CHRIS SUNNY THALIYATH

Quiz NLP

Data :

- doc1, doc2, doc3, doc4, doc5
  these are CORPUS

- W1, W2, W3, W4, W5, W6, W7
  words
- W6, W7 — STOP words

Represent Corpus

└→ 1 Binary Matrix
    └→ where

"1 if word present in document, else its 0"

        └→ Row : document is represented
        └→ Column : Represents words

$$TF-IDF(t,d) = TF(t,d) \times IDF(t)$$

IDF - TF (Inverse Document Frequency)

  └→ Histogram of words in all documents
     ie the count of different words

CBOW (continous Bag of words / Skip-Gram)

    └→ we create pair of words
       ( Target words & context words)

    └→ Skip Gram : predicts the context based on current word & surrounding words

# Binary Matrix

| Document | W1 | W2 | W3 | W4 | W5 | W6 | W7 |
|----------|----|----|----|----|----|----|----|
| doc 1 | 1 | 1 | 1 | 0 | 0 | 0 | 1 |
| doc 2 | 1 | 1 | 1 | 0 | 0 | 0 | 1 |
| doc 3 | 1 | 0 | 0 | 0 | 1 | 0 | 1 |
| doc 4 | 0 | 0 | 0 | 1 | 1 | 1 | 0 |
| do 5 | 0 | 0 | 0 | 0 | 0 | 1 | 0 |

$$IDF(t) = \log\left(N/d\,(f(t))\right)$$

N —: Total # of documents

$d.f.(t)$ : # of documents with term t

$df(w1) = 3$

$IDF(w1) = \log(5/3)$

$TF - IDF(w1, doc1) = 1 * IDF(w1)$

$= $

$df(w2) = 2$

$IDF(w2) = $

# Calculations

TF (doc4)

| | TF(doc1) | TF(doc2) | TF(doc3) | TF(doc4) | DF | IDF |
|---|---|---|---|---|---|---|
| W1 | 1 | 1 | 1 | 0   0 | 3 | $\log(5/3)$ |
| W2 | 1 | 1 | 0 | 0   0 | 2 | $\log(5/2)$ |
| W3 | 1 | 1 | 0 | 0   0 | 2 | $\log(5/2)$ |
| W4 | 0 | 0 | 6 | 1   0 | 1 | $\log(5/1)$ |
| W5 | 0 | 0 | 1 | 1   0 | 2 | $\log(5/2)$ |
| W6 | 0 | 0 | 0 | 1   1 | 2 | $\log(5/2)$ |
| W7 | 1 | 1 | 1 | 0   0 | 3 | $\log(5/3)$ |

## TF - IDF

| Document | W1 | W2 | W3 | W4 | W5 | W6 | W7 |
|---|---|---|---|---|---|---|---|
| doc 1 | 0.51 | 0.90 | 0.90 | 0 | 0 | 0 | 0.51 |
| doc 2 | 0.51 | 0.90 | 0.90 | 0 | 0 | 0 | 0.51 |
| " 3 | 0.51 | 0 | 0 | 0 | 0.90 | 0 | 0.51 |
| " 4 | 0 | 0 | 0 | 1.61 | 0.90 | 0.90 | 0 |
| " 5 | 0 | 0 | 0 | 0 | 0 | 0.90 | 0 |

## CBOW - Skip Gram

a) Tokenization : Break document into words
Remove stop words

b) Windowing : window size (2-3 words)

CBOW (continous Bag of words)
↳ input : words within the window
↳ output : Target words
↳ Model : Surrounding Context

Skip Gram :

- Input : Target word
- Output : Context words

Example with doc 2 with window size 2

Document 2 : ( $w_1$ , $w_2$ , $w_3$ , $w_7$ )

Training Examples

( $w_1$ , $w_2$ ) :     Predict $w_2$ given $w_1$

( $w_2$ , $w_3$ ) :          "

( $w_2$ $w_3$ ) :          "

( $w_3$ , $w_7$ ) :          "

( $w_3$ , $w_7$ ) :          "