

# Vector Space Models

Prof. Sameer Singh

---

CS 295: STATISTICAL NLP

WINTER 2017

January 19, 2017

# Outline

---

Latent Semantic Analysis

Vector Models for Words

Reducing the Dimensions

Direct Embeddings

# Outline

---

Latent Semantic Analysis

Vector Models for Words

Reducing the Dimensions

Direct Embeddings

# Example: Documents

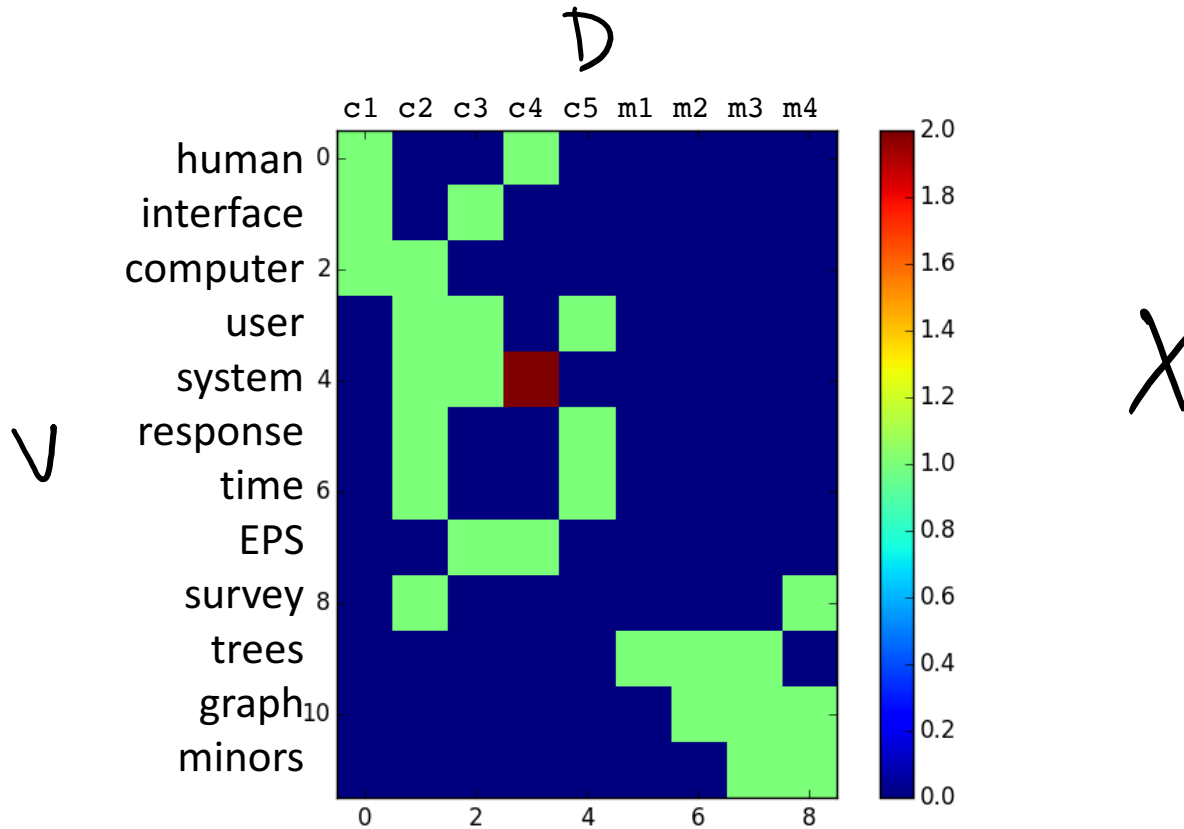
---

- c1: Human machine interface for ABC computer applications
- c2: A survey of user opinion of computer system response time
- c3: The EPS user interface management system
- c4: System and human system engineering testing of EPS
- c5: Relation of user perceived response time to error measurement

- m1: The generation of random, binary, ordered trees
- m2: The intersection graph of paths in trees
- m3: Graph minors IV: Widths of trees and well-quasi-ordering
- m4: Graph minors: A survey

From <http://lsa.colorado.edu/papers/dp1.LSAintro.pdf>

# Example: Term-Doc Matrix



# Problems with Sparse Vectors

---

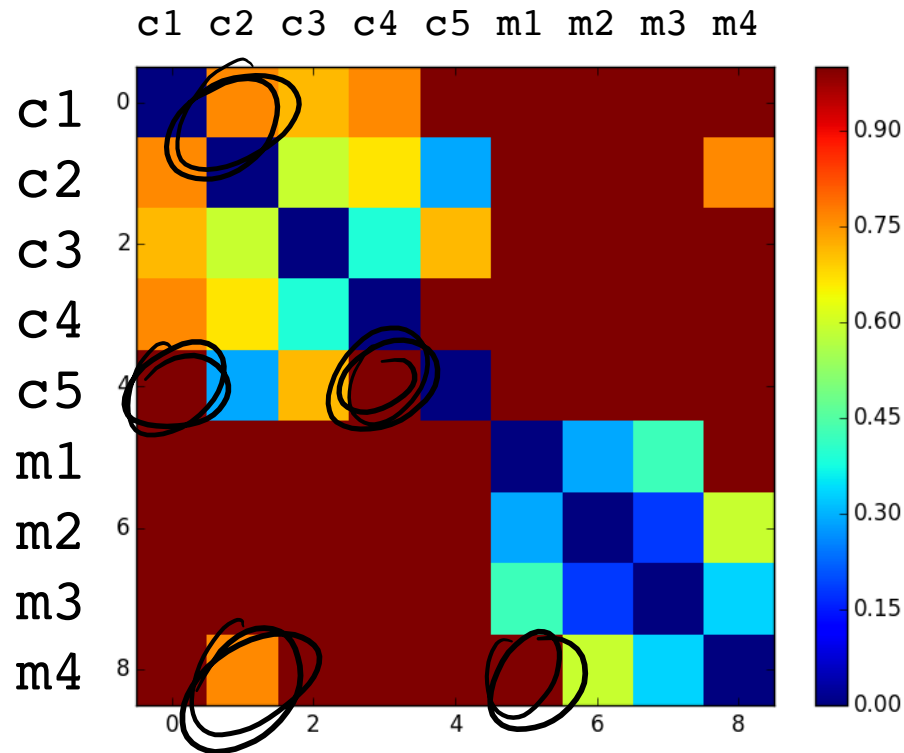
c2: A survey of user opinion of computer system response time

m4: Graph minors: A survey

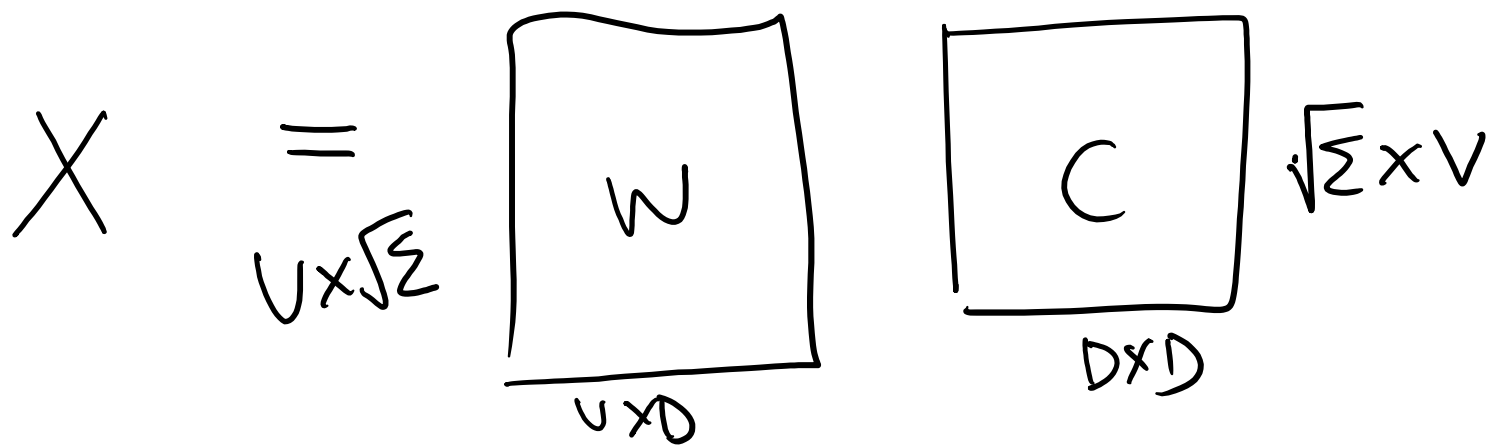
c1: Human machine interface  
for ABC computer applications

