

NLP - by Dr. Ajay Debnath

Day 5, 17th Dec. 2023

Recap: TF-IDF

BOW	TF-IDF
weight = 1	weight = \log
sent1	sent1
sent2	sent2
text	text
No	No
No	No

$$TF = \frac{\text{No. of words in sent}}{\text{Total no. of words in corpus}}$$

$$IDF = \log \left(\frac{\text{Total no. of sent}}{\text{no. of sent. containing that word}} \right)$$

sent1: int boy
sent2: int girl
sent3: int boy & girl

TF	sent1	sent2	sent3
int	1/2	1/2	1/3
boy	1/2	0/2	1/3
girl	0	1/2	1/3

IDF

words	IDF
int	$\log(3/3) = 0$
boy	$\log(3/2)$
girl	$\log(3/2)$

$$TF * IDF = \checkmark$$

	f1	f2	f3
sent1	0	$1/2 * \log(3/2)$	0
sent2	0	0	$1/2 * \log(3/2)$

ML

tfidf vectorizer → sklearn

text → Pre-process → Process → ML → predict → performing

end-to-end project

TF-IDF Advantage

- 1) It weights the words → importance to word

I love NLP boy

0.9 0.3 0.44 0

→ context, order, semantic information

I NLP love

0.9 0.44 0.3

Disadvantage

- 1) context info missing
- 2) Order " "
- 3) sensitive to corpus size
- 4) biased to rare terms

corpus

I was born in France.

I can speak French fluently

w1	w2	w3
✓	✓	✓

Data

sent1 spam
sent2 ham / not spam

10K

label

1 - spam

0 - ham

msg

$$X = \begin{bmatrix} 0 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 \end{bmatrix}$$

$$X = \begin{bmatrix} 0 & 0.5 & 0.67 \\ 0 & 0.5 & 0.67 \end{bmatrix}$$

$$Y = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

$$X = \begin{bmatrix} 1 & 0 & 1 & 0 & 0 & 1 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \end{bmatrix}$$

$$\frac{X}{Y}$$

Training
test

acc = 80%
acc ≈ 100%
Error = (y_true - y_pred) ≈ 0

model/robot/child

apple - 0

banana - 1

model → y-pred
apple

Train Validation Test

train { 1, 70, 71, 100 }
test { 71, 80, 81, 100 }

70 train

10 Validation

20 test

prepare mock final

100 → 70

30 test

70 y.true

70 y.pred

1 30

98%

acc BOW

97%

acc TFIDF

model BOW

model TFIDF

test → spam/ham

test1 Error1 = y_true - y_pred

= 0 - 0 = 0 ✓

test2 Error2 = 1 - 1 = 0 ✓

= 1 - 0 = 1 → misclassification

accuracy = ✓