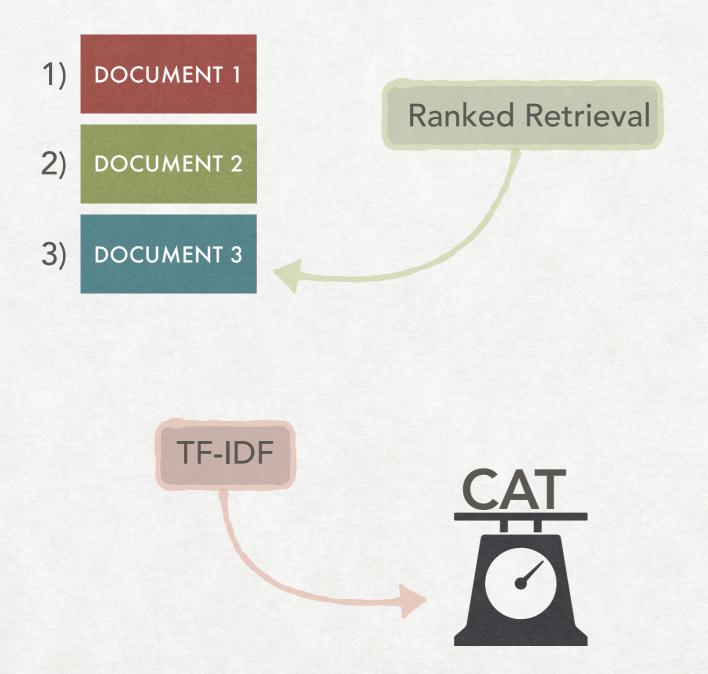
INFORMATION RETRIEVAL

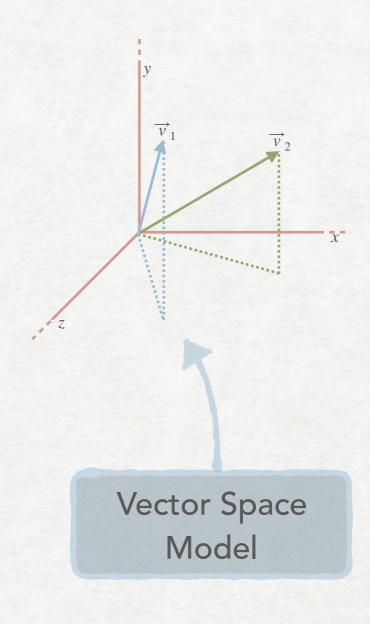
Luca Manzoni Imanzoni@units.it

Lecture 5

LECTURE OUTLINE

*MADE WITH ALIEN TECHNOLOGY





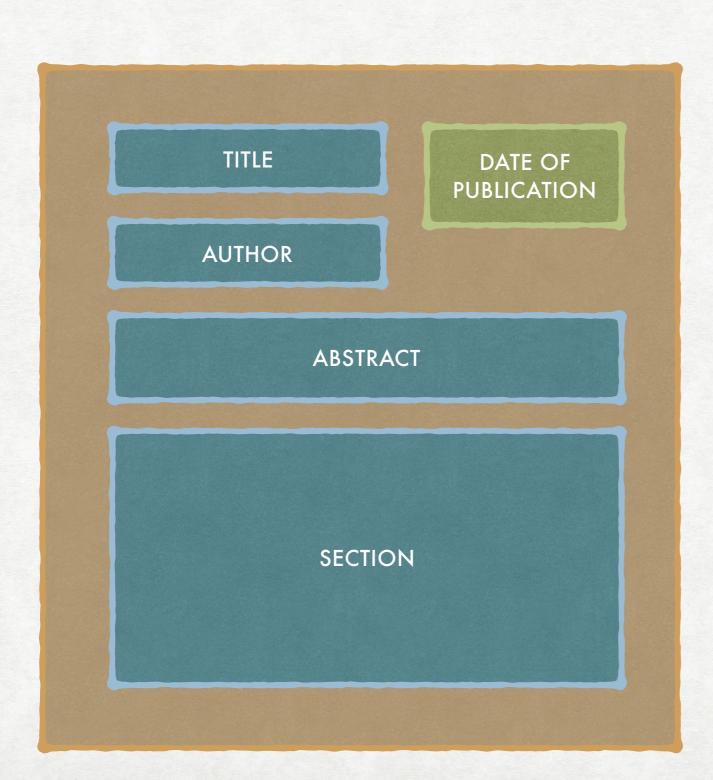
RANKED RETRIEVAL

MOTIVATIONS

- Until now we have returned all documents matching a Boolean query as a set.
- If many documents are returned then it might be important to rank them according to how relevant they are.
- A first way of ranking them is to "split" a document according to some structure and then weight different zones in different ways.
- We will then see how we can extend the idea of adding weights also to the terms of a document.

DOCUMENT STRUCTURE

METADATA, FIELDS, AND ZONES



- A text may have associated metadata.
- Some of them can be fields, with a set of values that can be finite, like publication dates.
- Others might be zones, arbitrary areas of free-form text (e.g., abstract, section, etc.).

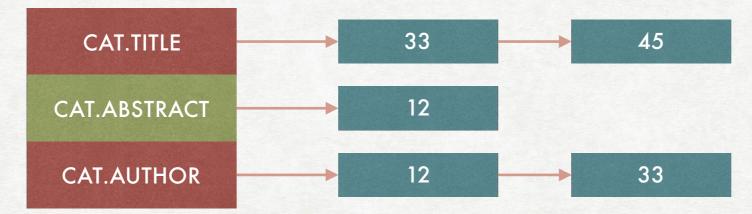
PARAMETRIC INDEXES

SEARCHING INSIDE FIELDS

- To allow for searching inside the fields we might want to build additional indexes, called parametric indexes.
- A parametric index can be thought as a standard index that only has information about a field (e.g., all the dates).
- If a query asks for "cat" in the title and "dog" inside the document we will retrieve the posting lists for dog from the "standard" index e "cat" from the parametric index for the title.
- The operations of union and intersections works as usual.

