# Data Mining and Machine Learning

# Vector Representation of Documents

Peter Jančovič

UNIVERSITY OF
BIRMINGHAM

# Objectives

- To explain vector representation of documents

- To understand cosine distance between vector representations of documents

Assignment Project Exam Help

https://powcoder.com

Add WeChat powcoder

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Vector Notation for Documents

- Suppose that we have a set of documents

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

think of this as the corpus for IR

- Suppose that the number of <u>different</u> words in the <u>whole corpus</u> is $V$ (vocabulary size)

- 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

Data Mining and Machine Learning

UNIVERSITY$^{OF}$ BIRMINGHAM

# Vector Notation

- The vector representation $vec(d)$ of $d$ is the $V$ dimensional vector:

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

$i(1)^{th}$ place

$i(2)^{th}$ place

$i(M)^{th}$ place

Notice that this is the weighting – i.e. the term frequency times the inverse document frequency $w_{i(1),d} = f_{i(1),d} \times IDF(i(1))$ from text IR

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Uniqueness

- Is the mapping between documents and vectors one-to-one?

- In other words:

  - if $d_1$, $d_2$ are documents, is it true 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 $d_2$?

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# 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 stayed at home → mouse stay home

- Vocabulary:

  - cat, chase, dog, home, mat, mouse, sat, stay

- To calculate the vector representations of these documents first calculate the TF-IDF weights

Data Mining and Machine Learning

UNIVERSITY OF
BIRMINGHAM

# 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 | |
| dog | | 1 | | 1 | 1.1 | | 1.1 | |
| home | | | 1 | 1 | 1.1 | | | 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 |

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Example (continued)

$$vec(d_1) = \begin{bmatrix} 0.81 \\ 0 \\ 0 \\ 0 \\ 1.1 \\ 0 \\ 1.1 \\ 0 \end{bmatrix} \qquad vec(d_2) = \begin{bmatrix} 0.41 \\ 1.1 \\ 1.1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \qquad vec(d_3) = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1.1 \\ 0 \\ 1.1 \\ 0 \\ 1.1 \end{bmatrix}$$

Data Mining and Machine Learning

UNIVERSITY OF
BIRMINGHAM

# Document length revisited

- Recall that the length of a vector

$$x = \left( x_1, \ldots, x_N \right)$$

is given by:

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

Data Mining and Machine Learning

UNIVERSITY$^{OF}$
BIRMINGHAM

# Document length

- In the case of a 'document vector'

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

$$\|vec(d)\| = \sqrt{w^2_{i(1)d} + w^2_{i(2)d} + ... + w^2_{i(M)d}} = \|d\|$$

Data Mining and Machine Learning

UNIVERSITY OF
BIRMINGHAM

# Document Similarity

- Suppose **d** is a document and **q** is a query

  - If **d** and **q** contain the <u>same words</u> in the <u>same proportions</u>, then *vec*(**d**) and *vec*(**q**) will point in the same direction

  - If **d** and **q** contain <u>different words</u>, then *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**

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Cosine similarity

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

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

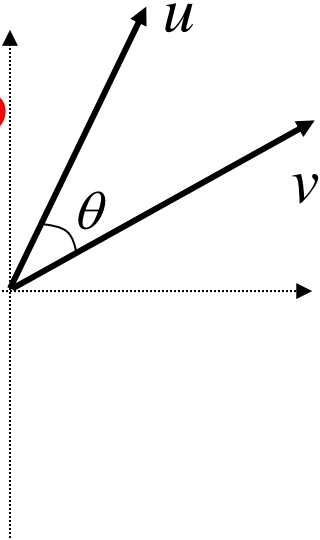where $\theta$ is the <u>angle</u> between $vec(q)$ and $vec(d)$

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

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

where $\theta$ is the <u>angle</u> between $vec(d_1)$ and $vec(d_2)$

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Cosine Similarity & Similarity

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

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

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

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# 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

https://powcoder.com

Add WeChat powcoder

Cosine similarity

$$CSim(q,d) = \cos(\theta) = \frac{vec(q) \cdot vec(d)}{\|q\|\|d\|} = \frac{\sum_{t \in q \cap d} w_{tq} \cdot w_{td}}{\|q\|\|d\|}$$

$$= Sim(q,d)$$

Similarity

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Summary

- Vector space representation of documents

Assignment Project Exam Help

- Cosine distance between vector representations of documents https://powcoder.com

Add WeChat powcoder

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM