

Word Meanings

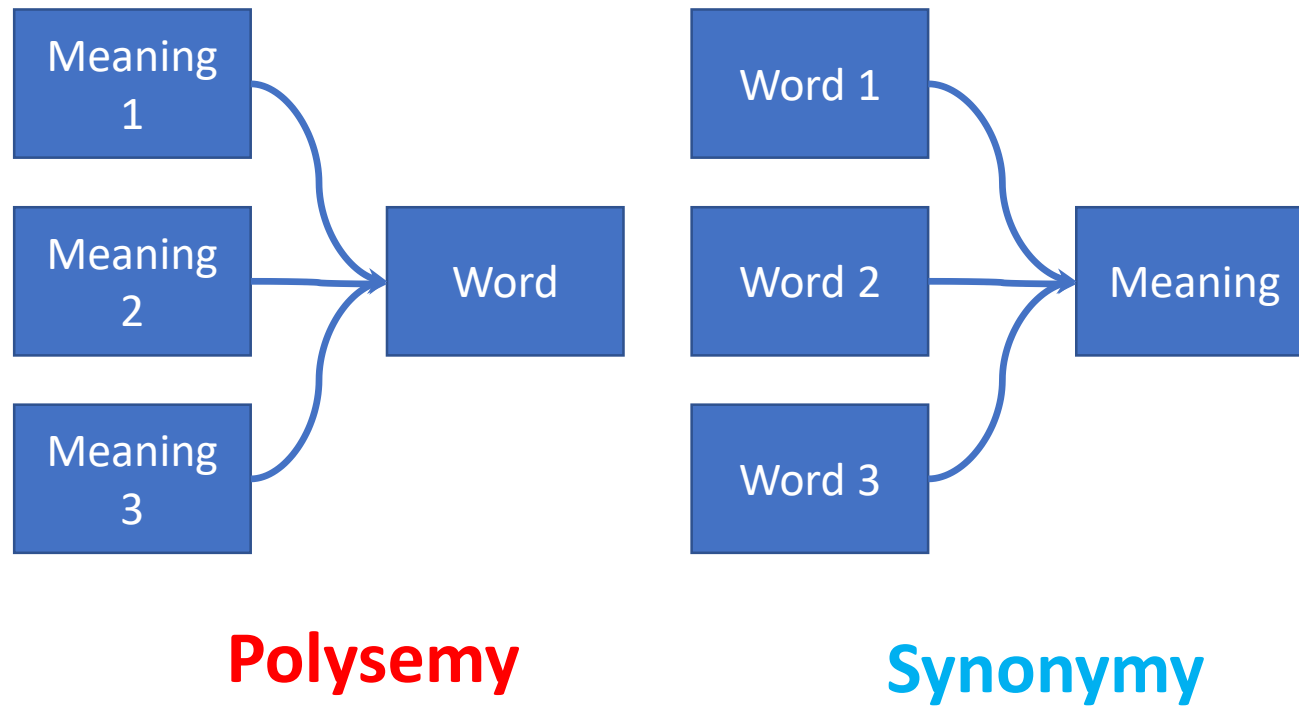
Latent Semantic Analytics

Dr. Uzair Ahmad

Program

- Word meanings and context
- The distributional hypothesis
- Design dimensions of word representations
- Latent semantic analysis
- Evaluation of representations

Word meanings: NLP Challenge



The distributional hypothesis

Acquire meaningful representations from unlabeled data

A bottle of _____ is on the table.

Everybody likes _____.

Don't have _____ before you drive.

We make _____ out of corn.

