⟶ AMBUJ NISHRA
⟶ 2021PCS1017

**Que② ⟶ Minimum loser tree ⟶**

ⓐ ⟶

6, 2, 1, 7, 3, 4, 5, 8

step① ⟶



6, 2 1, 7, 3, 4, 5, 8

step② ⟶



6, 2, 1, 7, 3, 4, 5, 8

step③ ⟶



6, 2, 1, 7, 3, 4, 5, 8
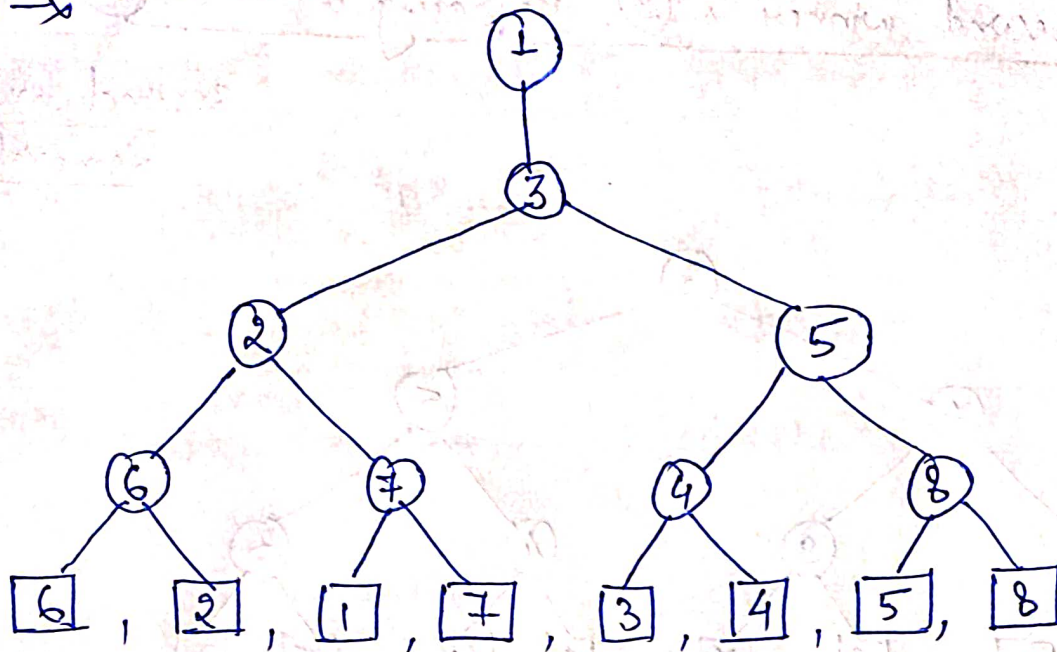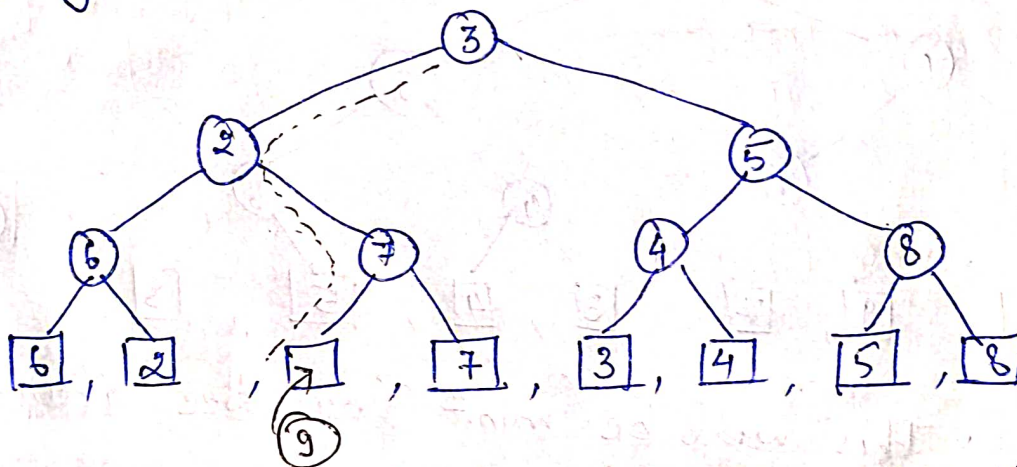
→ This is final min. loser tree.

