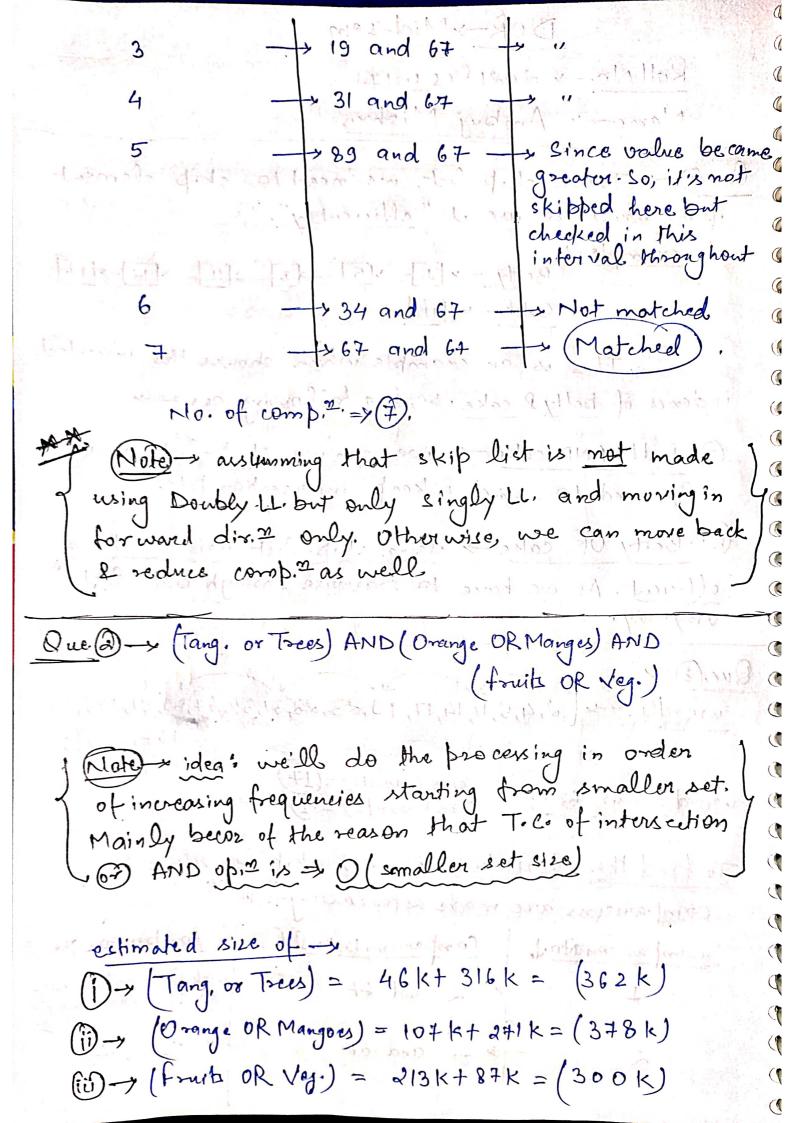Roll No.: → 2021 PCS 1017

Name → Ambuj Mishra

Que ① → To use skip-list, we need to skip element if we want to use it "efficiently".

for example →

Betty → $\boxed{1}$ → $\boxed{5}$ → $\boxed{8}$ → $\boxed{11}$ → $\boxed{12}$ → $\boxed{16}$

Cake → $\boxed{11}$.

→ This is an example which shows the inverted indexes of betty & cake. So, for following cases →

ⓐ. Betty AND cake → we can use skip-list as we only need to search & keep intersection pts.

ⓑ. Betty OR cake → Using skip-list will not be efficient. As we have to traverse through both lists anyway.

