

UNIVERSITY of DELAWARE

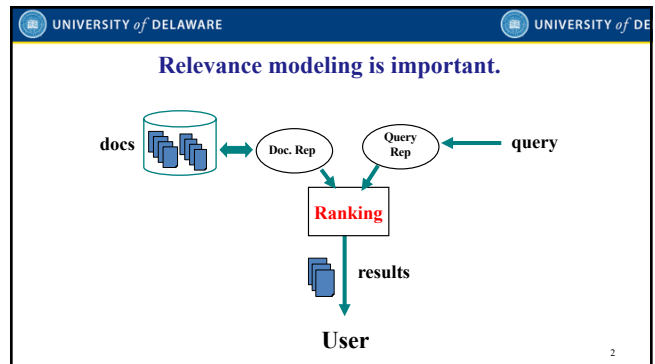
UNIVERSITY of DE

# Vector Space Retrieval Models

(Search and Data Mining)

Hui Fang  
Department of Electrical and Computer Engineering  
University of Delaware

1



UNIVERSITY of DELAWARE

UNIVERSITY of DE

## The Basic Question

Given a query, how do we know if document A is more relevant than B?

### One Possible Answer

**If document A uses more query words than document B**  
(Word usage in document A is more similar to that in query)

3

UNIVERSITY of DELAWARE

UNIVERSITY of DE

## Relevance = Similarity

- Assumptions
  - Query and document are represented similarly
  - A query can be regarded as a “document”
  - $\text{Relevance}(d,q) \propto \text{similarity}(d,q)$
- $R(q) = \{d \in C \mid f(d,q) > \theta\}$ ,  $f(q,d) = \Delta(\text{Rep}(q), \text{Rep}(d))$
- Key issues
  - How to represent query/document?
  - How to define the similarity measure  $\Delta$ ?

4

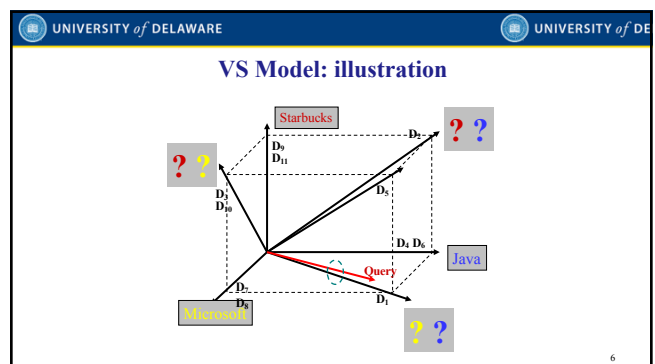
UNIVERSITY of DELAWARE

UNIVERSITY of DE

## Vector Space Model: Main Idea

- Documents and queries are represented as vectors in a high-dimensional space.
  - Each dimension corresponds to a term.
  - All the terms are assumed to be independent.
  - The value of an element in a vector corresponding to the weight of the term.
    - E.g.,  $d = (x_1, \dots, x_N)$ ,  $x_i$  is “importance” of term  $i$
- Relevance can be approximated by the similarity between the query vector and document vector in the vector space

5



UNIVERSITY of DELAWARE

UNIVERSITY of DE

### VS Models: Representations

- Query/Document: term vectors  
 $D = (w_{1,D}, \dots, w_{i,D})$      $Q = (w_{1,Q}, \dots, w_{i,Q})$
- Relevance between doc and query  
 → distance between two vectors

$$\text{similarity}(D, Q) = \sum w_{i,D} \times w_{i,Q}$$

7

UNIVERSITY of DELAWARE

UNIVERSITY of DE

### How to Measure Similarity?

$\vec{D}_i = (w_{1,i}, \dots, w_{N,i})$   
 $\vec{Q} = (w_{1,q}, \dots, w_{N,q})$

Dot product similarity:  $\text{sim}(\vec{Q}, \vec{D}_i) = \sum_{j=1}^N w_{qj} * w_{ij}$

Cosine: 
$$\text{sim}(\vec{Q}, \vec{D}_i) = \frac{\sum_{j=1}^N w_{qj} * w_{ij}}{\sqrt{\sum_{j=1}^N (w_{qj})^2} * \sqrt{\sum_{j=1}^N (w_{ij})^2}}$$

8

UNIVERSITY of DELAWARE

UNIVERSITY of DE

$$\text{Similarity}(D, Q) = \sum_{i \in D} w_{i,D} * w_{i,Q}$$

Q= "a dog"  
 D1="a dog walk dog animal dog cute"  
 D2="a cat walk cat cat dog"  
 D3="a book a book"

Assume  $w_{i,D} = 1$  when  $i \in D$   
 $w_{i,D} = 0$  when  $i \notin D$

	dog	a	walk	cute	animal	book	cat
Q=(							
D1=(							
D2=(							
D3=(							

Challenge: How to assign the weights to each term?