# Example: Distance Matrix

---

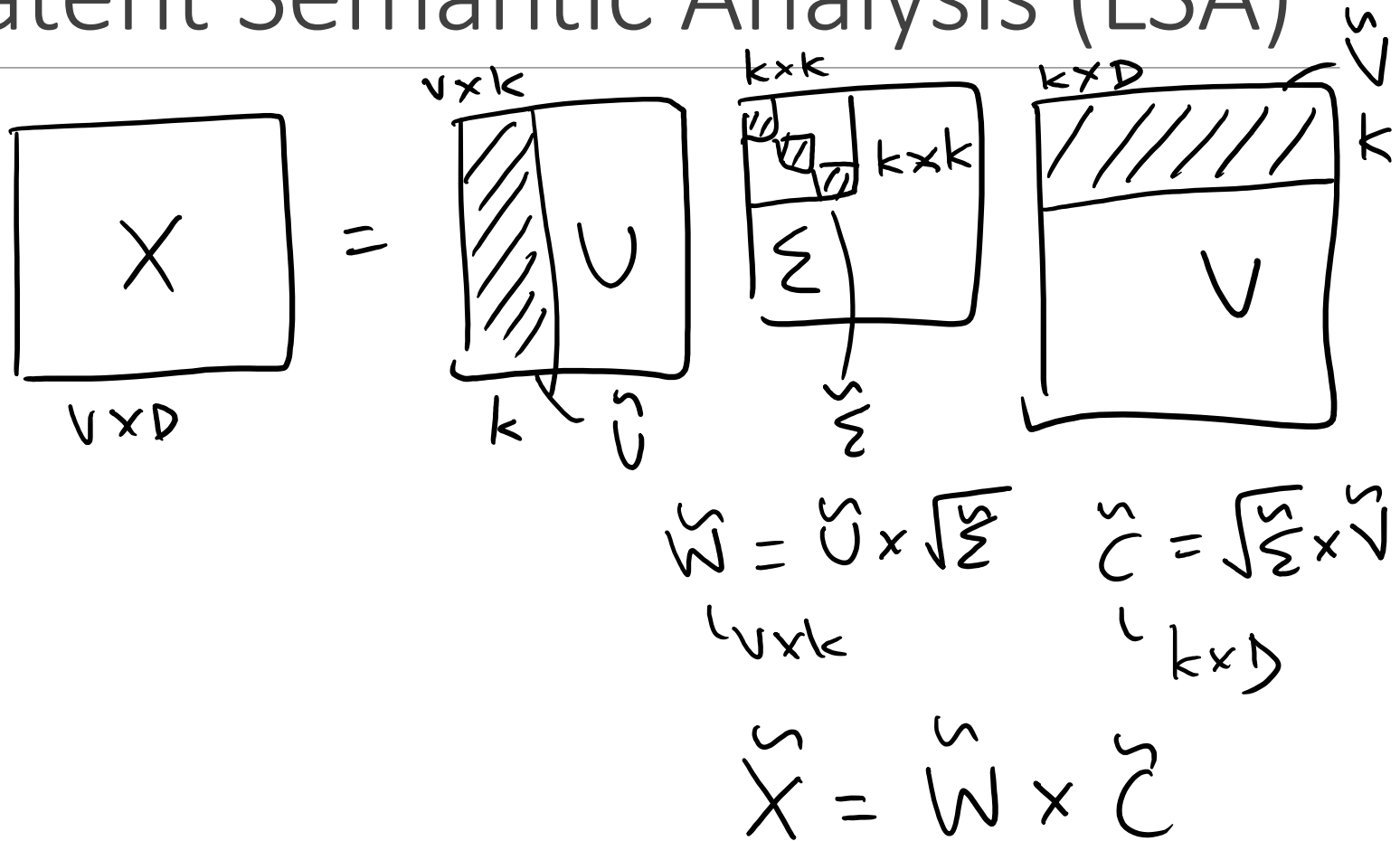


# Option 2: SVD $\hookrightarrow$ Singular Value Decomp.

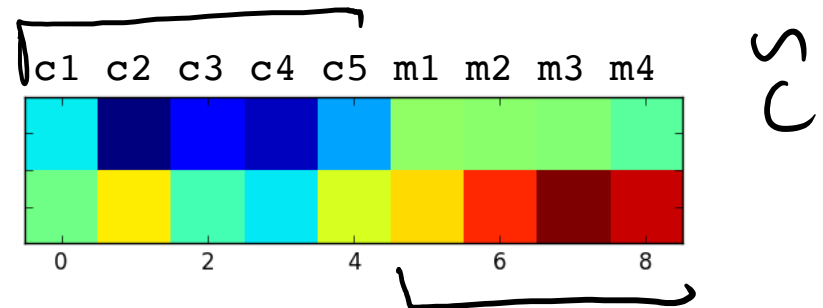
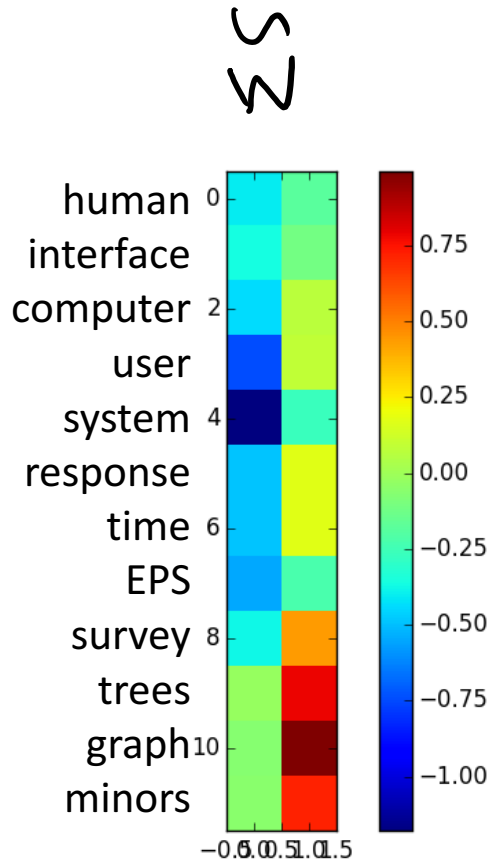




