

## Document representation

Word count document representation

Bag of words model

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

1	2	1	1	2	2	1
Carlos	the	Emily	soccer	calls	sport	but

"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

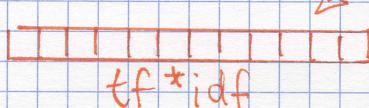
- Appears frequently in document (common locally)

Term frequency =  $\frac{\text{Word counts}}{\text{Total words}}$

- Appears rarely in corpus (rare globally)

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

Trade off: local frequency vs. global rarity



## Distance metrics: Euclidean and scaled Euclidean

Distance metric  
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

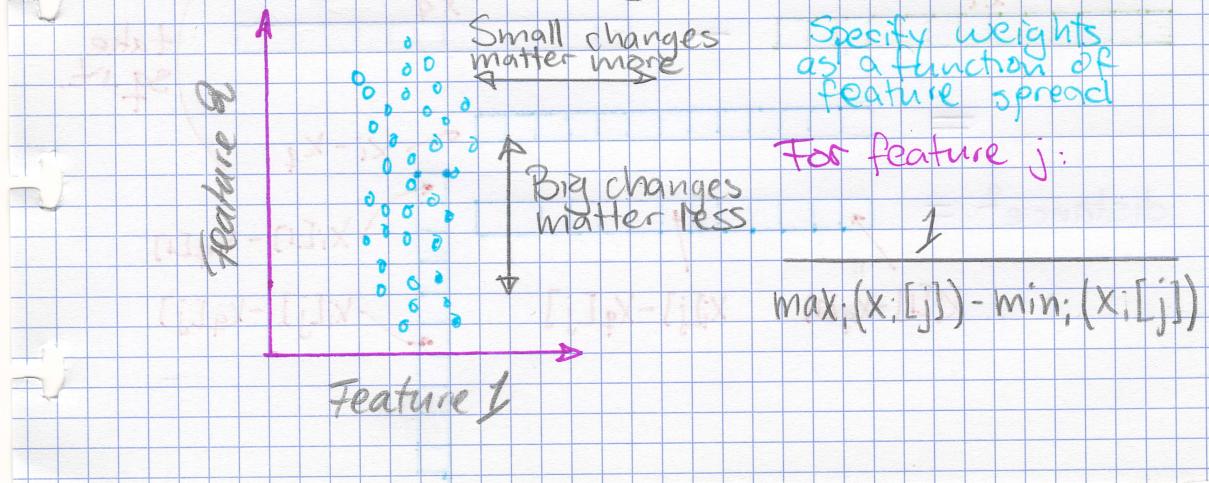
Weighting different features

Reasons:

- Some features are more relevant than others



- Some features vary more than others



## Scaled Euclidean distance

Formally, this is achieved via

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

$$\sqrt{a_1(x_i[1] - x_q[1])^2 + \dots + a_d(x_i[d] - x_q[d])^2}$$

weight on each feature  
(defining relative importance)

Setting weights as 0 or 1  
is equivalent to  
feature selection

{ Effect of  
binary  
weights

Feature engineering/  
selection is important  
but hard

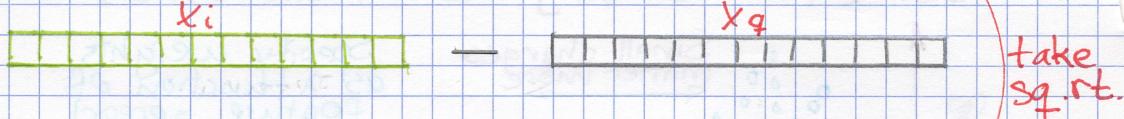
Writing (scaled) Euclidean distance using (weighted)  
inner products

(non-scaled) Euclidean distance

Defined in terms of inner product

$$\text{distance}(x_i, x_q) = \sqrt{(x_i - x_q)^T (x_i - x_q)} =$$

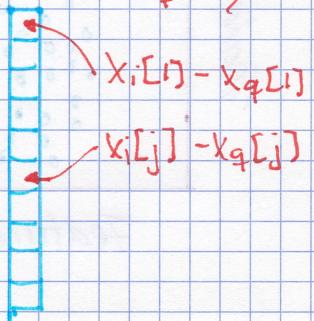
$$\sqrt{(x_i[1] - x_q[1])^2 + \dots + (x_i[d] - x_q[d])^2}$$