ZONE INDEXES POSSIBLE APPROACHES

Separate inverted index for each zone



Single inverted index in which the zones are part of the postings



WEIGHTED ZONE SCORING

AN ADDITIONAL USE FOR ZONES

- We now have a way of searching inside different parts of a document...
- ...but different parts might carry different importance: e.g., a title vs inside the main text.
- We can rank retrieved documents according to where the term is found inside the document.
- We can do this via weighted zone scoring (also called ranked Boolean retrieval).

SCORING FUNCTION

DEFINITION

- Consider a pair (q, d) of a query q and a document d.
- A scoring function associates a value in [0,1] to each pair (q,d).
- Higher scores are better.
- Suppose that a document has ℓ zones.
- Each zone has a weight $g_i \in [0,1]$ for $1 \le i \le \ell$.
- The weights sums to one:

$$\sum_{i=1}^{n} g_i = 1$$

SCORING FUNCTION

PART II

- Given a query q let s_i be defined as $s_i = \begin{cases} 1 & \text{if } q \text{ matches in zone } i \\ 0 & \text{otherwise} \end{cases}$
- Actually, s_i can also be defined to be any function that maps "how much" a query matches in the i-th zone.
- The weighted zone score in then defined as:

$$\sum_{i=1}^{\ell} g_i s_i$$

WEIGHTED ZONE SCORING

A SIMPLE EXAMPLE

Query: C	CAT
----------	-----

TITLE: LIFE OF A CAT AUTHOR: JAMES CAT ONCE THERE WAS A CAT...

TITLE: DOGS AND OTHER PETS
AUTHOR: ANONYMOUS
DOGS AND CATS ARE THE...

TITLE: ORCHARDS MANAGEMENT AUTHOR: JAMES CAT THE MANAGEMENT OF ORCHARDS...

