

Retrieving documents of interest

Document retrieval

- Currently reading article you like



Document retrieval

- Currently reading article you like
- Goal: Want to find similar article



Document retrieval



Challenges

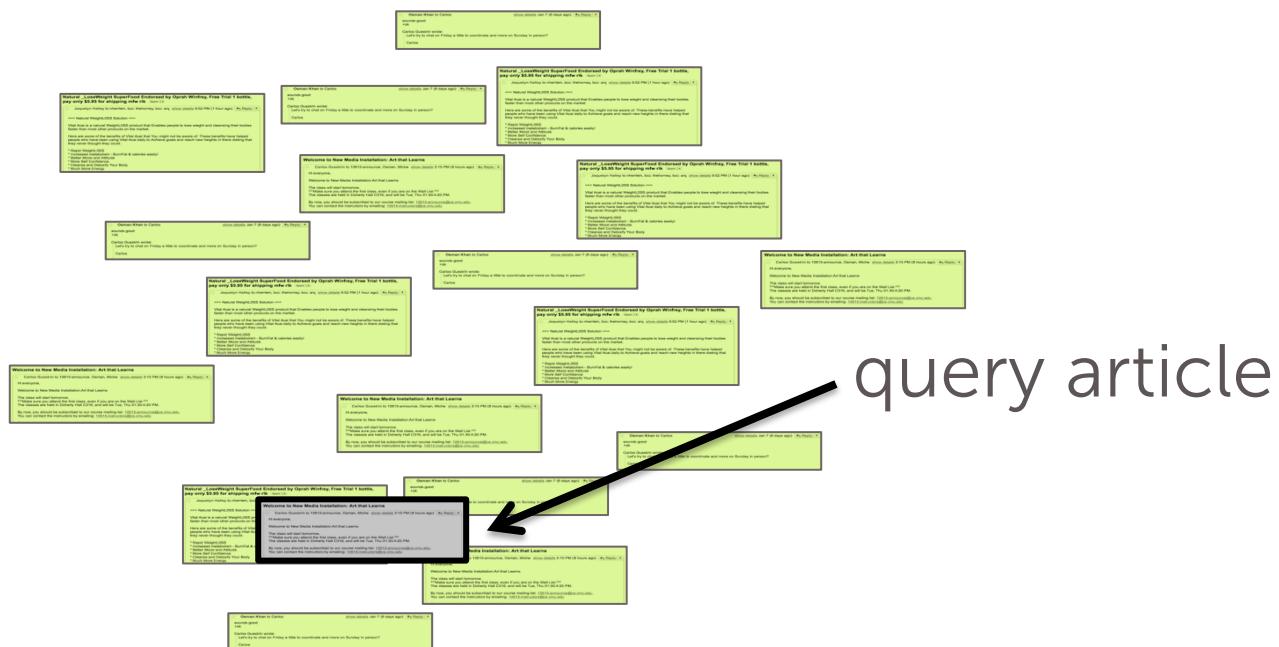
- How do we measure similarity?
- How do we search over articles?