Ⓑ → first winner ⇒ ①, so, replacing '1' with '9' &
then regeneration →
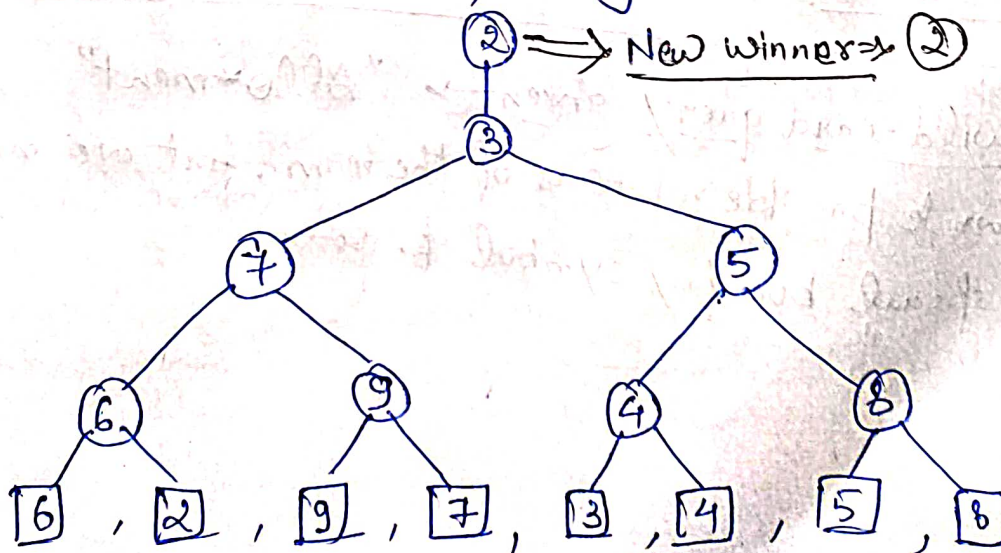
first winner
Removal.
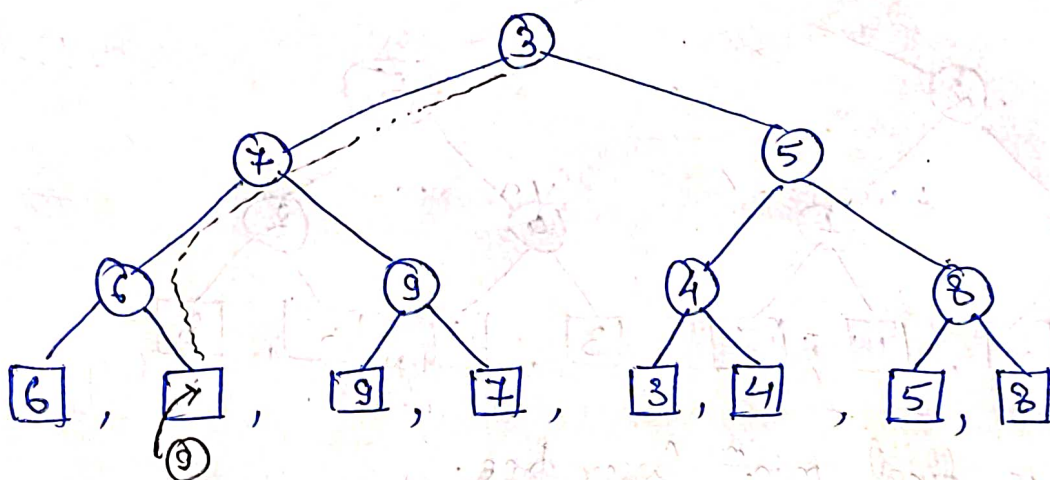


→ values may change in dotted path, so →

② ⟹ New winner ⇒ ②

Since second winner is ②, Replacing it with ⑨ →

→ values will change only in dotted path.



→ This will be min.$^m$ loser see after

2 winner Removals.

---

Que⑦ → wild-card query given → "Allo*ment"

→ To write possible rotations of the term, first we need
to use a special boundary symbol $, so →

⇒ Allo * ment $ ⟹ possible all rotations.

⟹ llo * ment $A , lo * ment $Al , o * ment$All ,

* ment$Allo, | ment $Allo * | , ent $Allo*m, nt $Allo*me,

t$Allo * men, $Allo*ment, Allo * ment$

→To find the fittest term, we'll take the rotation, in which * is at last. So→ | ment$Allo* |

ⓑ→ keys to look one→ from this we can conclude that we need terms, which start with | $Allo | and that ends with | ment$ |. We can create n-gram to find that.