	Body: 0.3	Author: 0.2	Title: 0.5	
1	0.3	0.2	0.5	
0.3	0.3	0	0	
0.2	0	0.2	0	

LEARNING WEIGHTS

OR SETTING THEM MANUALLY

- The new problem is now to find how to set the weights for the different scores.
- One possibility is to ask a domain expert.
- Another possibility is to have users label documents relevant or not with respect to a query...
- ...and trying to learn the weights using the training data.
- In addition to the binary classification (relevant or not) more nuanced classifications might be used.

THE TRAINING SET

Example	DocID	Query	In the title	In the body	Judgment
e1	43	LISP	1	1	Relevant
e2	43	BASIC	1	0	Relevant
e3	76	LISP	0	1	Non-relevant
e4	76	BASIC	0	1	Relevant
e5	87	SMALLTALK	1	1	Relevant
e6	87	APL	1	0	Non-relevant

COMPUTING THE ERROR HOW TO DECIDE IF OUR WEIGHTS WORKS

With only two zones, site score is computed as:

$$score(d, q) = g \cdot s_{title} + (1 - g) \cdot s_{body}$$

Since we know the queries and the real relevance of the documents in the training set we can compute the output that a weight g would give:

$$score(43, LISP) = g \cdot 1 + (1 - g) \cdot 1$$

$$score(43, BASIC) = g \cdot 1 + (1 - g) \cdot 0$$

$$score(76, LISP) = g \cdot 0 + (1 - g) \cdot 1$$

.

COMPUTING THE ERROR HOW TO DECIDE IF OUR WEIGHTS WORKS

If we decide that relevant is 1 and non-relevant is 0 we can compare the real score with the computed one and compute an error:

$$Err(g, e1) = (1 - score(43, LISP))^2$$

$$Err(g, e2) = (1 - score(43, BASIC))^2$$

$$Err(g, e3) = (0 - score(76, LISP)^2$$

•

MINIMISING THE ERROR

(AND MAYBE IT CANNOT BE ZERO)

We now want to minimise the sum of the errors:

$$\sum_{i=1}^{n} \operatorname{Err}(g, ei)$$

Notice that it might not be possible to reach an error of zero:

score(43,BASIC) =
$$g \cdot 1 + (1 - g) \cdot 0 = g$$

score(87,APL) = $g \cdot 1 + (1 - g) \cdot 0 = g$

But:

$$Err(g, e2) = (1 - g)^2$$

 $Err(g, e6) = g^2$

TF-IDF WEIGHTING

CHANGING SCORING

REFINING THE SCORING

- For now we have used a weight that is either 0 or 1 depending on wether a query term was present or not.
- We might want to assign different weight depending on the term and the number of times a term is present in the document.
- This works well with free-form text queries:
 - For each term in the query we compute a "match score"
 - The score of a document is the sum of the scores for each term