9

UNIVERSITY of DELAWARE

UNIVERSITY of DE

### Term Frequency (TF)

The relevance score of a document is related to the occurrence of a query term.

**dog**

The cat is a small carnivorous mammal. The cat is known to hunt over 1,000 species for food. The cat can be trained to obey simple commands. Cats, like dogs, are digitigrades. Cats are capable of walking very precisely.

Dogs were domesticated from wolves. Dogs were first domesticated in East Asia. Dogs, like humans, are highly social animals. This similarity has earned dogs a unique position in the realm of interspecies relationships.

10

UNIVERSITY of DELAWARE

UNIVERSITY of DE

$$\text{Similarity}(D, Q) = \sum_{i \in D} w_{i,D} * w_{i,Q}$$

Q= "a dog"  
 D1="a dog walk dog animal dog cute"  
 D2="a cat walk cat cat dog"  
 D3="a book a book"

$w_{i,D} = \text{TF}(t_i, D) = c(t_i, D)$

	dog	a	walk	cute	animal	book	cat
Q=(	1,	1,	0,	0,	0,	0,	0
D1=(	3,	1,	1,	1,	1,	0,	0
D2=(	1,	1,	1,	0,	0,	0,	3
D3=(	0,	2,	0,	0,	0,	2,	0

11

UNIVERSITY of DELAWARE

UNIVERSITY of DE

### Inverse Document Frequency (IDF)

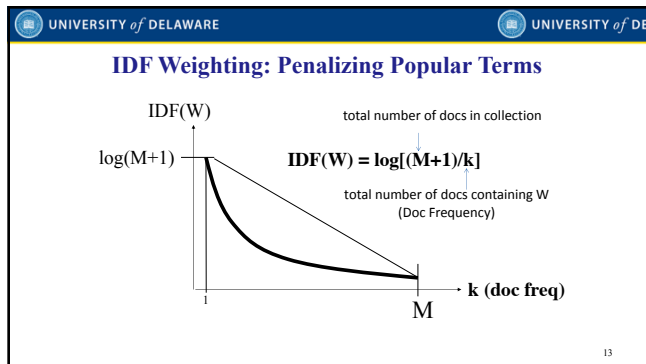
The importance of a term is related to document frequency of the term.

**the dog**

The earliest trees were tree ferns and horsetails, which grow in vast forests in the Carboniferous Period. Tree ferns still survive, but the only surviving horsetails are not of tree form. The listing below gives examples of many trees.

Dogs were domesticated from wolves. Dogs were first domesticated in East Asia. Dogs, like humans, are highly social animals. This similarity has earned dogs a unique position in the realm of interspecies relationships.

12



UNIVERSITY of DELAWARE

Similarity( $D, Q$ ) =  $\sum_{d \in Q} w_{d,D} \cdot w_{d,Q}$

$Q = \text{"a dog"}$   
 $D1 = \text{"a dog walk dog animal dog cute"}$   
 $D2 = \text{"a cat walk cat cat dog"}$   
 $D3 = \text{"a book a book"}$

$w_{d,D} = TF(t, D) \cdot IDF(t, D) = TF(t, D) \cdot \ln \frac{N+1}{df}$

$N=3, df(a)=3, df(dog)=df(walk)=2, df(cat)=df(book)=df(cute)=df(animal)=1$

	dog	a	walk	cute	animal	book	cat
$Q =$	$\log 2$	$\log 4/3$	0	0	0	0	0
$D1 =$	$3 \cdot \log 2$	$\log(4/3)$	$\log 2$	$\log 4$	$\log 4$	0	0
$D2 =$	$\log 2$	$\log(4/3)$	$\log 2$	0	0	0	$3 \cdot \log 4$
$D3 =$	0	$2 \cdot \log(4/3)$	0	0	0	$2 \cdot \log 4$	0

