

Information Retrieved

flavomile

Data
Sheet

kids → 20%, Smals → 40%, Assignments: 20%, Attendance: 20%

course material

First part: Manning's book

What is IR?

documents / web documents text (for this course)

→ Finding material of an unstructured nature that satisfies "information need" from within large collections of data.

query

Writs of documents

What are potential documents

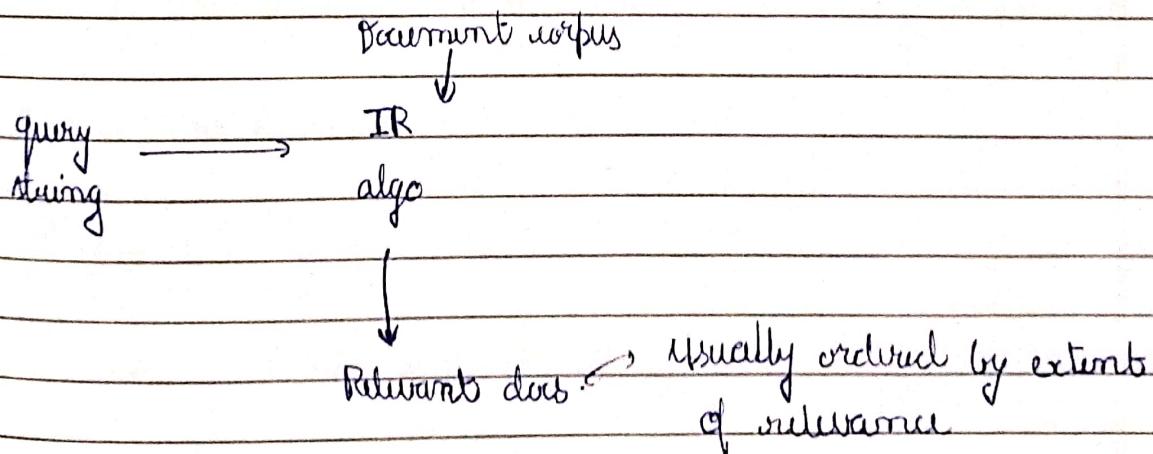
→ Legal documents / web pages / research papers

Challenges in IR:

→ Query & Data its outcome is unstructured & hence understanding the NL query.

→ Erroneous queries & understanding the intent of the query.

Components of query system



What does relevance mean?

- The response should be on the proper subject (topic).
- Response should be timely (i.e. not outdated response)
- Responses should be authoritative (should be of good value).
- Relevant to the information need ~~of~~ / intent of query.

Most vanilla form of IR:

Keyword search (string matching) : Match the query & document string (Moderately strict keyword matching)

variants of keyword matching

→ Word relations absent : Ex PRC & China ~~are the~~ the same things

→ Polysemy : (context based word meanings)

Ex: Apple \leftarrow fruit
 \leftarrow company.

: IR \rightarrow requires understanding of the intended meaning of word : requires context.

→ Not just syntax but semantics also important.

Boolean retrieval

Assumption:

- i) User interacting with an IR system.

- ii) Document : unit of retrieval

- iii) Doc ID : unique identifier of each document.

- iv) Corpus : collection of documents.

- v) User : has an info need

Expressed as

query $\xrightarrow{\text{Sent to}}$ IR system.

- vi) Term : Unit of information (word / phrase etc)

- vii) IR system : Tasked with retrieving relevant documents.

Running example for Boolean retrieval :

Corpus : All plays written by Shakespeare

Document : Individual play

Query : Boolean expression of terms connected by
Boolean operators (\wedge , \vee , \sim etc)

Ex: Query : Which all plays of Shakespeare contain the words
Brutus AND Caesar but not Calpurnia.

Possible solution: Get all Shakespearean plays containing Brutus
AND Caesar -

for line by line & omit all plays containing Calpurnia.

Cons of above solution

① Very slow

NOT Calpurnia is slowing things down.

ii) Queries like word A appearing near word B (Ex: Romans appears within 3 words of countrymen) is very hard to frame.

iii) Ranking not possible.

