

→ Document classification

eg- Naive Bayes, Bernoulli

→ Role of doc. classifn in ITR

↳ entertainment
↳ politics

↳ doc 1
↳ doc 2

↳ mapping of doc.

- we have training & testing.
- Repo. of doc. with labels available
(eg- doc 1- sports etc.)

• when new doc comes, we can automatically classify entertainment or politics / others

- By creating clusters, we can reduce our search space.
- more relevant queries
- reduction of search time. optimizⁿ of retrieval func.
- eg- in email :- spam or not spam
- more precision & recall.
- if image, audio, videos etc. already classified, then helpful for specific query like - cat image
- Naive Bayes most popular for text classifn.

→ Supervised learning

• A doc. d

$$\text{classes} = C = \{C_1, C_2, \dots, C_j\}$$

• A training set D of doc. with a label from r.

Bernoulli Trial \Rightarrow 1 variable, 2 values \Rightarrow over
multivariate Bernoulli trial \Rightarrow multiple variables with over
for each test.

$$P(C/d) = \frac{P(C) \times P(d|C)}{P(d)}$$

$P(d)$ \hookrightarrow for Naive Bayes

Determine

A learning classifier to classify class y .

we assign class y to d
 $y(d) \in C$

we train a model to classify on test object / data

\Rightarrow we use Naive Bayes

$$P(A/d) = P(B)$$

$$P(C/d) = \propto P(C) \prod P(t_k | c)$$

prob. of
 d belonging
to class c

\nearrow multiplication of term
presenting
for all terms

\downarrow division term
ignored assume
for all $\frac{2}{5} > \frac{2}{5}$

If class c is popular, then nc/d is highly probable for newest doc d .
(eg 90 of 100 are C , so C is probable for test d). $\frac{90}{100} > \frac{2}{5}$

$N_d = \text{length of doc. (no. of tokens)}$

including $\hookrightarrow P(t_k | c)$ is conditional prob. of term
multiple occurring in doc. of class c

occurrences $\hookrightarrow P(t_k | c)$ as a measure of how much
evidence $\&$ t_k contributes that c is
correct class

\hookrightarrow we check for each term t_k how many
times it appears in all docs. of c

- Our goal in naive Bayes is to find best class.
- Best class is most likely or max^m a posterior (MAP) class. (map is calculated)
- use argmax.

$P(C) = \frac{N_c}{N} \xrightarrow{\text{no. of class } c} \text{Prior prob.}$

conditional prob. = $P(+|c) = \frac{T_{ct}}{\sum T_{cv} T_{ct}}$

$\xrightarrow{\text{eg. - appears}}$
 $3, 2, 8, 1 \text{ in my}$
 index 23.
 $\therefore t_{ct} = 3, n = 6$

• we make naive bayes ~~post~~al independence assumption

$$P(t_{k_1}|c) = P(t_{k_2}|c)$$

1st place or any other position \Rightarrow same prob.

• we need to add 0's conditional probability

due to 1 term which may not be in training phase, prob. becomes 0
 \hookrightarrow so need for Laplacian correction.

eg class a, b, c

$$P(c|w/a) = 0 \quad (\text{if } w \text{ not in } a)$$

So, we add +1 to numerator, & no. of classes in denominator.

$P(c) \rightarrow c = \text{Chinese Chinese Chinese Tokyo}$
Japan.

3 classes - c

1 class - j

$$P(c) = \frac{3}{4}$$

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

$$P(\text{Chinese} | c) = \frac{5+1}{5+14} \rightarrow \begin{matrix} \text{5 Chinese in c} \\ \text{14 total words in corpus} \end{matrix}$$

$$\frac{8+6}{c+14} \rightarrow \begin{matrix} \text{vocabulary (entire corpus)} \\ \text{c = 4 total Chinese in dict of c} \\ \text{total words in c} \end{matrix}$$

$$P(\text{Tokyo} | c) = \frac{0+1}{8+6} = \frac{1}{14}$$

$$P(\text{Japan} | c) = \frac{0+1}{8+6} = \frac{1}{14}$$

\rightarrow class appears 3 times.