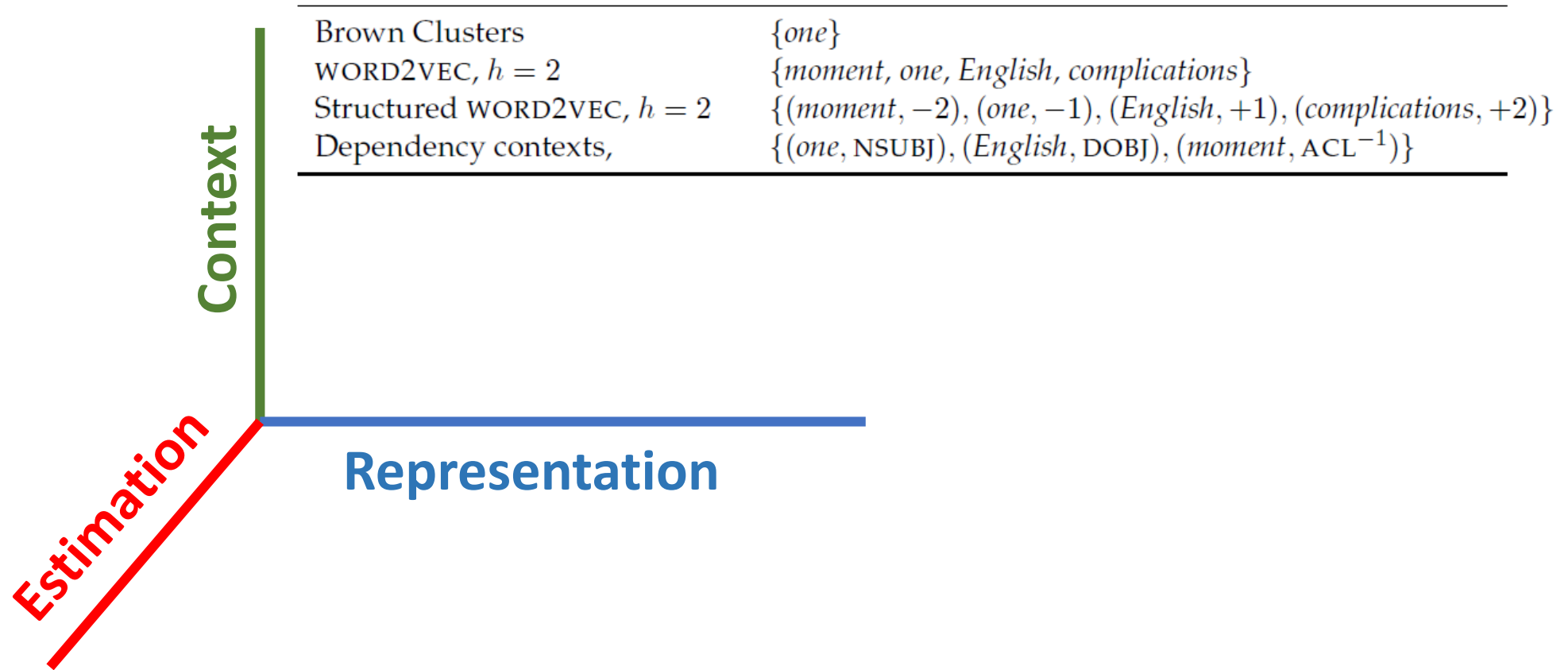
The distributional hypothesis

- (14.1) A bottle of ____ is on the table.
(14.2) Everybody likes _____.
(14.3) Don't have _____ before you drive.
(14.4) We make _____ out of corn.

	contextual properties				
	(14.1)	(14.2)	(14.3)	(14.4)	...
<i>tezgüino</i>	1	1	1	1	
<i>loud</i>	0	0	0	0	
<i>motor oil</i>	1	0	0	1	
<i>tortillas</i>	0	1	0	1	
<i>choices</i>	0	1	0	0	
<i>wine</i>	1	1	1	0	

YOU SHALL KNOW A WORD BY THE COMPANY IT KEEPS.
(FIRTH 1957)