TERM FREQUENCY

A SIMPLE SCORE

Term frequency: $tf_{t,d}$

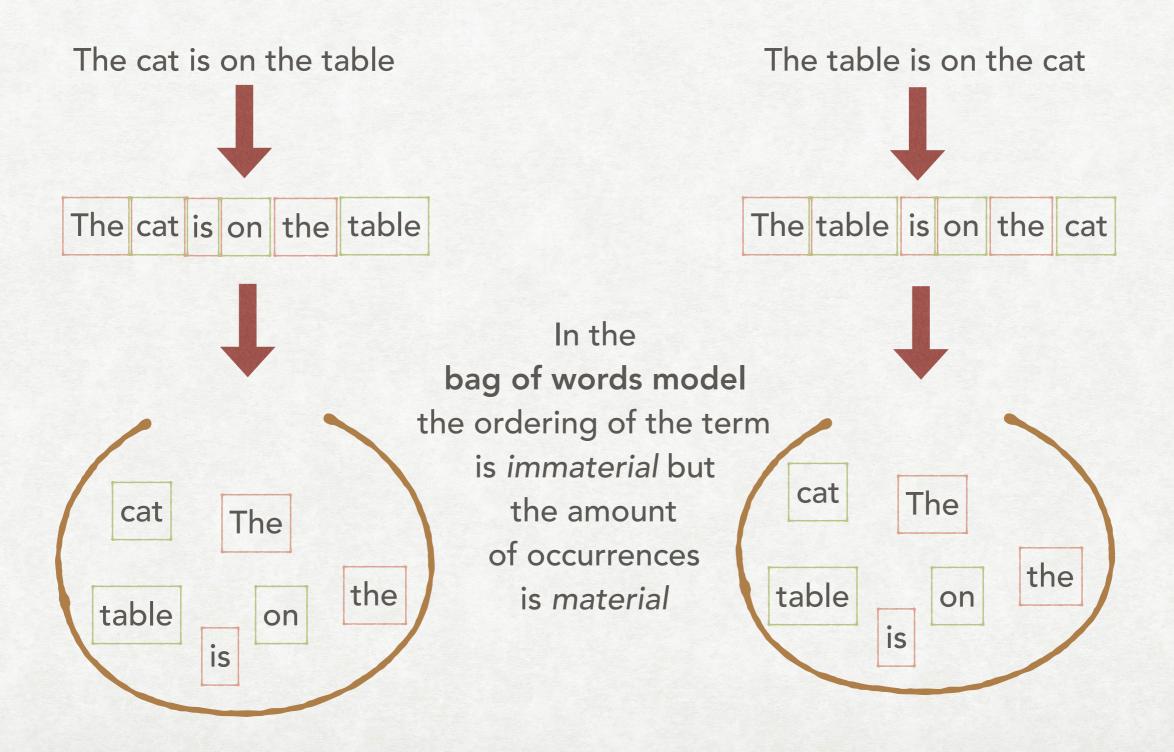
Number of occurrences of the term t inside the document d.

The main motivation is that the more a document is present inside a document the more we consider the document relevant with respect to that term.

But what about the order of the words?

BAG OF WORDS

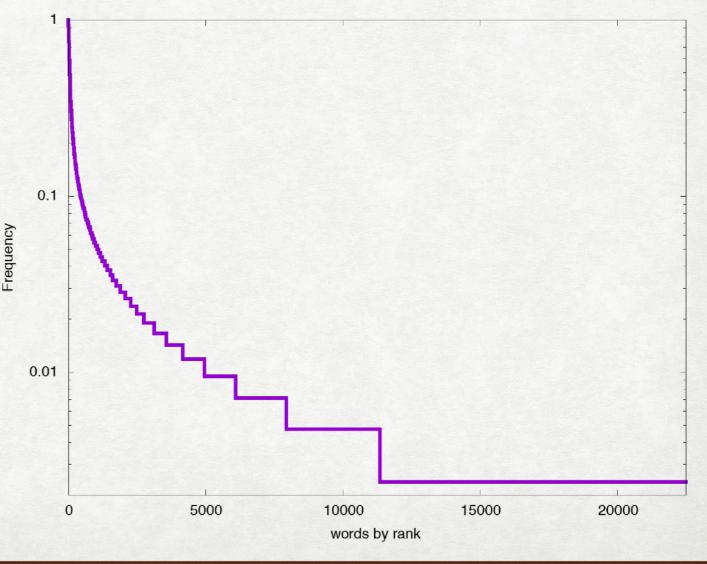
IGNORE THE ORDER!



TERM FREQUENCY

SOME LIMITATIONS

- Does the number of occurrences really represents the importance of a term?
- Which terms are more frequent?
- A small hint:
- Stop words!
- Not all terms carry the same weight in determining the relevancy of a document



COLLECTION AND DOCUMENT FREQUENCIES

RARE WORDS COUNT MORE

- The main characteristic of stop words is that they are present in most documents.
- Therefore, we might want to scale the importance of a word based on some measure of the frequency of the term:
- cf_t is the **collection frequency** of the term t: total number of occurrences of the term t in the collection.
- df_t is the document frequency of the term t: total number of document in which t appears in the collection.

