

**Lecture 3:
Similarity**

COMP90049
Knowledge
Technologies

Introduction and
Review

Comparing things
Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures

Lecture 3: Similarity

COMP90049 Knowledge Technologies

Lea Frermann and Karin Verspoor, CIS

Semester 2, 2019



THE UNIVERSITY OF
MELBOURNE

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and Review

Comparing things
Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures

Contact

Lea Frermann

Lecturer in Natural Language Processing (CIS)

lea.frermann@unimelb.edu.au

Please,

- Ask questions (!!)
- On-demand office hours (email me)
- Give feedback anytime (after lecture, email, ...)

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and Review

Comparing things
Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures

- When? Friday, August 30 2019, **8.30am – 9:40 am**
- Where? Wilson Hall and Kwong Lee Dow
- What? 50 minutes, with no reading time
Closed-book, no materials needed
All course content up to Friday, August 23rd (inclusive)
- Weight? 10 marks

A sample exam (with solutions) is available
in the course LMS → Assessment.

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and Review

Comparing things
Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures

Knowledge Technologies

- processing **unstructured data** (text)
- **extracting** and **analyzing** information
- Part 1: Basics
- Part 2: Information Retrieval
- Part 3: Machine learning

Last Week

- Regular expressions
- you know what you want to find **and** where to look

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and Review

Comparing things
Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures

Knowledge Technologies

- processing **unstructured data** (text)
- **extracting** and **analyzing** information
- Part 1: Basics
- Part 2: Information Retrieval
- Part 3: Machine learning

Last Week

- Regular expressions
- you know what you want to find **and** where to look
- But what if you don't? How do you judge what is **relevant enough**?

Similarity – So what?!

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and
Review

Comparing things

Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures

- humans can't help but group items into classes (categories)
- categories drive how we perceive and make decisions



Similarity – So what?!

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and
Review

Comparing things

Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures

- humans can't help but group items into classes (categories)
- categories drive how we perceive and make decisions



Similarity – So what?!

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and
Review

Comparing things

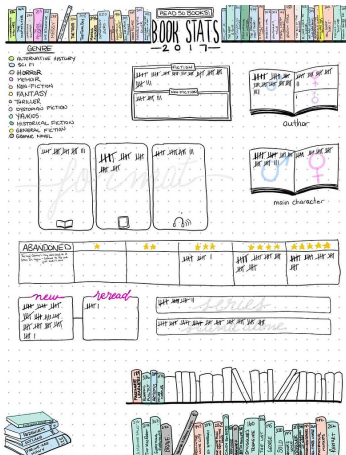
Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures

- humans can't help but group items into classes (categories)
- categories drive how we perceive and make decisions



Similarity – So what?!

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and
Review

Comparing things

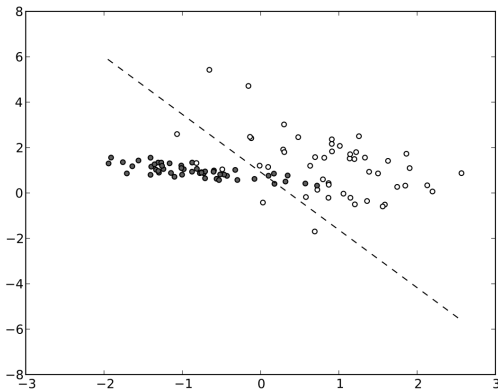
Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures

- humans can't help but group items into classes (categories)
- categories drive how we perceive and make decisions
- ...and we want to **develop knowledge technologies** to help us with this



