

Aaryan Bhagat abhag017 862468325

1

I have assumed constant b to be is x and 0.75 and for k1, since its a short collection, k1 is on the lower side which is 0.5.

First the formula

$$\text{score}(D, Q) = \sum_{i=1}^n \overbrace{\text{IDF}(q_i) \times f(q_i, D) \times (k_1 + 1)}^{\text{f}(q_i, D) + k_1 \times (1 - b + b \times \frac{|D|}{\text{avgdl}})}$$

For query 1

$R_1 = 0.5$   
 $b = 0.75$

$Q_1 = \text{university residence}$  (Assuming Case Insensitive)  
 $\uparrow \quad \uparrow$   
 $q_1 \quad q_2$

$$\text{IDF}(q_1) = \ln \left( \frac{10^4 - 500 + 0.5 + 1}{500 + 0.5} \right) \quad f(q_1, D) = 2$$

$$\text{IDF}(q_2) = \ln \left( \frac{10^4 - 100 + 0.5 + 1}{100 + 0.5} \right) \quad f(q_2, D) = 1$$

$$|D| = 48 \quad [\text{Calculated by tokenizing}]$$

$$\frac{\text{IDF}(q_1) \times \frac{2 \times (0.5 + 1)}{2 + 0.5 \times (1 - 0.75 + 0.75 \times 48)}}{\text{augll}} + \frac{\text{IDF}(q_2) \times \frac{1 \times (0.5 + 1)}{1 + 0.5 \times (1 - 0.75 + 0.75 \times 48)}}{\text{augll}}$$

$Q_2 = \text{diverse university}$   
 $\uparrow \quad \uparrow$   
 $q_1 \quad q_2$

$$\text{IDF}(q_1) = \ln \left( \frac{10^4 - 250 + 0.5 + 1}{250 + 0.5} \right) \quad f(q_1, D) = 1$$

$$\text{IDF}(q_2) = \ln \left( \frac{10^4 - 500 + 0.5 + 1}{500 + 0.5} \right) \quad f(q_2, D) = 2$$

$$|D| = 48$$

$$\frac{\text{IDF}(q_1) \times 1 \times (0.5 + 1)}{1 + 0.5 \times (1 - 0.75 + 0.75 \times 48)} + \frac{\text{IDF}(q_2) \times 1 \times (0.5 + 1)}{2 + 0.5 \times (1 - 0.75 + 0.75 \times 48)}$$

For Unigram it is very straightforward to compute

Unigram language Model

$Q_1: \text{university riverside}$

$$P(q_1|D) = \frac{1}{48} \quad P(q_2|D) = \frac{1}{48}$$

$Q_2: \text{diverse university}$

$$P(q_1|D) = \frac{1}{48} \quad P(q_2|D) = \frac{2}{48}$$

$\text{plex } \text{Unigram}(Q_1, D) = P(q_1|D) \times P(q_2|D) = \frac{1}{48} \times \frac{1}{48}$

$\text{plex } \text{Unigram}(Q_2, D) = P(q_1|D) \times P(q_2|D) = \frac{1}{48} \times \frac{2}{48}$

## 2

Program code for pagerank

```
import numpy as np

def pageRank(M, num_iter: int = 100, d: float = 0.85):
    N = M.shape[1]
    v = np.ones(N) / N
    M_hat = (d * M + (1 - d) / N)
    count = 0
    while True:
        count = count + 1
        v_new = M_hat @ v
        print("Iteration ", count)
        print(v)
```

```

        if(np.linalg.norm(v_new - v) <= 0.01):
            break
        v = np.copy(v_new)
    return v

M = np.array([[0, 1, 0, 0, 0],
              [1/3, 0, 1/2, 0, 1],
              [0, 0, 0, 0, 0],
              [1/3, 0, 0, 0, 0],
              [1/3, 0, 1/2, 0, 0]])
v = pageRank(M, 100, 0.85)
# print(v)