COLLECTION AND DOCUMENT FREQUENCIES

RARE WORDS COUNT MORE

- The document frequency df_t of a term is usually preferred.
- We prefer to use a document-based measure to weight documents.
- cf_t and df_t can behave quite differently. For example:
 - A single document with 1000 instances of a term t_1 in a collection of 1000 documents.
 - Each one of 1000 documents contains a term t_2 exactly once.

INVERSE DOCUMENT FREQUENCY

MODIFYING DOCUMENT FREQUENCY

 df_t is larger when we want the penalties to be larger

We use a modification of it:

Number of documents in the collection

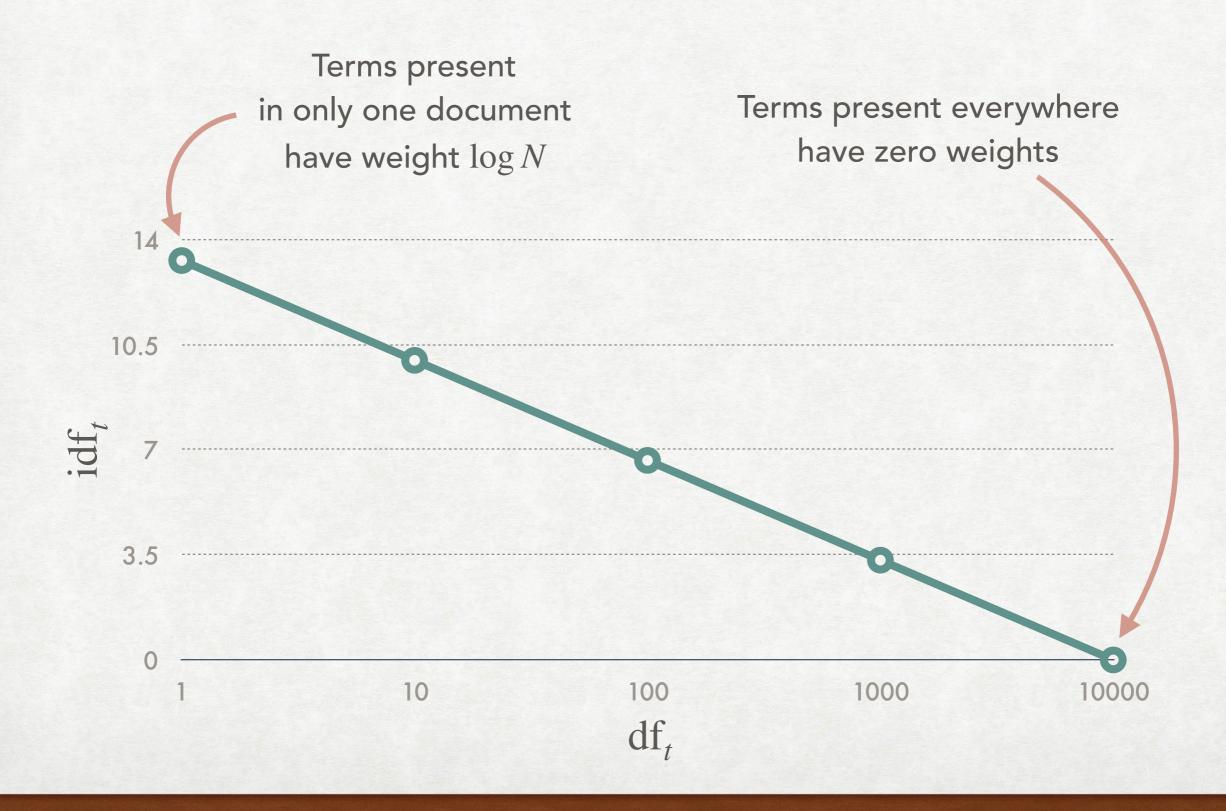
$$idf_t = \log \frac{N}{df_t}$$

Inverse document frequency

Document frequency

INVERSE DOCUMENT FREQUENCY

EFFECTS ON THE WEIGHTS



TF-IDF WEIGHTING

HOW TO COMBINE $tf_{t,d}$ AND idf_t

We now need to combine the two ideas:

$$tf\text{-}idf_{t,d} = tf_{t,d} \times idf_t$$

- When a rare term is present a many times in a document then the value is high
- When a frequent term is present many times or a rare term is present only a few time the value is low
- When a very frequent term is present only a few times then the value is the lowest

