

TEXT SIMILARITY

David Kauchak
CS159 Fall 2020

1

Admin

Assignment 4a

- Solutions posted
- If you're still unsure about questions 3 and 4, come talk to me.

Assignment 4b

Grading

Quiz #2 next Thursday covering material through 10/6

2

Course feedback

3

Text Similarity

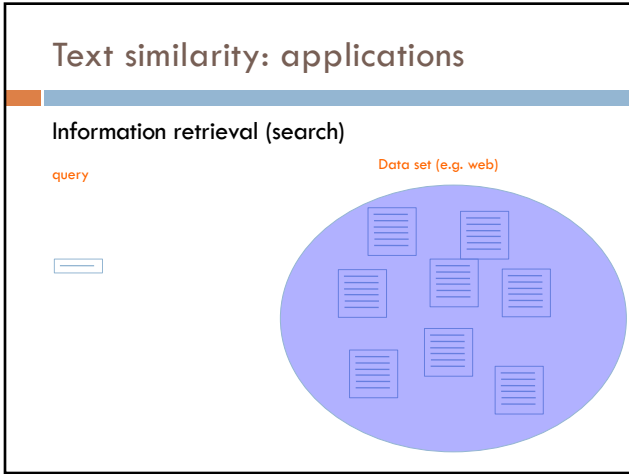
A common question in NLP is how similar are texts

score: $\text{sim}(\text{document}_1, \text{document}_2) = ?$

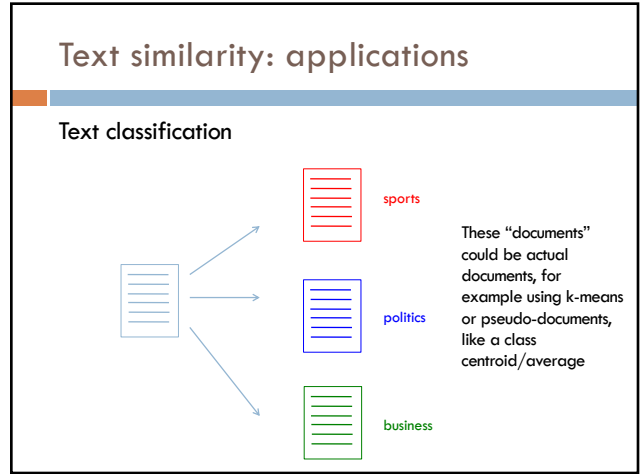
rank: $\text{document}_1 \text{ ? } \text{document}_2, \text{document}_3$

How could these be useful? Applications?

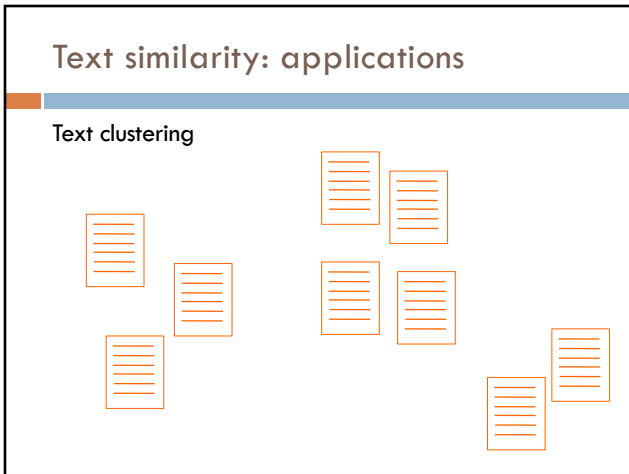
4



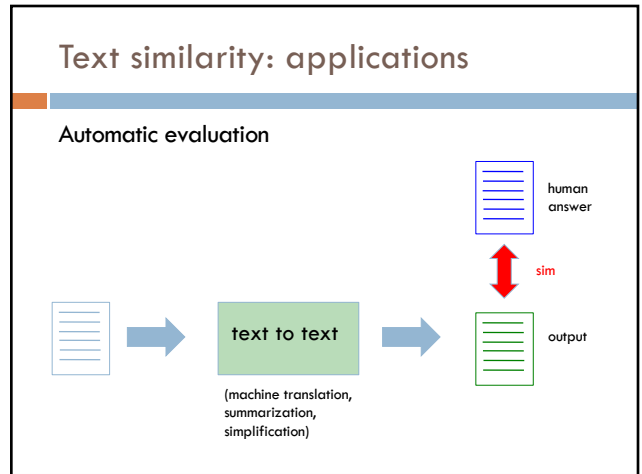
5



6



7



8

Text similarity: applications

Word similarity

$\text{sim}(\text{banana}, \text{apple}) = ?$

Word-sense disambiguation

I went to the *bank* to get some money.