Term incidence matrix

A matrix of terms * documents

Ex- Documents →

Antony & Antony & Julius & The Hamlet Macbeth
Cleopatra Cleopatra Caesar Tempest

Terms	Antony		0	0	1
-------	--------	--	---	---	---

J	Brutus	1	1	0	1	0
---	--------	---	---	---	---	---

	Cesar	1	1	0	1	0
--	-------	---	---	---	---	---

	Cleopatra	1	0	0	0	0
--	-----------	---	---	---	---	---

	Calpurnia	0	1	0	0	0
--	-----------	---	---	---	---	---

	Mary	1	0	1	1	1
--	------	---	---	---	---	---

$T_{ij} = 1$ indicates document j contains the term i.

Query of form Brutus AND Cesar AND NOT Calpurnia.

can be processed ~~using~~ using the AND of bit strings of appropriate terms.

Brutus AND Cesar AND NOT Calpurnia.

$$\equiv [11\ 010] \wedge [11010] \wedge [10\ 111]$$

$$\equiv [10010]$$

Problems with term incidence matrix

1) Size: Not scalable to something like the web.

Ex: 1M docs, each doc on avg: 1000 words & we need 6 bytes per word, then total document corpus size: 6GB.

If there are 500k distinct terms across documents, then size to store the table = $500k \times 1M \times \text{size/entry} = \frac{1}{2} \text{ trillion} \times \text{size/entry}$

Possible solution: Observe that the matrix is sparse (Each col has on average only 1000 1's). \therefore out of $\frac{1}{2}$ trillion entries, only 1B entries are informative

Inverted index

For each term T, list of all documents that contain T

called inverted as terms point to documents.

Inverted index example:

Brutus \rightarrow 1, 2, 4, 11, 31, 45, 173, 174.

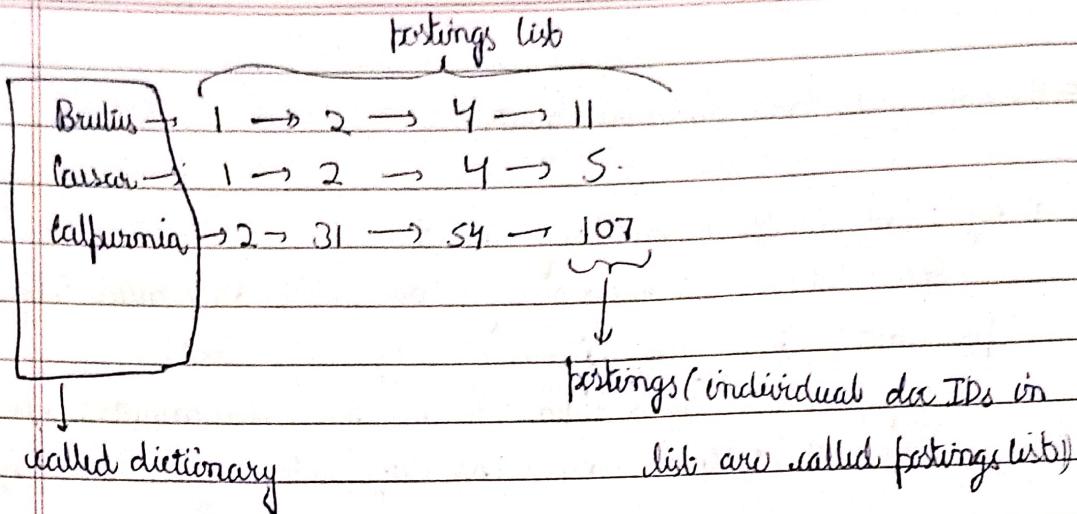
Caesar \rightarrow 1, 2, 4, 5, 6, 16, 57, 132.

Calpurnia \rightarrow 2, 31, 54, 101

Some questions:

i) Are the doc IDs contiguous or not?

In order to facilitate addition of new docs, linked list used.



ii) Where to store

Dictionary is usually stored in the main memory.
posting list in disk.

iii) How to go from Raw text doc to Inverted index

a) Tokeniser: Splits running text into individual tokens

(more about this later) (Normalisation)

b) Language models: Does some basic processing on words

(decapitalisation, prefix-suffix removal & so on)

c) Indexer: Builds index

e.g.: Friends, Romans, countrymen

↓
Tokeniser

Friends Romans Countrymen

↓ LM (process called stemming)

called stems ← Friends Roman countrymen

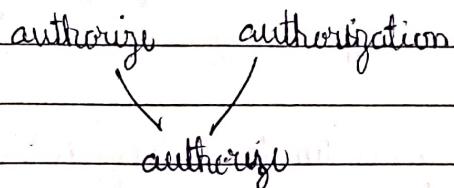
Ex:

Tokunization

state - of - the - arb → state of the arb

Normalization

U.S.A = USA

StemmingStopwords /# Stopwords / Function words

→ Words that do not have a meaning of their own (Ex: a, an, the etc). Occur very frequently.