SCORING A DOCUMENT TOWARDS THE VECTOR SPACE MODEL

The cat is on the table

We can see a document as a vector with a components for each term in the dictionary and having as elements the tf- $idf_{t,d}$ of the term t in the document

cat	is	on	table	the
tf-idf _{cat,d}	tf-idf _{is,d}	tf-idf _{on,d}	tf-idf _{table,d}	tf - $idf_{the,d}$

 $tf\text{-}idf_{t,d} = 0$ for all terms not in the document

SCORING A DOCUMENT TOWARDS THE VECTOR SPACE MODEL

To score a document for a query q we can simply sum the $\operatorname{tf-idf}_{t,d}$ values for all terms appearing in q:

$$Score(q, d) = \sum_{t \in q} tf\text{-}idf_{t,d}$$

Notice that in this way a document where a term does *not* appear might still have a positive score. The "penalty" will depend on which term is not present

VARIANTS OF TF-IDF

AND WHEN TO USE THEM

- There are some possible alternative in using directly tf-idf.
- One first consideration is that not all instances of a term inside a document carry the same weight.
- There is the idea of "diminishing returns": is a document with 20 occurrences really twice as important as one with 10 occurrences?
- Another observation is that we might be interested in the frequency of a term relative to the other terms in the document.

SUBLINEAR TF SCALING

We can scale the $tf_{t,d}$ value to have the influence of additional terms reduced:

$$wf_{f,d} = \begin{cases} 1 + \log tf_{t,d} & \text{if } tf_{t,d} > 0\\ 0 & \text{otherwise} \end{cases}$$

The new value can be replaced where $tf_{t,d}$ is used:

$$\operatorname{wf-idf}_{f,d} = \operatorname{wf}_{f,d} \times \operatorname{idf}_t$$

TF NORMALIZATION

We can scale the $\mathrm{tf}_{t,d}$ value to be dependant on the maximum term frequency in the document $\mathrm{tf}_{\mathrm{max}}(d)$:

$$\frac{\mathrm{tf}_{t,d}}{\mathrm{tf}_{\mathrm{max}}(d)}$$

Another possibility is to normalise according to the number of terms in the entire document:

$$\frac{\mathrm{tf}_{t,d}}{\sum_{t' \in d} \mathrm{tf}_{t',d}}$$

In both cases there are drawbacks and some smoothing might be applied to limit large swings in the normalised value

THE VECTOR SPACE MODEL

VERY BRIEF RECAP

JUST TO REFRESH SOME BASIC NOTION AND FIX NOTATION

• In \mathbb{R}^n the Euclidean length of a vector $\overrightarrow{v} = (v_1, v_2, ..., v_n)$ is

$$|\overrightarrow{v}| = \sqrt{\sum_{i=1}^{n} v_i^2}$$

- A vector is unit vector if its length is one.
- The inner products of two vectors

$$\overrightarrow{v} = (v_1, v_2, ..., v_n)$$
 and $\overrightarrow{u} = (u_1, u_2, ..., u_n)$ is defined as $\sum_{i=1}^n v_i u_i$