Retrieval as k-nearest neighbor search

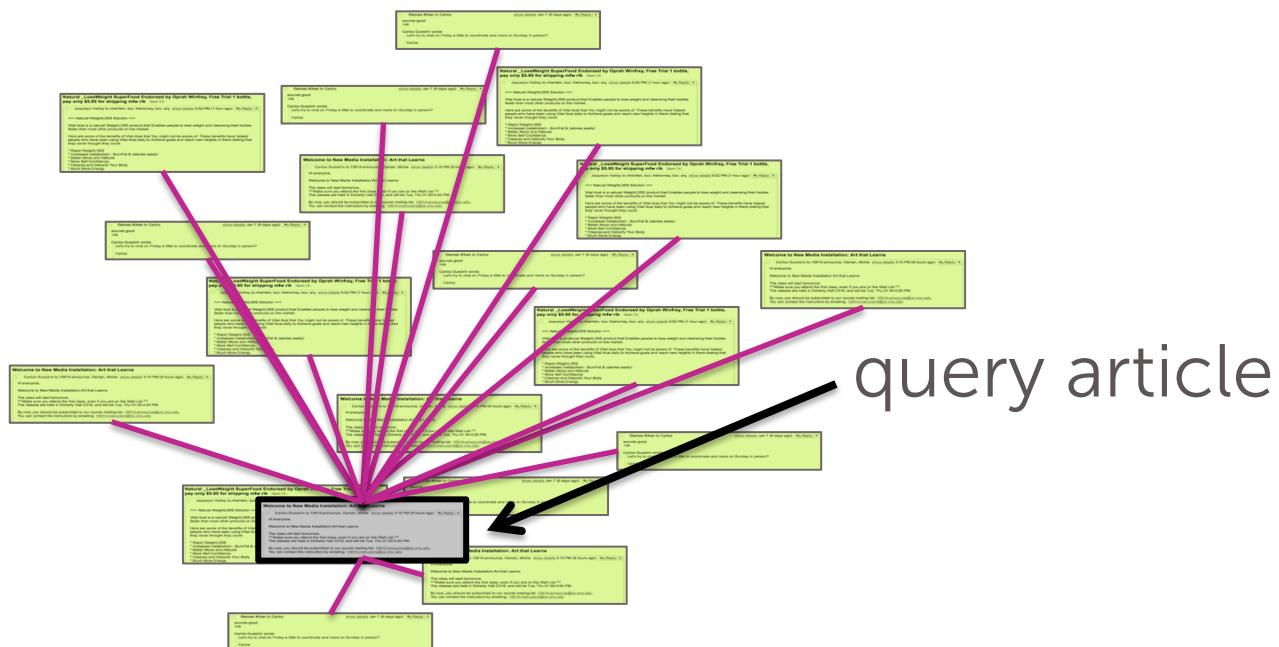
1-NN search for retrieval

Space of all articles,
organized by similarity of text



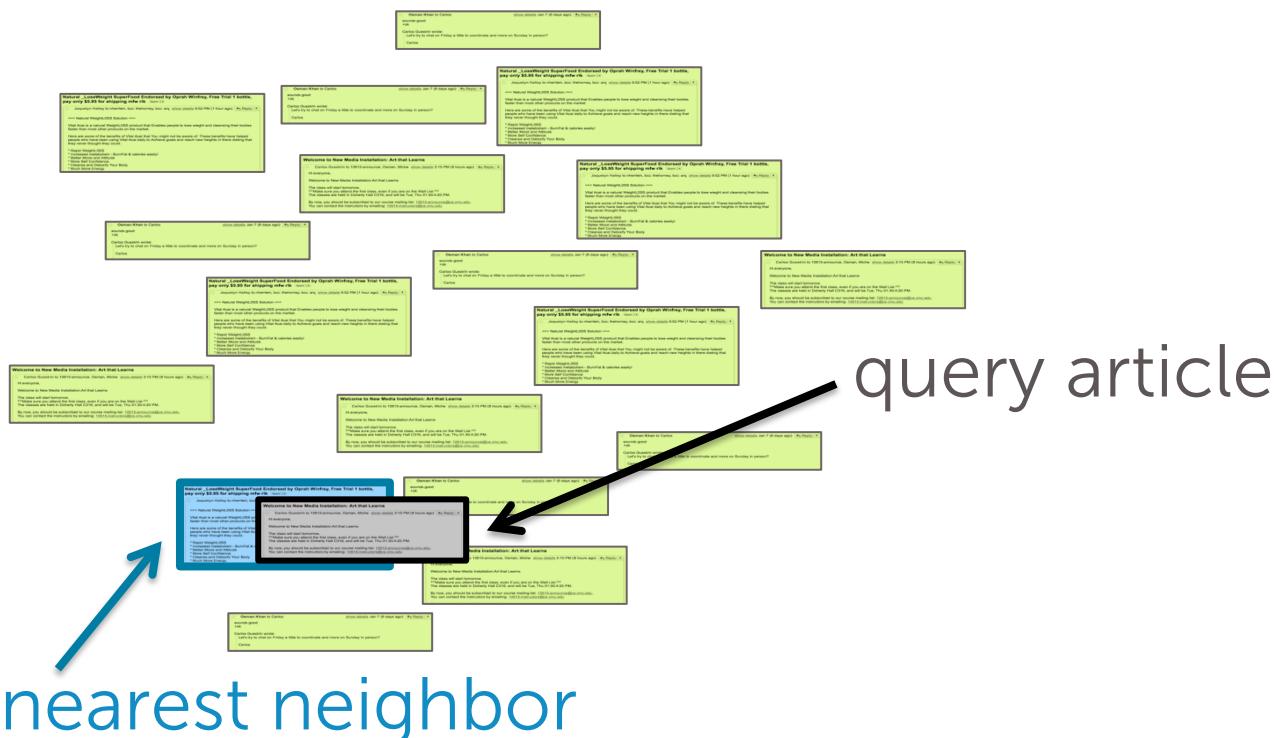
Compute distances to all docs

Space of all articles, organized by similarity of text



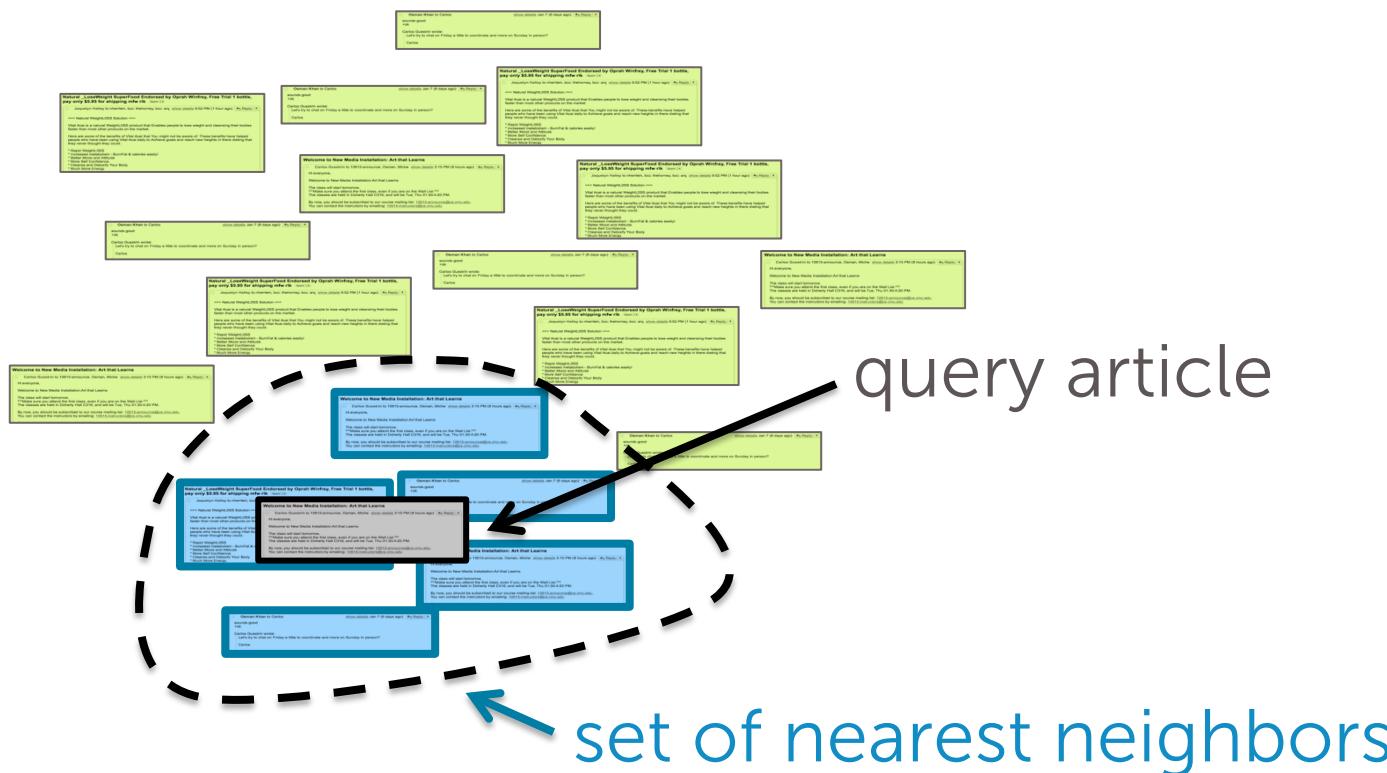
Retrieve “nearest neighbor”

Space of all articles,
organized by similarity of text



Or set of nearest neighbors

Space of all articles,