→ Usually not indexed

Indexer

We have the tokens & the doc ID.

Tokun	Da.ID
g	1
the	1
;	
caesar	2
was	2
cautious	2

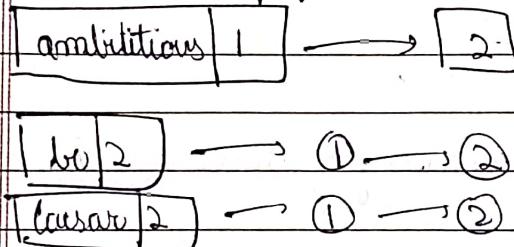
Two stages sorting is done
 (i) Lexicographically on terms
 ii) By doc ID

Firstly
 time sort
 :
 Caesar → 1 second,
 Caesar → 2 time
 :
 I → 1

ambitious → 2
 bc → 1
 bc → 2
 :
 Caesar → 1
 Caesar → 2
 Caesar → 2.

After 2 time sorting, we sort document frequency of term
 (no. of unique doc the term appears in)

Now this sorted list is converted to postings list.
 From Doc freq Postings list.



Processing query Brutus AND Caesar from inverted index

- i) Locate Brutus in dictionary
- ii) Retrieve postings list of Brutus from disk
- iii) Locate Caesar in dictionary
- iv) Retrieve postings list of Caesar from disk
- v) Find the AND of the 2 lists using the merge algorithm -

Time required = $O(|X| + |Y|)$ / works as the lists are stored in sorted fashion.

Merge (P_1, P_2)

ans = []

while $P_1 \neq \text{NULL}$ & $P_2 \neq \text{NULL}$:

if $*P_1 < *P_2$:

$P_1 = P_1 \rightarrow \text{next}$

else if $*P_1 == *P_2$

ans += $*P_1$

$P_1 = P_1 \rightarrow \text{next}$

$P_2 = P_2 \rightarrow \text{next}$

else

$P_2 = P_2 \rightarrow \text{next}$

Boolean retrieval system views any document as a set of terms
 → Still no ranking system, but still very popular.

Ex: What is the statute of limitations in cases involving
 Federal tort law acts?

Can be written as Boolean query:

LIMIT! /3 STATUTE ACTION/S FEDERAL /2 TORT /3 CLAIM

wildcard that

mimics all

max 3 words

STATUTE & ACTION

must be in same

variations of

sentences

words

#

Why doc frequency required

Ex: Mingo Brantes AND Califurmia AND Caesar

Brantes $\rightarrow 1 \rightarrow 4 \rightarrow 5 \rightarrow 16 \rightarrow 32 \rightarrow 64$

Caesar $\rightarrow 1 \rightarrow 2 \rightarrow 3 \rightarrow 6 \rightarrow 9 \rightarrow 45$

Califurmia $\rightarrow 13 \rightarrow 40$

Optimal sorting order merge order (Califurmia AND Brantes) AND Caesar.

This is known by the doc freq. (Merge the smallest 2 doc freq lists every time)

#

What should be order of OR queries?

Ex:

(madding or covid) AND (ignoble or strip)

Doc freq of OR is atomic sum of doc freqs of individual terms.

\therefore sort by sum of doc freqs & process in increasing order.

#

Google also uses boolean retrieval.

$[w_1, w_2, \dots, w_n] \equiv w_1 \text{ AND } w_2 \text{ AND } \dots \text{ AND } w_n$

Problems

- Very few docs with all words
- Word might be in anchor term & not doc
- Synonyms

Term vocabulary / Tokenisation & normalisation in IR.

In order to deal with documents, we need to process documents (parse documents).

Suppose to converting byte sequences into a linear sequence of characters.

Issues:

- i) Different formats & languages
 - ii) Different encoding
 - iii) Character sets in use
- } Each of these is a classification problem.

But generally, we use heuristics to deal with these problems rather than full fledged classification.

