

# CS 273A: Machine Learning

Winter 2021

## Lecture 16: Latent-Space Models

(cont.)

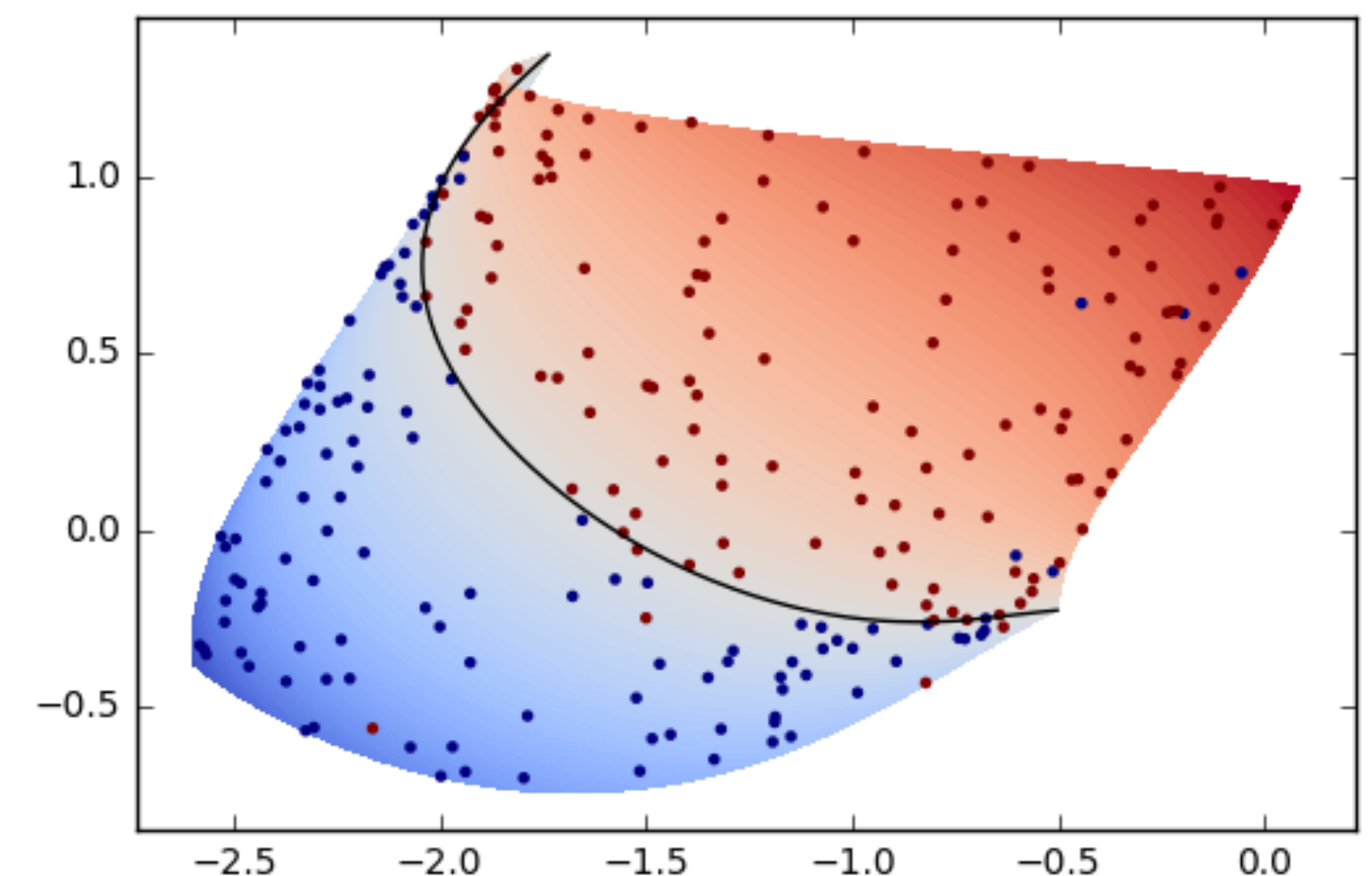
Roy Fox

Department of Computer Science

Bren School of Information and Computer Sciences

University of California, Irvine

All slides in this course adapted from Alex Ihler & Sameer Singh



# Logistics

---

assignments

- Assignment 5 **due Thursday**

project

- Final report **due next Thursday**

evaluations

- Evaluations **due end of next week**

final exam

- **Review:** next Thursday
- **Final:** Thursday, March 18, 1:30–3:30pm

# Today's lecture

---

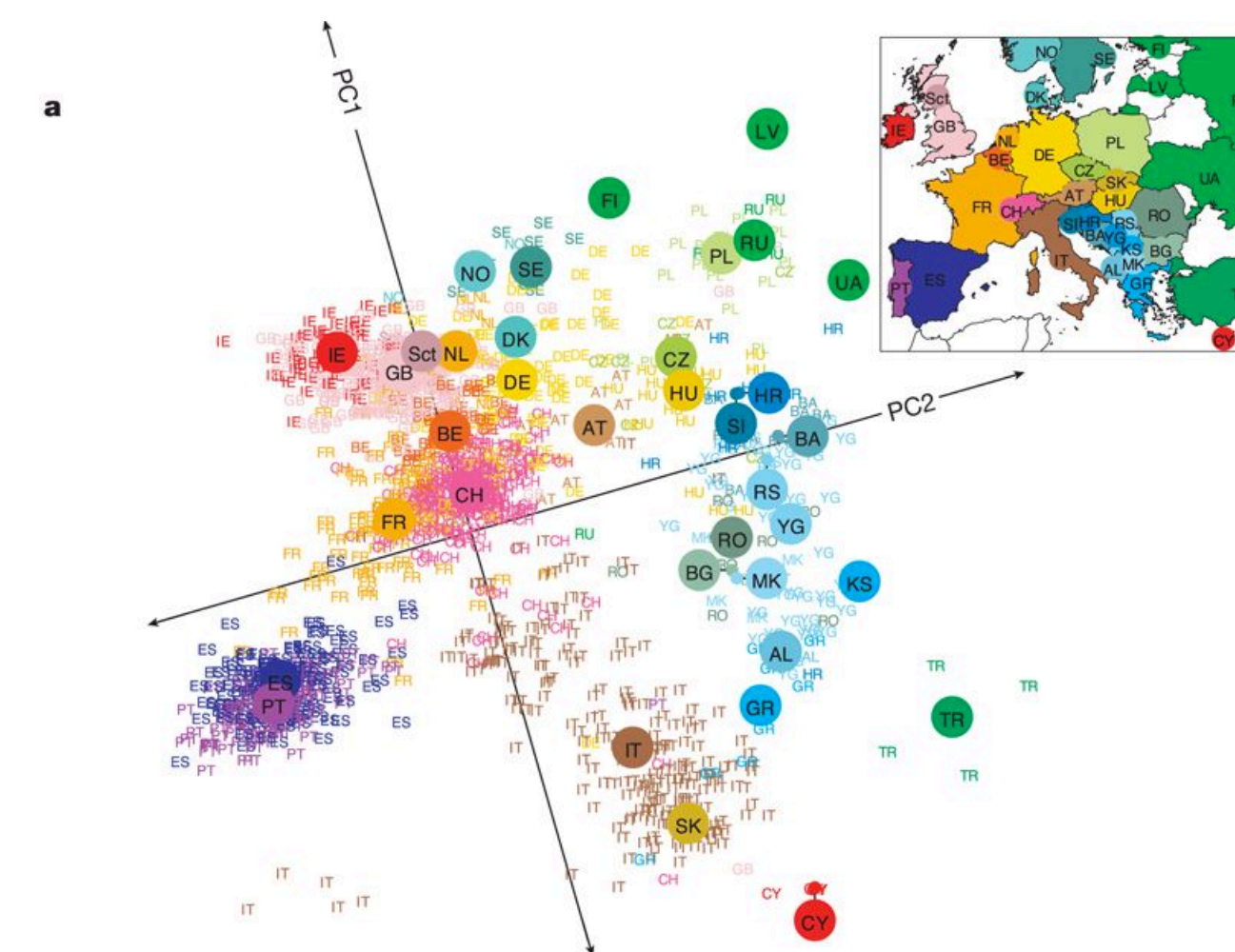
**Eigen-faces**

**Latent Semantic Analysis**

**Collaborative Filtering**

# Latent-space representations: uses

- **Remove** unneeded features
  - ▶ Features that add very little **information** (e.g. low **variability**, high **noise**)
  - ▶ Features that are **similar** to others (e.g. almost linearly dependent)
  - ▶ Reduce dimensionality for downstream application
    - **Supervised learning**: fewer parameters, need less data
    - **Compression**: less bandwidth
- Can also **add** features
  - ▶ **Summarize** multiple features into few cleaner / higher-level ones