organized by similarity of text



1-NN algorithm

1 – Nearest neighbor

- **Input:** Query article :  \underline{x}_q
Corpus of documents  $(N \text{ docs})$
- **Output:** *Most* similar article  $\leftarrow x^{NN}$

Formally:

$$x^{NN} = \min_{x_i} \text{distance}(x_q, x_i)$$

1-NN algorithm

Initialize $\text{Dist2NN} = \underline{\infty}$,

For $i=1,2,\dots,N$

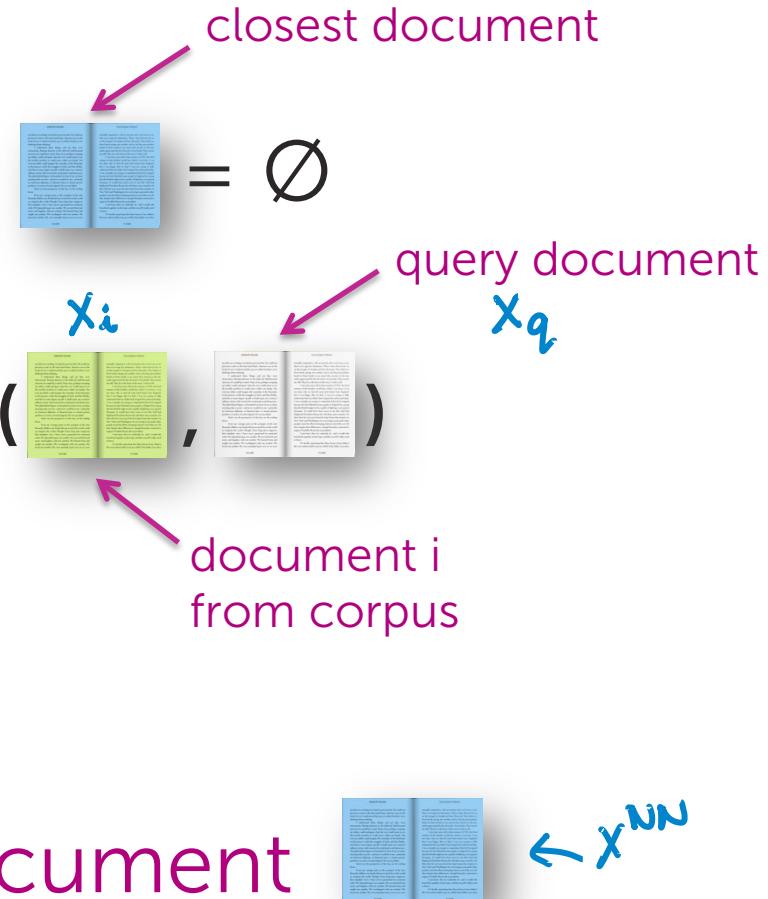
Compute: $\delta = \text{distance}(x_i, x_q)$

