Lecture 3:
Similarity

COMP90049
Knowledge
Technologies

Comparing things
Sets of descriptors
Features, Vectors

Comparing
Documents

Distance
Measures

**Lecture 3: Similarity**

Assignment Project Exam Help

https://eduassistpro.github.i

Sarah Erfani and Karin Verspoor and Jeremy Nicholson, CIS

Add WeChat edu_assist_pr

Semest

THE UNIVERSITY OF
MELBOURNE

Assignment Project Exam Help

https://eduassistpro.github.i

Add WeChat edu_assist_pr

Assignment Project Exam Help

https://eduassistpro.github.i

Add WeChat edu_assist_pr

# Lecture 3: Similarity
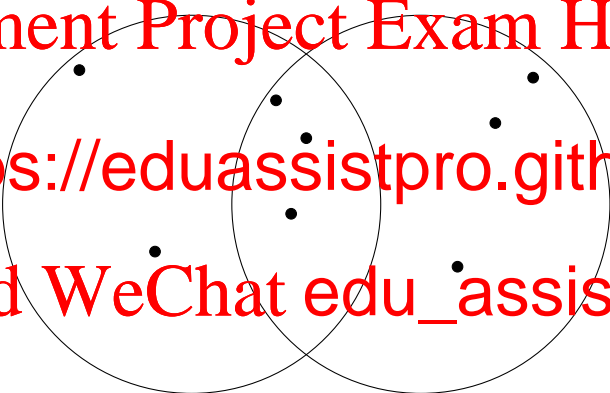
COMP90049
Knowledge
Technologies

**Comparing things**

**Sets of descriptors**

Features, Vectors

**Comparing Documents**

**Distance Measures**

Many similarity assessments can be framed as set intersection.

- Amazon: Book purchases
- 

- Rating sets (stars)
  - thresholding using ratings
  - different subsets for different ratings
- Categories of items
  - generalisation
  - book or movie genres

Assignment Project Exam Help

https://eduassistpro.github.i

Add WeChat edu_assist_pr

$$\frac{|A \cap B|}{|A \cup B|}$$

$$sim(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{3}{8}$$

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

$$sim(A, B) = \frac{2|A \cap B|}{|A| + |B|} = \frac{2*3}{5+6} = \frac{6}{11}$$

fruit above.

A feature vector is an $n$-dimensional vector of *features* that represent some object.

c of

- 
- Features may be ordinal (e.g. cool )
- Features may be numeric/continuo

A vector locates an object (document, pers
$n$-space. The angle of the vector in that space is determined by the relative weight of each term.

Assignment Project Exam Help

https://eduassistpro.github.i

Add WeChat edu_assist_pr

| age | income | credit |
| --- | --- | --- |
| 33 | 6 | low |
| 58 | 42 | low |
|  |  | low |
|  |  | low |
|  |  | hig |
|  |  | hig |

…

How should we compare documents to assess their similarity?

- String-level similarity (e.g., edit distance)
- ■                                                                    s)
- ■

How similar are these sentences?

1. Mary is quicker than John.
2. John is quicker than Mary.
3. Mary is slower than John.
4. Jane is quicker than Mary.

Assignment Project Exam Help

1 Mary is quicker than John.

https://eduassistpro.github.i

| Sentence | "Mary" | "John" | "Jane" | "quicker" | "slower" |
|----------|--------|--------|--------|-----------|----------|
| 1 | 1 | 1 | 0 | 1 | 0 |
| 2 | 1 | 0 | 1 | 1 | 0 |
| 3 | 1 | 0 | 0 | 0 | 1 |
| 4 | 1 | 0 | 1 | 1 | 0 |

Add WeChat edu_assist_pr

One of the earliest models proposed for retrieval of documents (information retrieval, in 1962) was the vector-space model.

Suppose there are $n$ distinct indexed terms in the collection. Then each document $d$ can be thought of as a vector

$$\langle w_{d,1}, w_{d,2}, \ldots, w_{d,t}, \ldots, w_{d,n}\rangle$$

(Most $w_{d,t}$ values will be zero, because m
tiny proportion of a collection's terms.)

Intuitively, if some other document $d'$

$$\langle w_{d',1}, w_{d',2}, \ldots, w_{d',t}, \ldots, w_{d',n}\rangle$$

where the weights are close to those of $d$ – in particular, if the non-zero $w$ values are for much the same set of terms – then $d$ and $d'$ are likely to be similar in topic.

We have discussed similarity at an intuitive and quantitative level.

$$sim_D(A, B) \quad = \quad \frac{2|A \quad B|}{|A|} \qquad \underline{2 \quad 3} \qquad \underline{6}$$

What is the relationship between similarity

A distance measure on a space is a function that takes two points in a space as arguments.

1. No negative distances.

2.

3. Distance is symmetric.

$$d(x, y) =$$

4. The *triangle inequality* typically holds.
   (Distance measures the length of the *shortest path* between two points.)

$$d(x, y) \leq d(x, z) + d(z, y)$$

---

Given two items $A$ and $B$, and their corresponding feature vectors $\vec{a}$ and $\vec{b}$, respectively, we can calculate their similarity via their distance $d$ in euclidean space:

Assignment Project Exam Help

https://eduassistpro.github.i

Add WeChat edu_assist_pr

In n-dimensional space:

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

Given two items $A$ and $B$, and their corresponding feature vectors $\vec{a}$ and $\vec{b}$, respectively, we can calculate their similarity via their *vector cosine* (the cosine of the angle $\theta$ between the two vectors):

Assignment Project Exam Help

https://eduassistpro.github.i

Add WeChat edu_assist_pr

$$sim(A, B) = \frac{\vec{a} \cdot \vec{b}}{|\vec{a}||\vec{b}|} = \frac{\sum_i a_i b_i}{\sqrt{\sum_i a_i^2}\sqrt{\sum_i b_i^2}}$$

Assignment Project Exam Help
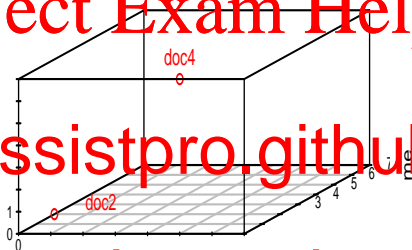
https://eduassistpro.github.i

Add WeChat edu_assist_pr



- Doc4, like Doc1, is all about "tea" and "two".
- But because it is longer, it is in a space by itself.

["City block" distance or "Taxicab geometry" or "$L_1$ distance"]

Given two items $A$ and $B$, and their corresponding feature vectors $\vec{a}$ and $\vec{b}$, respectively, we can calculate their similarity via their distance $d$ based on the absolute differences of their cartesian coordinates.

In n-dimensional space:

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

Relative entropy:

$$D(x \parallel y) = \sum_i x_i (\log_2 x_i - \log_2 y_i)$$

or *Jensen-Shannon divergence*:

$$JSD(x \parallel y) = \frac{1}{2} D(x \parallel m) + \frac{1}{2} D(y \parallel m)$$

where $m = \frac{1}{2}(x + y)$

NB: Probability will be reviewed next lecture!

- How can we represent a set of objects?
- What are some methods for measuring similarity between objects?

Reading

`http://infolab.stanford.edu/~ullman/mmds.html`

Or document representation
Chapter 6
*Information Retrieval*, Manning *et al.*
`http://nlp.stanford.edu/IR-book/html/htmledition/`
`scoring-term-weighting-and-the-vector-space-model-1.html`