$$P(c | d_s) \propto \frac{3}{4} \times \left(\frac{3}{7}\right)^3 \times \frac{1}{4} \times \frac{1}{4}$$

Similarly for $P(\text{Chinese} | j), P(\text{Tokyo} | j)$
 $\& P(\text{Japan} | j)$

Multiplication $\Rightarrow P(j | d_s)$

which ever greater is class.

$$P(w | c) = \frac{\text{count}(w, c) + 1}{\text{count}(c) + \text{Vocab}}$$

\rightarrow count(c) + Vocab

→ Bernoulli

- Doesn't see term freq. only presence or absence.
- $P(t/c)$ estimated as fraction of doc. of class containing term t .
- Probability of non occurrence is also calculated.
- Binary values put - 0, 1.

$$\text{eg. } P(\text{chinese}/c) = \frac{3+1}{3+2}$$

we calculate probability of each unique term in corpus to check if term may not appear.

$$\Rightarrow P(\text{chinese}/c) = \frac{\text{in how many it appears}}{\text{total classes}} = \frac{3+1}{3+2} \rightarrow \text{aparior}$$

$\text{total classes} = (C_{ij})$
 $\text{does } C_{ij}$
 $\text{of } c$

$$P(\text{Japan}/c) = P(\text{Tokyo}/c) = \frac{0+1}{3+2} = \frac{1}{5}$$

$$P(\text{chinese}/J) = \frac{1+1}{1+2}$$

similarly for others.

$$\therefore P(\text{city}/c) \propto P(c) \times P(\text{chinese}/c) \times P(\text{Japan}/c)$$

$$\times (-P(\text{macau}/c)) \times P(\text{Japan Tokyo}/c) \\ \times (1-P(\text{shanghai}/c)) \times (1-P(c) \times P(\text{Beijing}/c))$$

multinomial

Naive Bayes question

training set	doc words	in English?
1	Chinese Beijing Chinese	yes
2	Chinese Chinese Shanghai	yes
3	Chinese man	yes
4	Tokyo Japanese Chinese	no
test \Rightarrow	d 5 Chinese Chinese Tokyo Chinese Tokyo Japan	?

$$\rightarrow P(\text{Chinese}/c) = \frac{5+1}{8+6} = \frac{6}{14} = \frac{3}{7}$$

$$P(\text{Tokyo}/c) = \frac{0+1}{8+6} = \frac{1}{14}$$

$$P(c) = \frac{9}{14}$$

$$\therefore P(\text{Chinese}/d_s) \propto \frac{3}{4} \times \left(\frac{3}{7}\right)^3 \times \frac{1}{14} \times \frac{1}{14} \approx 0.0003$$

$$\rightarrow P(\text{Chinese}/\bar{c}) = \frac{1+1}{3+6} = \frac{2}{9}$$

$$P(\text{Tokyo}/\bar{c}) = \frac{1+1}{3+6} = \frac{2}{9}$$

$$P(\text{Japan}/c) = \frac{1+1}{3+6} = \frac{2}{9}, P(\bar{c}) = \frac{1}{4}$$

$$P(\text{Chinese}/d_s) \propto \frac{1}{4} \times \left(\frac{2}{9}\right)^3 \times \frac{2}{9} \times \frac{2}{9} \approx 0.0001$$

$\therefore c$ is more probable

• Laplacian correction is for smoothening.

|| Page No. _____
Date : _____ ||

steps

From training corpus extract vocabulary calculate req. $P(c_j)$ & $P(x_k | c_j)$ terms.

For each class c_j in C_{cls} :

- $\rightarrow d_{cls} \leftarrow$ subset of docs. for which the target class is c_j

$$P(c_j) = \frac{(d_{cls})}{\text{total docs}}$$

\Rightarrow class of class c_j
total docs

$$P(x_k | c_j) = \frac{n_j + 1}{n_j + \alpha}$$

where

$K = 1$ to all unique words
in query

\Rightarrow in how many docs.