$$= \sqrt{\sum_{j=1}^d (x_i[j] - x_q[j])^2}$$

$$\text{distance}^2 = \sum_{j=1}^d (x_i[j] - x_q[j])^2$$

$$= (x_i[1] - x_q[1])^2 + (x_i[2] - x_q[2])^2 + \dots + (x_i[d] - x_q[d])^2$$



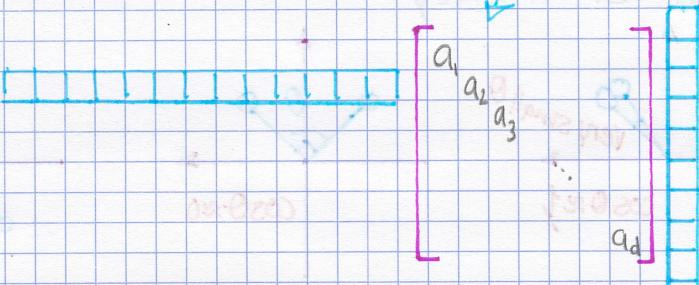
## Scaled Euclidean distance

Defined in terms of inner product

$$\text{distance}(x_i, x_q) = \sqrt{(x_i - x_q)^T A (x_i - x_q)}$$

$$\sqrt{a_1(x_i[1] - x_q[1])^2 + \dots + a_d(x_i[d] - x_q[d])^2}$$

$$\text{distance}^2 =$$



## Distance metrics: cosine similarity

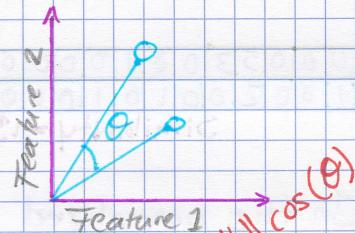
## Another natural inner product measure

$$\begin{matrix} 1 & 0 & 0 & 0 & 5 & 3 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 3 & 0 & 0 & 0 & 2 & 0 & 0 & 1 & 0 & 1 & 0 & 0 & 0 \end{matrix} \xrightarrow{\text{X}_q} \text{Similarity} = \mathbf{x}_i^T \mathbf{x}_q$$

$$k_i = \sum_{j=1}^d x_i[j] x_q[j]$$

$$= 13$$

$\begin{array}{ccccccccc} 1 & 0 & 0 & 0 & 5 & 3 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 0 & 0 & 9 & 0 & 0 \end{array} \begin{array}{ccccccccc} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 6 & 0 & 4 & 0 \end{array}$   $\rightarrow$  Similarity = 0



Cosine Similarity — normalize

$$\text{Similarity} = \sum_i x_i[j] x_q[j]$$

- Not a proper distance metric

- Efficient to compute for sparse vecs

$$\sqrt{\sum_{j=1}^d (x_i[j])^2} \sqrt{\sum_{j=1}^d (x_d[j])^2}$$

$$= \frac{x_i^T x_d}{\|x_i\| \|x_d\|} = \cos \theta$$

$$\mathbf{a}^T \mathbf{b} = \| \mathbf{a} \|$$

first  
normalize

$$\left( \frac{\mathbf{x}_i}{\| \mathbf{x}_i \|} \right)^T \left( \frac{\mathbf{x}_q}{\| \mathbf{x}_q \|} \right)$$

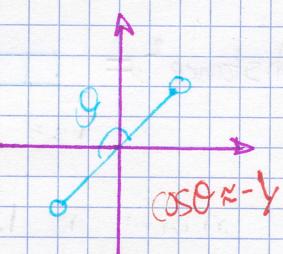
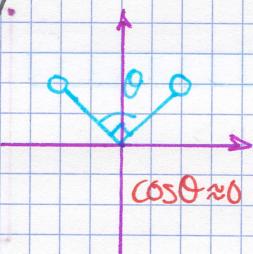
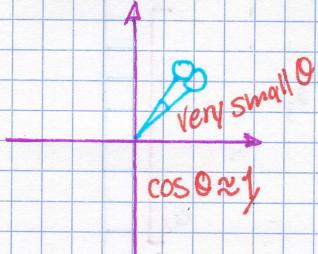
Normalize

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

$$\sqrt{1^2 + 5^2 + 3^2 + 7^2}$$

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

Cosine similarity.



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

For positive features (like tf-idf) } Our focus

$$0 \leq \text{similarity} \leq 1$$

Define distance = 1 - similarity

To normalize or not and other distance considerations

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 1 0 6 0 0 2 0 0 0 0  
6 2 0 0 4 0 0 2 0 2 0 0 0  
Similarity = 52

In the normalized case

1 5 3 1  
0 0 0 1 0 0 0 0  
6 6 6 6  
3 1 2 1 1  
1 0 0 0 1 0 0 0  
4 4 4 4  
Similarity = 13/24

But not always desired . . .



long document

Normalizing can make  
dissimilar objects  
appear more similar



long document long document

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]