

Q1.

(1).

According to the question, word appears at least once in each of all N documents.

Lets Assume that posting IDs are always start from 1 to N .

For uncompressed encoding. each posting ID requires 4 bytes.

Therefore, the total encoding is $4N$ bytes

By applying Elias-r encoding. and gap,
we know that the first posting ID is 1, and all the following gaps are 1 as well.

Since Elias-r encoding for 1 is a bit '0'.

Therefore, the total encoding is $\lceil \frac{N}{8} \rceil$ bytes

By applying Elias-s encoding. and gap,

Similar to Elias-r, and Elias-s encoding for 1 is a bit '0' as well.

Therefore, the total encoding is $\lceil \frac{N}{8} \rceil$ bytes

Since, compression ratio = $\frac{\text{Uncompressed Size}}{\text{Compressed Size}}$

Elias-r compression ratio = Elias-s compression ratio = $\frac{4N}{\lceil \frac{N}{8} \rceil}$

(2). According to the question, we are using gaps and term frequency in our postings; and term frequency for each doc is known to be 10.

For uncompressed encoding, we need 4 bytes for first posting ID and 4 bytes for frequency; and all the following gaps use 4 bytes as well.

Therefore, the total encoding is $\lceil \frac{8}{5}N \rceil$ bytes for $\frac{N}{5}$ docs.

By applying Elias-r encoding.

Assume the id of the first doc is R , $k \in [1, N]$

Then, bits for first posting ID is $2\lfloor \log_2 k \rfloor + 1$ bits

bits for term frequency 10 is 7 bits

For the following postings, we are using gap 5.

Then bits for gap is 5 bits

$$\begin{aligned} \text{Therefore, in total } & 2\lfloor \log_2 k \rfloor + 1 + 7 + (5+7) \cdot \left(\frac{N}{5}-1\right) \\ &= 2\lfloor \log_2 k \rfloor + \frac{12}{5}N - 4 \text{ bits} \\ &= \boxed{\left[\frac{2\lfloor \log_2 k \rfloor + \frac{12}{5}N - 4}{8} \right] \text{ bytes}} \end{aligned}$$

By applying Elias-S encoding.

Assume the id of the first doc is R , $k \in [1, N]$

Then, bits for first posting ID is $2\log_2 \log_2 k + \log_2 k$ bits

bits for term frequency 10 is 8 bits

For the following postings, we are using gap 5.

Then bits for gap is 5 bits

$$\text{Therefore, in total } 2\log_2 \log_2 k + \log_2 k + 8 + (5+8) \cdot \left(\frac{N}{5}-1\right)$$

$$= \left\lceil \frac{2\log_2 \log_2 k + \log_2 k + \frac{13}{5}N - 5}{8} \right\rceil \text{ bytes}$$

Since, compression ratio = $\frac{\text{Uncompressed Size}}{\text{Compressed Size}}$

$$\text{Elias-r compression ratio} = \frac{\lceil \frac{8}{5}N \rceil \text{ bytes}}{\left[\frac{2\lfloor \log_2 k \rfloor + \frac{12}{5}N - 4}{8} \right] \text{ bytes}}$$

$$\text{Elias-s compression ratio} = \frac{\lceil \frac{8}{5}N \rceil \text{ bytes}}{\left\lceil \frac{2\log_2 \log_2 k + \log_2 k + \frac{13}{5}N - 5}{8} \right\rceil \text{ bytes}}$$

Q2.

- (1). Assume after all logarithmic merges have finished, there will be n number of sub-indexes formed in disk.

Therefore, we can have $2^n = t$.

Because for each merge operation, we merge 2 into 1.

And there are t sub-indexes in total.

Since $2^n = t$,

$n = \lceil \log_2 t \rceil$ for n to be an integer.

For special case $t=1$, we say that $n \leq \lceil \log_2 t \rceil + 1$

(2). From the proof of (1), we know that for t sub-indexes, there will be at most $\log_2 t$ sub-indexes in disk.

Therefore, for each sub-index in t , it will be called at most $\log_2 t$ times for merge operation.

Since for each sub-index, it has M pages. The read/write cost for it would be M . for each merge operation.

There are t sub-indexes in total.

Therefore the I/O cost for the whole merge operation will be bounded by $2 \cdot t \cdot M \cdot \log_2 t$ for read and write.

∴ the total I/O cost of the logarithmic merge is $O(t \cdot M \cdot \log_2 t)$.

Q3.

(1). The formula for recall is $\frac{TP}{TP+FN}$.

By using stemming, the original TP will not reduce because the stemming results must include the original correct words.

The original FN may reduce because we may include actual positives. For example.

When the query term is "automate". and when using Boolean query retrieval, docs having "automates" will not be considered; But it should consider.

By applying stemming, both "automate" and "automates" are "automa". Therefore "automates" docs are considered.

Therefore FN decreases.

Therefore, we can conclude that stemming will not hurt recall.

(2). The formula for F_1 is

$$F_1 = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}} = \frac{TP}{TP + \frac{1}{2}(FP + FN)}$$

We know that when FN decrease by x , TP would increase by x .
when FP increase by y , TN would decrease by y .

From (1), we know that stemming will always decrease FN by an amount. Assume it is x , and $x \geq 0$.

And we know that stemming will always increase FP by an amount. Assume it is y , and $y \geq 0$.

$$\therefore \text{new } F_1 = \frac{TP + x}{TP + x + \frac{1}{2}(FP + y + FN - x)} = \frac{TP + x}{TP + \frac{1}{2}(FP + FN) + \frac{1}{2}(x+y)}$$

$$\text{Let } TP = a, \quad TP + \frac{1}{2}(FP + FN) = b, \quad b \geq a$$

$$\therefore F_1 = \frac{a}{b}, \quad \text{new } F_1 = \frac{a+x}{b+\frac{1}{2}(x+y)}$$

If $F_1 = \text{new } F_1$,

$$\text{Then. } \frac{a}{b} = \frac{a+x}{b+\frac{1}{2}(x+y)}$$

$$\frac{a}{b}(x+y) = b x$$

$$ax+ay = 2bx$$

$$ay = (2b-a)x$$

$$\frac{y}{x} = \frac{2b-a}{a}$$

Therefore. when $\frac{y}{x} > \frac{2b-a}{a}$, F_1 will decrease after stemming
when $\frac{y}{x} < \frac{2b-a}{a}$, F_1 will increase after stemming.

In conclusion, it is wrong that stemming always helps or hurts the F_1 score.

Intuitively, when the stemming word includes a lot of words that are not "same" as query word, F_1 score will be dominant by increasing in FP and will hurt F_1 score.

In opposite, if the stemming word include a lot of words that are "same" as query word, F_1 score will be dominant by decreasing in FN and will help F_1 score.