Compare and Contrast

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and
Review

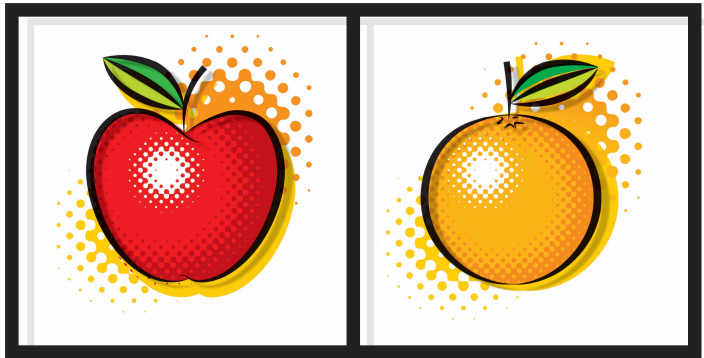
Comparing things

Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures



Compare and Contrast

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and
Review

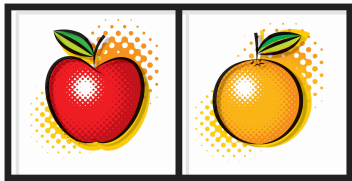
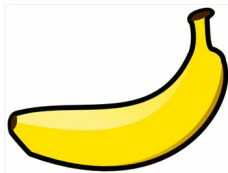
Comparing things

Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures



Venn Diagram

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

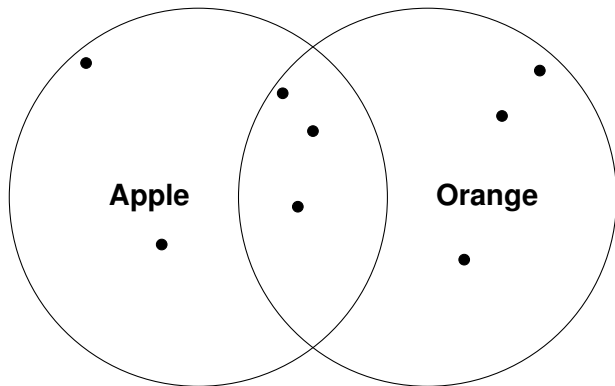
Introduction and
Review

Comparing things
Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures



Venn Diagram

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

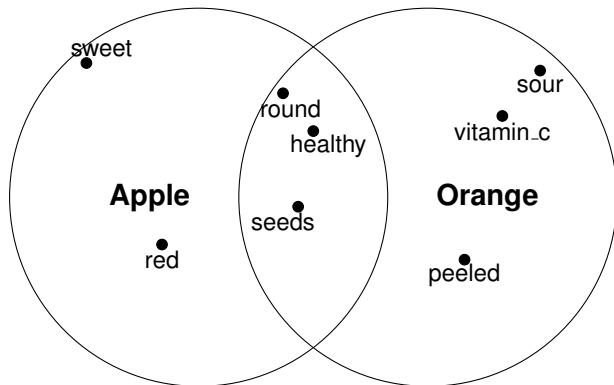
Introduction and
Review

Comparing things
Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures



Venn Diagram

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and
Review

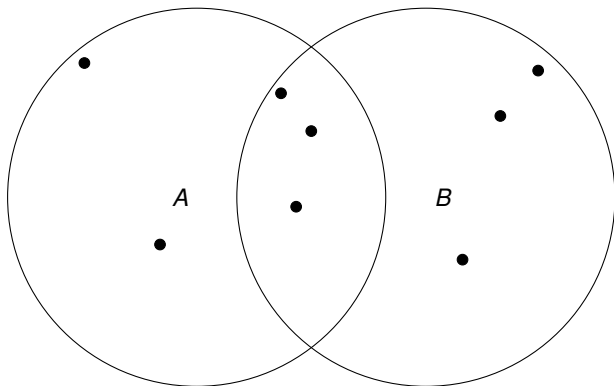
Comparing things

Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures



Venn Diagram

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

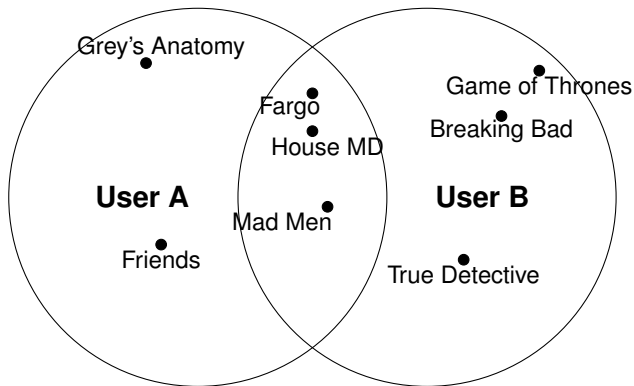
Introduction and
Review

Comparing things
Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures



How similar is User A to User B?

