + \cancel{\frac{3}{5}}} - \cancel{\frac{18}{30}} \end{aligned}$$

$$= \frac{5+4+3}{17} = \frac{12}{17}$$

Rank Index

T_p :- two similar docs (objects) assigned to same cluster

T_N :- two dissimilar docs (objects) assigned to different cluster

F_p :- two dissimilar docs (objects) assigned to same cluster

F_{NT} :- two similar docs (objects) assigned to different cluster

| | similar | dissimilar |
|------------|----------|------------|
| similar | $T_p(a)$ | $F_N(c)$ |
| dissimilar | $F_p(b)$ | $T_p(a)$ |

$$P = \frac{a+d}{a+b+c+d}, R = \frac{a}{a+b}, F = \frac{a}{a+b+c}$$

| class / cluster | v_1 | v_2 | \dots | v_c | SUM |
|-----------------|----------|----------|----------|----------|----------|
| U_1 | U_{11} | U_{12} | \vdots | \dots | U_{1c} |
| U_2 | U_{21} | U_{22} | \vdots | \dots | U_{2c} |
| \vdots | \vdots | \vdots | \vdots | \vdots | \vdots |
| U_R | U_{R1} | U_{R2} | \vdots | \dots | U_{Rc} |
| SUM. | n_1 | n_2 | \vdots | \dots | n_c |

| class / cluster | v_1 | v_2 | v_3 | Sum |
|-----------------|----------|----------|----------|--------|
| U_1 | 1 | 1 | 0 | 2 |
| U_2 | 1 | 2 | 1 | 4 |
| U_3 | 0 | 0 | 4 | 4 |
| Sum | 2 | 3 | 5 | 10 [n] |
| | v_{01} | v_{02} | v_{03} | |

$$0.5(18+22+4) = 20$$

$$a = \sum_{i,j} \binom{n_{ij}}{2}$$

$a \rightarrow$ no. of pairs of objects in the same class $U \cap V$ in the same cluster in V

$b \rightarrow$ no. of pairs of objects in the same class $U \cap V$ not in same cluster in V

$c \rightarrow$ no. of pairs of objects in the different class $U \setminus V$ in the same cluster in V

$d \rightarrow$ no. of pairs of objects in the different class $U \setminus V$ in different cluster in V .

$$a = \sum_{i,j} \binom{n_{ij}}{2}$$

$$= \binom{2}{2} + \binom{4}{2} = 1 + 6 = 7$$

$$b = \sum_{i,j} \binom{n_{ij}}{2} - \sum_{i,j} \binom{n_{ij}}{2}$$

$$= \cancel{\binom{2}{2}} + \cancel{\binom{4}{2}} + \binom{4}{2} - \cancel{\binom{2}{2}} - \cancel{\binom{4}{2}}$$

$$= 6$$

$$c = \binom{2}{2} + \binom{3}{2} + \binom{5}{2} \rightarrow \left[\sum_j \binom{n_{ij}}{2} - \sum_{i,j} \binom{n_{ij}}{2} \right]$$

$$= 1 + 3 + 10 - 7$$

$$= 7.$$

$$d = \binom{n}{2} - (a+b+c)$$

$$= \binom{10}{2} - \binom{20}{2} = \frac{45-20}{2} = 25$$

$$P(D) = \frac{a+d}{a+b+c+d} = \frac{7+25}{2+4+2+5} = \frac{32}{70}.$$

| class cluster | 1 | 2 | 3 | |
|-----------------|--------|--------|--------|----------|
| R | 5 | 1 | 2 | 8 $n_1.$ |
| G | 0 | 1 | 3 | 4 $n_2.$ |
| B | 1 | 4 | 0 | 5 $n_3.$ |
| | 6 | 6 | 5 | 17 |
| | $n_1.$ | $n_2.$ | $n_3.$ | |

$$a = \binom{5}{2} = 10.$$

$$b = \binom{8}{2} + \binom{4}{2} + \binom{5}{2} - 10$$

$$= 28 + 6 + 10 - 10 = 34$$

$$c = \binom{6}{2} + \binom{6}{2} + \binom{5}{2} - 10$$

$$= 15 + 15 + 10 - 10 = 30.$$

$$d = \binom{17}{2} - (10 + 34 + 30)$$

$$= 136 - 74 = 62$$

$$RI = \frac{10 + 62}{10 + 34 + 30 + 62} = \frac{72}{136} = 0.52$$

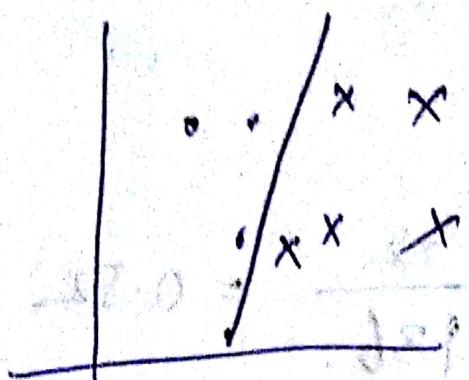
nMI (normalized mutual information)