# Latent Semantic Analysis (LSA)



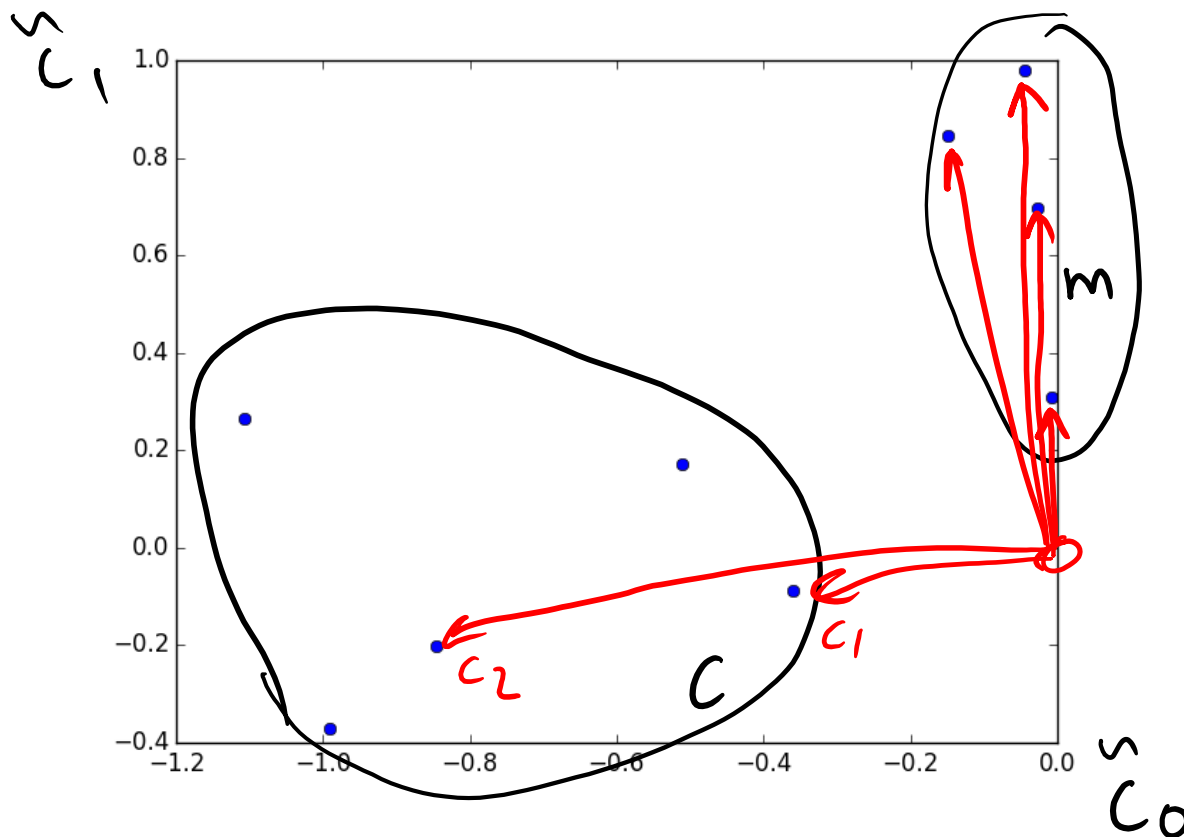
# Example: Decomposition, $k=2$



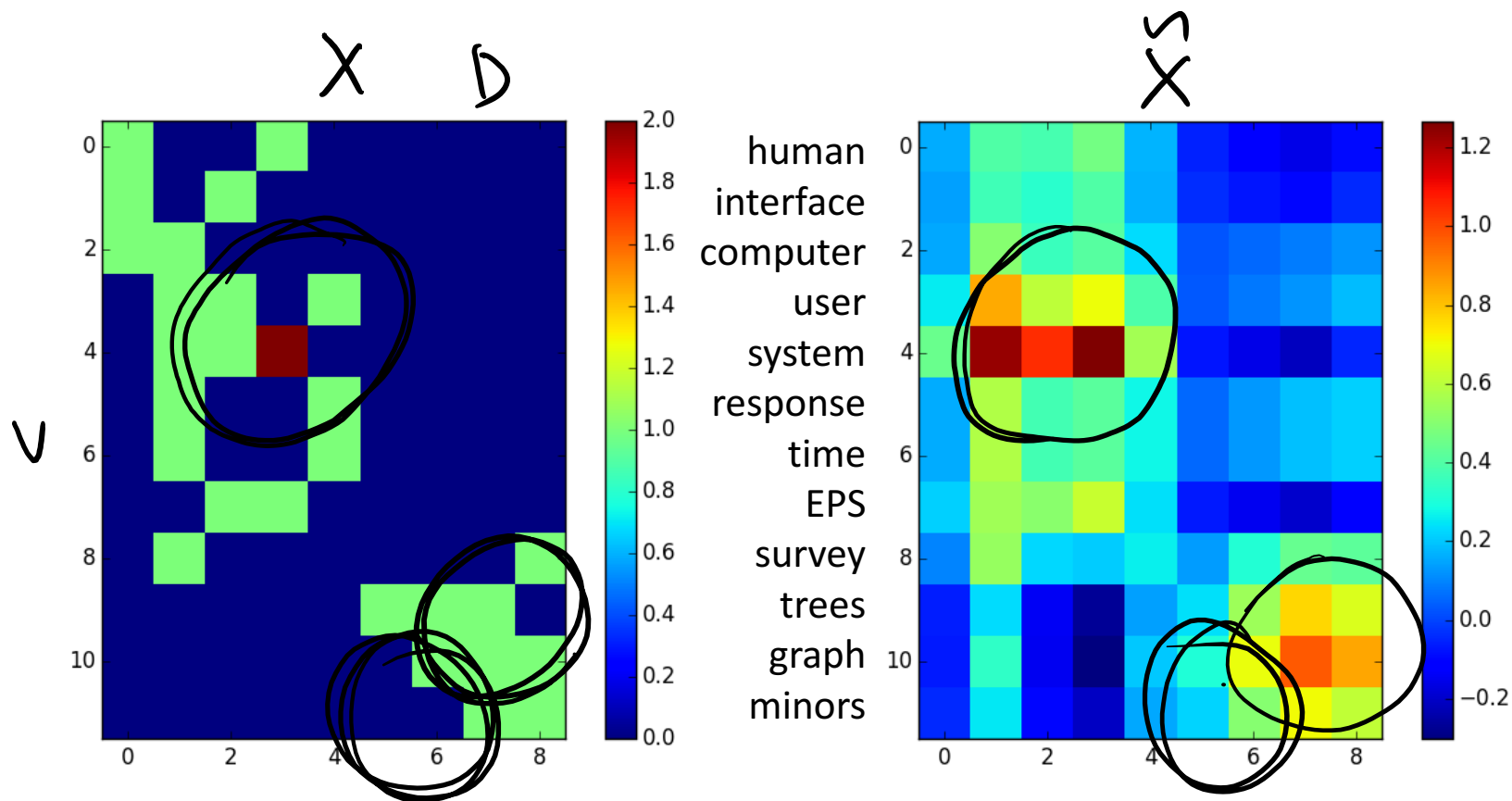
$$\phi(x_i) = \begin{array}{|c|c|c|c|c|c|c|c|c|c|} \hline \text{||} & \text{///} & \text{///} & \text{///} & \text{///} & \text{///} & \text{///} & \text{///} & \text{///} & \text{///} \\ \hline \end{array}$$