# PCA: applications

---

- Eigen-faces
  - Represent image data (e.g. faces) using PCA
- Latent-Space Analysis (topic models)
  - Represent text data (e.g. bag of words) using PCA
- Collaborative Filtering for Recommendation Systems
  - Represent sentiment data (e.g. ratings) using PCA

# Singular Value Decomposition (SVD)

- Alternative method for finding covariance **eigenvectors**
  - Has many other uses
- **Singular Value Decomposition (SVD):**  $X = UDV^T$ 
  - $U$  and  $V$  (left- and right **singular vectors**) are orthogonal:  $U^T U = I$ ,  $V^T V = I$
  - $D$  (**singular values**) is rectangular-diagonal
  - $\Sigma = X^T X = VD^T U^T U D V^T = V(D^T D) V^T$
- $UD$  matrix gives **coefficients** to reconstruct data:  $x_i = U_{i,1} D_{1,1} v_1 + U_{i,2} D_{2,2} v_2 + \dots$ 
  - We can truncate this after **top  $k$  singular values** (square root of eigenvalues)

$$\begin{array}{|c|} \hline X \\ \hline m \times n \\ \hline \end{array} \approx \begin{array}{|c|} \hline U \\ \hline m \times k \\ \hline \end{array} \cdot \begin{array}{|c|} \hline D \\ \hline k \times k \\ \hline \end{array} \cdot \begin{array}{|c|} \hline V^T \\ \hline k \times n \\ \hline \end{array}$$

