Hemin Patel - Hw 8

① 

A) Consider table of term frequencies for 3 docs. Compute tf-idf weights for each terms. Corpus size = 100K documents.

| term | df | idf | |
|------|------|------|------|
| Car | 18165 | 1.74 | $1 + \log\left(\frac{100K}{18165}\right)$ |
| auto | 6723 | 2.17 | $1 + \log\left(\frac{100K}{6723}\right)$ |
| insurance | 19241 | 1.72 | $1 + \log\left(\frac{100K}{19241}\right)$ |
| best | 25235 | 1.60 | $1 + \log\left(\frac{100K}{25235}\right)$ |

weights

| | Doc 1 | Doc 2 | Doc 3 |
|------|-------|-------|-------|
| Car | 46.98 | 6.96 | 41.76 |
| auto | 6.51 | 71.61 | 0 |
| insurance | 0 | 56.76 | 49.88 |
| best | 22.4 | 0 | 27.2 |

$$tf\text{-}idf = tf_{t,d} \cdot idf_t$$

$\begin{bmatrix} 27 \cdot 1.74 & 4 \cdot 1.74 & 24 \cdot 1.74 \end{bmatrix}$

B) Compute normalized doc vectors for each doc
- Doc 1: 46.98, 6.51, 0, 22.4
  
  $|doc1|: \sqrt{46.98^2 + 6.51^2 + 22.4^2} = 52.45$
  
  Doc1/|doc1|: $\boxed{0.90, \ 0.12, \ 0, \ 0.43}$

- Doc2: 6.96, 71.61, 56.76, 0
  
  $|doc2|: \sqrt{6.96^2 + 71.61^2 + 56.76^2} = 91.64$
  
  Doc2/|doc2|: $\boxed{0.076, \ 0.78, \ 0.62, \ 0}$

- Doc 3: 41.76, 0, 49.88, 27.2
  
  $|Doc3|: \sqrt{41.76^2 + 49.88^2 + 27.2^2} = 70.51$
  
  Doc3/|Doc3|: $\boxed{0.59, \ 0, \ 0.71, \ 0.40}$

| | Doc1 | Doc2 | Doc3 |
|------|------|-------|------|
| Car | 0.90 | 0.076 | 0.59 |
| auto | 0.12 | 0.78 | 0 |
| insur | 0 | 0.62 | 0.71 |
| best | 0.43 | 0 | 0.40 |

2) normal: url 1: normal.org
　　　　 2: merrian-webster.com
　　　　 3: dictionary.com

　　form: url 1: google.um/forms/abod
　　　　　 2: merrian-webster.com
　　　　　 3: form.com

　　　　　 cmd + F　w/o space

① c) compute cosine similarity between each pair of doc

- Doc1 · Doc 2: $(46.98 \cdot 6.96) + (6.51 \cdot 71.61) + / + / = 793.1619$

$|Doc1| = 52.45$
$|Doc2| = 91.64$　$\left.\right]$ $\cos(Doc1, Doc2) = \dfrac{793.1619}{(52.45 \cdot 91.64)} = \boxed{0.1650}$

- Doc2 · Doc 3: $(6.96 \cdot 41.76) + / + (56.76 \cdot 49.88) + / = 3121.8384$

$|Doc2| = 91.64$
$|Doc3| = 70.51$　$\left.\right]$ $\cos(Doc2, Doc3) = \dfrac{3121.8384}{(91.64 \cdot 70.51)} = \boxed{0.4831}$

- Doc1 · Doc3: $(46.98 \cdot 41.76) + / + / + (22.4 \cdot 27.2) = 2571.1648$

$|Doc1| = 52.45$
$|Doc3| = 70.51$　$\left.\right]$ $\cos(Doc1, Doc3) = \dfrac{2571.1648}{(52.45 \cdot 70.51)} = \boxed{0.6952}$

② Let google query Q = normal form | 3 URLs, Rank approx cos similarity

N = 1 trillion

| term | df | idf | | query Q |
|---|---|---|---|---|
| normal | 2,440,000,000 | 5.61261 | $1 + \log \dfrac{1 \text{ tril}}{2,440 \text{ mil}}$ | Norm(idf) |
| form | 4,970,000,000 | 5.30364 | | 7.7220 |

Term Frequency

| | url 1 | url 2 | url 3 |
|---|---|---|---|
| normal | 28 | 119 | 38 |
| form | 29 | 151 | 13 |

weights

| | url 1 | url 2 | url 3 | $tf\text{-}idf = tf_{t,d} \cdot idf_t$ |
|---|---|---|---|---|
| normal | 157.1531 | 667.9006 | 213.2792 | |
| form | 153.8056 | 800.8496 | 68.9473 | |

vectors + Normalization

| | url 1 | url 2 | url 3 |
|---|---|---|---|
| normal | 17.9397 | 21.4666 | 18.6841 |
| form | 16.9026 | 20.7030 | 15.0545 |
| NORM | 24.6481 | 29.8233 | 23.9944 |

prefers to weights

$(1 + \log(tf)) \cdot \left(1 + \log\left(\dfrac{N}{df}\right)\right)$ → or idf

$\llcorner$ $1 + \log(157.1531) \cdot \left(1 + \log \dfrac{1 \text{ tril}}{2,440 \text{ mil}}\right)$

EX: $\sqrt{17.9397^2 + 16.9026^2}$

Hemin    HW 8 - cont.

②  Cosine  Similarity  w/  query Q  &  3  web  pages

• URL 1 · Q: $(17.9397 \cdot 5.61261) + (16.9026 \cdot 5.30364) = 190.3338$

$\left.\begin{array}{l} |URL\ 1| = 24.6481 \\ |Q| = 7.7220 \end{array}\right]$ $cos(URL\ 1, Q) = \dfrac{190.3338}{(24.6481 \cdot 7.7220)} = 1.000006$

• URL 2 · Q: $(21.4666 \cdot 5.61261) + (20.7030 \cdot 5.30364) = 230.2849$

$\left.\begin{array}{l} |URL\ 2| = 29.8233 \\ |Q| = 7.7220 \end{array}\right]$ $cos(URL\ 2, Q) = \dfrac{230.2849}{(29.8233 \cdot 7.7220)} = 0.999954$

• URL 3 · Q: $(18.6841 \cdot 5.61261) + (15.0545 \cdot 5.30364) = 184.7102$

$\left.\begin{array}{l} |URL\ 3| = 23.9944 \\ |Q| = 7.7220 \end{array}\right]$ $cos(URL\ 3, Q) = \dfrac{184.7102}{(23.9944 \cdot 7.7220)} = 0.996899$