$\uparrow$   $\uparrow$   
 $c_{i0}$   $c_{i1}$

# New Document Vectors

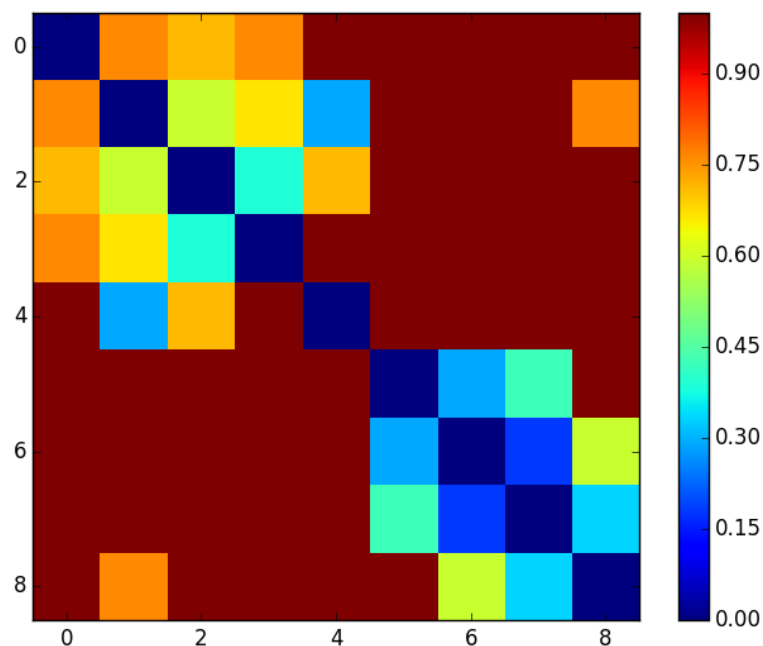


# Example: Reconstruction

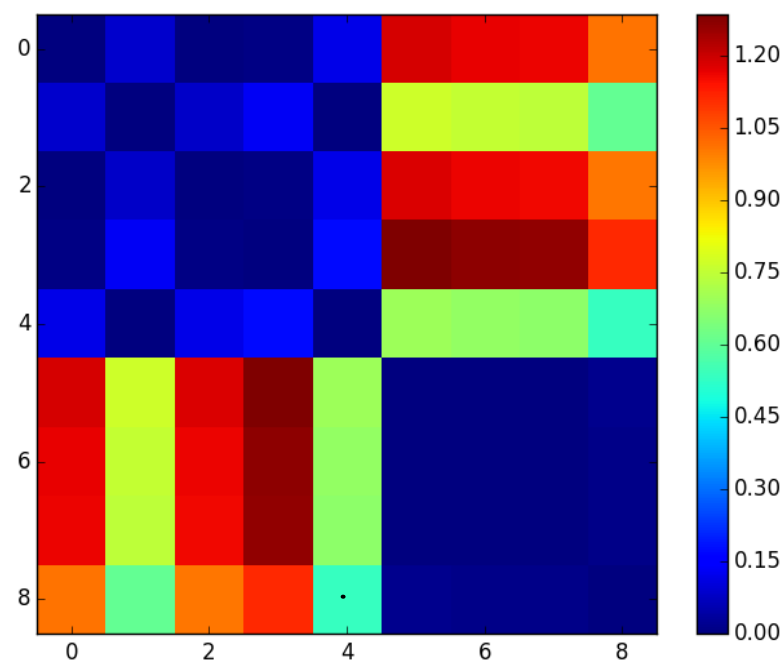


# Example: Distance Matrix

$\text{dist}(x_i, x_j)$



$\text{dist}(\tilde{x}_i, \tilde{x}_j)$



# Outline

---

Latent Semantic Analysis

Vector Models for Words

Reducing the Dimensions

Direct Embeddings

# Let's look at words

---

A bottle of **tezguino** is on the table.  
Everybody likes **tezguino**.  
**Tezguino** makes you drunk.  
We make **tezguino** out of corn.

What does **tezguino** mean?

Loud, motor oil, tortillas, choices, wine



You shall know a word by the company keeps.  
(Firth, 1957)

# Term-Context Matrix

---

C1: A bottle of \_\_\_\_\_ is on the table.

C2: Everybody likes \_\_\_\_\_.

C3: \_\_\_\_\_ makes you drunk.

C4: We make \_\_\_\_\_ out of corn.

	C1	C2	C3	C4
tezguino	1	1	1	1
loud	0	0	0	0
motor oil	1	0	0	0
tortillas	0	1	0	1
choices	0	1	1	0
wine	1	1	1	0



# What is a “Context”?

---

Can be anything you want!

- Entire contents of the sentence
- One word before and after
- Words in the same sentence
- Document it appears in
- Many other variations...

A bottle of **tezguino** is on the table.  
**Tezguino** makes you drunk.

...

I had a fancy bottle of **wine** and  
got drunk last night!  
The terrible **wine** is on the table.

# What is a “Context”?

Can be anything you want!

- Entire contents of the sentence
  - Unlikely to occur again!
- One word before and after
- Words in the same sentence
- Document ID it appears in
- Many other variations...

C<sub>1</sub> A bottle of **tezguino** is on the table.

C<sub>2</sub> **Tezguino** makes you drunk.

...

C<sub>3</sub> I had a fancy bottle of **wine** and got drunk last night!

C<sub>4</sub> The terrible **wine** is on the table.

	C1	C2	C3	C4
tezguino	1	1	0	0
wine	0	0	1	1

# What is a “Context”?

Can be anything you want!

- Entire contents of the sentence
- One word before and after
  - Or n-words
- Words in the same sentence
- Document it appears in
- Many other variations...

A bottle of **tezguino** is on the table.  
**Tezguino** makes you drunk.

...

I had a fancy bottle of **wine** and  
got drunk last night!

The terrible **wine** is on the table.

bottle-of <sup>on</sup> is-of ~~is on~~ makes-you and-got the-terrible ~~is on~~

tezguino	1	1	1	0	0
wine	1	1	0	1	1

# What is a “Context”?

Can be anything you want!

- Entire contents of the sentence
- One word before and after
- Words in the same sentence
  - Filter: nouns and verbs?
  - Bag of words in a window
- Document it appears in
- Many other variations...

A bottle of **tezguino** is on the table.  
**Tezguino** makes you drunk.

...

I had a fancy bottle of **wine** and  
got drunk last night!

The terrible **wine** is on the table.

	bottle	table	you	drunk	fancy	night	terrible
tezguino					0	0	0
wine			0				

# What is a “Context”?

Can be anything you want!

- Entire contents of the sentence
- One word before and after
- Words in the same sentence
- Document it appears in
  - Term-document matrix!
  - Latent Semantic Analysis
- Many other variations...

D<sub>1</sub> A bottle of **tezguino** is on the table.

D<sub>2</sub> **Tezguino** makes you drunk.

...

D<sub>3</sub> I had a fancy bottle of **wine** and got drunk last night!

D<sub>4</sub> The terrible **wine** is on the table.

	D1	D2	D3	D4
tezguino				
table				
bottle				
drunk				
wine				

# Pointwise Mutual Information

## Raw counts are not good

- Skewed towards common words/contexts
- Many of them are not *informative*
  - is, the, it, they, ...

## PMI(w,c)

- How much more likely is  $w$  to occur in  $c$ , than just randomly?