The term document itself is well defined:

- i) Is document a single file or collection of files.
 - ii) Is document a single email, chain of emails, emails with attachment
- } Design choices that depend upon downstream applications.

Some definitions we'll use

i) Word:

Delimited string of characters as it occurs in running text (corpus)

ii) Term:

Normalised word.

iii) Token

An instance of a word.

~~Input~~

An instance of distinctiveness of a word (or token)

In the inverted index construction

Input:

Text -> [Friends, Norman, Countryman...]

Text -> [It is the with lesson...]

The needed terms
to build posting list.

Output: tokens

[Friends] [Norman] [It] [is]

Steps to go from input to output:

Tokenization: splitting of string of characters into tokens/words.

Usually done with splitting using whitespace delimiter

Given:

i) One word or many:

Hewlett-Packard

State-of-the-art.

San Francisco.

Things like: San Francisco - Los Angeles

★ Design issue: consider these as single token (will

have to build a dictionary with these multi-word tokens)

or break them into multiple tokens

ii) Numbers

- a) Differentiating b/w various formats of dates.
- b) Differentiating b/w phone numbers & IPs.

iii) Languages without delimiters (like urdu)

iv) Language specific names

v) Morphology of words

vi) New words

vii) Accents

Do we store Universität & Universitaet differently? (Hiragami)

→ handled usually using careful observation of user: Users usually drop accents while querying. ∴ drop accents while indexing.

viii) Case folding: Usually indexing is done after blindly converting everything to lowercase.

Ex: MIT & mit.

feel vs fed:

Normalisation: we want to identify query terms that have the same form.

Ex: U.S.A & USA.

~~Automobiles~~ cars & automobiles } Ideally we would want these in the same index

Requires generation of equivalence classes among terms. Can be done using handwritten rules or asymmetric expansion.

Asymmetric expansion:

window → window, windows

windows → windows, Windows

But Windows → X restricted to offend. (-: asymmetric)

Asymmetric expansion used to replace query words with equivalents

Stopwords & their removal

→ Stopwords: words that do not have a meaning of their own but rather confer meaning to other words.

Ex: Articles, prepositions, conjunctions, auxiliary verbs.

→ Occur very frequently in a document. Some of them are frequently occurring.

→ Typical design policy is to not index these words.

Caveats: Phrase queries like "Joan of Arc" can't be removed

because provides a greater purpose than just stopword

Stemming & Lemmatization (Not Morphological processing)

→ Converting the word to its root form.

↳ called stem / lemma.

Lemmatization: proper reduction of the word to its lemma.

Usually by matching the headword from a dictionary.

Stemming: just chops the end of the word & assumes the stem has been obtained. Usually done using heuristics.

is, am, are $\xrightarrow{\text{lemmatize}}$ be
car, cars, car's, cars $\xrightarrow{\text{stemming}}$ car.

automation, automatic... $\xrightarrow{\text{stemming}}$ all automat

↓
automat is a good approximation for automation

Stemming usually good

Enough for practical purposes.

Inflectional morphology

Word gets an inflection. POS remains same.

Ex - car → cars
 ↓ |
 noun noun.

Derivational morphology

POS changes.

Ex - automobile → automation
 ↗ verb noun

Stemmer: Stemming algorithms are called stemmers.# Porter stemmer

Only 1 rule applies at each stage. Rule decided by longest suffix match.

Stage 1

sses → ss (caresses → caress)

ies → i (ponies → pony)

ss → ss (ss → carress → caress)

s → φ (days → day)

caress matches both s & ss - ss is chosen.

Stage 2

(★ v ★) ing → φ (walking → walk)
(★ v ★) ed → φ (played → play)

Stage 3

atmational → ati

inger → ige

ator → atu

:

Stage 4

ati → φ (survival → survi)

atu → φ (activitate → activ)

able → φ (adjustable → adjust)

:

These are just a few examples of such rules.

These are not the only rules. (these are just examples).

There are many more rules that aren't mentioned here.

#

Skip pointers

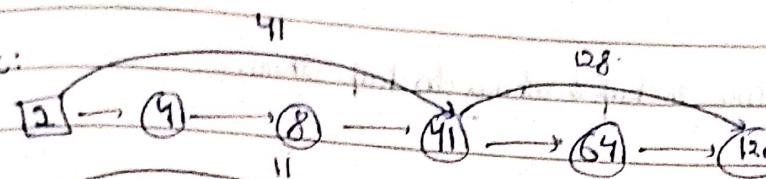
Used to make fastening lists operations efficient.