If $\delta < \text{Dist2NN}$

set $x_i =$

set $\text{Dist2NN} = \delta$

Return most similar document



closest document in
corpus to query article

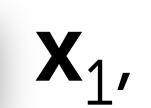
k-NN algorithm

k – Nearest neighbor

- **Input:** Query article : \mathbf{x}_q



Corpus of documents



: $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$

- **Output:** *List of k* similar articles

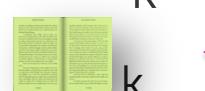


Formally:

$$X^{NN} = \{x^{NN_1}, \dots, x^{NN_k}\}$$

for all x_i not in X^{NN} , $\text{distance}(x_i, x_q) \geq \max_{j=1, \dots, k} \text{distance}(x^{NN_j}, x_q)$

k-NN algorithm

Initialize $\text{Dist2kNN} = \text{sort}(\delta_1, \dots, \delta_k)$ ← list of sorted distances
=  \dots  δ_1  \dots  δ_k ← list of sorted docs

For $i=k+1, \dots, N$

Compute: $\delta = \text{distance}(\text{book}_i, \text{book}_q)$ ← query doc

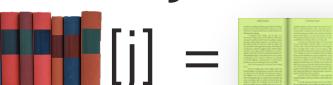
If $\delta < \text{Dist2kNN}[k]$ ← distance to k^{th} NN (furthest NN in set)

find j such that $\delta > \text{Dist2kNN}[j-1]$ but $\delta < \text{Dist2kNN}[j]$

remove furthest house and shift queue:

 $[1:k] =$  -1

$\text{Dist2kNN}[j+1:k] = \text{Dist2kNN}[j:k-1]$

set $\text{Dist2kNN}[j] = \delta$ and  $[j] = \text{book}_i$ ← closest k docs to query doc

Return k most similar articles



Critical elements of NN search

Item (e.g., doc) representation

$$\mathbf{x}_q \leftarrow$$



Measure of **distance** between items:

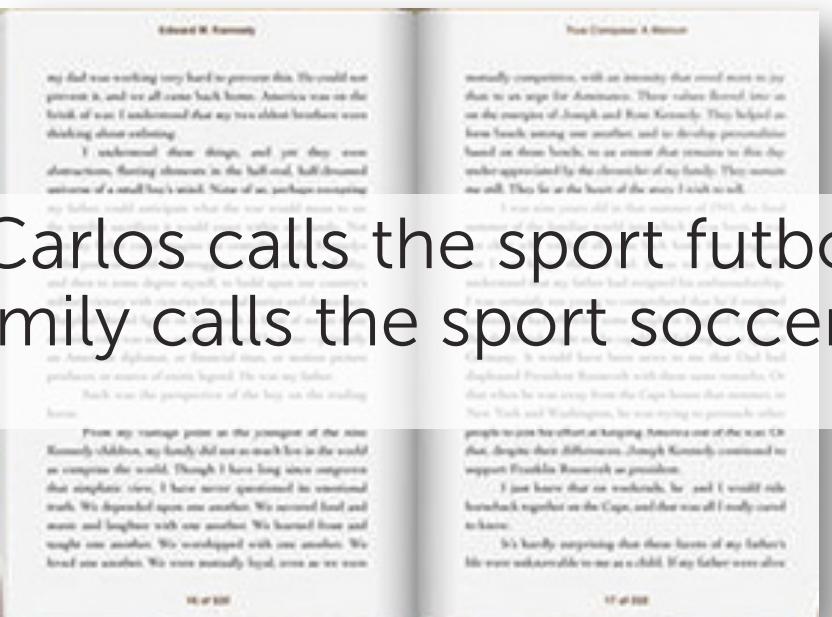
$$\delta = \text{distance}(\mathbf{x}_i, \mathbf{x}_q)$$

Document representation

Word count document representation

Bag of words model

- Ignore order of words
 - Count # of instances of each word in vocabulary



“Carlos calls the sport futbol.
Emily calls the sport soccer.”