positive

$$\text{PMI}(w, c) = \log_2 \frac{P(w, c)}{P(w)P(c)}$$

$P(w, c)$

$\frac{f_{w,c}}{N}$

$\frac{\sum_j f_{w,j}}{N}$

$\frac{\sum_i f_{i,c}}{N}$

$\rightarrow P(w)P(c)$

# Outline

---

Latent Semantic Analysis

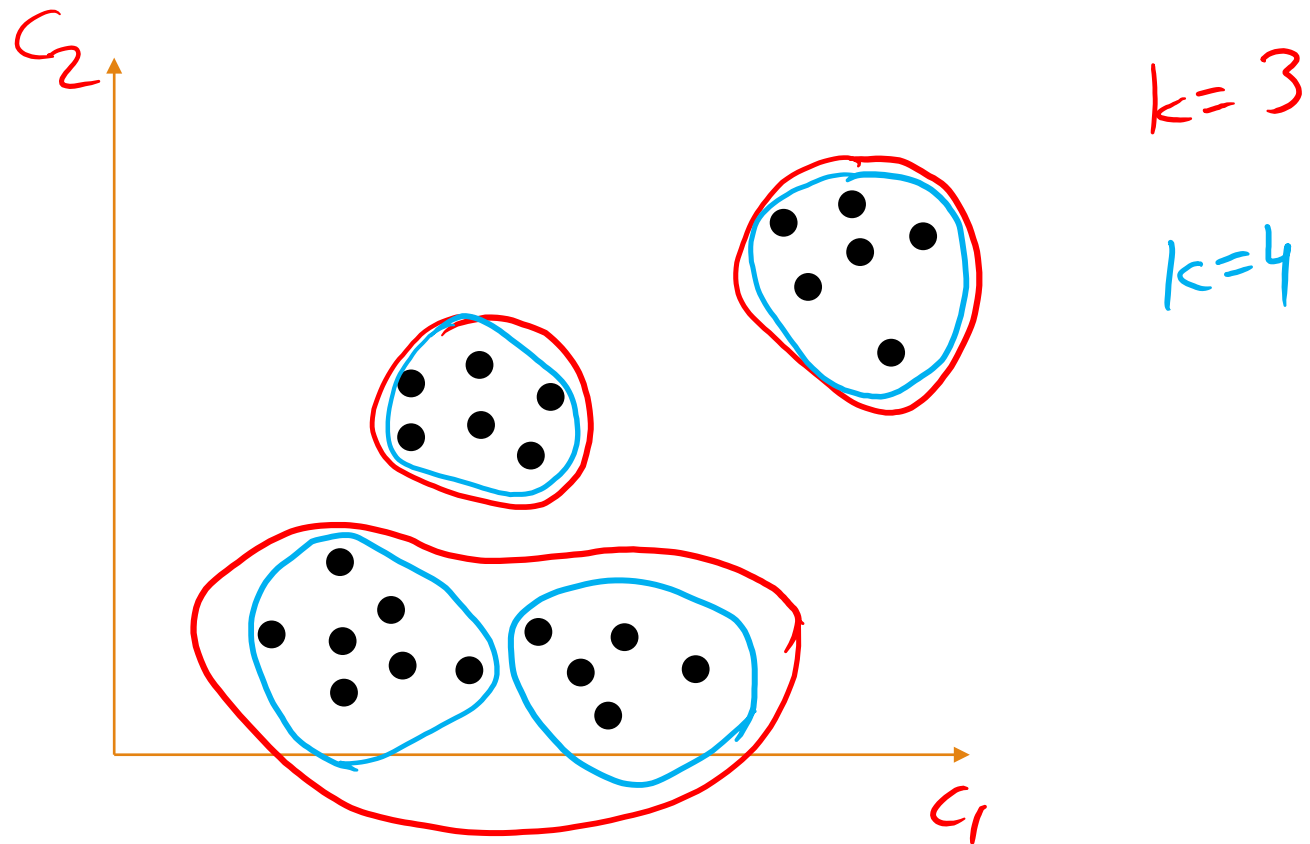
Vector Models for Words

Reducing the Dimensions

Direct Embeddings

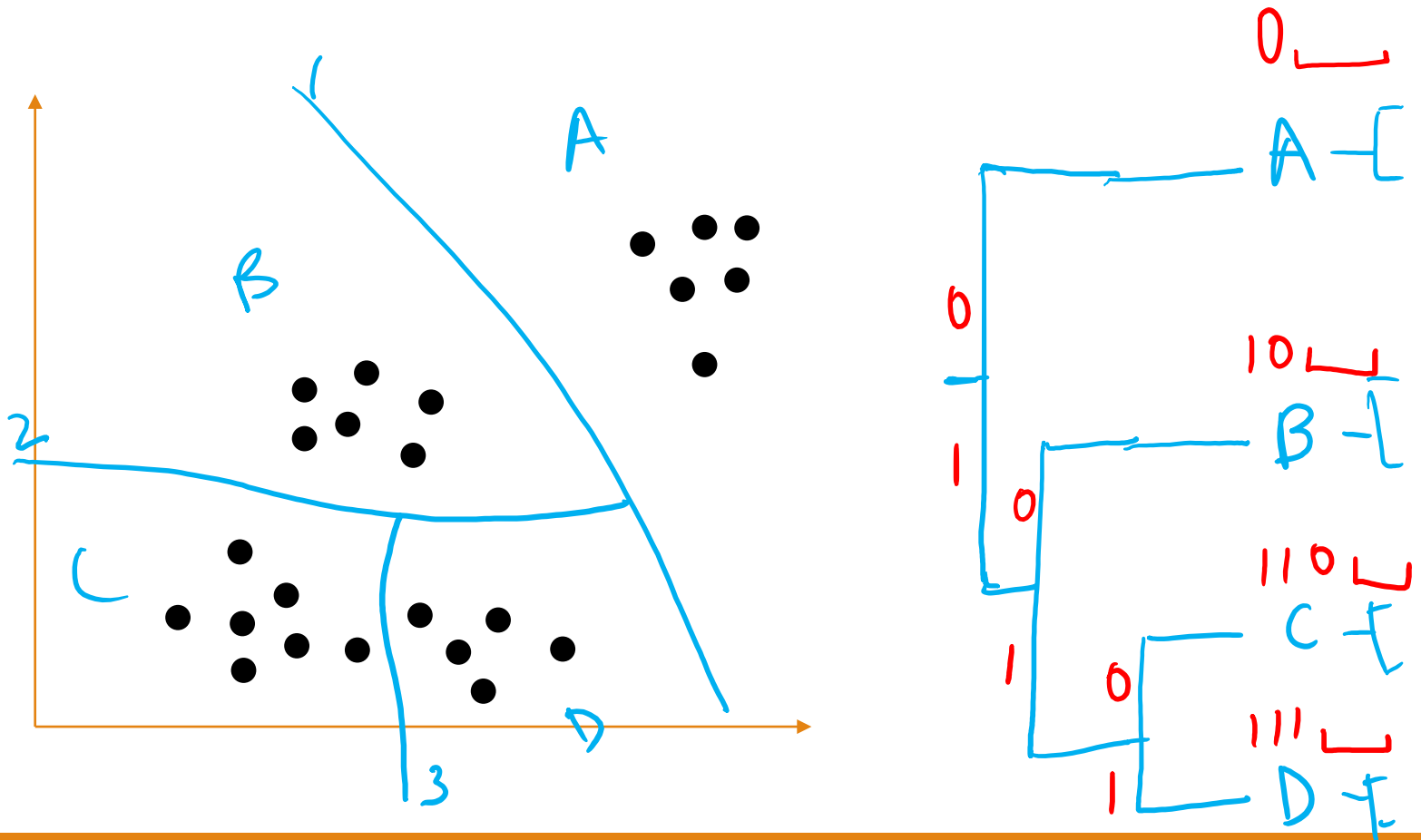
# Option 1: Revisiting Clustering

---

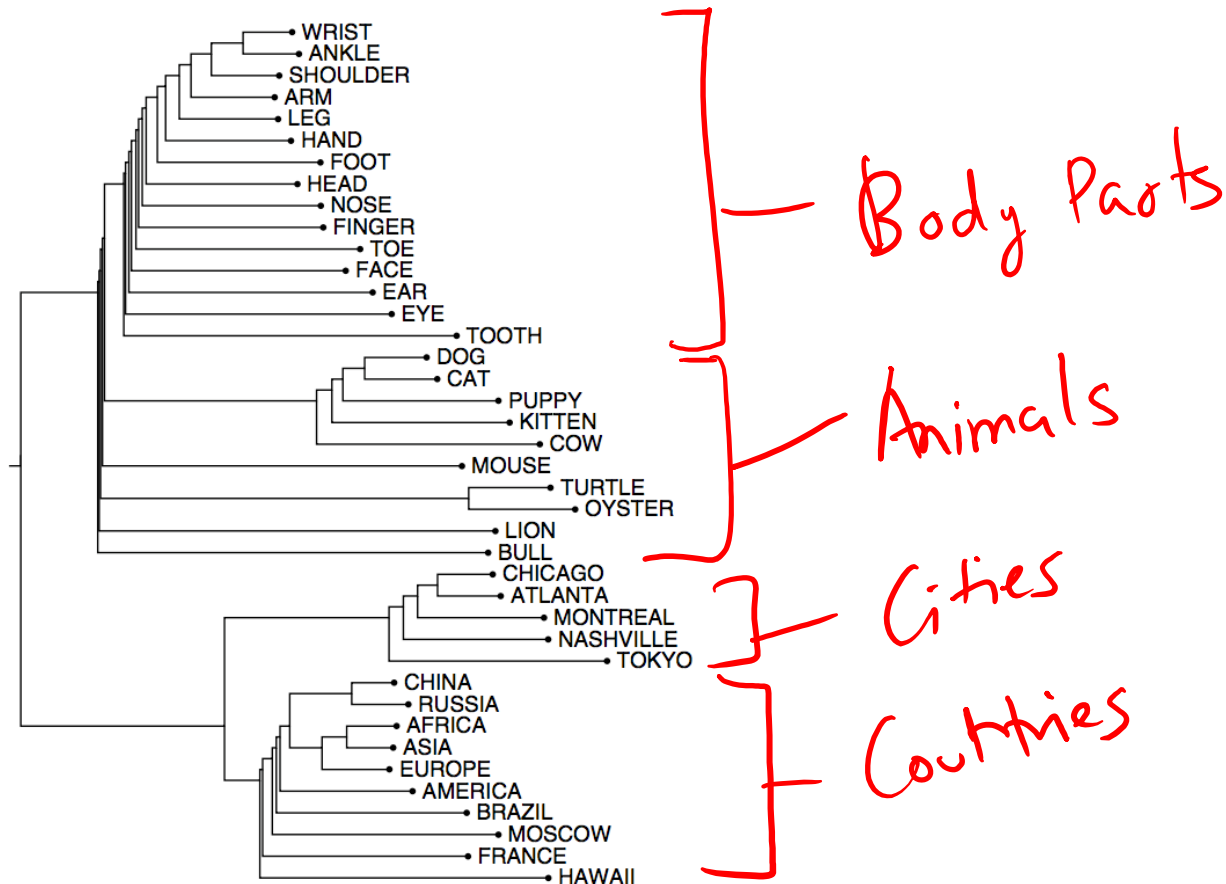




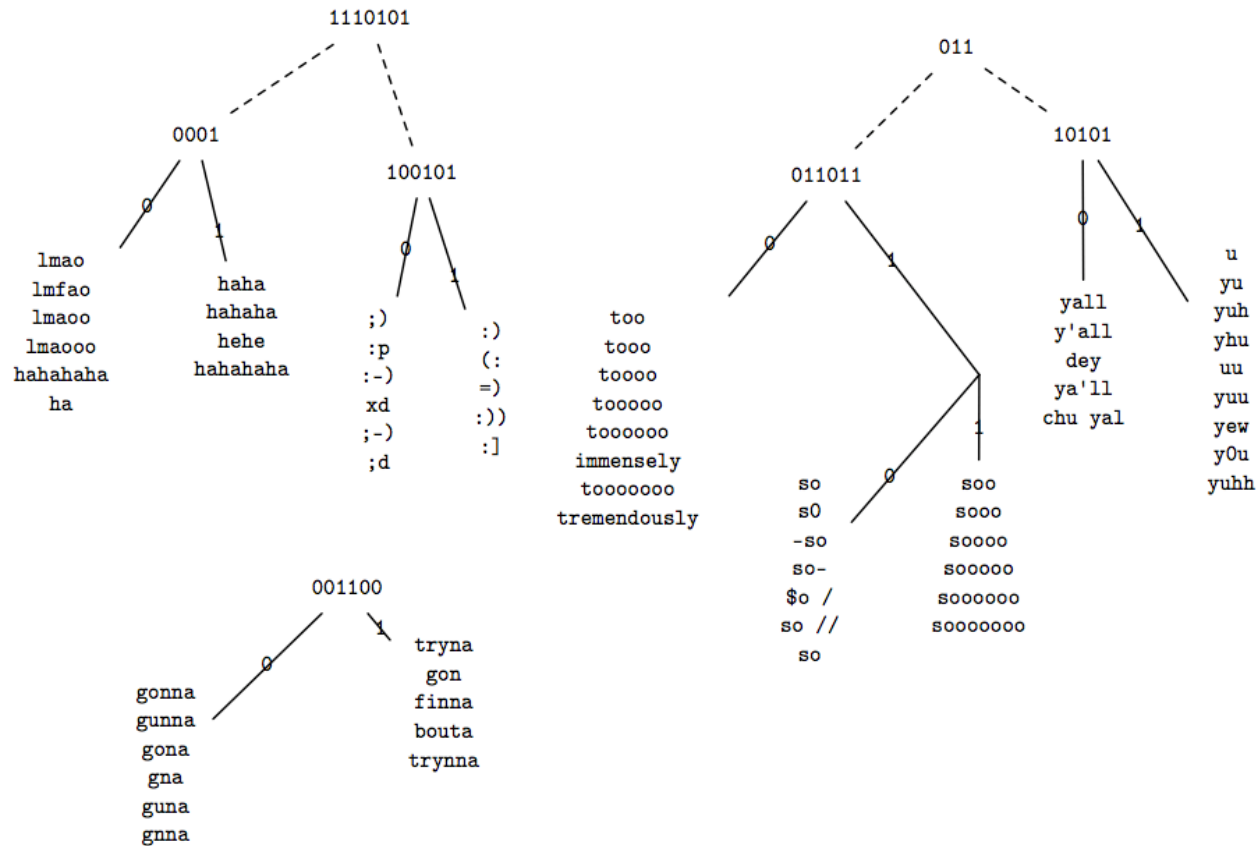
# Hierarchical Clustering



# Example



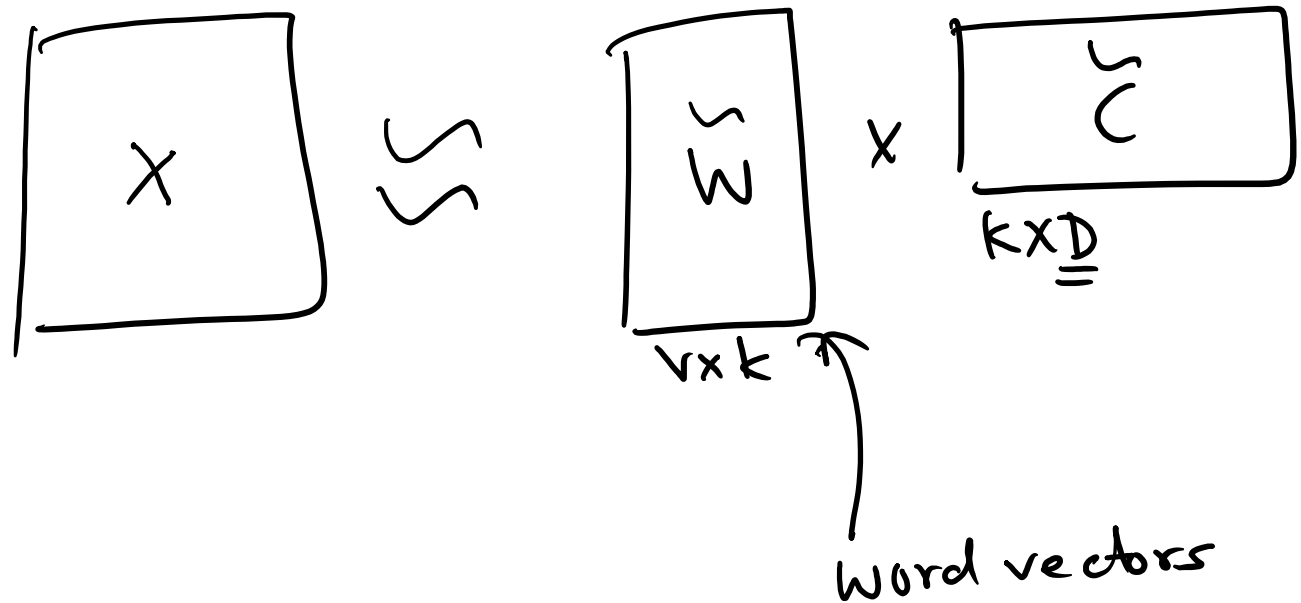