Word representations



Latent Semantic Analysis

Contexts |C|

		c1	c2	c3	c4	c5	m1	m2	m3	m4
V Words	<i>human</i>	1	0	0	1	0	0	0	0	0
	<i>interface</i>	1	0	1	0	0	0	0	0	0
	<i>computer</i>	1	1	0	0	0	0	0	0	0
	<i>user</i>	0	1	1	0	1	0	0	0	0
	<i>system</i>	0	1	1	2	0	0	0	0	0
	<i>response</i>	0	1	0	0	1	0	0	0	0
	<i>time</i>	0	1	0	0	1	0	0	0	0
	<i>EPS</i>	0	0	1	1	0	0	0	0	0
	<i>survey</i>	0	1	0	0	0	0	0	0	1
	<i>trees</i>	0	0	0	0	0	1	1	1	0
	<i>graph</i>	0	0	0	0	0	0	1	1	1
	<i>minors</i>	0	0	0	0	0	0	0	1	1

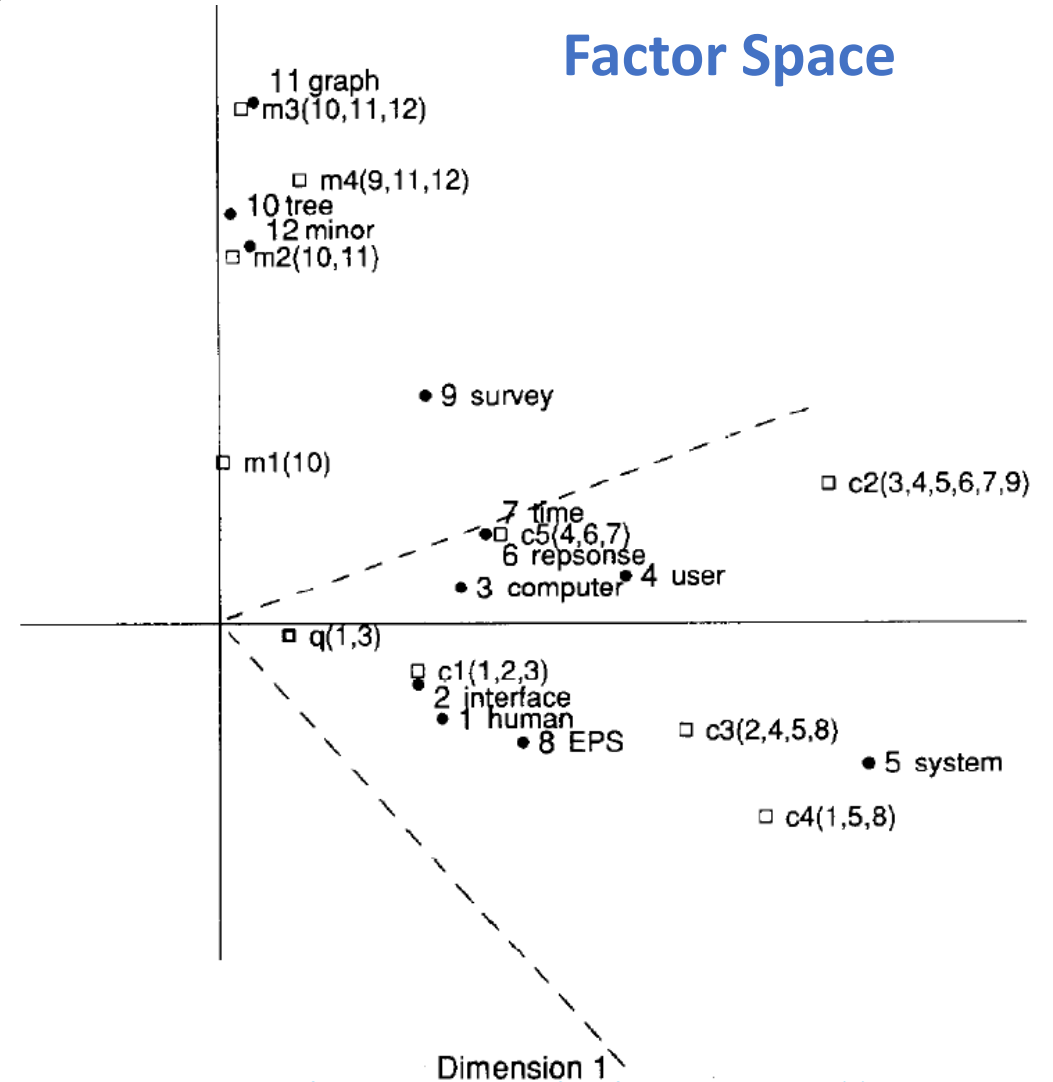
**Term-Document
Matrix**

Latent Semantic Analysis

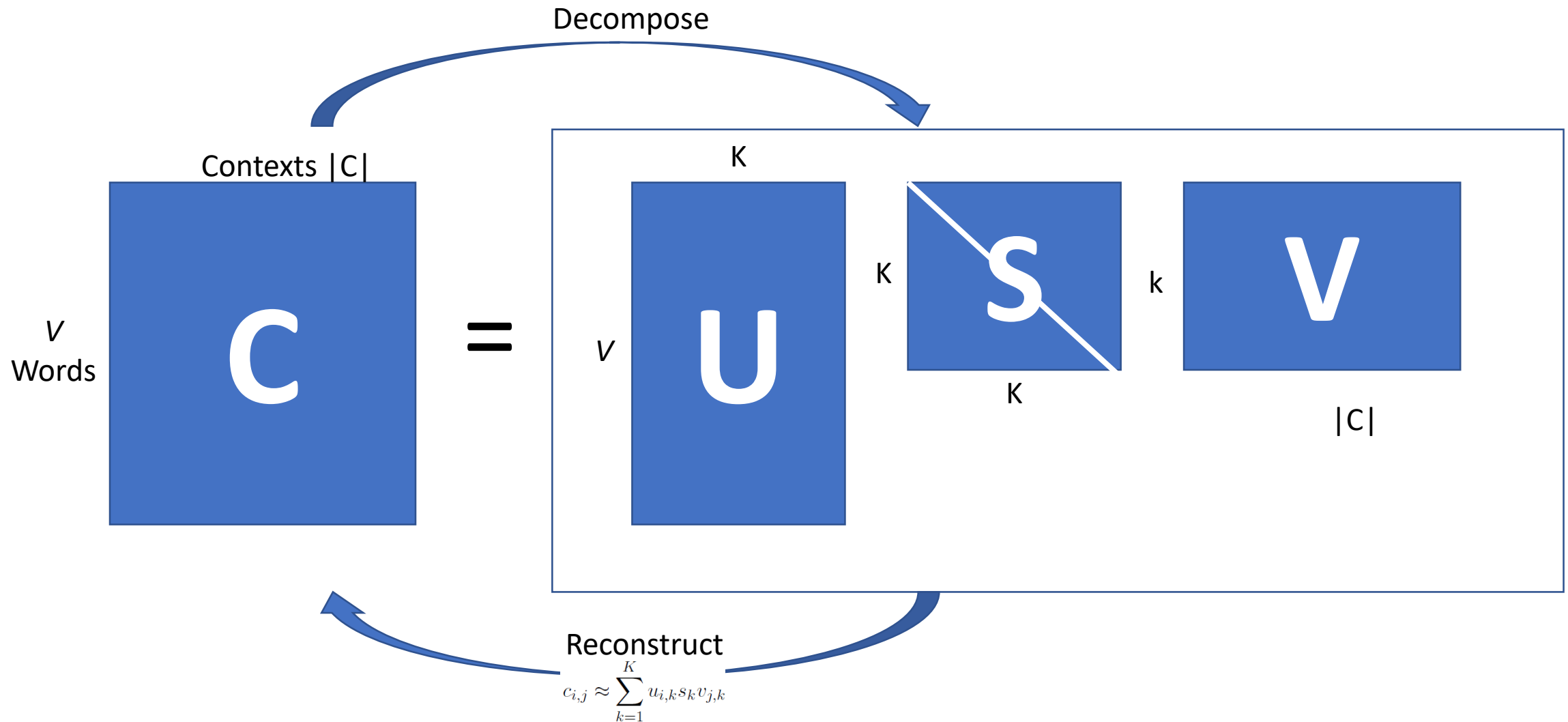
		c1	c2	c3	c4	c5	m1	m2	m3	m4
1	human	1	0	0	1	0	0	0	0	0
2	interface	1	0	1	0	0	0	0	0	0
3	computer	1	1	0	0	0	0	0	0	0
4	user	0	1	1	0	1	0	0	0	0
5	system	0	1	1	2	0	0	0	0	0
6	response	0	1	0	0	1	0	0	0	0
7	time	0	1	0	0	1	0	0	0	0
8	EPS	0	0	1	1	0	0	0	0	0
9	survey	0	1	0	0	0	0	0	0	1
10	trees	0	0	0	0	0	1	1	1	0
11	graph	0	0	0	0	0	0	1	1	1
12	minors	0	0	0	0	0	0	0	1	1