Issues with word counts – Rare words



Common words in doc: “the”, “player”, “field”, “goal”

Dominate rare words like: “futbol”, “Messi”

TF-IDF document representation

Emphasizes **important words**

- Appears frequently in document (**common locally**)
- Appears rarely in corpus (**rare globally**)

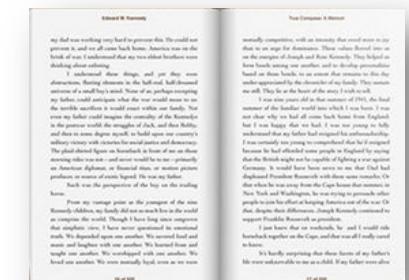
TF-IDF document representation

Emphasizes **important words**

- Appears frequently in document (**common locally**)

Term frequency =  word counts

- Appears rarely in corpus (**rare globally**)



TF-IDF document representation

Emphasizes **important words**

- Appears frequently in document (**common locally**)

Term frequency =  word counts

- Appears rarely in corpus (**rare globally**)

Inverse doc freq. = $\log \frac{\# \text{ docs}}{1 + \# \text{ docs using word}}$



TF-IDF document representation

Emphasizes **important words**

- Appears frequently in document (**common locally**)

Term frequency =  word counts

- Appears rarely in corpus (**rare globally**)

Inverse doc freq. = $\log \frac{\# \text{ docs}}{1 + \# \text{ docs using word}}$

Trade off: **local frequency vs. global rarity**

$tf * idf$

Distance metrics

Distance metrics: Defining notion of “closest”

In 1D, just Euclidean distance:

$$\text{distance}(x_i, x_q) = |x_i - x_q|$$

In multiple dimensions:

- can define many interesting distance functions
- most straightforwardly, might want to weight different dimensions differently

Weighting different features

Reasons:

- Some features are more relevant than others



bedrooms
bathrooms
sq.ft. living
sq.ft. lot
floors
year built
year renovated
waterfront



Weighting different features

Reasons:

- Some features are more relevant than others



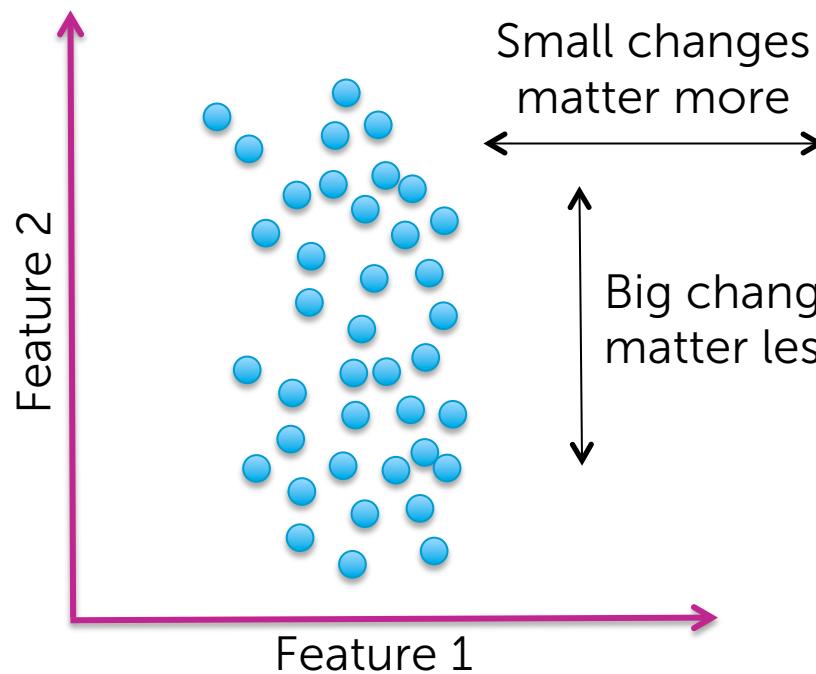
**title
abstract
main body
conclusion**



Weighting different features

Reasons:

- Some features are more relevant than others
- Some features vary more than others



Small changes
matter more

Big changes
matter less

Specify weights
as a function of
feature spread

For feature j:

$$\frac{1}{\max_i(\mathbf{x}_i[j]) - \min_i(\mathbf{x}_i[j])}$$

Scaled Euclidean distance

Formally, this is achieved via

$$\text{distance}(\mathbf{x}_i, \mathbf{x}_q) = \sqrt{a_1(\mathbf{x}_i[1]-\mathbf{x}_q[1])^2 + \dots + a_d(\mathbf{x}_i[d]-\mathbf{x}_q[d])^2}$$

weight on each feature
(defining relative importance)

Effect of binary weights

$$\text{distance}(\mathbf{x}_i, \mathbf{x}_q) = \sqrt{a_1(\mathbf{x}_i[1]-\mathbf{x}_q[1])^2 + \dots + a_d(\mathbf{x}_i[d]-\mathbf{x}_q[d])^2}$$