ⓒ. Different Trigram→ $Al, All, llo, men, ent, nt$

**Que④** →



so, structural terms→

| check Next page |

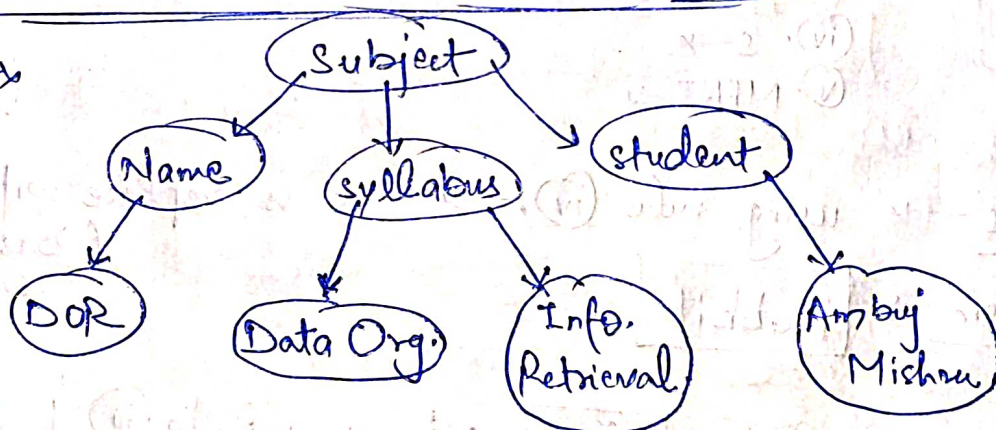(DOR), (Data Org.), (Info. Petrieval), (Ambaj Mishra)

(Name) $\rightarrow$ (DOR) , (Syllabus) $\rightarrow$ (Data Org.) , (Syllabus) $\rightarrow$ (Info. Petrieval.) , (Student) $\rightarrow$ (Ambaj Mishra)

(Subject) $\rightarrow$ (Name) $\rightarrow$ (DOR) , (Subject) $\rightarrow$ (Syllabus) $\rightarrow$ (Data Org.) , (Subject) $\rightarrow$ (Syllabus) $\rightarrow$ (Info. Petrieval) , (Subject) $\rightarrow$ (Student) $\rightarrow$ (Ambaj Mishra)

---

Que. ⑥ $\rightarrow$ Rules $\rightarrow$ ⓘ. ION $\rightarrow$ E

⓪ ⇒

ⓙ. Rabbits $\rightarrow$ using rule ⓔ. $\rightarrow$ 's' is replace with blank. (basically Removed)

stem $\rightarrow$ | Rabbit |

ⓙ Flying $\rightarrow$ stem $\rightarrow$ | flye | ; { using Rule ⓙ }.

ⓙ. Cement $\rightarrow$ stem $\rightarrow$ | ce | ; { using Rule ⓥ }.

ⓥ. Ration $\rightarrow$ stem $\rightarrow$ | Rate | ; { using Rule ⓘ }.

(ii). ING $\rightarrow$ E
(iii). SS $\rightarrow$ SS
(iv). S $\rightarrow$
Ⓥ. MENT $\rightarrow$

(B)→ These stemming rules are case-specific. They'll generate meaningful words few types of words like in case of "Rabbits" but not meaningful words in case of "Cement".

(i)→ Meaningful,
(ii)→ Not meaningful (03 may be semi-meaningful depending upon how we use it further).
(iii)→ Not meaningful
(iv)→ Not meaningful (context is changed completely).

→ We can say, in general that these rule may ~~work~~ work if these suffixes are added extra on the words. Like if plurals are given. But wouldn't ~~work~~, if these are part of word meaning. 'ING' suffix case wouldn't work if it is not added by removing ie.

for example→ "Caution" will become → "Cautie" {which is meaningless}.

(C). Modification→ Using "lemmatization" would be an appropriate sol^n: But, if we want to change these rules to improve accuracy, then → we need to try to convert all possible form in one common "stem" form. So →

ING →                   {"ing" can be completely removed}
MENT→ MENT            {leave it as it is}.

⚡ → we can also put cond.^n to have a chek on which vowel can occur in case of "ING".

**Que ①** → No. of Documents → 1 M.

Vocabulary → 1000 terms.

So → term Document incidence matrix →

Document id vectors corresponding to terms.

Terms. ↓



Issues → ① any term will be present in only a limited no. of Documents. So, matrix will be a sparse matrix. It means that space will be wasted.

② → entries in matrix ⟹ $10^3 \times 10^6$ ⟹ ⑩⁹ : so, size of will be very large of matrix. and memory can get full.

Alternative → Instead of term-document incidence matrix, inverted indexing method should be used.

---

**Que ②** → RDBMS can't be used for structured and unstructured IR because →

① Users generally don't know about the structure of the database system.