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and
Review

Comparing things
Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures

Many similarity assessments can be framed as set intersection.

- Amazon: Book purchases
- Netflix: Movies that you have watched

Refinements

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

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and
Review

Comparing things
Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures

How should we compare documents to assess their similarity?

- String-level similarity (e.g., edit distance)
- Sets of common substrings (sentences, phrases, words, n-grams)
- “bag of words”
- Meaning (whatever that is?)

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and
Review

Comparing things
Sets of descriptors
Documents

Features, Vectors

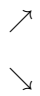
Documents,
revisited

Distance
Measures

How should we compare documents to assess their similarity?

- String-level similarity (e.g., edit distance)
- Sets of common substrings (sentences, phrases, words, n-grams)
- “bag of words”
- Meaning (whatever that is?)

How similar are these sentences?

- a. Mary is quicker than John.
- b. Mary is slower than John.
- c. John is quicker than Mary.
- 

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and
Review

Comparing things

Sets of descriptors

Documents

Features, Vectors

Documents,
revisited

Distance
Measures

How similar are these sentences?

a. Mary is quicker than John.



b. Mary is slower than John.



c. John is quicker than Mary.

Sentence	"Mary"	"John"	"quicker"	"slower"
a.	1	1	1	0
b.	1	1	0	1
c.	1	1	1	0

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and
Review

Comparing things

Sets of descriptors

Documents

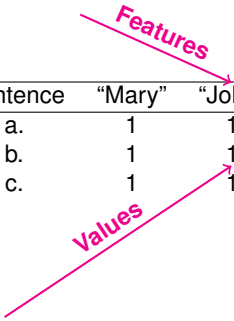
Features, Vectors

Documents,
revisited

Distance
Measures

How similar are these sentences?

- a. Mary is quicker than John.
- b. Mary is slower than John.
- c. John is quicker than Mary.



Sentence	"Mary"	"John"	"quicker"	"slower"
a.	1	1	1	0
b.	1	1	0	1
c.	1	1	1	0

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and
Review

Comparing things
Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures

A **feature** or *attribute* is any distinct aspect, quality, or characteristic of that object

Features may be:

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and
Review

Comparing things
Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures

A **feature** or *attribute* is any distinct aspect, quality, or characteristic of that object

Features may be:

- symbolic/categorical/discrete (e.g. colour, gender)

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and
Review

Comparing things
Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures

A **feature** or *attribute* is any distinct aspect, quality, or characteristic of that object

Features may be:

- symbolic/categorical/discrete (e.g. colour, gender)
- ordinal (e.g. cool < mild < hot [temperature])

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and
Review

Comparing things
Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures

A **feature** or *attribute* is any distinct aspect, quality, or characteristic of that object

Features may be:

- symbolic/categorical/discrete (e.g. colour, gender)
- ordinal (e.g. cool < mild < hot [temperature])
- numeric/continuous (e.g., height, age)

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and
Review

Comparing things
Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures

A **feature** or *attribute* is any distinct aspect, quality, or characteristic of that object

Features may be:

- symbolic/categorical/discrete (e.g. colour, gender)
- ordinal (e.g. cool < mild < hot [temperature])
- numeric/continuous (e.g., height, age)
- binary (e.g., has_feathers, is_round)

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and
Review

Comparing things
Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures

A **feature** or *attribute* is any distinct aspect, quality, or characteristic of that object

Features may be:

- symbolic/categorical/discrete (e.g. colour, gender)
- ordinal (e.g. cool < mild < hot [temperature])
- numeric/continuous (e.g., height, age)
- binary (e.g., has_feathers, is_round)