DOCUMENTS AS VECTORS

THE START OF THE VECTOR SPACE REPRESENTATION

$$e_{\mathsf{bart}} = (1,0,0,0,0)$$

$$e_{\text{box}} = (0,1,0,0,0)$$

$$e_{cat} = (0,0,1,0,0)$$

$$e_{dog} = (0,0,0,1,0)$$

$$e_{\text{drone}} = (0,0,0,0,1)$$

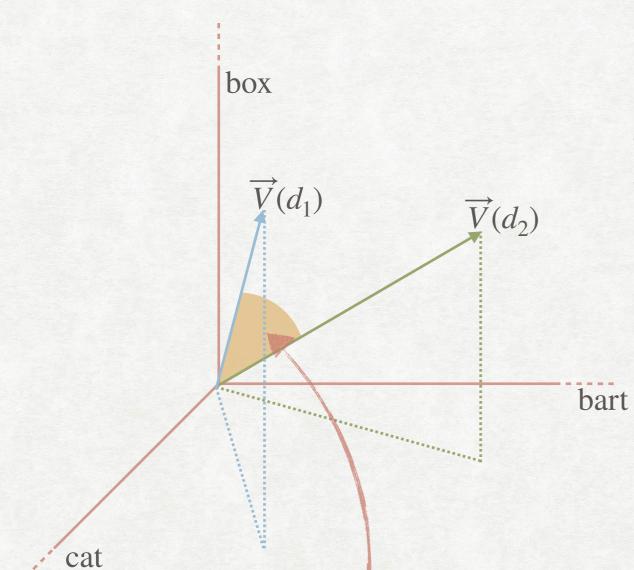
Each term is an element of the canonical base of \mathbb{R}^n with n the number of terms in the dictionary.

A document is a point in this n-dimensional space:

$$\overrightarrow{V}(d) = (0.6, 0.5, 0.1, 0, 0.9)$$
 tf-idf_{cat,d}

COSINE SIMILARITY

HOW TO COMPARE DOCUMENTS



We can compute the similarity of two documents by computing the *cosine similarity* between the two corresponding vectors:

$$\operatorname{sim}(d_1, d_2) = \frac{\overrightarrow{V}(d_1) \cdot \overrightarrow{V}(d_2)}{|\overrightarrow{V}(d_1)||\overrightarrow{V}(d_2)|}$$

Which represents the cosine of the angle formed by the two vectors

The similarity is the cosine of this angle

NORMALISING VECTORS LOOKING AGAIN AT COSINE SIMILARITY

If we look again at cosine similarity we can see that we can replace a vector $\overrightarrow{V}(d)$ with the unit vector $\overrightarrow{v}(d)$:

$$\overrightarrow{v}(d) = \frac{\overrightarrow{V}(d)}{|\overrightarrow{V}(d)|}$$

In fact, since the angle formed by the vectors does not depend on the magnitude of the vectors, we can assume, without loss of generality, each document vector to be a unit vector.

QUERIES AS VECTORS THE MISSING HALF OF THE REPRESENTATION

In addition to documents, also queries can be represented as vectors

Query: CAT Vector: (0,0,1,0,0)

Query: CAT DOG Vector: $(0,0,1/\sqrt{2},1/\sqrt{2},0)$

Each query is a unit vector with the non-zero components corresponding to the query terms

ANSWERING QUERIES

COSINE SIMILARITY (AGAIN)

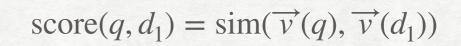
 $\overrightarrow{v}(d_2)$

 $\overrightarrow{v}(q)$

dog

 $\overrightarrow{v}(d_1)$

The answer to the query can be computed using (again) the cosine similarity:



$$score(q, d_2) = sim(\overrightarrow{v}(q), \overrightarrow{v}(d_2))$$

Since all vectors are unit vectors this is equivalent to:

$$score(q, d_1) = \overrightarrow{v}(q) \cdot \overrightarrow{v}(d_1)$$

$$score(q, d_2) = \overrightarrow{v}(q) \cdot \overrightarrow{v}(d_2)$$

VECTOR SPACE MODEL

CONSIDERATIONS

- The fact that we compute a similarity score means that we have a ranking of documents; we can retrieve the K most relevant documents.
- A document might have a non-zero similarity score even if not all terms are present: the matching is not exact like in the Boolean model.
- Even if we have used tf-idf to define the document vectors, any other measure might be used.
- Notice that we cannot exclude (for now) the computation of the cosine similarity for each document in the collection!