# Eigen-faces

- “Eigen- $X$ ” = represent  $X$  using its principal components
- Viola Jones dataset:  $24 \times 24$  images  $\in \mathbb{R}^{576}$ 
  - Can represent vector as image

$$\begin{array}{|c|} \hline X \\ \hline m \times n \\ \hline \end{array} \approx \begin{array}{|c|} \hline U \\ \hline m \times k \\ \hline \end{array} \cdot \begin{array}{|c|} \hline D \\ \hline k \times k \\ \hline \end{array} \cdot \begin{array}{|c|} \hline V^T \\ \hline k \times n \\ \hline \end{array}$$



⋮

⋮

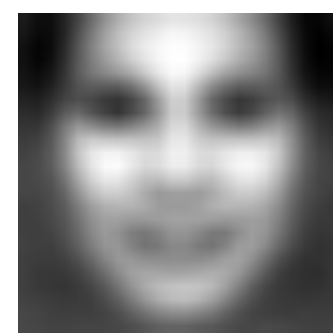
# Eigen-faces

- “Eigen- $X$ ” = represent  $X$  using its principal components

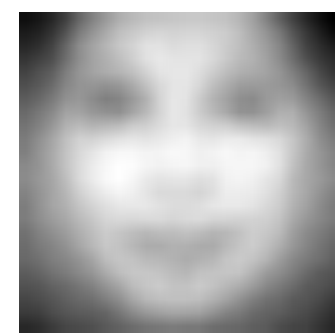
$$\begin{matrix} X \\ m \times n \end{matrix} \approx \begin{matrix} U \\ m \times k \end{matrix} \cdot \begin{matrix} D \\ k \times k \end{matrix} \cdot \begin{matrix} V^T \\ k \times n \end{matrix}$$

- Viola Jones dataset:  $24 \times 24$  images  $\in \mathbb{R}^{576}$

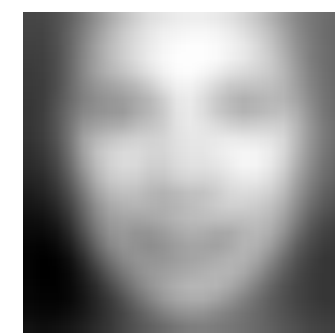
- ▶ Can represent vector as image



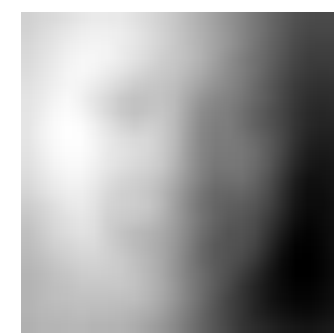
mean



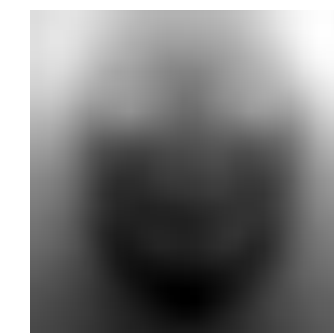
$v_1$



$v_2$



$v_3$



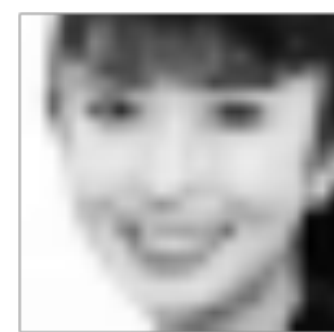
$v_4$

...

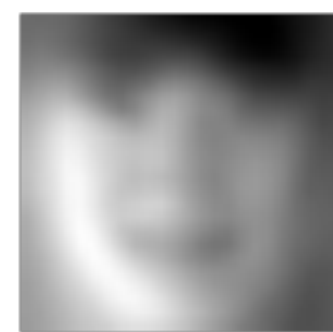
somewhat interpretable

- ▶ Project data on  $k$

principal components



$x_i$



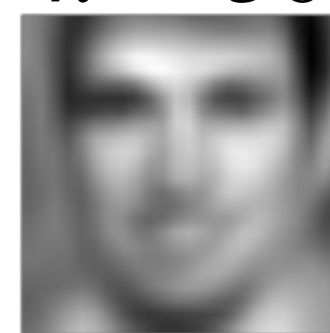
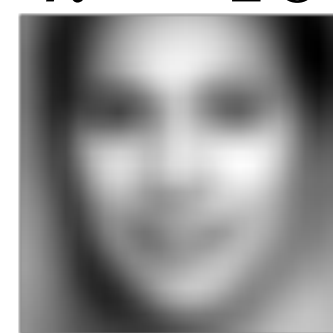
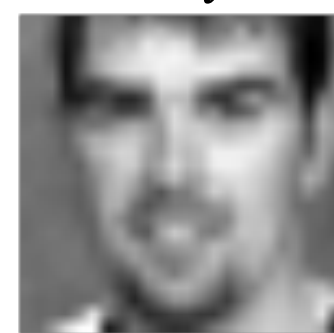
$k = 5$