Que ③ →

word① → [2, 4, 5, 11, 14, 17, 19, 23, 28, 31, 34, 67, 89, 91, 145, 167, 180];

word② → [67];  size (word 1) = ⑰
                size (word 2) = ①

To find the intersection → with skip length = 3, comparisons are made as following →

| Comparison No.↓ | Comparison Detail ↓ | Comparison Result |
|---|---|---|
| 1 | 2 and 67 | skip 3 values bcoz lesser no. |
| 2 | 11 and 67 | " " |

| | | | |
|---|---|---|---|
| 3 | $\longrightarrow$ | 19 and 67 | $\rightarrow$ " |
| 4 | $\longrightarrow$ | 31 and 67 | $\rightarrow$ " |
| 5 | $\longrightarrow$ | 89 and 67 | $\rightarrow$ Since value became greater. So, it's not skipped here but checked in this interval throughout |
| 6 | $\longrightarrow$ | 34 and 67 | $\rightarrow$ Not matched. |
| 7 | $\longrightarrow$ | 67 and 67 | $\rightarrow$ (Matched). |

No. of comp.$^n$ $\Rightarrow$ ⑦.

**(Note)** $\rightarrow$ assuuming that skip list is _not_ made using Doubly LL. but only singly LL. and moving in forward dir.$^n$ only. Otherwise, we can move back & reduce comp.$^n$ as well.

---

**Que ②** $\longrightarrow$ (Tang. or Trees) AND (Orange OR Manges) AND
(fruits OR Veg.)

**(Note)** $\rightarrow$ idea: we'll do the processing in order of increasing frequencies starting from smaller set. Mainly becoz of the reason that T.C. of intersection or AND op.$^n$ is $\Rightarrow$ $O$(smaller set size)

estimated size of $\longrightarrow$

① $\rightarrow$ (Tang, or Trees) = 46k + 316k = (362k)

② $\rightarrow$ (Orange OR Mangoes) = 107k + 271k = (378k)

③ $\rightarrow$ (fruits OR Veg.) = 213k + 87k = (300k)

. we need to calculate → ① AND ⓘⓘ AND ⓘⓘⓘ.

Query processing order → { ⓘⓘⓘ. AND ① } AND ⓘⓘ.

→ we'll take intersection of smallest set ⓘⓘⓘ. with ①. & then intersect it with ⓘⓘ.

Modifications → ①. We can use term frequency as well. which will give directly idea about it's posting size

ⓘⓘ. We ~~have because~~ have to assume that

$$[X \ OR \ Y = size(X) + size(Y)] ;$$ although it's not always true. As 2 sets having 9 values & 990 values will have higher probability of resulting in larger union sets than 2 sets with 500 values, Although $size(X) + size(Y)$ is smaller. (Depends upon similarity).

---

Que ⑤ → step ① → Tokenization → And [Normalization] →

Doc① → ["Betty", "bought", "a", "butterscotchcake"]

Doc② → ["The", "cake", "was", "very", "bitter"]

Doc③ → ["Betty", "returned", "the", "cakebos", "the", "cake", "was", "bitter"]

Doc④ → ["Betty", "got", "a", "new", "cake", "andcowtd", "the", "party"]

{ Note → using Normalization we Removed symbols. (", -, ') & replaced with NULL }

step.② → stop-words →

   [ "a", "The", "was", "very", "and," ] will
                              "new"
be removed from all the documents. if
they are present.

step③ → stemming →
   { "Cake cos" → "cake"
     "Bitter" → "Bitt"
     "Returned" → "Return"
     "andcontd" → "and cont" } → These
     "Grat" → "Gro"
   changes are made throughout the documents

After that posting list →

   "Betty" → ① → ② → ③
   "Bought" → ①
   "Butterscotch cake" → ①
   "cake" → ② → ③ → ④
   "Bitt" → ② → ③
   "return" → ③
   "Gro" → ④
   "and cont" → ④
   "party" → ④

   ‿‿‿‿       ‿‿‿‿
   Dictionary   posting.

→ Now sort, on
the basis of
Dictionary Key
     ↓
That is final
result of
postings.
(sorry, No time).