→ We saw the merge operation took $O(X+Y)$ time.

→ We can make it faster by augmenting the list with skip pointers.

→ Skip pointers can be used to skip over all the nodes that do not feature in the search.

Ex:



Modified merge algorithm using skip pointers:

Intuitively, (p_1, p_2)

answer = []

while $p_1 \neq \text{NULL}$ & $p_2 \neq \text{NULL}$

if $p_1 \rightarrow \text{desc.ID} == p_2 \rightarrow \text{desc.ID}$

 answer.append($p_1 \rightarrow \text{desc.ID}$)

$p_1 = p_1 \rightarrow \text{next}$

$p_2 = p_2 \rightarrow \text{next}$

else:

 if $p_1 \rightarrow \text{desc.ID} < p_2 \rightarrow \text{desc.ID}$:

 if hasskip(p_1) & & $p_1 \rightarrow \text{skip} \rightarrow \text{desc.ID} \leq p_2 \rightarrow \text{desc.ID}$:

$p_1 = p_1 \rightarrow \text{skip}$

 else

$p_1 = p_1 \rightarrow \text{next}$

 else:

 if hasskip(p_2) & & $p_2 \rightarrow \text{skip} \rightarrow \text{desc.ID} \leq p_1 \rightarrow \text{desc.ID}$:

$p_2 = p_2 \rightarrow \text{skip}$

 else

$p_2 = p_2 \rightarrow \text{next}$

How many skip pointers do we keep & where do we keep them?

Depends on the designer.

Tradeoffs involved in no. of skip pointers:

More skips

Less skips

- shorter skip spans (so will usually choose skips)
- longer skip spans (so will usually not choose skips)
- more comparisons
- fewer pointer comparisons.

Where to place the skip pointers?

If the list size is L then keeping \sqrt{L} skip pointers evenly spaced gives empirically the best performance (Moffat & Zobel 1996).

Easy if L is not growing frequently (would require pointer redistribution otherwise).

Note

Skip pointers are generally only used when partitions don't change frequently.

More memory means more skip pointers means more memory required.

Problems when loading the list onto memory for merging. I/O cost might surpass the benefit of having skip pointers.

#

Phrase queries

Queries that enable us to do searches like [Stanford University] together as a unit (instead of searching Stanford & University) AND

Boolean queries won't work here because Stanford AND University might bring up documents about a guy named Stanford that went to University & documents regarding other universities in Stanford.

Ways to resolve this problem

i) Biword indexing: Indexes bigrams (2 consecutive words)

Ex: "Friends, Romans, Countrymen" will have biword index entries "friend roman", "roman countrymen".

Each of these biwords are added to this dictionary.

Queries like "Stanford University Palo Alto" can be broken as "Stanford university" AND "Palo Alto".

Layouts: Phrases like "abolition of slavery" cannot still be processed. (as this cannot be broken into bigrams).

Stemmed biwords: Requires POS tagging of query & document.

Buckets all the nouns (N) and all articles & prepositions (X)

Extended biwords: $N \times^* N$

Ex: catharsis in the city
N X X N

King of Biwords
N X N

Now the vocabulary contains all extended by biwords, triwords & single words.

Drawbacks: → Index would easily blow up.

→ Queries with greater than 3 nouns can't be processed.

Note

If we're using extended biword ~~for~~ indexing, we can't skip our own stopwords.

i) Positional index (reconstruction of postings list)

Postings list in a non-positional index is just a list of doc.ID.

But in positional index, each postings is a docID + a list of positions.

Eg: We have a query "to be or not to be". → (to, 1) is different than (to, 5).

Postings list changes as: "to be or not to be" → doc freq.

to : 993427; → doc freq.

1 : < 7, 18, 33, 52, 86, 231 >; → doc freq.

2 : < 1, 17, 212 >; → doc freq.

4 : < 8, 16, 90, 400, 429, 433 >; → doc freq.

5 : < 363, 367 >; → doc freq.

7 : < 13, 23 >; → doc freq.

be : 138239; → doc freq.

1 : < 17, 25 >; → doc freq.

4 : < 17, 191, 291, 400, 429, 434 >; → doc freq.

