

## DATA SCIENCE - II Homework - 3

Q1)

Name : Abhishek. Milind. Patwardhan  
ID : 81127139

Data Science - II - Homework - 3

1) Text Retrieval

- Today the Dawgs Won!
- Dawgs have won today
- Dawgs the champion!
- The Dawgs news today

1) Preprocessing

- 1) Today Dawgs Won
- 2) Dawgs have Won today
- 3) Dawgs champion
- 4) Dawgs news Today

2) Term-Document incidence Matrix

| Document →<br>Term ↓ | D <sub>1</sub> | D <sub>2</sub> | D <sub>3</sub> | D <sub>4</sub> |
|----------------------|----------------|----------------|----------------|----------------|
| Today                | 1              | 1              | 0              | 1              |
| Dawgs                | 1              | 1              | 1              | 1              |
| Won                  | 1              | 1              | 0              | 0              |
| Champion             | 0              | 0              | 1              | 0              |
| news                 | 0              | 0              | 0              | 1              |

3) Implement TF-IDF

→ Term Frequency (TF) =  $\frac{\text{No. of times term 't' appears in doc 'd'}}{\text{Total No. of items in doc 'd'}}$

Document 1 :

TF (Today) =  $\frac{1}{3}$

TF (Dawgs) =  $\frac{1}{3}$

TF (Won) =  $\frac{1}{3}$

TF (Champion) = 0

TF (news) = 0

(2)

Document 2 :

$$\begin{aligned} \text{TF}(\text{Today}) &= \frac{1}{3} \\ \text{TF}(\text{Dawgs}) &= \frac{1}{3} \\ \text{TF}(\text{won}) &= \frac{1}{3} \\ \text{TF}(\text{News}) &= 0 \\ \text{TF}(\text{Champion}) &= 0 \end{aligned}$$

Document 3 :

$$\begin{aligned} \text{TF}(\text{Today}) &= 0 \\ \text{TF}(\text{Dawgs}) &= \frac{1}{2} \\ \text{TF}(\text{won}) &= 0 \\ \text{TF}(\text{News}) &= 0 \\ \text{TF}(\text{Champion}) &= \frac{1}{2} \end{aligned}$$

Document 4 :

$$\begin{aligned} \text{TF}(\text{Today}) &= \frac{1}{3} \\ \text{TF}(\text{Dawgs}) &= \frac{1}{3} \\ \text{TF}(\text{won}) &= 0 \\ \text{TF}(\text{News}) &= \frac{1}{3} \\ \text{TF}(\text{Champion}) &= 0 \end{aligned}$$

• Inverse Document Frequency (IDF)

$$\text{IDF}(t) = \log_{10} \left( \frac{\text{Total No. of documents}}{\text{No of documents with term 't'}} \right)$$

$$\text{IDF}(\text{Today}) = \log(4/3) = \log(4/3)$$

$$\text{IDF}(\text{Dawgs}) = \log(4/4) = 0$$

$$\text{IDF}(\text{won}) = \log(4/2) = \log(2)$$

$$\text{IDF}(\text{News}) = \log(4/1) = \log(4)$$

$$\text{IDF}(\text{Champion}) = \log(4/1) = \log(4)$$



(3)

∴ TF-IDF Matrix is as follows:

| Document →<br>Term ↓ | D <sub>1</sub>                | D <sub>2</sub>                | D <sub>3</sub>              | D <sub>4</sub>                |
|----------------------|-------------------------------|-------------------------------|-----------------------------|-------------------------------|
| Today                | $\frac{1}{3} \cdot \log(4/3)$ | $\frac{1}{3} \cdot \log(4/3)$ | 0                           | $\frac{1}{3} \cdot \log(4/3)$ |
| Dawgs                | 0                             | 0                             | 0                           | 0                             |
| won                  | $\frac{1}{3} \cdot \log(2)$   | $\frac{1}{3} \cdot \log(2)$   | 0                           | 0                             |
| champion             | 0                             | 0                             | $\frac{1}{2} \cdot \log(4)$ | 0                             |
| News                 | 0                             | 0                             | 0                           | $\frac{1}{3} \cdot \log(4)$   |

4) Inverted Index for given query "Dawgs won".