\Rightarrow for Laplacian correct

$n_j + \alpha$ Vocabulary

\Rightarrow no of unique words in corpus.

positions = all word posn in current doc. which contain tokens found in vocab.

we calculate $\Rightarrow C_{NB} = \operatorname{argmax}_{C_j \in C} (P(c_j)) \times \prod_{i \in positions} P(x_i | c_j)$

$\{ C_j \in C \}$

\Rightarrow for class

having

maxth value.



Bernoulli Naïve Bayes classifier

• considers presence & absence of words (0 or 1)

• 2 possible values for a word \Rightarrow 0 or 1.

$$P(x_i | c_j) = \frac{N(\text{or } x_i = 1, c = c_j) + 1}{N(c = c_j) + k}$$

↗ no. of time word appears
 ↗ index containing label c
 ↗ naptorian
 ↗ correction

↗ no. of values for
 ↗ a word
 ↗ more 2 - 0 or 1)

↗ no. of
 ↗ does with
 ↗ label c

	doc IP	words	in $c = \text{China}$?
training set	1	Chinese Beijing Chinese	Yes
	2	chinese chinese Shanghai	yes
	3	chinese Macao	yes
	4	Tokyo Japan Chinese	no
test set	5	Chinese Chinese Chinese Tokyo Japan	?

$$P(\text{Chinese} | c) = \frac{3+1}{3+2} = \frac{4}{5}$$

↗ Chinese appears in 3 documents
 ↗ total
 ↗ appears or not
 ↗ doc of c (0, 1)

$$P(\text{Tokyo} | c) = \frac{0+1}{3+2} = \frac{1}{5}$$

$$P(\text{Japan} | c) = \frac{0+1}{3+2} = \frac{1}{5}$$

$$P(\text{Beijing} | c) = \frac{1+1}{3+2} = \frac{2}{5}$$

$$P(\text{Macao} | c) = \frac{1+1}{3+2} = \frac{2}{5} = P(\text{Shang...})$$

$$P(c) = \frac{3}{9} \therefore P(c | d_5) \propto \frac{3}{9} \times \frac{4}{5} \times \frac{1}{5} \times \frac{1}{5} \times \left(1 - \frac{2}{5}\right) \times \left(1 - \frac{2}{5}\right) \times \left(1 - \frac{2}{5}\right) \approx 0.005$$

$$P(\text{Chinese}/\bar{c}) = \frac{1+1}{1+2} = \frac{2}{3}$$

$$P(\text{Tokyo}/\bar{c}) = \frac{1+1}{1+2} = \frac{2}{3} = P(\text{Japan}/\bar{c})$$

$$P(\text{Beijing}/\bar{c}) = P(\text{Macau}/\bar{c}) = P(\text{Shanghai}/\bar{c}) \\ = \frac{0+1}{1+2} = \frac{1}{3}$$

$$P(c) = \frac{1}{9}$$

$$\therefore P(\bar{c}/d_5) = \frac{1}{4} \times \frac{2}{3} \times \frac{2}{3} \times \frac{2}{3} \times \left(1 - \frac{1}{3}\right) \times \left(1 - \frac{1}{3}\right)$$

$$\approx 0.022$$

$\therefore [\bar{c}]$ is predicted for d_5

→ How to improve relevance through feedback

- Objective is to provide relevant docs to user
- we want to improve relevance of docs \Rightarrow goal
- eg. purchase implies high rating
 - if we click on retrieved result, like us
 - click 4th doc of out of 5 retrieved, then 1st 3 are not relevant, but 4 was.
- This is a kind of feedback
- we do query expansion
- we want to provide better relevance to users
- Improved relevance model, can be provided to same or other users, when same query was asked.
 - L) all this is query expansion.

eg - we ask "plane"

- results :-

= flight
= aircraft

= by air
= aeroplane

= by air

= fly

= fight figt

= arcrft.