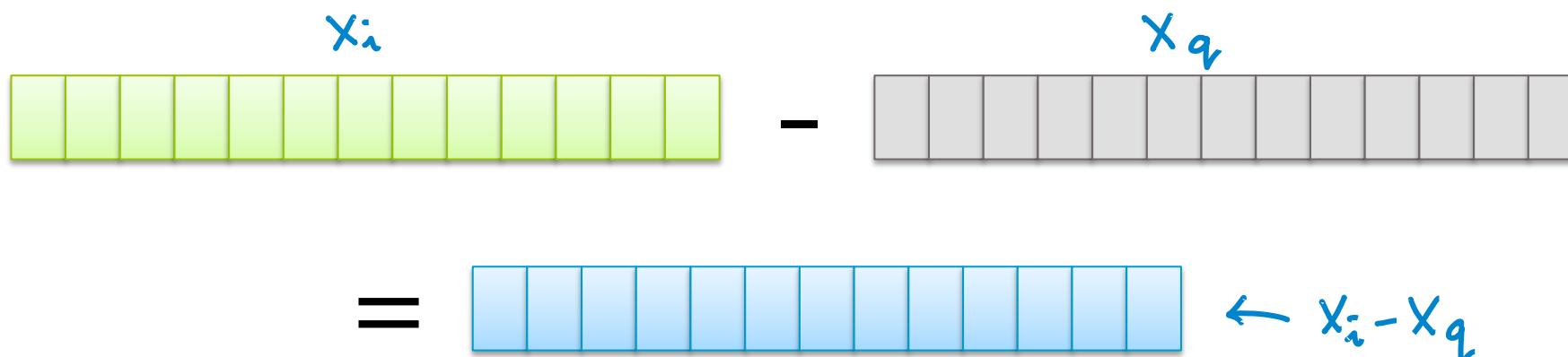
Setting weights as 0 or 1
is equivalent to
feature selection

Feature engineering/
selection is
important, but hard

(non-scaled) Euclidean distance

Defined in terms of inner product

$$\text{distance}(\mathbf{x}_i, \mathbf{x}_q) = \sqrt{(\mathbf{x}_i - \mathbf{x}_q)^T (\mathbf{x}_i - \mathbf{x}_q)}$$
$$= \sqrt{(\mathbf{x}_i[1] - \mathbf{x}_q[1])^2 + \dots + (\mathbf{x}_i[d] - \mathbf{x}_q[d])^2}$$



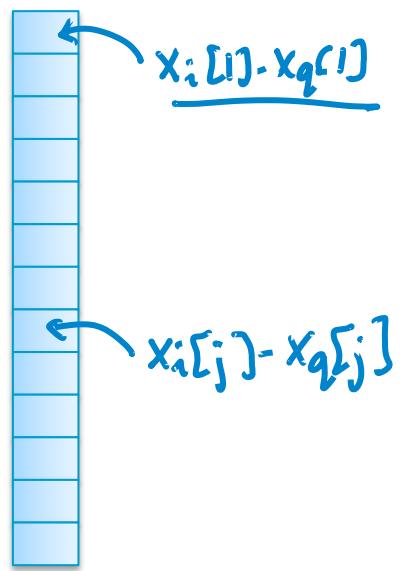
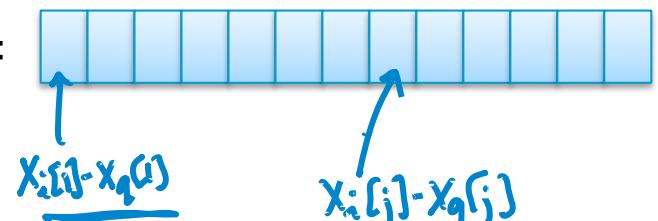
(non-scaled) Euclidean distance

Defined in terms of inner product

$$\text{distance}(\mathbf{x}_i, \mathbf{x}_q) = \sqrt{(\mathbf{x}_i - \mathbf{x}_q)^\top (\mathbf{x}_i - \mathbf{x}_q)}$$

↙

$$= \sqrt{(\mathbf{x}_i[1] - \mathbf{x}_q[1])^2 + \dots + (\mathbf{x}_i[d] - \mathbf{x}_q[d])^2}$$

