

## In class activity week 3

$$1. A. v_1 \text{ Doc1} = [1, 1, 0, 0, 1]$$

$$v_1 \text{ Doc2} = [0, 0, 1, 1, 1]$$

$$v_1 \text{ Doc3} = [1, 0, 1, 0, 1]$$

## 1.A. TF-IDF

$$v_1 \text{ doc1} = [\frac{1}{3}, \frac{1}{3}, 0, 0, \frac{1}{3}]$$

$$\text{TF}[\frac{1}{3}, \frac{1}{3}, \phi, \phi, \frac{1}{3}] \times \text{IDF}[\ln(\frac{3}{2}), \ln(\frac{3}{1}), \ln(\frac{3}{3})]$$

$$\text{TF-IDF}_{\text{doc1}} = (0, 13, 0, 366, 0)$$

free price in now

$$(0, 13, 0, 366, 0, 0, 0) \text{ TF-IDF Doc1}$$

## 1A. DOC2 TF-IDF

TF

$$[0, 0, \frac{1}{3}, \frac{1}{3}, \frac{1}{3}] \times \text{IDF}[\ln(\frac{3}{2}), \ln(\frac{3}{1}), \ln(1)]$$

$$0.405, 1.09,$$

$$[\text{TF-IDF} \quad 0, 0, 0.135, 0.366, \phi] \text{ Doc2}$$

I.A. &amp; TF IDF doc 3

$$\text{TF} \left[ \frac{1}{3}, 0, \frac{1}{3}, 0, \frac{1}{3} \right] \times \text{IDF} \left[ \ln\left(\frac{3}{2}\right), 0, \ln\left(\frac{3}{1}\right), 0, \ln\left(\frac{3}{1}\right) \right]$$

$\frac{1}{13}$

0.405 1.09

$[6, 135, \phi; 0.366, 0, 0]$  TF-IDF doc 3

Prob 1.B

$$\begin{aligned} \text{Bow Doc1} &= [1, 1, 0, 0, 0] \\ \text{Bow Doc2} &= [0, 0, 1, 1, 0] \\ \text{Bow Doc3} &= [0, 0, 0, 0, 1] \end{aligned}$$

1.B TF IDF

$$\text{TF} \left[ \frac{1}{2}, \frac{1}{2}, 0, 0, 0 \right] \times \text{IDF} \left[ \ln\left(\frac{3}{1}\right), \ln\left(\frac{3}{1}\right), 0, 0, 0 \right]$$

$0.5, 0.5$

1.09, 1.09

$[0.545, 0.545, 0, 0, 0]$  Doc 1 TF-IDF

1.B TFIDF

$$[0, 0, \frac{1}{2}, \frac{1}{2}, 0] \times [0, 0, \ln\left(\frac{3}{1}\right), \ln\left(\frac{3}{1}\right), 0]$$

DOC 2 TFIDF  $[0, 0, 0.545, 0.545, 0]$

1.B TPFDF

$$\begin{aligned} &[0, 0, 0, 0, \frac{1}{2}] \times [0, 0, 0, \ln\left(\frac{3}{1}\right)] \\ &[0, 0, 0, 0, 0.545] \quad \text{Doc 3 TPFDF} \end{aligned}$$

2. 100 spam 1000 ham

2 correct spam and correct ham

$$2, A, \text{TP} = \frac{\text{spam correctly classified}}{\text{total # emails}} = \frac{2}{1100} = 0.018$$

~~TPR~~

0.0018	0.081
0.0009	0.908

$$\frac{98}{1100} = 0.081 = FN$$

$$TN = \frac{99}{1100} = 0.1 - 0.08$$

$$\frac{1}{1100} = 0.0009 = FP$$

$$B, \text{Accuracy} = \frac{TP + TN}{TP + FP + TN + FN} = \frac{0.0018 + 0.99}{0.0018 + 0.0009 + 0.081 + 0.1} = 0.9098$$

$$\text{Accuracy} = 0.9098$$

$$\text{precision} = \frac{TP}{TP + FP} = \frac{0.0018}{0.0018 + 0.0009} = \frac{0.0018}{0.0027} = 0.66$$

$$\text{recall} = \frac{TP}{TP + FN} = \frac{0.0018}{0.0018 + 0.081} = \frac{0.0018}{0.0908}$$

$$\text{Recall} = 0.0198$$

$$F1 = 2 \times \frac{0.66 \times 0.018}{0.66 + 0.0198} = \frac{0.026}{0.6798} = 0.0382$$

$$F1 = 0.0382$$

2. C. The evaluation metric that best reflects performance in this situation is the F1 score. This is because it is the average of precision and recall. A high number shows the filter is both good at minimizing false positives and false negatives. This is very important for spam filter. This has a low F2 score showing it is not precise or has high recall. This means overall the positive is not accurate most of the time which is the most important part.

3. A. Snippet B is technically correct

3. B. In snippet A since we vectorize then split the model learns on all the data. This makes it so that the test and train data will both show perfect F1 scores when tested on, not giving an accurate test of the model. We need to split first so the model doesn't learn on all the data and can be tested properly on parts of the dataset.