$k = 10$



$k = 50$





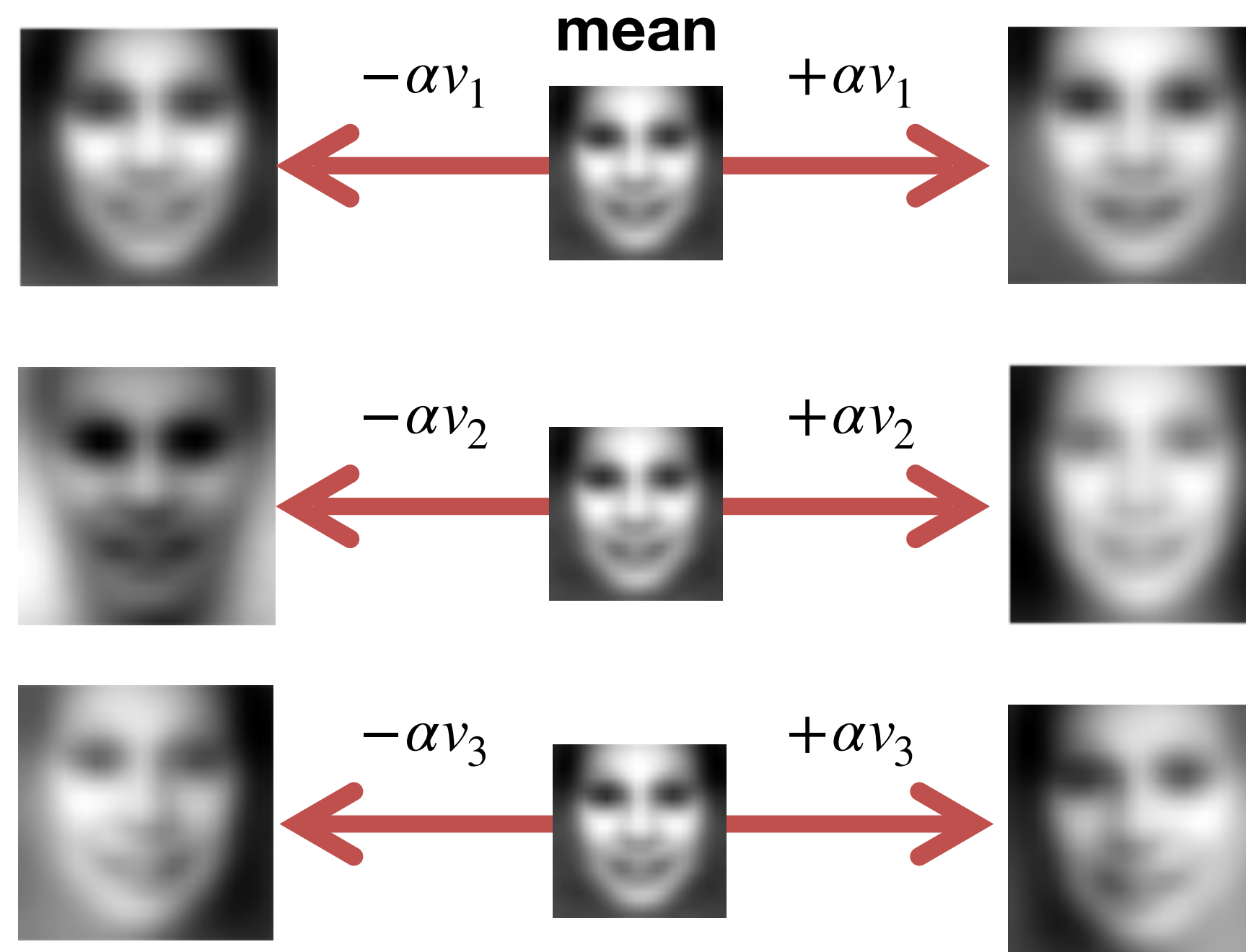
# Eigen-faces

- “Eigen- $X$ ” = represent  $X$  using its principal components

$$\begin{matrix} X \\ m \times n \end{matrix} \approx \begin{matrix} U \\ m \times k \end{matrix} \cdot \begin{matrix} D \\ k \times k \end{matrix} \cdot \begin{matrix} V^T \\ k \times n \end{matrix}$$