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and
Review

Comparing things
Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures

A **feature** or *attribute* is any distinct aspect, quality, or characteristic of that object

Features may be:

- symbolic/categorical/discrete (e.g. colour, gender)
- ordinal (e.g. cool < mild < hot [temperature])
- numeric/continuous (e.g., height, age)
- binary (e.g., has_feathters, is_round)

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

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and
Review

Comparing things
Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures

A **feature** or *attribute* is any distinct aspect, quality, or characteristic of that object

Features may be:

- symbolic/categorical/discrete (e.g. colour, gender)
- ordinal (e.g. cool < mild < hot [temperature])
- numeric/continuous (e.g., height, age)
- binary (e.g., has_feathters, is_round)

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

A vector **locates** an object (document, person, ...) as a point in n -dimensional **space**. The **angle** of the vector in that space is determined by the relative **weight** of each term.

Feature vectors and vector space

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and
Review

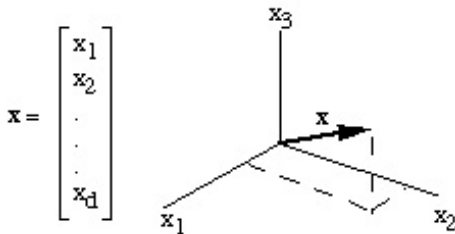
Comparing things

Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures



Credit as a function of age and income

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and
Review

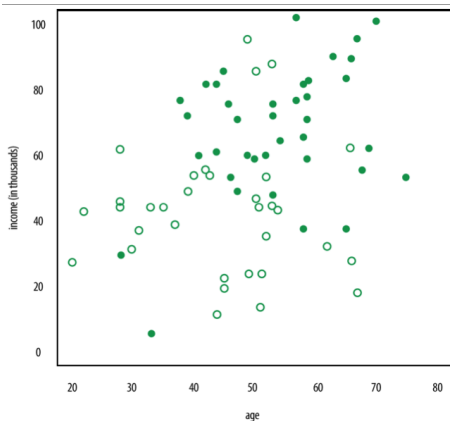
Comparing things
Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures

age	income	credit
33	8	low
58	42	low
49	79	low
49	17	low
58	26	high
44	71	high
...		



Activity – features for spam prediction

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and
Review

Comparing things
Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures

In pairs, think of (at least) one **conituous**, one **categorical** and one **binary** feature predictive for spam vs non-spam e-mails.

title	re [8] : dear friend
body	ORDER CONFIRMATION!!! your order should be shipped by january, via fedex. your federal express tracking number is 45954036. thank you for registering. your userid is: 56075519 learn to make a fortune with ebay! complete turnkey system software - videos - tutorials clk here for information cililings.

title	vacation payback
body	I just received a call from payroll - they miscalculated your vacation payback. Rather than the \$1,113.66, you actually only owe \$957.04 (they forgot to take away your or state taxes in the first calculation). Can you please write another check and we will void the first one?

Sorry for the inconvenience. Please let me know if you have any questions.
Amy

Examples from Enron4 spam data: http://nlp.cs.aueb.gr/software_and_datasets/Enron-Spam/index.html

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and
Review

Comparing things
Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures

In pairs, think of (at least) one **continuous**, one **categorical** and one **binary** feature predictive for spam vs non-spam e-mails.

Continuous features

- average word length (body)
- average word length (title)
- time of sending (continuous)

Categorical features

- title length
- message length
- topic
- time of sending (binned)

Binary features

- word occurrence
- non-empty title
- contains capitalized words
- contains URL
- contains repeated punctuation

Vector space model for document retrieval

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and
Review

Comparing things
Sets of descriptors
Documents

Features, Vectors

**Documents,
revisited**

Distance
Measures

- proposed in 1962, one of the earliest document **retrieval** models

Vector space model for document retrieval

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and
Review

Comparing things
Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures

- proposed in 1962, one of the earliest document **retrieval** models
- n distinct indexed terms (aka. terms of interest) in the collection
- term importance weight $w_{d,t}$ for each document d
e.g., count, binary, importance indicator (tf-idf, next)