# Brown Clusters for Twitter



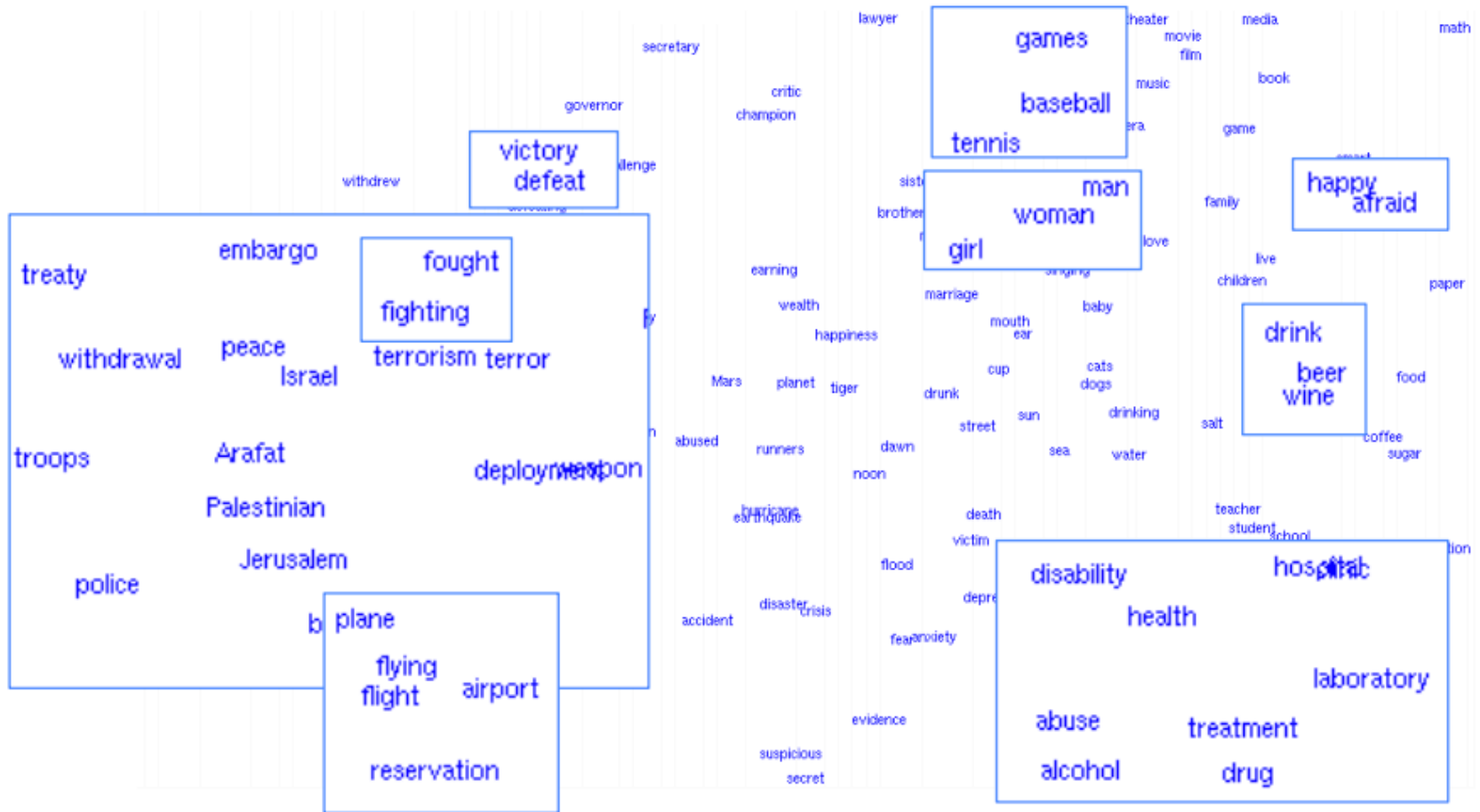
[http://www.cs.cmu.edu/~ark/TweetNLP/cluster\\_viewer.html](http://www.cs.cmu.edu/~ark/TweetNLP/cluster_viewer.html)

# Option 2: SVD

---



# Example Word Projection



# Problem with SVD & Clustering

## Computational Complexity

- SVD:  $O(mn^2)$
- Clustering:  $O(knm)$  per iteration, or  $O(n^3)$
- But,  $n$  can be 100,000!

## "One shot"

- Difficult to add new documents or words
- Cannot work with streaming data

$C \sim 100k$   