- Viola Jones dataset:  $24 \times 24$  images  $\in \mathbb{R}^{576}$

- ▶ Can represent vector as image



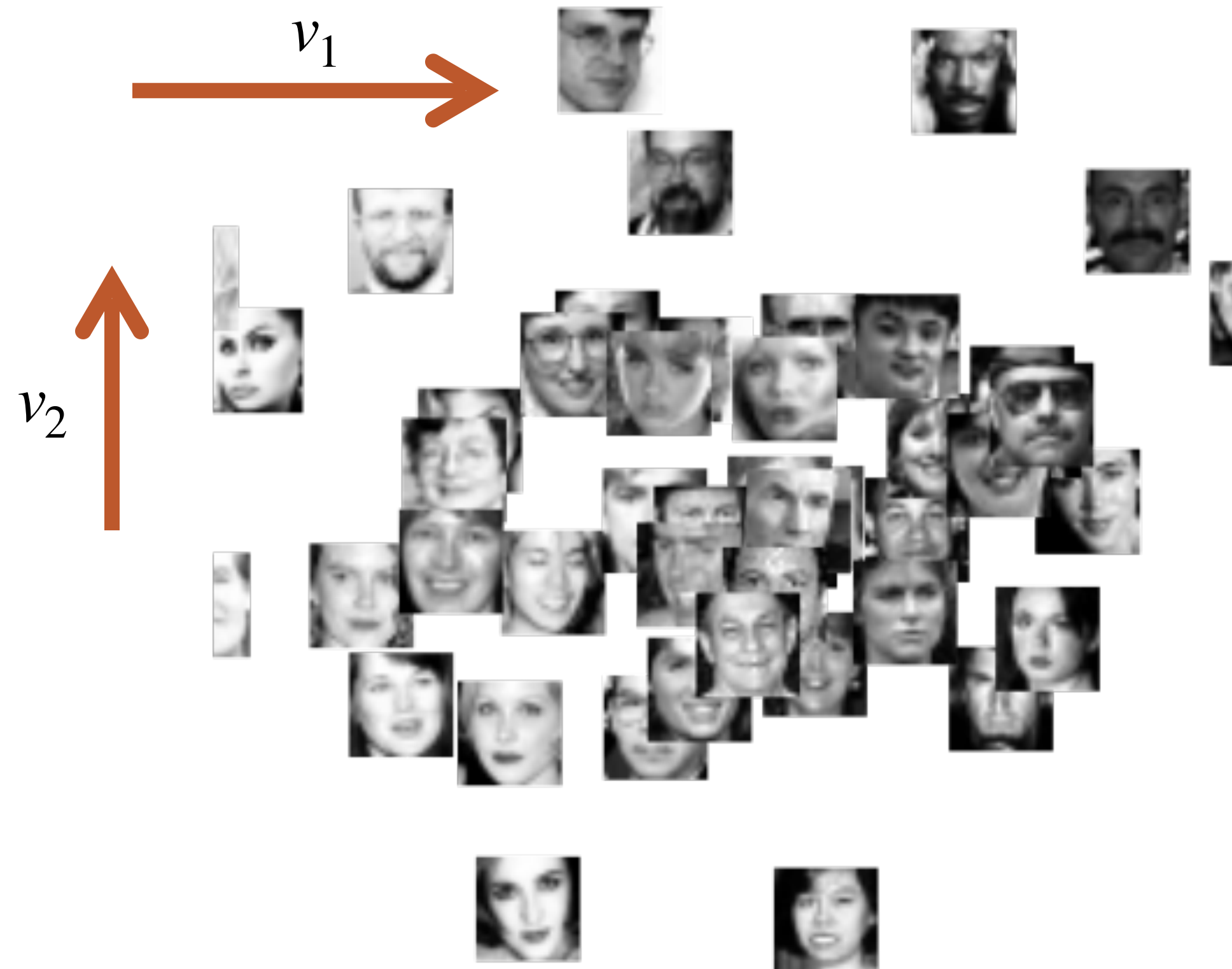
- ▶ Visualize basis vectors  $v_i$

as  $\mu \pm \alpha v_i$

# Eigen-faces

- “Eigen- $X$ ” = represent  $X$  using its principal components
- Viola Jones dataset:  $24 \times 24$  images  $\in \mathbb{R}^{576}$