1) Token Sequence

| Term     | Document ID |
|----------|-------------|
| Today    | 1           |
| Dawgs    | 1           |
| Won      | 1           |
| Dawgs    | 2           |
| Won      | 2           |
| Today    | 2           |
| Dawgs    | 3           |
| champion | 3           |
| Dawgs    | 4           |
| News     | 4           |
| Today    | 4           |

(4)

2) Sorting by terms :

| Term     | Document ID |
|----------|-------------|
| Champion | 3           |
| Dawgs    | 1           |
| Dawgs    | 2           |
| Dawgs    | 3           |
| Dawgs    | 4           |
| news     | 4           |
| Today    | 1           |
| Today    | 2           |
| Today    | 4           |
| Won      | 1           |
| Won      | 2           |

3) Dictionary &amp; Posting :

| Term     | Frequency | → | Posting Lists |
|----------|-----------|---|---------------|
| Dawgs    | 4         | → | 1 → 2 → 3 → 4 |
| Champion | 1         | → | 3             |
| news     | 1         | → | 4             |
| today    | 3         | → | 1 → 2 → 4     |
| won      | 2         | → | 1 → 2         |

4) Query : "Dawgs Won"

→ We will locate "Dawgs" & "Won" separately & retrieve their respective postings.

|       |   |   |               |
|-------|---|---|---------------|
| Dawgs | 4 | → | 1 → 2 → 3 → 4 |
| Won   | 2 | → | 1 → 2         |



(5)

Finding the Intersection of Document ID's for all terms in query :

$$[1, 2, 3, 4] \cap [1, 2] = [1, 2].$$

Hence Document 1 & 2 are retrieved documents for given query "Dawgs won".

5) Explain advantages of inverted index over term-document incidence matrices for organizing texts?

→ 1) Efficient Retrieval :

- Inverted index allows fast retrieval of documents containing specific terms because index directly maps terms to the documents they appear in.

2) Scalability :

- Inverted indexes are generally more scalable than term-document incidence matrices due to their distributed and parallel nature.

3) Dynamic Updation :

- Inverted index can be easily updated when new documents are added. Whereas term-document incidence matrices may require recalculation or resizing operations.

4) Space Efficiency :

- Inverted index stores only the unique terms and their associated document references. Hence it is more space efficient than term-document incidence matrices.

Q 2)

```
PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS SEARCH ERROR
PS C:\Users\abhis\OneDrive - University of Georgia\Desktop\Data Science\Iris_clf_mask_code> python -u "c:\Use

Model: LogReg
-----
Train Accuracy=0.9429 | F1=0.9444 | AUC=0.9812
Test Accuracy=0.9333 | F1=0.9231 | AUC=0.9864

Model: MLP
-----
Train Accuracy=0.9857 | F1=0.9867 | AUC=0.9992
Test Accuracy=0.9667 | F1=0.9630 | AUC=0.9955

Model: NaiveBayesian
-----
Train Accuracy=0.9286 | F1=0.9296 | AUC=0.9308
Test Accuracy=0.9333 | F1=0.9231 | AUC=0.9321

Model: DecisionTree
-----
Train Accuracy=0.9857 | F1=0.9867 | AUC=0.9848
Test Accuracy=0.7667 | F1=0.7586 | AUC=0.7760
PS C:\Users\abhis\OneDrive - University of Georgia\Desktop\Data Science\Iris_clf_mask_code> |
```

## **Conclusions :**

The results indicate that all models perform well on the training set, with Logistic Regression, MLP, and Decision Tree achieving high accuracy, F1 score, and AUC. However, on the test set, MLP stands out as the best-performing model with the highest accuracy, F1 score, and AUC, indicating its superior generalization ability. Logistic Regression also performs reasonably well on the test set, showing robustness in its predictions. Naive Bayesian and Decision Tree models exhibit moderate performance, with Decision Tree showing signs of overfitting as evidenced by the significant gap between training and test accuracies. In conclusion, while MLP emerges as the top-performing model, Logistic Regression demonstrates competitive performance with good generalization capabilities, making it a reliable choice when interpretability and simplicity are valued.