5 : < 14, 19, 101 >; → doc freq.

Most likely matches ab docID: 4; word: 429.

Merge method

- i) Get postings list of to, by, or, not.
- ii) Progressively intersect the postings list of words starting with to, be.

(Position list could be merged by a slight modification of merge algorithm (we check for difference instead of equality))

Biword indexing used when bigrams are extremely frequent, otherwise we use positional index.

Multiword phrases with bigram indexing:

State of the art = "State of;" AND "of the" AND "the art"

Document structure:

→ Some fields in the document need to be given special treatment

Ex: Dates, titles: Main titles/ section titles.

→ These fields are usually indexed separately.

One of the ways: Bitmap lists

Format: Bitmap

Precomputed scores

Postings list now contains Doc ID & also relevance score of that word for the document.

Ex: fish → [1, 3.6] → [3, 2.]

How does this help?

While merging, ignore postings with relevance less than a certain threshold. Reconstruct the postings list after removing these entries \rightarrow smaller postings list \Rightarrow faster merge.

Suppose we query "tropical fish aquarium".

Some questions that come up are:

- i) What is the size of the result set for this query?
i.e. for a query "a b c" how many docs contain all 3 of these terms?

Let f_a = frequency of "a" in the corpus. f_b & f_c .
Assuming a, b & c are independent & N is the number of tokens in the corpus,

Expected frequency of a, b & c together is:

Probability of a's occurrence = f_a , probability for b & c is

$$\frac{f_b}{N} \text{ & } \frac{f_c}{N}$$

$$\text{Probability of occurrence of a b c} = \frac{f_a f_b f_c}{N^3}$$

$$\text{Expected no of documents} = \frac{N(f_a f_b f_c)}{N^3}$$

$$= \frac{f_a f_b f_c}{N^2}$$

* Independence assumption leads to very poor results.

One solution is to use document frequency of biword queries.

$$P(a \cap b \cap c) = P(a \cap b) \cdot P(c | a \cap b)$$

$$P(a \cap b \cap c) = P(a \cap b) \cdot \min(P(c|b), P(c|a))$$

$$P(a \cap b) = \frac{\text{docs}(\{a, b\})}{N}$$

\hookrightarrow no of documents.

$$\begin{aligned} E(f_{\text{biword}}) &= N \times P(a \cap b) \times \min(P(c|b), P(c|a)) \\ &= f_{\text{biword}} \times \min\left(\frac{f_{\text{c} \cap \text{b}}}{f_b}, \frac{f_{\text{c} \cap \text{a}}}{f_a}\right). \end{aligned}$$

$$\therefore f_{\text{aq, tropical} \cap \text{fish}} = f_{\text{aq, tropical}} \times \min\left(\frac{f_{\text{aq} \cap \text{fish}}}{f_{\text{aq}}}, \frac{f_{\text{tropical} \cap \text{fish}}}{f_{\text{tropical}}}\right).$$

Still tedious as requires bigram frequencies.

A simpler estimate

If 26480 documents have the word aquarium. Select some random 3000 docs out of it. Let 258 of them have tropical, fish & aquarium all 3 phrases. Then we estimate about $\frac{258}{3000} \times 26480$ documents will have tropical, fish & aquarium.

Retrieved documents must not only be relevant to the query but the intents of the user (even if the query may not accurately represent the intent) (i.e. information needs of the user)

We desire relevant retrieval & not exact retrieval.

documents matching the query & documents matching perturbations in the query.

HashTables not very useful because they don't help in finding minor variants / prefix matches which are necessary for relevant retrieval.

Tree algorithms

Take log-n time complexity for operations, which is not good enough.

Retrieving terms matching pro* & work* & interest*

Naive idea:

Query pro* & * work & interest - do an OR operation

Problem: Intersection is costly

How bigram indices are space efficiently compared to permutation indices?

Consider a word like among. In permutation index, it creates 5 entries but in bigram index, we just store it once (the inverted index just stores pointer to the word). (Extra space for the bigrams is atmost 27^2)

Bigram indices faster than permutation indices because we can use hashable instead of trees.

- Intersection of bigram index disk faster than intersection of words ~~is~~ retrieved from trees (the case where we retrieve ~~from & with & intersect) because of the sorted structure of lists~~ ($\propto O(m+n)$ v/s $O(m^2n)$)