- ▶ Can represent vector as image

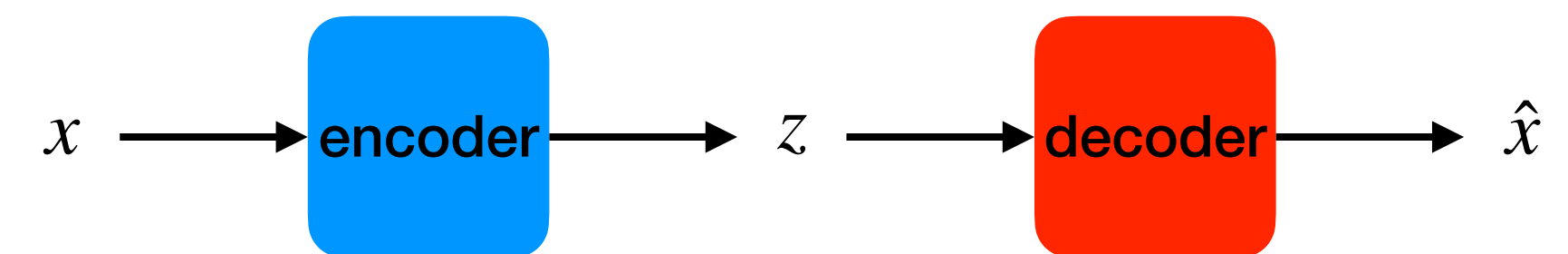


- ▶ Visualize data by projecting onto 2 principal components

# Nonlinear latent spaces

- Latent-space **representation** = represent  $x_i$  as  $z_i$ 
  - Usually more **succinct**, less noisy
  - Preserves most (interesting) information on  $x_i \implies$  can **reconstruct**  $\hat{x}_i \approx x_i$

▸ **Auto-encoder** = encode  $x \rightarrow z$ , decode  $z \rightarrow \hat{x}$

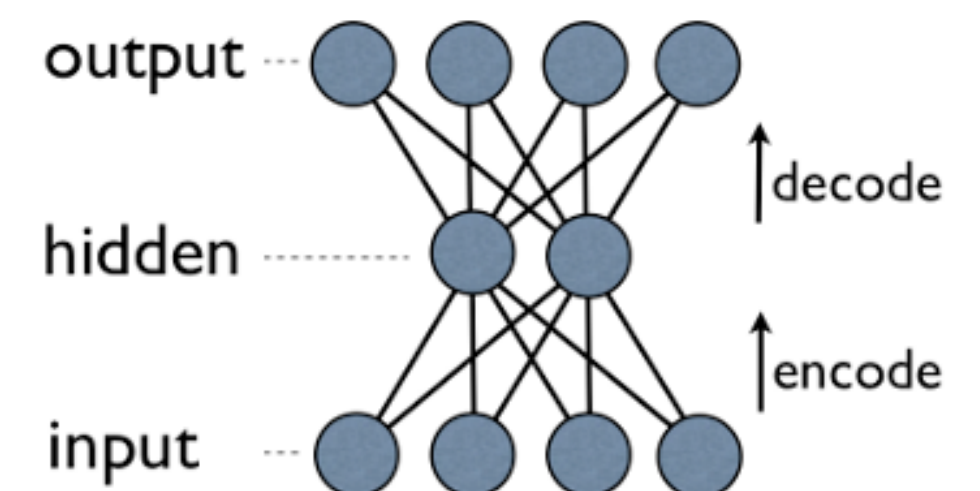


- **Linear** latent-space representation:

▸ **Encode:**  $Z = XV_{\leq k} = (UDV^T V)_{\leq k} = U_{\leq k} D_{\leq k}$ ; **Decode:**  $X \approx ZV_{\leq k}^T$