Vector space model for document retrieval

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and
Review

Comparing things
Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures

- proposed in 1962, one of the earliest document **retrieval** models
- n distinct indexed terms (aka. terms of interest) in the collection
- term importance weight $w_{d,t}$ for each document d
e.g., count, binary, importance indicator (tf-idf, next)
- documents d_i are represented as a vector

$$d_1 = \langle w_{d_1,1}, w_{d_1,2}, \dots, w_{d_1,t}, \dots, w_{d_1,n} \rangle$$

$$d_2 = \langle w_{d_2,1}, w_{d_2,2}, \dots, w_{d_2,t}, \dots, w_{d_2,n} \rangle$$

Vector space model for document retrieval

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and
Review

Comparing things
Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures

- proposed in 1962, one of the earliest document **retrieval** models
- n distinct indexed terms (aka. terms of interest) in the collection
- term importance weight $w_{d,t}$ for each document d
e.g., count, binary, importance indicator (tf-idf, next)
- documents d_i are represented as a vector

e.g., d_1 could
be a query and
 d_2 a website

$$d_1 = \langle w_{d_1,1}, w_{d_1,2}, \dots, w_{d_1,t}, \dots, w_{d_1,n} \rangle$$

$$d_2 = \langle w_{d_2,1}, w_{d_2,2}, \dots, w_{d_2,t}, \dots, w_{d_2,n} \rangle$$

Vector space model for document retrieval

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and
Review

Comparing things
Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures

- proposed in 1962, one of the earliest document **retrieval** models
- n distinct indexed terms (aka. terms of interest) in the collection
- term importance weight $w_{d,t}$ for each document d
e.g., count, binary, importance indicator (tf-idf, next)

- documents d_i are represented as a vector

e.g., d_1 could
be a query and
 d_2 a website

$$d_1 = \langle w_{d_1,1}, w_{d_1,2}, \dots, w_{d_1,t}, \dots, w_{d_1,n} \rangle$$

$$d_2 = \langle w_{d_2,1}, w_{d_2,2}, \dots, w_{d_2,t}, \dots, w_{d_2,n} \rangle$$

- if the weights for d_1 and d_2 are similar, then they likely to be **similar in topic**
- the content of individual document d is limited and focussed, so that most $w_{d,t}$ values will be zero (**sparsity**).

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and
Review

Comparing things
Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures

Document 1	Document 2	Document 3
Yesterday, Mary went to the store and bought a bunch of apples.	Yesterday, Peter talked to the agents and filed a bunch of complaints.	Agents typically ignore and delete the complaints.

- which document is most relevant to **Document 2**?
- which words are **informative**?
- what distributional properties should these words have?

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and
Review

Comparing things
Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures

Document 1	Document 2	Document 3
Yesterday , Mary went to the store and bought a bunch of apples.	Yesterday , Peter talked to the agents and filed a bunch of complaints.	Agents typically ignore and delete the complaints.

- which document is most relevant to **Document 2**?
- which words are **informative**?
- what distributional properties should these words have?

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and
Review

Comparing things
Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures

Document 1	Document 2	Document 3
Yesterday, Mary went to the store and bought a bunch of apples.	Yesterday, Peter talked to the agents and filed a bunch of complaints .	Agents typically ignore and delete the complaints .

- which document is most relevant to **Document 2**?
- which words are **informative**?
- what distributional properties should these words have?