words  $\sim 50k$   
Truncated SVD

# Outline

---

Latent Semantic Analysis

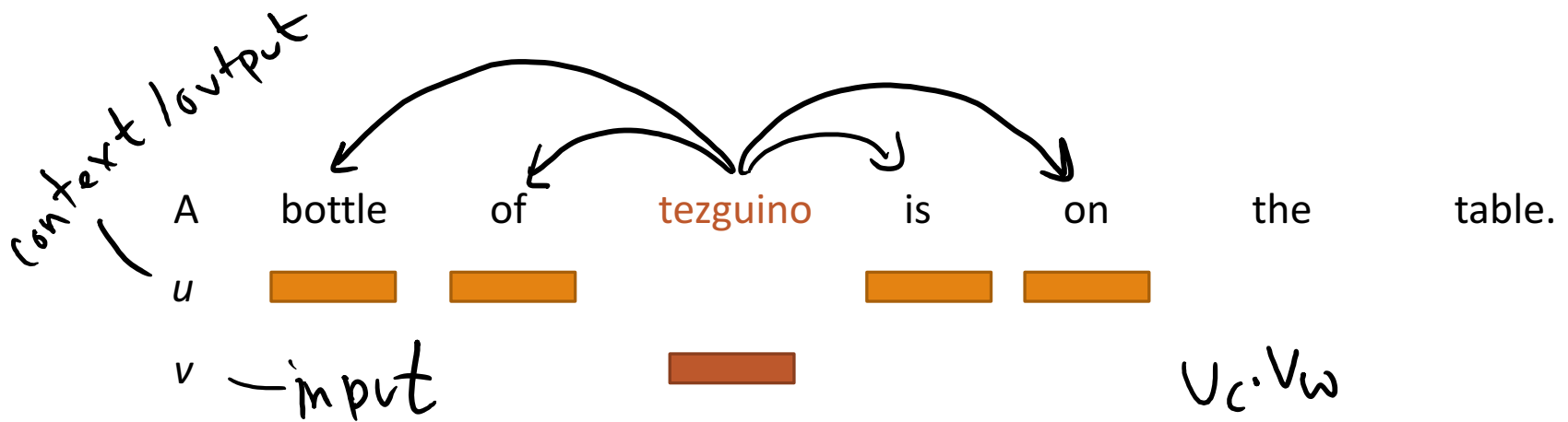
Vector Models for Words

Reducing the Dimensions

Direct Embeddings

# Predict surrounding words

$$\forall t \forall j, -m \dots m \quad j \neq 0 \quad \underbrace{P(w_{t+j} | w_t)}$$



$$P(c|w) = \frac{e^{u_c \cdot v_w}}{\sum_i e^{u_i \cdot v_w}}$$



# Estimating the Word Vectors

$$\operatorname{argmax}_{u,v} \underbrace{\sum_{t=1}^T \sum_{\substack{-m \dots m \\ j \neq 0}} \log P(w_{t+j} | w_t)}_{J(u,v)}$$