- **Nonlinear:** e.g., encoder + decoder are neural networks

▸ Restrict  $z$  to be shorter than  $x \implies$  requires succinctness



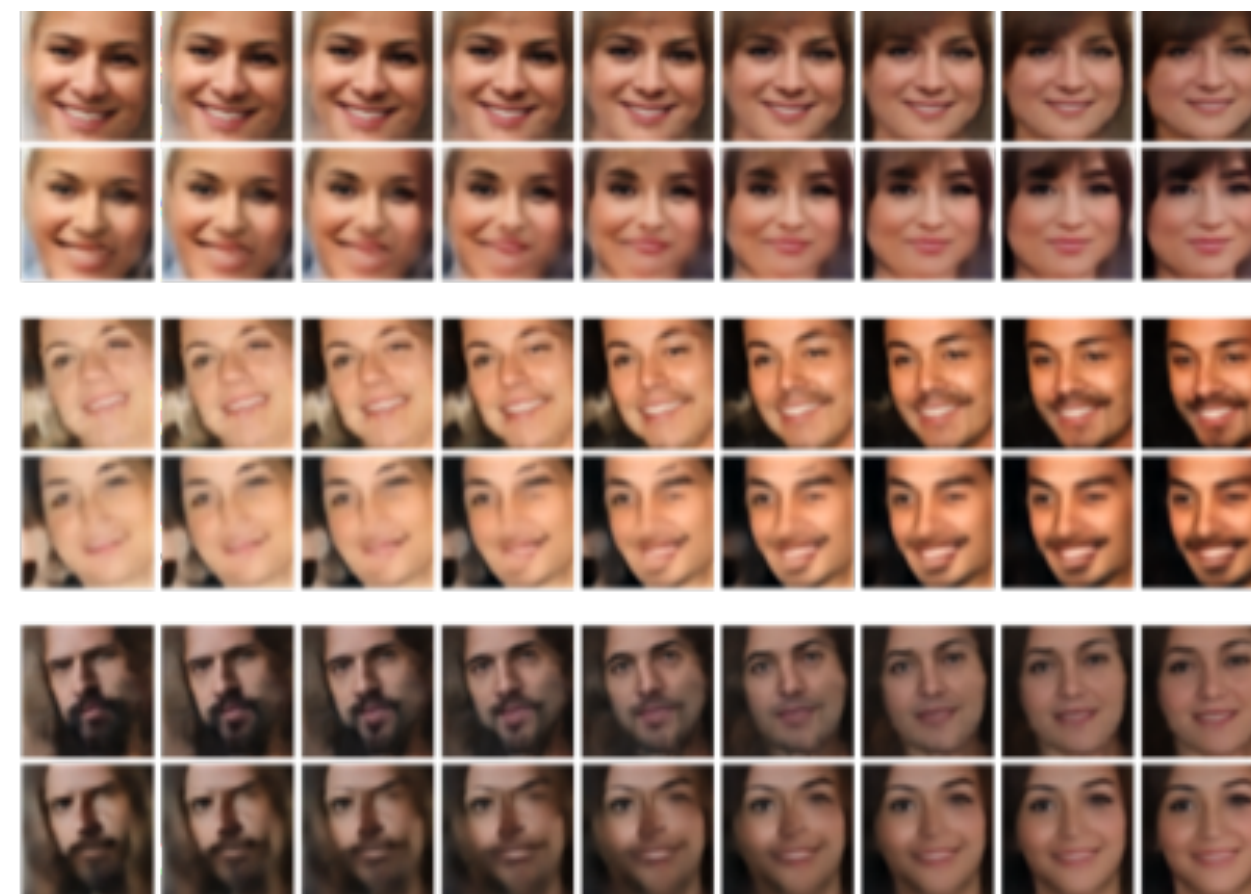
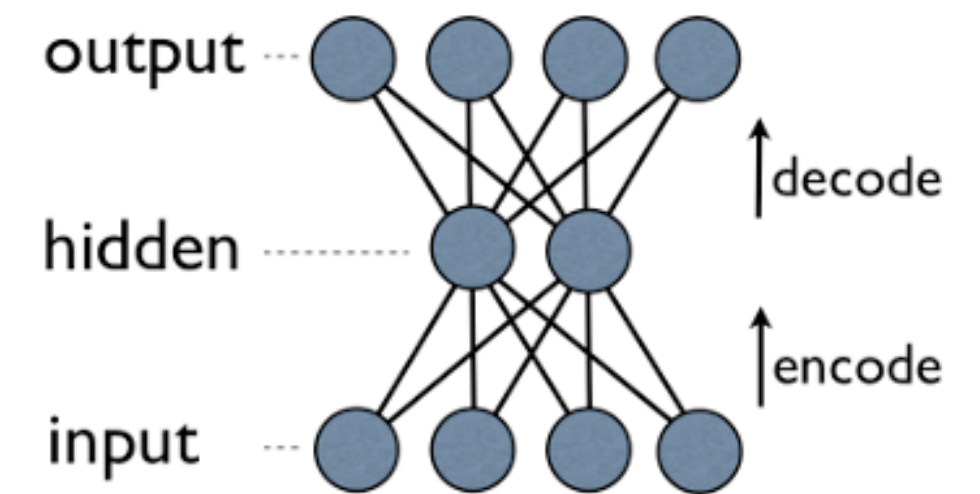
# Variational Auto-Encoders (VAE)

- Probabilistic model:

- ▶ Simple **prior** over latent space  $p(z)$  (e.g. Gaussian)

- ▶ **Decoder = generator**  $p_{\theta}(x | z)$ , tries to match **data** distribution  $p_{\theta}(x) \approx \mathcal{D}$

- ▶ **Encoder = inference**  $q_{\phi}(z | x)$ , tries to match **posterior**  $q_{\phi}(z | x) \approx \frac{p(z)p_{\theta}(x | z)}{p_{\theta}(x)}$



# Today's lecture

---

Eigen-faces

**Latent Semantic Analysis**

Collaborative Filtering

# Textual features

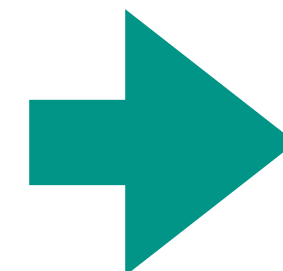
- How to **extract features** from text for model inputs?
  - Idea: **bag of words** = indicate word count, not order

Rain and chilly weather didn't keep thousands of paradegoers from camping out Friday night for the 111th Tournament of Roses.

Spirits were high among the street party crowd as they set up for curbside seats for today's parade.

"I want to party all night," said Tyne Gaudielle, 15, of Glendale, who spent the last night of the year along Colorado Boulevard with a group of friends.

Whether they came for the partying or the parade, campers were in for a long night. Rain continued into the evening and temperatures were expected to dip down into the low 40s.



## Observed Data (text docs):

DOC #	WORD #	COUNT
1	29	1
1	56	1
1	127	1
1	166	1
1	176	1
1	187	1
1	192	1
1	198	2
1	356	1
1	374	1
1	381	2

...

## VOCABULARY:

0001 ability  
0002 able  
0003 accept  
0004 accepted  
0005 according  
0006 account  
0007 accounts  
0008 accused  
0009 act  
0010 acting  
0011 action  
0012 active  
....