Term-Document
Matrix

Dimension 2



Latent Semantic Analysis



Latent Semantic Analysis

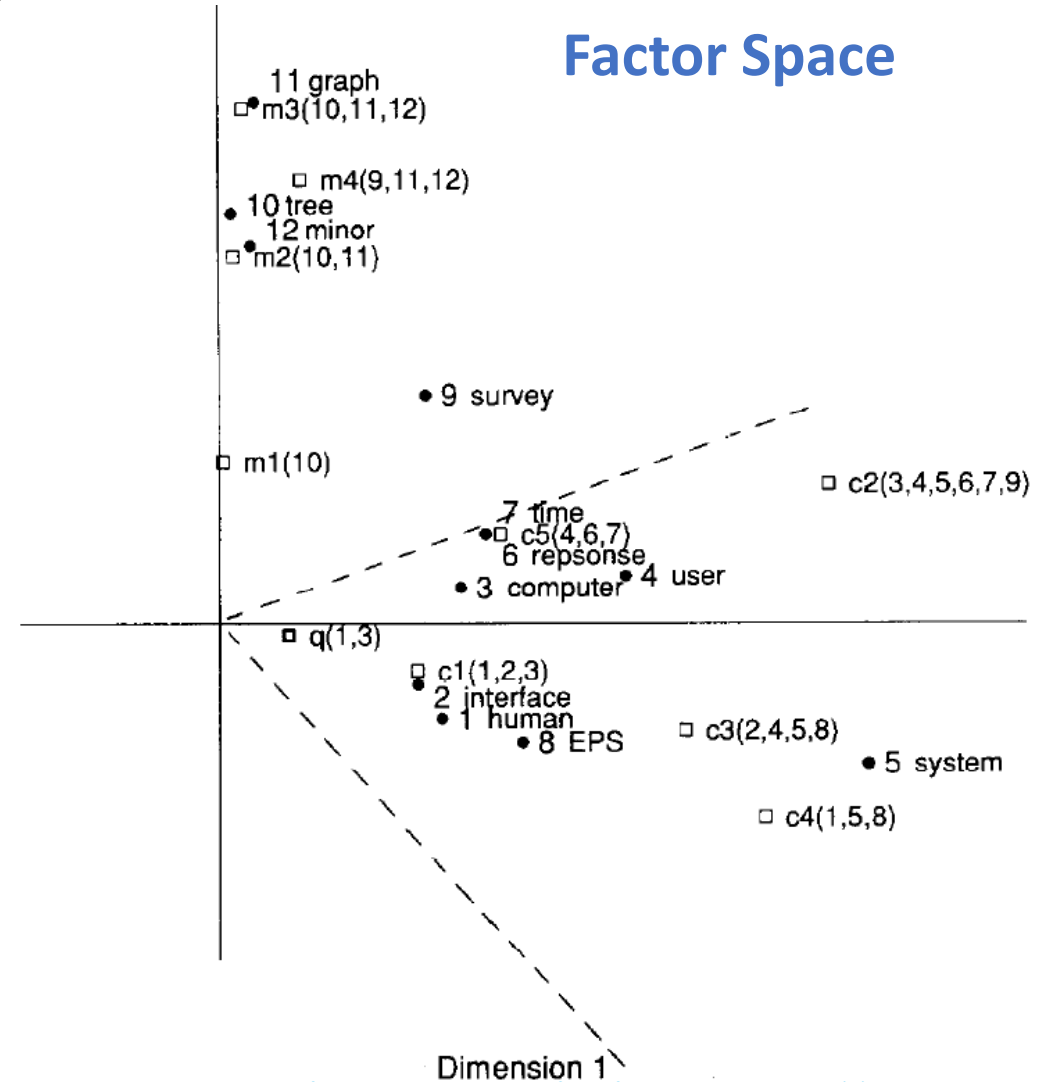
$$\begin{aligned} \min_{\mathbf{U} \in \mathbb{R}^{V \times K}, \mathbf{S} \in \mathbb{R}^{K \times K}, \mathbf{V} \in \mathbb{R}^{|\mathcal{C}| \times K}} \quad & ||\mathbf{C} - \mathbf{USV}^\top||_F \\ \text{s.t.} \quad & \mathbf{U}^\top \mathbf{U} = \mathbb{I} \\ & \mathbf{V}^\top \mathbf{V} = \mathbb{I} \\ & \forall i \neq j, \mathbf{S}_{i,j} = 0, \end{aligned}$$