9

Text similarity: application

Automatic grader

Question: what is a variable?

Answer: a location in memory that can store a value

How good are:

- a variable is a location in memory where a value can be stored
- a named object that can hold a numerical or letter value
- it is a location in the computer's memory where it can be stored for use by a program
- a variable is the memory address for a specific type of stored data or from a mathematical perspective a symbol representing a fixed definition with changing values
- a location in memory where data can be stored and retrieved

10

Text similarity

There are many different notions of similarity depending on the domain and the application

Today, we'll look at some different tools

There is no one single tool that works in all domains

11

Text similarity approaches

$\text{sim}(\text{document 1}, \text{document 2}) = ?$

A: When the defendant and his lawyer walked into the court, some of the victim supporters turned their backs to him.

B: When the defendant walked into the courthouse with his attorney, the crowd turned their backs on him.

How can we do this?

12

The basics: text overlap

Texts that have overlapping words are more similar

A: When the defendant and his lawyer walked into the court, some of the victim supporters turned their backs to him.

B: When the defendant walked into the courthouse with his attorney, the crowd turned their backs on him.

13

Word overlap: a numerical score

Idea 1: number of overlapping words

A: When the defendant and his lawyer walked into the court, some of the victim supporters turned their backs to him.

B: When the defendant walked into the courthouse with his attorney, the crowd turned their backs on him.

$\text{sim}(T_1, T_2) = 11$ problems?

14

Word overlap problems

- Doesn't take into account word order
- Related: doesn't reward longer overlapping sequences

A: defendant his the When lawyer into walked backs him the court, of supporters and some the victim turned their backs him to.

B: When the defendant walked into the courthouse with his attorney, the crowd turned their backs on him.

$\text{sim}(T_1, T_2) = 11$

15

Word overlap problems

Doesn't take into account length

A: When the defendant and his lawyer walked into the court, some of the victim supporters turned their backs to him.

B: When the defendant walked into the courthouse with his attorney, the crowd turned their backs on him. *I ate a large banana at work today and thought it was great!*

$\text{sim}(T_1, T_2) = 11$

16

Word overlap problems

Doesn't take into account synonyms

A: When the defendant and his lawyer walked into the court, some of the victim supporters turned their backs to him.

B: When the defendant walked into the courthouse with his attorney, the crowd truned their backs on him.

$$\text{sim}(T1, T2) = 11$$

17

Word overlap problems

Doesn't take into account spelling mistakes

A: When the defendant and his lawyer walked into the court, some of the victim supporters turned their backs to him.

B: When the defendant walked into the courthouse with his attorney, the crowd truned their backs on him.

$$\text{sim}(T1, T2) = 11$$

18

Word overlap problems

Treats all words the same

A: When the defendant and his lawyer walked into the court, some of the victim supporters turned their backs to him.

B: When the defendant walked into the courthouse with his attorney, the crowd truned their backs on him.

19

Word overlap problems

May not handle frequency properly

A: When the defendant and his lawyer walked into the court, some of the victim supporters turned their backs to him. I ate a banana and then another banana and it was good!

B: When the defendant walked into the courthouse with his attorney, the crowd truned their backs on him. I ate a large banana at work today and thought it was great!

20

Word overlap: sets

A: When the defendant and his lawyer walked into the court, some of the victim supporters turned their backs to him.



A
and
backs
court
defendant
him
...

B: When the defendant walked into the courthouse with his attorney, the crowd turned their backs on him.



B
and
backs
courthouse
defendant
him
...

21

Word overlap: sets

What is the overlap, using set notation?

- $|A \cap B|$ the size of the intersection

How can we incorporate length/size into this measure?

22

Word overlap: sets

What is the overlap, using set notation?

- $|A \cap B|$ the size of the intersection

How can we incorporate length/size into this measure?

Jaccard index (Jaccard similarity coefficient)

$$J(A,B) = \frac{|A \cap B|}{|A \cup B|}$$

Dice's coefficient

$$Dice(A,B) = \frac{2 |A \cap B|}{|A| + |B|}$$

23

Word overlap: sets

$$J(A,B) = \frac{|A \cap B|}{|A \cup B|} \quad Dice(A,B) = \frac{2 |A \cap B|}{|A| + |B|}$$

How are these related?