(synonyms as well as wrongly spelled ones are retrieved)

→ How to ensure good results

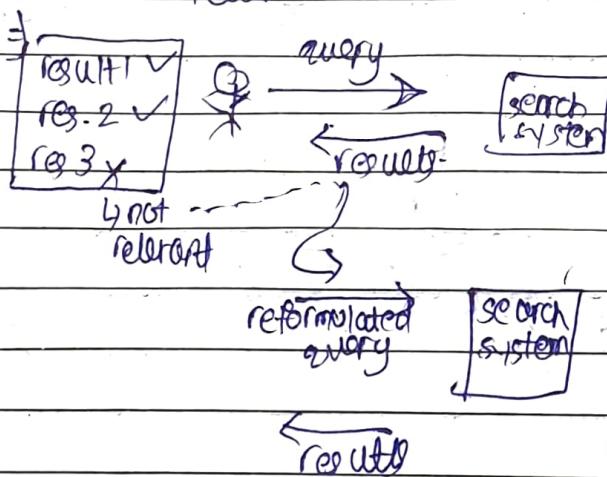
automatic
query
refinement

- we use query or refined results to reformulate & improve every time for relevance eg - pseudorelevance
- indirect relevance feedback
- relevance feedback

global

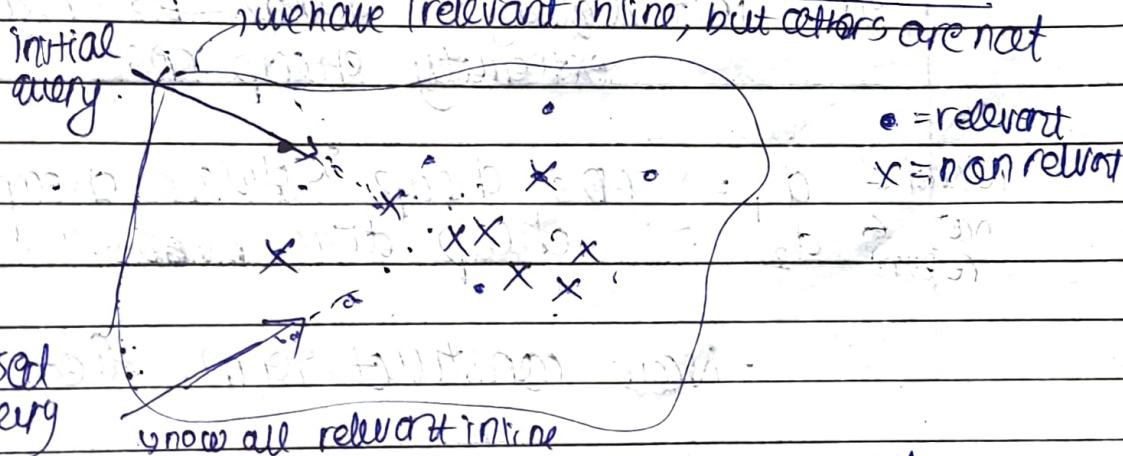
- we do not use query or results to improve & reformulate every time for relevance

eg - the query spelling correct



- Refmet takes explicit feedback from user
- It asks if result is relevant or not.
- But its bad
- user doesn't have time & patience to provide explicit feedback
- But, still it is one of the way
- Indexed content is unknown to user
 - ↳ we don't know, if other relevant docs. are present but not shown
- "Inform" need changes after looking at the results.
 - ↳ eg- we search for a specific song. But looking at song results, we change to our mind & search other song.
 - ↳ so need changes & its necessary to do query expansion.

Rocchio algorithm (Relevance feedback (local))



can we shift our query from non relevant to relevant docs? → do query expansion

- In practice, we only know few relevant & non-relevant.

\vec{q}_0 = initial query vector
 D_r = set of relevant docs $\beta \uparrow$ its weight
 D_{nr} = set of non rel. docs $\beta \downarrow$ its weight

$$\vec{q}_m = \alpha \vec{q}_0 + \beta \frac{1}{|D_r|} \sum_{d_j \in D_r} d_j - \gamma \frac{1}{|D_{nr}|} \sum_{d_j \in D_{nr}} d_j$$

α, β, γ are weights.