Latent Semantic Analysis

		c1	c2	c3	c4	c5	m1	m2	m3	m4
1	human	1	0	0	1	0	0	0	0	0
2	interface	1	0	1	0	0	0	0	0	0
3	computer	1	1	0	0	0	0	0	0	0
4	user	0	1	1	0	1	0	0	0	0
5	system	0	1	1	2	0	0	0	0	0
6	response	0	1	0	0	1	0	0	0	0
7	time	0	1	0	0	1	0	0	0	0
8	EPS	0	0	1	1	0	0	0	0	0
9	survey	0	1	0	0	0	0	0	0	1
10	trees	0	0	0	0	0	1	1	1	0
11	graph	0	0	0	0	0	0	1	1	1
12	minors	0	0	0	0	0	0	0	1	1

Term-Document Count Matrix

Dimension 2



Latent Semantic Analysis

- Transforming the count matrix **C**
 - Pointwise mutual information

$$\text{PMI}(i, j) = \log \frac{p(i, j)}{p(i)p(j)} = \log \frac{p(i | j)p(j)}{p(i)p(j)} = \log \frac{p(i | j)}{p(i)}$$

Latent Semantic Analysis

- Transforming the count matrix **C**
 - Positive Pointwise mutual information

$$\text{PPMI}(i, j) = \begin{cases} \text{PMI}(i, j), & p(i | j) > p(i) \\ 0, & \text{otherwise.} \end{cases}$$

PMI

	computer	data	result	pie	sugar	count(w)
cherry	2	8	9	442	25	486
strawberry	0	0	1	60	19	80
digital	1670	1683	85	5	4	3447
information	3325	3982	378	5	13	7703
count(context)	4997	5673	473	512	61	11716

$$P(w=\text{information}, c=\text{data}) = \frac{3982}{11716} = .3399$$

$$P(w=\text{information}) = \frac{7703}{11716} = .6575$$

$$P(c=\text{data}) = \frac{5673}{11716} = .4842$$

$$\text{ppmi}(\text{information}, \text{data}) = \log_2(.3399 / (.6575 * .4842)) = .0944$$

Summary

- Distributed representations
 - Latent Semantic Analysis
 - Information retrieval
- Evaluation methods
 - Intrinsic
 - Extrinsic