```

## Output

```

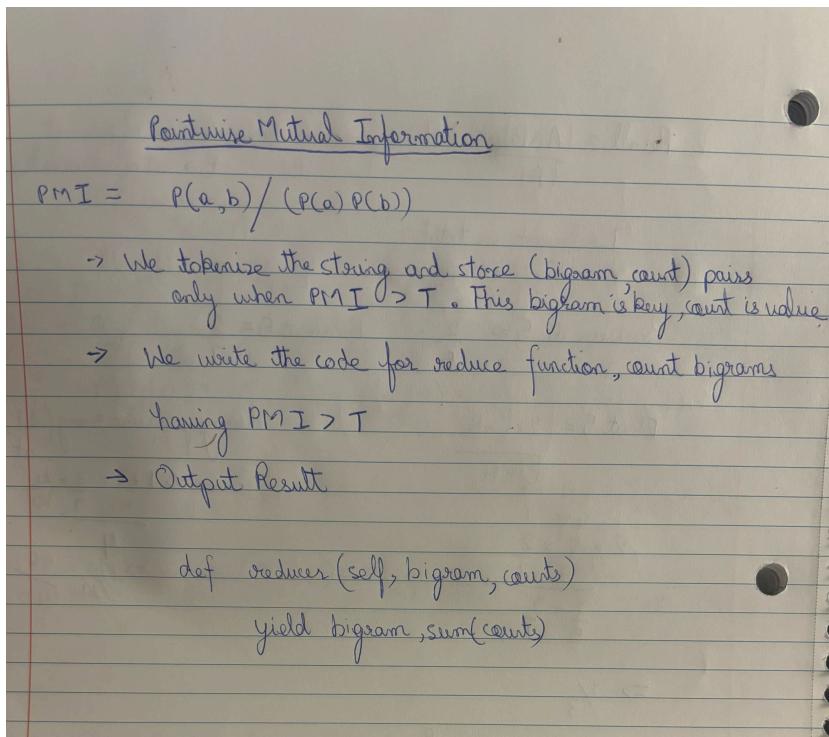
Iteration 1
[0.2 0.2 0.2 0.2 0.2]
Iteration 2
[0.2 0.34166667 0.03 0.08666667 0.17166667]
Iteration 3
[0.31531667 0.24023333 0.0249 0.08156667 0.09431667]
Iteration 4
[0.22688833 0.20278139 0.02269 0.11202972 0.12261222]
Iteration 5
[0.19297423 0.19875872 0.02061005 0.08489508 0.09453833]
Iteration 6
[0.1866982 0.16154617 0.01775329 0.07242932 0.0811886 ]
Iteration 7
[0.15290272 0.14504175 0.01558847 0.06848629 0.07603144]
Iteration 8
[0.137027 0.12831578 0.01374152 0.05706396 0.06368905]
Iteration 9
[0.12106353 0.11079528 0.01199512 0.05081944 0.05665958]
Iteration 10
[0.10471598 0.09809989 0.01053999 0.04484132 0.04993925]
Iteration 11
[0.092629 0.08584148 0.00924409 0.03891362 0.04339311]
Iteration 12
[0.08106589 0.07515841 0.00810064 0.03434552 0.03827426]
Iteration 13
[0.07099299 0.06605291 0.00710834 0.03007701 0.03351978]
Iteration 14
[0.0623775 0.05786007 0.00623253 0.02634721 0.02936826]
Iteration 15
[0.05464663 0.05075104 0.00546557 0.02313919 0.02578802]
Iteration 16
[0.04793209 0.04451961 0.00479371 0.02027692 0.02259979]

```

So converging after 16 iterations

## 3

For map reduce this is the steps of how we will devise an algorithm to solve a bigram having a condition of pointwise mutual information



If we use Combiner here, it can optimize the reduce function by doing certain data preprocessing in different shards where data is stored.

## 4

All the formulas are applied below, the example mentioned for high recall is google map which will certainly give the right answers as well as other answers.

For high precision it is the normal google search

$$\text{Recall} = \frac{|A \cap B|}{|A|} \quad \text{Precision} = \frac{|A \cap B|}{|B|} \quad A \rightarrow \text{set of relevant} \quad B \rightarrow \text{set of retrieved}$$

7 relevant in total

~~r r x x x x r x x r~~

$$\text{Precision at } 5 \Rightarrow \frac{2}{5} \quad \text{Recall at } 5 = \frac{2}{7}$$

$$F1 \text{ at } 5 = \frac{2 \times P_s \times R_s}{P_s + R_s}$$

$$\text{Average precision} = \overline{\text{avg}(\text{at } 1, \text{at } 2, \dots, \text{at } 7) \text{ at } 8}$$

$$\Rightarrow \overline{1 + 1 + \frac{3+4}{7}} = \overline{1 + 1 + \frac{7}{7}} = \overline{1 + 1} = \overline{2}$$

$$\frac{2 \times 2 \times 2}{5 \cdot 7} = \frac{2 \cdot 2}{5 \cdot 7}$$

$$\Rightarrow \frac{8}{35} = \frac{8}{35} = \frac{99}{140}$$

$$\Rightarrow \frac{1}{3}$$

$$DCG = \frac{1}{p} \sum_{i=2}^p \frac{\text{rel}_i}{\log i}$$

$$p = 5$$

$$1 + \sum_{i=2}^5 \frac{\text{rel}_i}{\log i} \Rightarrow 1 + \frac{1}{\log 2} + 0 + 0 + 0$$

$$\Rightarrow 2$$

Higher Recall when we need all correct relevant pages to be present  
Example is google maps which will show full correct + more

Higher Precision when we have more relevant in output  
Example is google search engine