Hint: break them down in terms of

$|A - B|$ words in A but not B
 $|B - A|$ words in B but not A
 $|A \cap B|$ words in both A and B

24

Word overlap: sets

$$J(A,B) = \frac{|A \cap B|}{|A \cup B|}$$

$$= \frac{|A \cap B|}{|A - B| + |B - A| + |A \cap B|}$$

↑ in A but not B
↑ in B but not A

$$Dice(A,B) = \frac{2 |A \cap B|}{|A| + |B|}$$

$$= \frac{2 |A \cap B|}{|A - B| + |B - A| + 2 |A \cap B|}$$

Dice's coefficient gives twice the weight to overlapping words

25

Set overlap

Our problems:

- ▣ word order
- ▣ length
- ▣ synonym
- ▣ spelling mistakes
- ▣ word importance
- ▣ word frequency

Set overlap measures can be good in some situations, but often we need more general tools

26


Bag of words representation

When the defendant and his lawyer walked into the court, some of the victim supporters turned their backs to him.

↓

When the defendant and his lawyer walked into the court, some of the victim supporters turned their backs to him.

↓



What information do we lose?

27

Bag of words representation

For now, let's ignore word order:

Obama said banana repeatedly last week on tv, "banana, banana, banana"

(4, 1, 1, 0, 0, 1, 0, 0, ...)

banana

obama

said

california

across

tv

wrong

capital

Frequency of word occurrence

"Bag of words representation": multi-dimensional vector, one dimension per word in our vocabulary

28

Bag of words representation



<http://membercentral.aaas.org/blogs/member-spotlight/ron-mitchell-studies-human-language-both-man-and-machine>

29

Vector based word

A

| | |
|----------------|---|
| a1: When | 1 |
| a2: the | 2 |
| a3: defendant | 1 |
| a4: and | 1 |
| a5: courthouse | 0 |
| ... | |

B

| | |
|----------------|---|
| b1: When | 1 |
| b2: the | 2 |
| b3: defendant | 1 |
| b4: and | 0 |
| b5: courthouse | 1 |
| ... | |

Multi-dimensional vectors,
one dimension per word in
our vocabulary

How do we calculate the
similarity based on these
vectors?

30

Vector based similarity

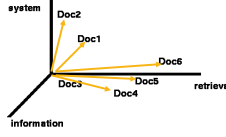
We have a $|V|$ -dimensional vector space

Terms are axes of the space

Documents are points or vectors in this space

Very high-dimensional

This is a very sparse vector - most entries are zero



What question are we asking in this space for similarity?

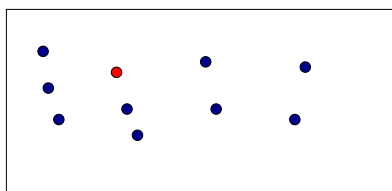
31

Vector based similarity

Similarity relates to distance

We'd like to measure the similarity of documents in the $|V|$ dimensional space

What are some distance measures?



32

Distance measures

Euclidean (L2)

$$\text{dist}(A, B) = \sqrt{\sum_{i=1}^n (a_i - b_i)^2}$$

Manhattan (L1)

$$\text{dist}(A, B) = \sum_{i=1}^n |a_i - b_i|$$

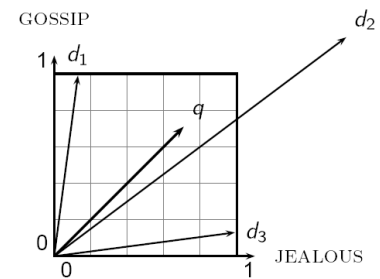
What do these mean for our bag of word vectors?

33

Distance can be problematic

Which d is closest to q using one of the previous distance measures?

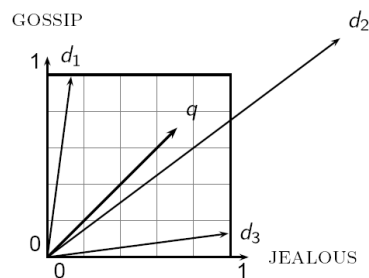
Which do you think should be closer?



34

Distance can be problematic

The Euclidean (or L1) distance between q and d_2 is large even though the distribution of words is similar



35

Use angle instead of distance

Thought experiment:

- take a document d
- make a new document d' by concatenating two copies of d
- "Semantically" d and d' have the same content

What is the Euclidean distance between d and d' ?
What is the angle between them?