② Using RDBMS would result is larger no. of returned result. which will unnecessarily increase the processing time

③ Users may not know how to write queries in RDBMS query languages. To search, users should have a proper knowledge about query languages.

(iv). RBBMS system is an unranked system. So, in case of multiple results, it wouldn't be able to present most relatable and highest ranked results.

(v). On unstructured data, applying RDBMS is not good as the ~~result~~ query can be wrong because of unstructured info while trying to access RDBMS.

---

Que③ → |false|. ⇒ Because if we don't applying stemming on query then stem of query term and stem of it's correct index term (if it exists) __may be__ different

for example→   while indexing, suppose we have converted word "flying" into "fly" & then saved it in index, then if "flying" is present in query, then before searching, we need to convert this {word also in its stem form "fly".
query

---

⇒Que⑤ → In map-reduce algo., parser is used in map phase to distribute data & inverter are used in reduce phase to collect the data and bind it together.

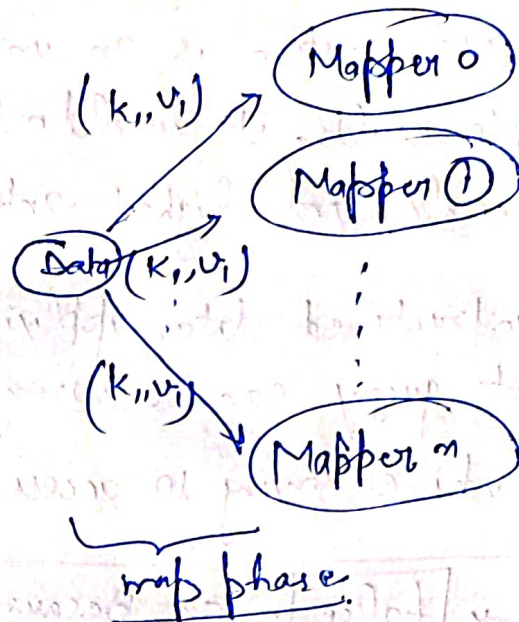structure→      {term → Doc. id, term freq.}

Pseudocode→ first input all data.
   ⓐ Map() fn→    →Take terms & T.f as input.
           input ( list (k,v)) ;  { list of all terms}.
           Partionar( list (k,v)) ;  { partitioners are used to distribute data}.

→ partitioner takes keys & dictibute to Mappers.

{ every 'Kth term is provided/ dictributed to 'n' mappers }.

(k₁,v₁) → Mapper 0

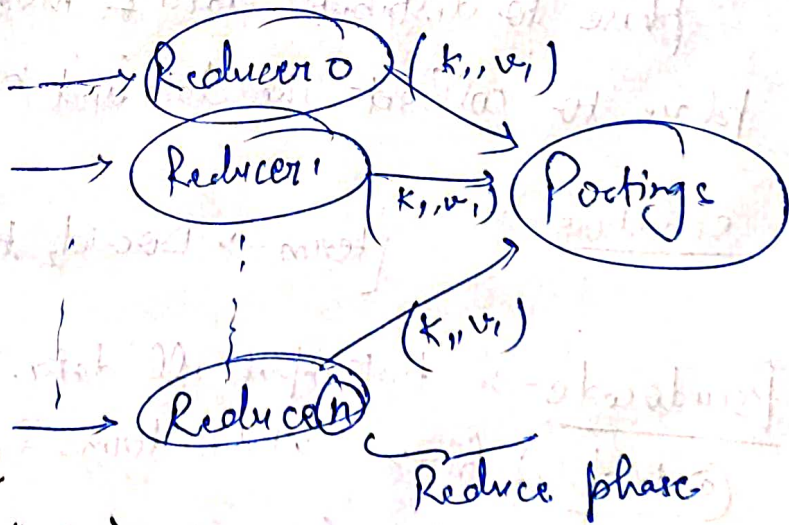Data (k₁,v₁) → Mapper ①

⋮

(k₁,v₁) → Mapper n

map phase

pseudocode for Reduce() → Combiner or Reducers are used to collect data.

→ Take input from all mappers.
→ for {k:all term }:
  Take kth term and it's partions.
  Combiner (-k, list(v));

→ Output (k, list(v));

→ Reducer 0 (k₁,v₁)
→ Reducer 1 (k₁,v₁) → Poctings
⋮
Reducer n (k₁,v₁)

Reduce phase

for example →

("Hello", 1)
("Hello", 1)   } combiner/ Reducer.
("Hello", 1)

("Hello", 3)

input taken from map phase