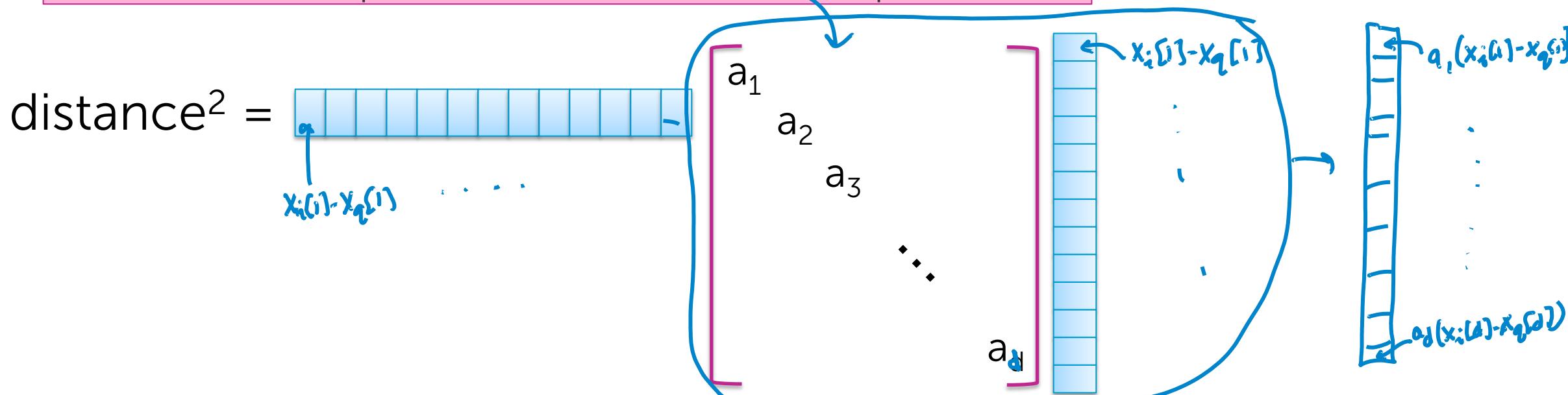


*take
sq.rt.*

Scaled Euclidean distance

Defined in terms of inner product

$$\text{distance}(\mathbf{x}_i, \mathbf{x}_q) = \sqrt{(\mathbf{x}_i - \mathbf{x}_q)^T \mathbf{A} (\mathbf{x}_i - \mathbf{x}_q)}$$
$$= \sqrt{a_1(x_i[1] - x_q[1])^2 + \dots + a_d(x_i[d] - x_q[d])^2}$$



Another natural inner product measure



x_q



x_i



Similarity

$$= \mathbf{x}_i^T \mathbf{x}_q$$

$$= \sum_{j=1}^d \mathbf{x}_i[j] \mathbf{x}_q[j]$$

$$= 13$$

Another natural inner product measure



Similarity

= 0



Cosine similarity – normalize

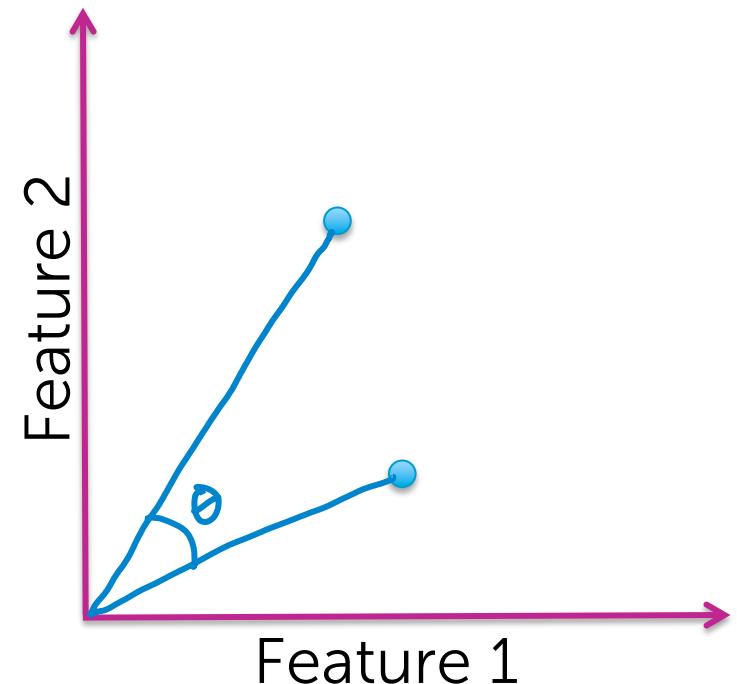
Similarity =

$$\frac{\sum_{j=1}^d \mathbf{x}_i[j] \mathbf{x}_q[j]}{\sqrt{\sum_{j=1}^d (\mathbf{x}_i[j])^2} \sqrt{\sum_{j=1}^d (\mathbf{x}_q[j])^2}}$$
$$\mathbf{x}_i^\top \mathbf{x}_q = \cos(\theta)$$
$$= \left(\frac{\mathbf{x}_i}{\|\mathbf{x}_i\|} \right)^\top \left(\frac{\mathbf{x}_q}{\|\mathbf{x}_q\|} \right)$$