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and
Review

Comparing things
Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures

We will formalize **two basic intuitions**: term weights are:

- 1 terms that occur frequently in a given document have high utility:

$$w_{d,t} \propto f_{d,t}$$

$f_{d,t}$: frequency of
term t in doc d

Term Frequency-Inverse Document Frequency

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and
Review

Comparing things
Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures

We will formalize **two basic intuitions**: term weights are:

- 1 terms that occur frequently in a given document have high utility:

$$w_{d,t} \propto f_{d,t}$$

$f_{d,t}$: frequency of
term t in doc d

- 2 terms that occur in a wide variety of documents have low utility:

$$w_t \propto \frac{1}{f_t}$$

f_t : # documents
containing t

Term Frequency-Inverse Document Frequency

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and
Review

Comparing things
Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures

We will formalize **two basic intuitions**: term weights are:

- 1 terms that occur frequently in a given document have high utility:

$$w_{d,t} \propto f_{d,t} \quad \leftarrow \begin{array}{l} f_{d,t}: \text{frequency of} \\ \text{term } t \text{ in doc } d \end{array}$$

- 2 terms that occur in a wide variety of documents have low utility:

$$w_t \propto \frac{1}{f_t} \quad \leftarrow \begin{array}{l} f_t: \# \text{ documents} \\ \text{containing } t \end{array}$$

Models which weigh up these two are referred to as **TF-IDF** (term frequency–inverse document frequency) models

The “classic” TF-IDF formulation is:

$$w_{d,t} = f_{d,t} \times \log \frac{N}{f_t} \quad \leftarrow \begin{array}{l} N: \# \text{ documents} \\ \text{in collection} \end{array}$$

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and
Review

Comparing things

Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures

- ✓ We have discussed similarity at an intuitive level.
- ✓ We have formalized feature (term) weighting.
- ✗ How do we measure similarity quantitatively?