Adjusted Ranked Index

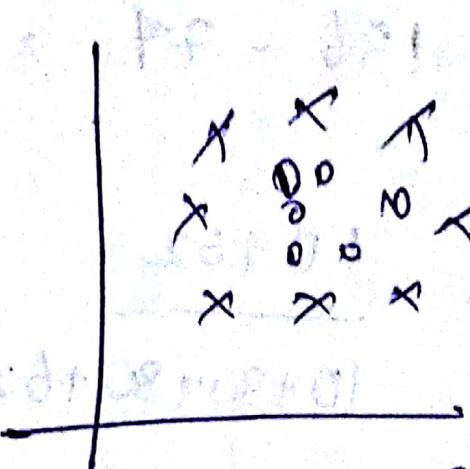
Maximum Adjusted Ranked Index

\hat{v} \hat{b}

Support Vector Machine



linearly Separable



non linearly Separable

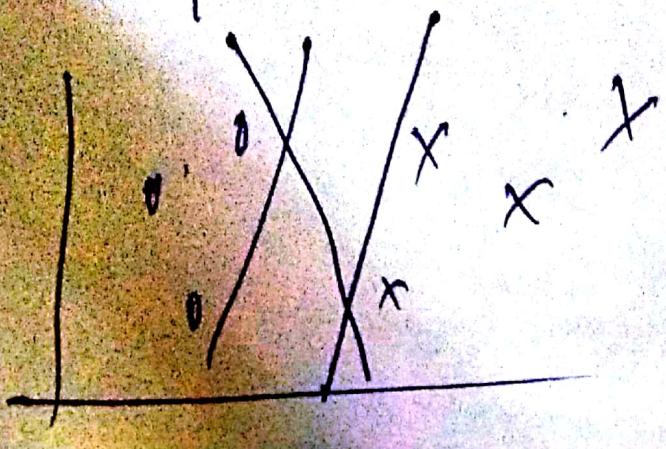
linearly Separable in

1D plane — line

3D plane — plane

\geq 3D plane — hyperplane

for linear Separable



many lines can be fit,
hence SVC gives best
fit line.

It is a kind of optimization problem.

Objective function

constraint

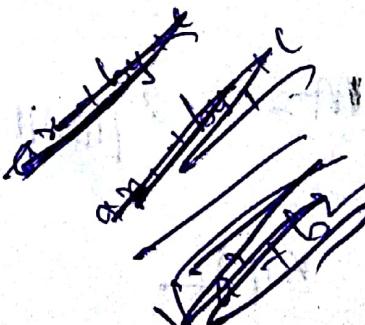
constraints

non negative
restrictions

simple
~~constant~~
constraints

Objective is to maximize the margin.

two arbitrary line in
equations || to each other



$$w^T x + b \geq 1$$

$$w^T x + b \leq -1$$

distance two lines =

$$\frac{|c|}{\sqrt{a^2+b^2}}$$

$$1^r \text{ dist} = \frac{\|b + w\|}{\|w\|}$$

$$\|w\| = \sqrt{w_1^2 + w_2^2}$$

$$1^r \text{ dist} = \frac{\|b - w\|}{\|w\|}$$

$$\text{distance} = \frac{\|b + w\| - \|b - w\|}{\|w\|}$$

$$= \frac{2}{\|w\|}$$

$m = \frac{2}{\|w\|}$, we need to maximize m .

or $m/2$ or minimize $\frac{2}{\|w\|}$

$$\text{minimize} = \frac{1}{2} \|w\|^2$$

$$\text{Subject to } y_i(w^T x_i + b) \geq 1$$

always > 1
because either
both should +ve
or -ve

y_i Class label of i^{th} instance
1 x_i i^{th} instance of data

Data = $\{x_1, \dots, x_n\}$, Class labels = {-1, 1}

↓
instances of Data. (Vector)

After solving optimization, then we get value of w & b. (use Lagrange Multiplier)

objective function — non linear
constraints — linear. } overall non linear.

Those α_i 's \rightarrow Lagrange coefficient ≥ 0 , corresponding examples are

called Support Vectors

↳ points on ~~w^T x + b = -1~~ $w^T x + b = +1$

$\dim(w) = \dim(x_i) = \text{no. of attributes}$