$$P(w|c) = \frac{e^{u_w \cdot v_c}}{\sum_{\text{all words}} e^{u_w \cdot v_c}}$$

$$\forall w, c \quad P(w|c) \neq 1$$

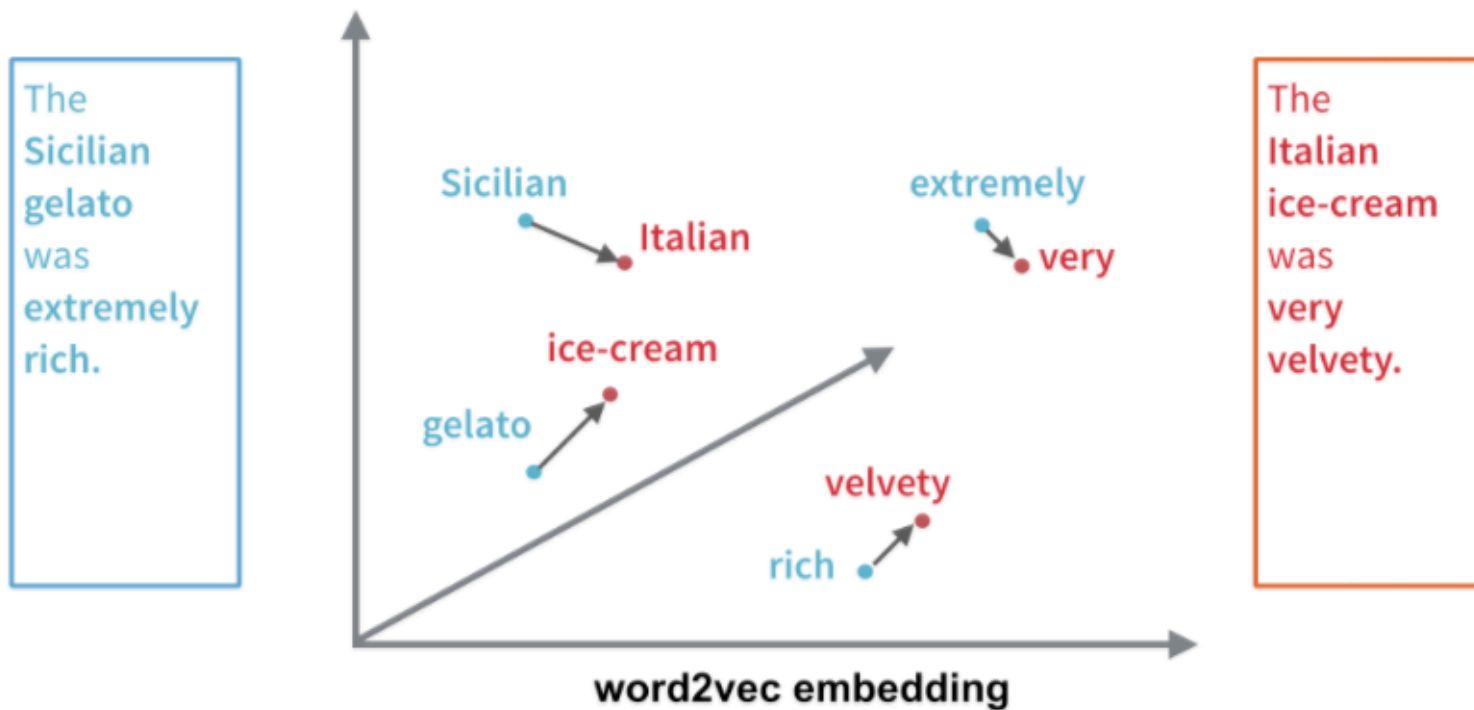
$$\sum_w P(w|c) = 1$$

# Similar Meaning = Close

---

Target Word	BoW5	BoW2	Target Word	BoW5	BoW2
batman	nightwing aquaman catwoman superman manhunter	superman superboy aquaman catwoman batgirl	florida	gainesville fla jacksonville tampa lauderdale	fla alabama gainesville tallahassee texas
hogwarts	dumbledore hallows half-blood malfoy snape	evernight sunnydale garderobe blandings collinwood	object-oriented	aspect-oriented smalltalk event-driven prolog domain-specific	aspect-oriented event-driven objective-c dataflow 4gl
turing	nondeterministic non-deterministic computability deterministic finite-state	non-deterministic finite-state nondeterministic buchi primality	dancing	singing dance dances dancers tap-dancing	singing dance dances breakdancing clowning

# Similar Meaning = Close



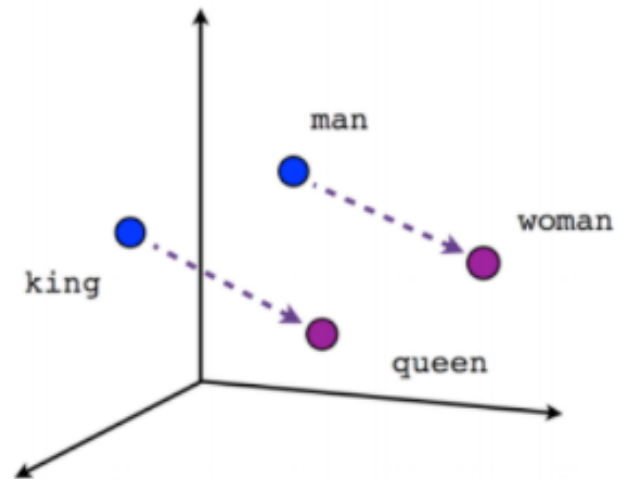
<https://siddhant7.github.io/Vector-Representation-of-Words/>

# Vectors “know” Gender

---

male : female :: King : **queen**

King - male + female **queen**



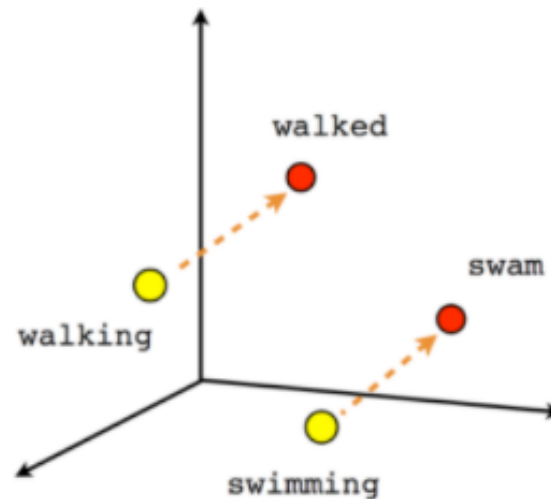
<https://siddhant7.github.io/Vector-Representation-of-Words/>

# They “know” Tenses!

---

walking : walked :: swimming : swam

swimming – walking + walked = swam



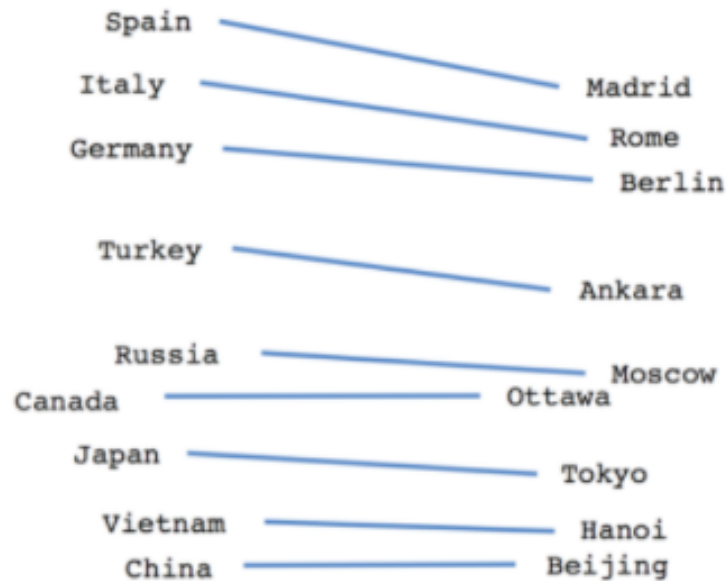
<https://siddhant7.github.io/Vector-Representation-of-Words/>

# They “know” Facts!

---

Country – Capital + Spain

Madrid



<https://siddhant7.github.io/Vector-Representation-of-Words/>

# Upcoming...

---

## Homework

- Homework 1 is up!
- No more material will be covered
- Due: **January 26, 2017**

## Project

- Project pitch is due **January 23, 2017!**
- Start assembling teams now
- Tons of datasets on the “projects” page on website