

Find the two closest documents based on their Freq feature vectors (use Euclidean distance), after having removed the stopwords {I, like, in, the, my, has, a, was}.

D1 = I like swimming in the sea

D2 = My car has a flat tire

D3 = A fish was swimming in the sea

$W' = \{ \text{I, like, swimming, in, the, sea, my, car, has, a, flat, tire, fish, was} \}$

$W = W' - \text{stopwords} = \{ \text{swimming, sea, car, flat, tire, fish} \}$

Frequency

$F_{D1} = (1, 1, 0, 0, 0, 0)$ $F_{D2} = (0, 0, 1, 1, 1, 0)$

$F_{D3} = (1, 1, 0, 0, 0, 1)$

Normalized term-frequency

$TF_{D1} = (\frac{1}{2}, \frac{1}{2}, 0, 0, 0, 0)$ $TF_{D2} = (0, 0, \frac{1}{3}, \frac{1}{3}, \frac{1}{3}, 0)$

$TF_{D3} = (\frac{1}{3}, \frac{1}{3}, 0, 0, 0, \frac{1}{3})$

Term frequency inverse-document-frequency

$TFIDF_{D1} = (\frac{1}{2} \ln(\frac{3}{2}), \frac{1}{2} \ln(\frac{3}{2}), 0, 0, 0, 0)$

$TFIDF_{D2} = (0, 0, \frac{1}{3} \ln(\frac{3}{1}), \frac{1}{3} \ln(\frac{3}{1}), \frac{1}{3} \ln(\frac{3}{1}), 0)$

$TFIDF_{D3} = (\frac{1}{3} \ln(\frac{3}{2}), \frac{1}{3} \ln(\frac{3}{2}), 0, 0, 0, \frac{1}{3} \ln(\frac{3}{1}))$