

Intelligent Data Analysis 2020

Lecture 4

Vector Representation of Documents

Martin Russell

Objectives

- To explain **vector representation** of documents
- To understand **cosine distance** between vector representations of documents
- To understand, intuitively, how **Latent Semantic Analysis (LSA)** can
 - Discover latent **topics** in a corpus and represent them in terms of words
 - Achieve dimension reduction for document vectors
 - Represent words in terms of topics

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 **different words** in the **whole corpus** 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

Vector Notation

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

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

\uparrow $i(1)^{th}$ place \uparrow $i(2)^{th}$ place \uparrow $i(M)^{th}$ place

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

- $vec(d)$ is the **document vector** for d

freq (i) · [ln(ND) - ln(ND_i)]

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 $d_1 = d_2$? *False*
- If λ is a scalar and $vec(d_1) = \lambda vec(d_2)$ what does this tell you about d_1 and d_2 ?

*$vec(d_1) = 2 vec(d_2)$
twice*

Example

- d_1 = the cat sat on the cat's mat \rightarrow cat sat mat cat
- d_2 = the dog chased the cat \rightarrow dog chase cat
- d_3 = the mouse stayed at home \rightarrow mouse stay home
- **Vocabulary:** *Corpus* $M = 8$.
 - cat, chase, dog, home, mat, mouse, sat, stay
- To calculate the vector representations of these documents first calculate the **TF-IDF weights**

Example (continued)

	d1	d2	d3	^{# doc} Nd	IDF	^{TF-IDF} w(t,d1)	w(t,d2)	w(t,d3)
cat	2	1		3 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

Example (Continued)

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

Document length revisited

- Recall that the **length (norm)** of a vector

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

is given by:

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

Document length

- In the case of a **document vector**

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

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

- Where $\|d\|$ is the length of the document d from last week's lecture

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**, 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

$\angle \uparrow$, $sim(\) \downarrow$

Cosine similarity

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

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

where θ is the **angle** between $vec(q)$ and $vec(d)$

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

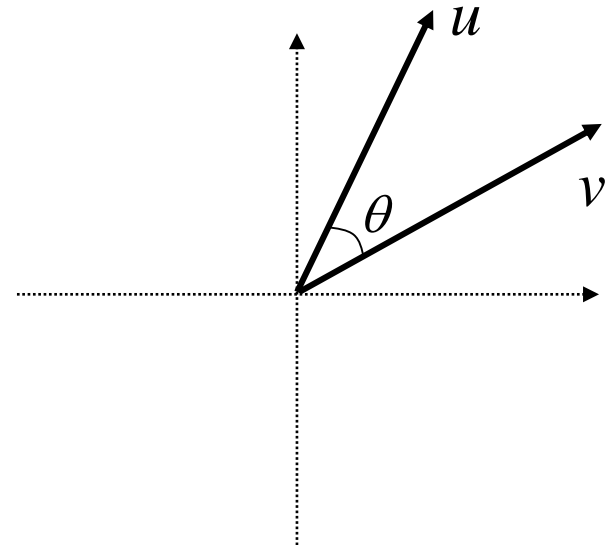
$$CSim(d_1, d_2) = \cos \theta$$

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

Cosine Similarity & Similarity

- Recall that if $u=(x_1, y_1)$ and $v=(x_2, y_2)$ are 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

Cosine Similarity & Similarity

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

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

Summary

- Vectorisation of documents
- Cosine similarity is equivalent to TF-IDF similarity
- Document length revisited