

# Data Mining and Machine Learning

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Vector Re Documents

Add WeChat edu\_assist\_pro

Peter Jančovič

# Objectives

- To explain vector representation of documents

- To understand **Assignment Project Exam Help** **een** vector representation **https://eduassistpro.github.io/**

**Add WeChat edu\_assist\_pro**

# Vector Notation for Documents

- Suppose that we have a set of documents

$$D = \{d_1, d_2, \dots, d_N\}$$

think of this

- Suppose that the whole corpus is  $V$  (vocabulary) words in the whole corpus
- Now suppose a document  $d$  in  $D$  contains  $M$  different terms:  $\{t_{i(1)}, t_{i(2)}, \dots, t_{i(M)}\}$
- Finally, suppose term  $t_{i(m)}$  occurs  $f_{i(m)}$  times

# Vector Notation

- The vector representation  $\text{vec}(d)$  of  $d$  is the  $V$  dimensional vector:

Assignment Project Exam Help

$$(0, \dots, 0, w_{i(1),d}, 0, \dots, 0, w_{i(2),d}, 0, \dots, 0, w_{i(M),d}, 0, \dots, 0)$$

Add WeChat edu\_assist\_pro

$i(1)^{\text{th}}$   
place

$i(2)^{\text{th}}$   
place

$i(M)^{\text{th}}$   
place

Notice that this is the weighting – i.e. the term frequency times the inverse document frequency

$$w_{i(1),d} = f_{i(1),d} \times \text{IDF}(i(1)) \text{ from text IR}$$

# Uniqueness

- Is the mapping between documents and vectors one-to-one?
- In other words:
  - if  $d_1, d_2$  such that  $vec(d_1) = vec(d_2)$  if and only if  $d_1 = d_2$ ?
- If  $\lambda$  is a scalar and  $vec(d_1) = \lambda vec(d_2)$  what does this tell you about  $d_1$  and  $d_2$ ?

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# Example

- $d_1$  = the cat sat on the cat's mat → cat sat cat mat
- $d_2$  = the dog chased the cat → dog chase cat
- $d_3$  = the mouse stay home
- Vocabulary: <https://eduassistpro.github.io/>
  - cat, chase, dog, home, mat, stay
- To calculate the vector representations of these documents first calculate the TF-IDF weights

# Example (continued)

	d1	d2	d3	Nd	IDF	w(t,d1)	w(t,d2)	w(t,d3)
cat	2	1		2	0.41	0.81	0.41	
chase	1	1		1	1.1	1.1		
dog							1.1	
home								1.1
mat	1			1	1.1	1.1		
mouse			1	1	1.1			1.1
sat	1			1	1.1	1.1		
stay			1	1	1.1			1.1

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# Example (continued)

$$vec(d_1) = \begin{bmatrix} 0.81 \\ 0 \\ 0 \\ 0 \\ 1.1 \\ 0 \\ 1.1 \\ 0 \end{bmatrix} \quad \begin{bmatrix} 0.41 \\ 1.1 \\ 1.1 \end{bmatrix} \quad vec(d_3) = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1.1 \\ 0 \\ 1.1 \\ 0 \\ 1.1 \end{bmatrix}$$

Assignment Project Exam Help  
<https://eduassistpro.github.io/>  
Add WeChat edu\_assist\_pro



# Document length revisited

- Recall that the length of a vector

$$x = (x_1, \dots, x_N)$$

is given by: <https://eduassistpro.github.io/>

$$\|x\| = \sqrt{x_1^2 + x_2^2 + \dots + x_N^2}$$

# Document length

- In the case of a ‘document vector’

$$vec(d) = (0, \dots, 0, w_{i(1)d}, 0, \dots, 0, w_{i(2)d}, \dots, w_{i(M)d}, 0, \dots, 0)$$

$$\|vec(d)\| = \sqrt{w_{i(1)d}^2 + w_{i(2)d}^2 + \dots + w_{i(M)d}^2} = \|d\|$$

# Document Similarity

- Suppose  $d$  is a document and  $q$  is a query
  - If  $d$  and  $q$  contain the same words in the same proportions, then  $vec(d)$  and  $vec(q)$  will point in the same direction
  - If  $d$  and  $q$  contain different words,  $vec(d)$  and  $vec(q)$  will point in different directions
  - Intuitively, the greater the angle between  $vec(d)$  and  $vec(q)$  the less similar the document  $d$  is with the query  $q$

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# Cosine similarity

- Define the **Cosine Similarity** between document  $d$  and query  $q$  by:

$$CSim(q, d) = \cos \theta$$

where  $\theta$  is the angle between  $vec(q)$  and  $vec(d)$

- Similarly, define the Cosine Similarity between documents  $d_1$  and  $d_2$  by:

$$CSim(d_1, d_2) = \cos \theta$$

where  $\theta$  is the angle between  $vec(d_1)$  and  $vec(d_2)$

# Cosine Similarity & Similarity

- Let  $u=(x_1,y_1)$  and  $v=(x_2,y_2)$  be vectors in 2 dimensions, then

Assignment Project Exam Help  
<https://eduassistpro.github.io/>  
Add WeChat edu\_assist\_pro

$$\cos(\theta) = \frac{x_1 x_2 + y_1 y_2}{\|u\| \|v\|} = \frac{u \cdot v}{\|u\| \|v\|}$$


- In fact, this result holds for vectors in any  $N$  dimensional space

# Cosine Similarity & Similarity

- Hence, if  $q$  is a query,  $d$  is a document, and  $\theta$  is the angle between  $vec(q)$  and  $vec(d)$ , then:

Assignment Project Exam Help

Cosine  
similarity

<https://eduassistpro.github.io/>

Add WeChat: edu\_assist\_pro

$$CSim(q, d) = \cos(\theta) = \frac{vec(q) \cdot vec(d)}{\|q\| \|d\|}$$
$$= Sim(q, d)$$

Similarity

# Summary

- Vector space representation of documents
  - Cosine distance representations of documents
- Assignment Project Exam Help**  
**<https://eduassistpro.github.io/>**  
**Add WeChat edu\_assist\_pro**