14

UNIVERSITY of DELAWARE

### VS Example: Raw TF & Dot Product

doc1: information, retrieval, search, engine, information  
 doc2: travel, information, map, travel  
 doc3: government, president, congress

query="information retrieval"

$Sim(q, doc1) = 4.8 \cdot 2.4 + 4.5 \cdot 4.5$   
 $Sim(q, doc2) = 2.4 \cdot 2.4$   
 $Sim(q, doc3) = 0$

	info	retrieval	travel	map	search	engine	govern	president	congress
$IDF(faked)$	2.4	4.5	2.8	3.3	2.1	5.4	2.2	3.2	4.3
doc1	2(4.8)	1(4.5)			1(2.1)				
doc2	1(2.4)		2(5.6)	1(3.3)					
doc3							1(2.2)	1(3.2)	1(4.3)
query	1(2.4)	1(4.5)							

15

UNIVERSITY of DELAWARE

### Pivoted Normalization Formula

$$S(Q, D) = \sum_{t \in Q \cap D} \frac{1 + \ln(1 + \ln(c(t, D)))}{(1-b) + b \frac{|D|}{avdl}} \cdot c(t, Q) \cdot \ln \frac{N+1}{df(t)}$$

IDF weighting

16

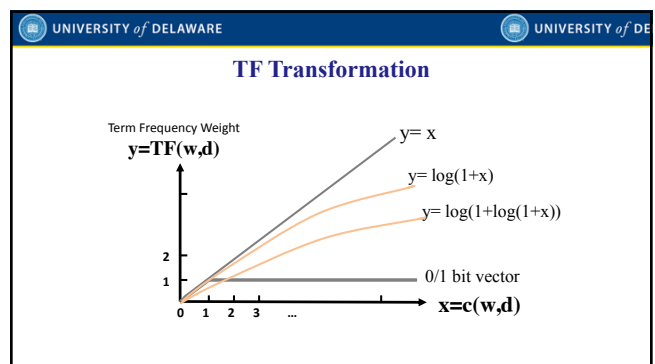
UNIVERSITY of DELAWARE

### Pivoted Normalization Formula

$$S(Q, D) = \sum_{t \in Q \cap D} \frac{1 + \ln(1 + \ln(c(t, D)))}{(1-b) + b \frac{|D|}{avdl}} \cdot c(t, Q) \cdot \ln \frac{N+1}{df(t)}$$

TF weighting

17



UNIVERSITY of DELAWARE

UNIVERSITY of DE

### Pivoted Normalization Formula

$$S(Q, D) = \sum_{t \in Q, D} \frac{1 + \ln(1 + \ln(c(t, D)))}{(1 - b) + b \frac{|D|}{avdl}} \cdot c(t, Q) \cdot \ln \frac{N + 1}{df(t)}$$

Length normalization

19

UNIVERSITY of DELAWARE

UNIVERSITY of DE

### Document Length Normalization (LN)

The relevance score of a document is related to the document length.

dog

The **cat** is a small carnivorous mammal. It is valued by humans for its companionship and its ability to hunt vermin, and has been associated with humans for at least 9,500 years. A skilled predator, the **cat** is known to hunt over 1,000 species for food. The **cat** is intelligent and can be trained to obey simple command. **Cats**, like **dogs**, walk directly on their toes. **Cats** are capable of walking very precisely, because **cats** ...

**Dogs**, like humans, are highly social animals and this similarity in their overall behavioral pattern accounts for their trainability.

20

UNIVERSITY of DELAWARE

UNIVERSITY of DE

### Document Length Normalization

- Penalize a long document with a document length normalizer.
  - Long documents have a better chance to match any query.
  - Need to avoid over-penalization.
- A document is long because
  - It uses more words → more penalization
  - It has more content → less penalization

21

UNIVERSITY of DELAWARE

UNIVERSITY of DE

### Pivoted Length Normalization

- Use average document length as pivot

$b \in [0, 1]$

$normalizer = 1 - b + b \frac{|d|}{avdl}$

1.0

0 1 2 ... avdl ... |d|

Shorter than avdl Longer than avdl

$b \gg 0$   
 $b > 0$   
 $b = 0$

Reward  
 Penalization

22