- we represent query as weighted vector
- use ranking formula to get new query vector
- Now calculate cosine similarities
- Reranking model increases recall
- +ve feedback is much more valuable than -ve feedback

eg - initial query = " cheap CDs cheap DVDs
extremely cheap CDs "

relevant d_1 = " CDs cheap Software cheap DVDs "
 not relevant d_2 = " cheap drills DVDS "

- Now construct term. if query vector

	Cheap	CDS	+ -
q_0	3	-	
d_1	2	-	
d_2	1	-	

eg - to move $(2, 2)$ to $(5, 2)$, add 3 to x.

$$q_{\text{new}} = 1 \times q_0 + 0.75 \times d_1 - 0.25 \times d_2$$

do this for each term.

$$\begin{aligned} \text{eg- } q_{\text{new}} &= 1 \times 3 + 0.75 \times 2 - 0.25 \times 1 \\ &\quad (\downarrow) \text{for } d_2 \\ &= 4.25 \end{aligned}$$

$q_{\text{new}} =$	4.25	-	+ -	
--------------------	------	---	-----	--

we can't have $q_{\text{new}} = -0.5$, this is wrong, so -ve values are put as 0.
 $\therefore q_{\text{new}} = 0$

→ Pseudo Relevance (or Ind.) Relevance Feedback

- No user judgement
- we completely trust our web search engine
- we assume top k given by model are relevant to query
- now we can use ranking algorithm to change the query.

→

problem

- we search "Dhoni retires from 1st class cricket". now, southafrica tests, etc. are also relevant hence retrieved.
- so matrix includes southafrica, newzealand etc.
- The next query w vector will thus drift. so, query drift problem occurs, as more weightage to terms which are not relevant.

→ Indirect (implicit) relevance Feedback

- we do not ask explicit feedback.
 - we consider clickstream mining: userclicks on nth doc, then top $n-1$ not relevant, & n is relevant
- query expansion would be done based on implicit user preference (clicks)

→ In global approaches, we do not consider either query or retrieved results.

Global (User) result - independent

- automatic thesaurus gener'n.

eg - fast = rapid

four = boundary \Rightarrow cricket

sound = noise

tall = height

- how to improve relevance?

- slangs?

\hookrightarrow brother \approx bro

- we can do co-occurrence analysis for the doc.

eg - mainframes
(doc)

mainframe
(doc?)

primarily, we
are ...

usually, we
refer

primarily & usually
are used like
synonyms.

eg - term incident matrx :-

words	BOOK 1	B2	B3	B4
w1	100	1	110	2
w2	50	169	47	8
w3	25	3	24	4
avg	75	6	73	7

Book 1 & 3,
Book 2 & 4
are similar.

- So, we can do co-occurrence analysis
- 2 terms are similar if term vectors are similar.

- Context is important
 - ↳ same video may not be relevant on weekend with but relevant on week.

eg - term-doc matrix

term-doc matrix • term doc matrix (Transpose)
= term-term matrix

from term-term matrix,
we can see results :-

	t_1	t_2	t_3	t_4	t_5	t_6
f_1	3	-	1	2	1	2
f_2	1	-				1
f_3	2	1	-	2		2
f_4	1	2	1	-	1	1
f_5	2	1	2	1	-	2
f_6	3	1	2	1	2	3

similar

so, $f_1 \approx f_6$

→ PRECISION & RECALL

- Precision - fraction of retrieved docs. are relevant.

$$P(\text{relevant}) = \frac{\text{relevant}}{\text{retrieved}}$$
- Recall - fraction of relevant docs that are retrieved.

$$P(\text{retrieved}) = \frac{\text{retrieved}}{\text{relevant}}$$

		true		false +ve → confusion matrix
		relevant	non-relevant	
retrieved	relevant	tp	fp	→ false +ve
	not-retrieved	fn	tn	→ false -ve
				→ confusion matrix

• Precision $P = \frac{tp}{tp+fp} \rightarrow$ correct prediction
 $\frac{tp}{tp+fp} \rightarrow$ all retrieved docs.

• Recall = $\frac{tp}{tp+fn} \rightarrow$ correctly predicted
 $\frac{tp}{tp+fn} \rightarrow$ total relevant in corpus

• There has to be balance b/w precision & recall

• If we try to increase one, other may decrease

e.g. if we retrieve more docs for more recall, but precision may decrease, & vice versa

• Diff. b/w recall & precision must not be more than 5%.

e.g. if precision = 95%, recall = 10%, then diff. $> 5%$

- For recommender systems,

	Item 1	i2	i3	in
User 1	~	.	.	.
U2	~	~	~	.
U3	~	~	~	.
:	~	~	~	.
Un	~	~	~	~

→ we try to predict each cell rating.