first normalize

- Not a proper distance metric
- Efficient to compute for sparse vecs

$$\mathbf{a}^\top \mathbf{b} = \|\mathbf{a}\| \|\mathbf{b}\| \cos(\theta)$$



Normalize



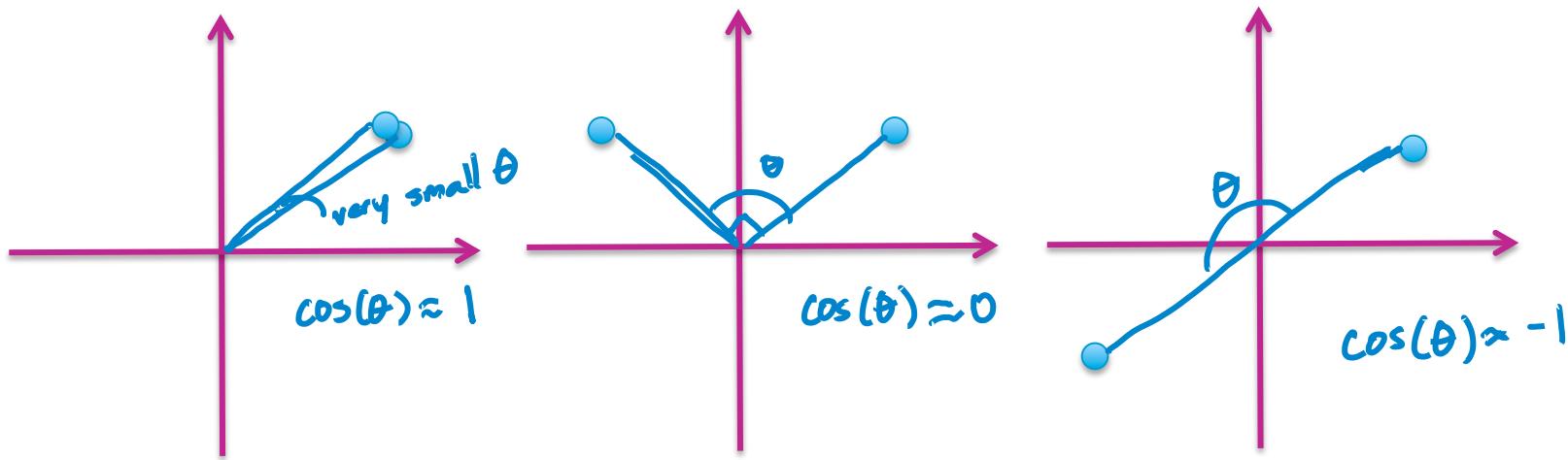
1	0	0	0	5	3	0	0	1	0	0	0	0
---	---	---	---	---	---	---	---	---	---	---	---	---

x_i

$$\sqrt{(1^2 + 5^2 + 3^2 + 1^2)} \quad \leftarrow \|x_i\| = \sum_{j=1}^d x_i[j]^2$$

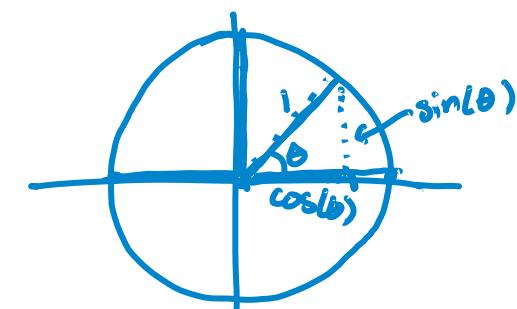
1					5	3		1				
/	0	0	0	/	/	0	0	/	0	0	0	0
6				6	6		6					

Cosine similarity



In general, $-1 < \text{similarity} < 1$

For positive features (like tf-idf)
 $-1 < \text{similarity} < 1$



Define **distance** = **1-similarity**

To normalize or not?



1 0 0 0 5 3 0 0 1 0 0 0 0

3 1 0 0 2 0 0 1 0 1 0 0 0
Similarity = 13



2 0 0 0 10 6 0 0 2 0 0 0 0

6 2 0 0 4 0 0 2 0 2 0 0 0
Similarity = 52



Normalize



1	0	0	0	5	3	0	0	1	0	0	0	0
---	---	---	---	---	---	---	---	---	---	---	---	---

$$\sqrt{1^2 + 5^2 + 3^2 + 1^2}$$

1					5	3		1				
/	0	0	0	/	/	0	0	/	0	0	0	0
6				6	6		6					

In the normalized case

1				5	3			1					
/	0	0	0	/	/	0	0	/	0	0	0	0	0
6				6	6			6					

3	1			1	/	0	0	1	/	0	1	0	0	0
/	/	0	0	/	0	0	0	/	0	0	/	0	0	0
4	4			2				4			4			

Similarity

= 13/24

1	0	0	0	5	3	0	0	1	/	0	0	0	0
/	0	0	0	/	/	0	0	6	/	0	0	0	0
6				6	6				6				

3	1			1	/	0	0	1	/	0	1	/	0	0	0
/	/	0	0	/	0	0		/	0		/	0	0	0	0
4	4			2				4			4				

Similarity
= 13/24

But not always desired...

long document

long document

He was working very hard to prove this. He could not prove it, and so we came back home. America was on the road of progress, but we were still far behind. The old brother was still silent about this.

I understood him better than the hell had told him, and I was not surprised. Of course, no one among us could have been surprised. We were all aware that the terrible socialist was to come within our family. Not even the old brother could have been surprised. He was not surprised because he was fully aware of the terrible power the socialist of czar, and he was fully aware of the terrible power the socialist of America.

my military service with the regulars and sailors and Marines. I was a member of the 1st Marine Division, which was the most active division in the Pacific during the war. We were in the Philippines, Okinawa, and Iwo Jima. We were also involved in the Korean War. I served as a gunner on a tank and was assigned to the 1st Marine Division. I was promoted to sergeant and became a platoon leader. I was awarded the Purple Heart for my injuries in the Korean War. I was honorably discharged from the military in 1952.

long document long document

Normalizing can
make dissimilar
objects appear
more similar

Common compromise: Just cap maximum word counts

Other distance metrics

- Mahalanobis
- rank-based
- correlation-based
- Manhattan
- Jaccard
- Hamming
- ...

Combining distance metrics

Example of document features:

1. Text of document
 - Distance metric: Cosine similarity
2. # of reads of doc
 - Distance metric: Euclidean distance

Add together with user-specified weights