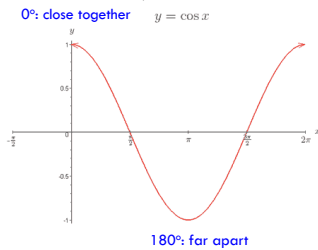
- The Euclidean distance can be large
- The angle between the two documents is 0

36

From angles to cosines

Cosine is a monotonically decreasing function for the interval $[0^\circ, 180^\circ]$

decreasing angle is equivalent to increasing cosine of that angle
(larger cosine means more similar)



37

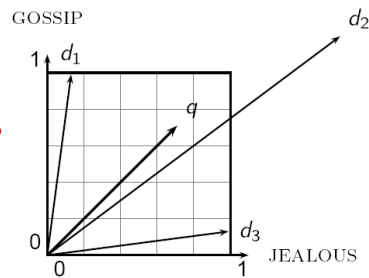
Near and far

<https://www.youtube.com/watch?v=iZhEcRrMA-M>

38

cosine

How do we calculate the cosine between two vectors?



39

Cosine of two vectors

Dot product

$$A \cdot B = \|A\| \|B\| \cos \theta$$

$$\cos \theta = \frac{A \cdot B}{\|A\| \|B\|} = \frac{A}{\|A\|} \cdot \frac{B}{\|B\|}$$

Dot product between unit length vectors

40

Cosine as a similarity

$$sim_{cos}(A, B) = A \cdot B = \sum_{i=1}^n a_i b_i \quad \text{ignoring length normalization}$$

Just another distance measure, like the others:

$$dist_{L_2}(A, B) = \sqrt{\sum_{i=1}^n (a_i - b_i)^2}$$

$$dist_{L_1}(A, B) = \sum_{i=1}^n |a_i - b_i|$$

41

Cosine as a similarity

$$sim_{cos}(A, B) = A \cdot B = \sum_{i=1}^n a_i b_i \quad \text{ignoring length normalization}$$

For bag of word vectors, what does this do?

42

Cosine as a similarity

$$sim_{cos}(A, B) = A \cdot B = \sum_{i=1}^n a_i b_i \quad \text{ignoring length normalization}$$

Only words that occur in both documents count towards similarity

Words that occur more frequently in both receive more weight

43

Length normalization

A vector can be length-normalized by dividing each of its components by its length

Often, we'll use L_2 norm (could also normalize by other norms):

$$\|\vec{x}\|_2 = \sqrt{\sum_i x_i^2}$$

Dividing a vector by its L_2 norm makes it a unit (length) vector

44

Unit length vectors

In many situations, normalization improves similarity, but not in all situations

45

Normalized distance measures

Cosine

$$sim_{cos}(A,B) = A \cdot B = \sum_{i=1}^n a_i b_i = \frac{\sum_{i=1}^n a_i b_i}{\sqrt{\sum_{i=1}^n a_i^2} \sqrt{\sum_{i=1}^n b_i^2}}$$

L2

$$dist_{L2}(A,B) = \sqrt{\sum_{i=1}^n (a_i - b_i)^2}$$

L1

$$dist_{L1}(A,B) = \sum_{i=1}^n |a_i - b_i|$$

a' and b' are length normalized versions of the vectors

46

Distance measures

Cosine

$$sim_{cos}(A,B) = A \cdot B = \sum_{i=1}^n a_i b_i$$

L2

$$dist_{L2}(A,B) = \sqrt{\sum_{i=1}^n (a_i - b_i)^2}$$

L1

$$dist_{L1}(A,B) = \sum_{i=1}^n |a_i - b_i|$$

Cosine is the most common measure. Why do you think?

47

Distance measures

Cosine

$$sim_{cos}(A,B) = A \cdot B = \sum_{i=1}^n a_i b_i$$

L2

$$dist_{L2}(A,B) = \sqrt{\sum_{i=1}^n (a_i - b_i)^2}$$

L1

$$dist_{L1}(A,B) = \sum_{i=1}^n |a_i - b_i|$$

- L1 and L2 penalize sentences for not having words, i.e. if a has it but b doesn't
- Cosine can be significantly faster since it only calculates over the intersection

48

Our problems

Which of these have we addressed?

- word order
- length
- synonym
- spelling mistakes
- word importance
- word frequency

49

Our problems

Which of these have we addressed?

- word order
- length
- synonym
- spelling mistakes
- word importance
- word frequency

50