- Then we calculate error by $\sqrt{\sum (\text{actual rating} - \text{predicted rating})^2}$ (Root mean square error)
- This is our accuracy parameters.

e.g. IRS retrieved 20 docs.

there are 100 relevant in corpus.
8 are correct & relevant

| R |
| T |
→ correct

$$\text{Precision} = \frac{8}{20}$$

similarly
20 cells

$$\text{Recall} = \frac{8}{100}$$

| R |
R
-
→ here
all the
R are
correct

8 Retrieved docs → RRNNNNNNN RNRNNNR
NNNNR

$$\text{precision} = \frac{6}{20}$$

recommended

• let 8 relevant docs (total)

$$\therefore \text{recall} = \frac{6}{8}$$

→ better matching & better recommendation.

→ better recommendation.

F-measure

• Takes both precision & recall into account

$$\therefore F = \frac{2 \cdot P \cdot R}{P + R} = \frac{2}{\frac{1}{P} + \frac{1}{R}}$$

Q - why F-measures take harmonic mean?
(not arith. or geometric)?

A - arithmetic - biased towards high values.
(tail)

the geometric mean - ?

→ eg -

R	R	R	P
R	R	R	P

 case 1

R	R	R	R
R	R	R	R

 case 2

↳ Both have same precision & recall, but case 2 is better,
as 1 to 8 are ranked
relevant and ^{case 2} correctly
retrieves ranking wise.

→ Solution - precision at 8th index

1 2 3 4 5 6 7 8 9 10



$$P@1 = \frac{0}{1} = 0$$

retained
ICell

← At 1st index / location

$$P@2 = \frac{1}{2}$$

I | R

← At 2nd index

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

I | R | R

$$P@4 = \frac{2}{4}$$

I | R | R |

→ Interpolated precision

cutoff at k^{th} relevance level.

1st R

$$P@1 = \frac{1}{2}$$

I | R

eg - if total
5 relevant,

1st locatn
where R appears

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

2nd R

then 5
levels.

Avg. precisⁿ =

$$\frac{0.5 + 0.66}{2}$$

1st locatn
= where
first R
appears,

similarly
2nd locatn
where 2nd
R appears

& likewise.

* case 1 → To calculate Avg. precisⁿ

I | R | R | R | R | R

$$\text{Avg. precis}^n = \frac{\frac{1}{2} + \frac{1}{2} + \frac{1}{2} + \frac{1}{2} + \frac{1}{2}}{5}$$

$\therefore \text{If 10 docs,}$

$$\text{avg} = \frac{\frac{1}{2} + \frac{1}{2} + \dots + \frac{1}{10}}{10} \text{ times}$$

case 2 \rightarrow For Avg. precⁿ (AP) at level 4,

$$\frac{P@1 + P@2 + P@3 + P@4}{4}$$

If multiple queries we can average MAP of the 2 queries
 \rightarrow MAP (mean avg. precision) is important - calculate every precⁿ for all queries
 eg - Q1 Q2 Q3 Q4 Q5 Q6 Q7 Q8 Q9 Q10
 5 levels of relevance

$$\therefore \frac{P@1 + P@2 + P@3 + P@4 + P@5}{5}$$

$$= \frac{\frac{1}{1} + \frac{2}{3} + \frac{3}{6} + \frac{4}{9} + \frac{5}{10}}{5}$$

$$= \frac{1 + 0.66 + 0.5 + 0.5 + 0.44}{5} = 0.62$$

query-2) Q1 Q2 Q3 Q4 Q5 Q6 Q7 Q8 Q9 Q10
 73 levels of relevance

$$\frac{P@1 + P@2 + P@3}{3}$$

$$> \frac{\frac{1}{2} + \frac{2}{5} + \frac{3}{7}}{3} = \frac{0.5 + 0.4 + 0.428}{3} = 0.428$$

$$\therefore \text{final MAP} = \frac{0.67 + 0.44}{2} = 0.53$$

Q- case 1: P@+

Q- every / =