Jaccard Similarity

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and
Review

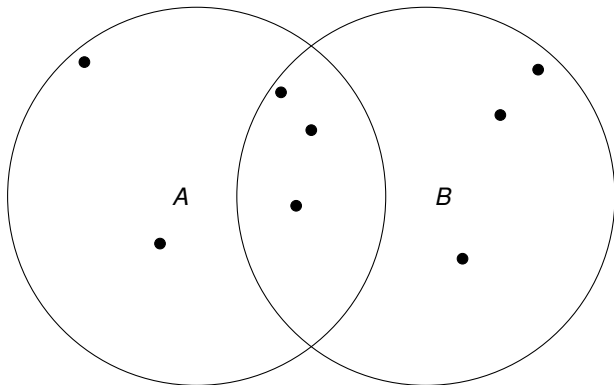
Comparing things

Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures



$$\text{sim}(A, B) = \frac{|A \cap B|}{|A \cup B|} = ?$$

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

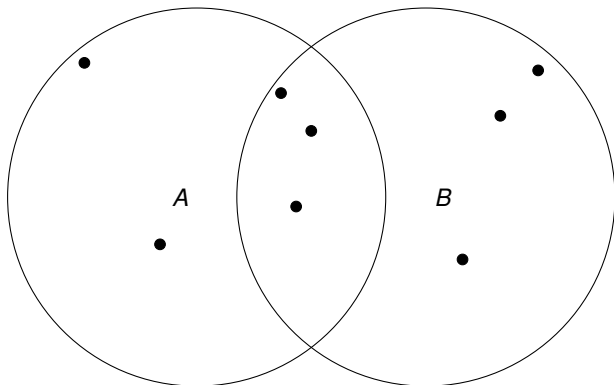
Introduction and
Review

Comparing things
Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures



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

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and
Review

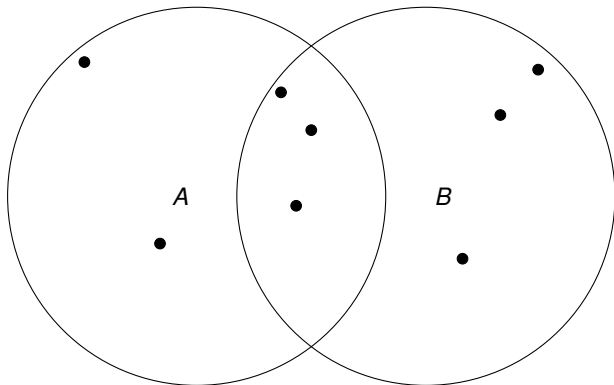
Comparing things

Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures



$$\text{sim}(A, B) = \frac{2|A \cap B|}{|A| + |B|} = ?$$

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

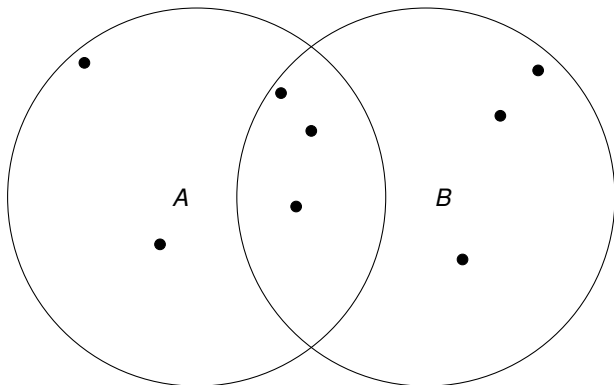
Introduction and
Review

Comparing things
Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures



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

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and
Review

Comparing things
Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

**Distance
Measures**

What is the relationship between similarity and distance?

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and
Review

Comparing things
Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures

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

- 1 No negative distances.

$$d(x, y) \geq 0$$

- 2 Distances are positive, except for the distance from a point to itself.

$$d(x, y) = 0 \text{ if and only if } x = y$$

- 3 Distance is symmetric.

$$d(x, y) = d(y, x)$$

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

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

Manhattan Distance (aka L1 Distance)

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and
Review

Comparing things

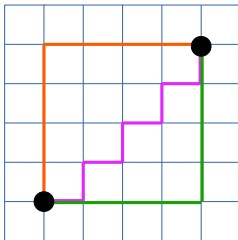
Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures

- Given two items A and B , as their feature vectors \vec{a} and \vec{b}
- Absolute differences of their cartesian coordinates
- Travel coordinate by coordinate. No short cuts!



In n -dimensional space:

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

Euclidean Distance (aka L2 Distance)

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and
Review

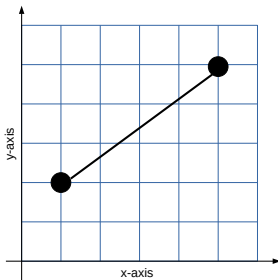
Comparing things
Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures

- Given two items A and B , as their feature vectors \vec{a} and \vec{b}
- Straight line between the two points



In n -dimensional space:

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

Cosine Distance

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and
Review

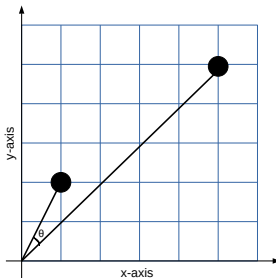
Comparing things
Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures

- Given two items A and B , as their feature vectors \vec{a} and \vec{b}
- Similarity as the cosine of the angle θ between the two vectors



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

“Long” documents & Euclidean distance

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and
Review

Comparing things

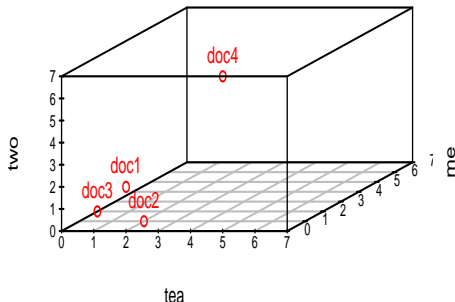
Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures

	tea	me	two
doc1	2	0	2
doc2	2	1	0
doc3	0	2	0
doc4	5	0	7



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

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and
Review

Comparing things
Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures

Relative entropy:

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

or alternatively *skew divergence*:

$$s_\alpha(x, y) = D(x \parallel \alpha y + (1 - \alpha)x)$$

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!

Lecture 3: Similarity

COMP90049
Knowledge
Technologies

Introduction and
Review

Comparing things
Sets of descriptors
Documents

Features, Vectors

Documents,
revisited

Distance
Measures

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

Reading

- On distance measures:

Chapter 3, especially Section 3.5

Mining of Massive Datasets

<http://infolab.stanford.edu/~ullman/mmds.html>

- On 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>