# Topic models

- Word distribution is different for different document **topics**
  - Represent the bag-of-words feature vector in a **latent space**
  - Such that representation **keep topic** information

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

**Human-Computer Interfaces**

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

**Graph Theory**

# Latent Semantic Analysis (LSA)

- Represent dataset as matrix of **word counts**:  $X_{ij}$  = # of word  $j$  in document  $i$ 
  - **Sparse matrix** = mostly 0s
  - Typically **normalize** rows:  $X_{ij}$  = probability that random word in doc  $j$  is word  $i$ 
    - Make documents of **different length** comparable
  - Typically don't shift or scale columns to have mean 0, var 1
- Perform **PCA** on  $X$  to get latent representation = top  $k$  components
  - Allows **fuzzy search** = documents of given topic, rather than word matching



# Word-doc matrix: example

- Observation: many words have **little overlap** between the topics

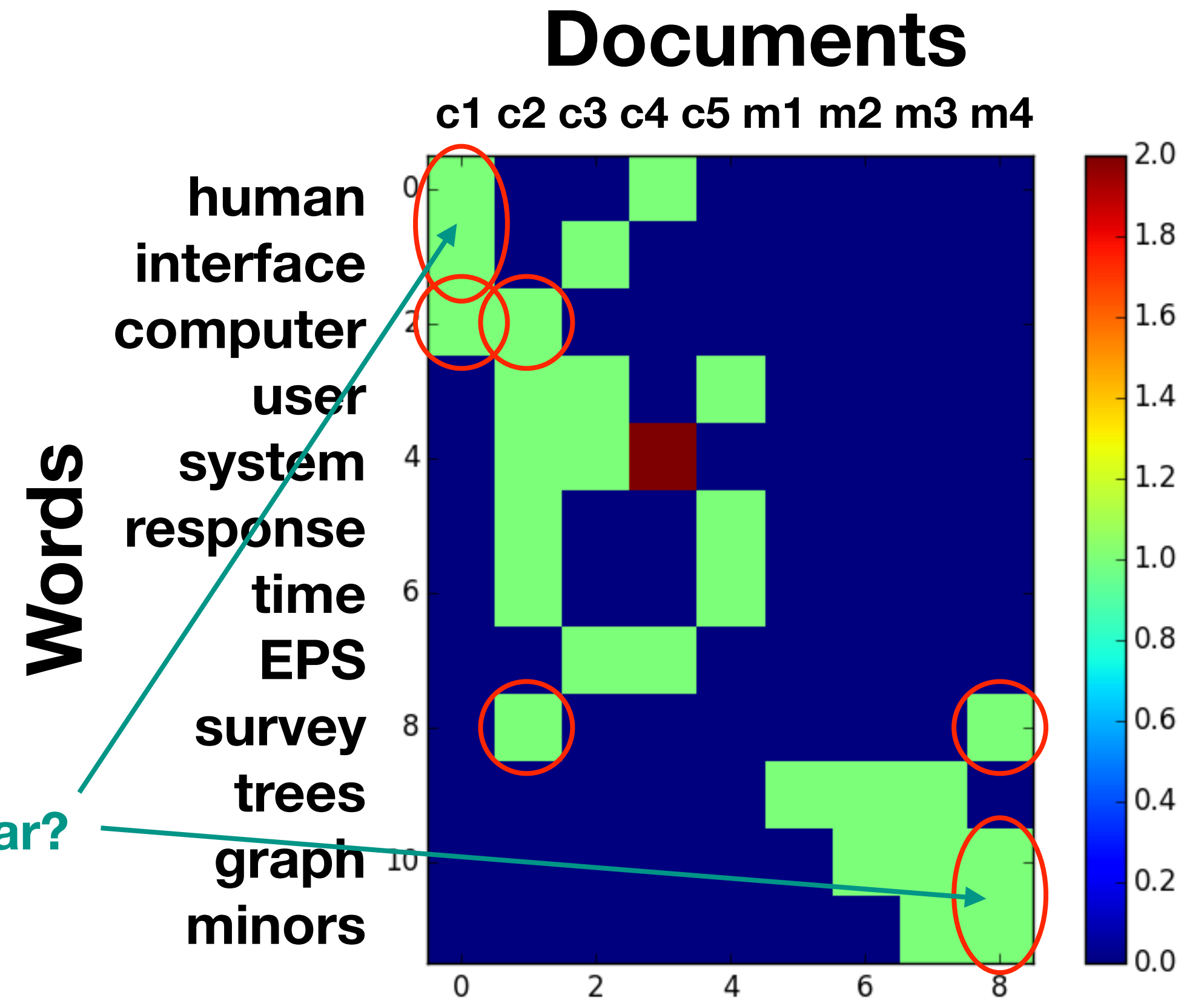
- Typical sizes:

- ▶ #docs =  $D \sim 10^6$

- ▶ #words in vocabulary =  $W \sim 10^5$

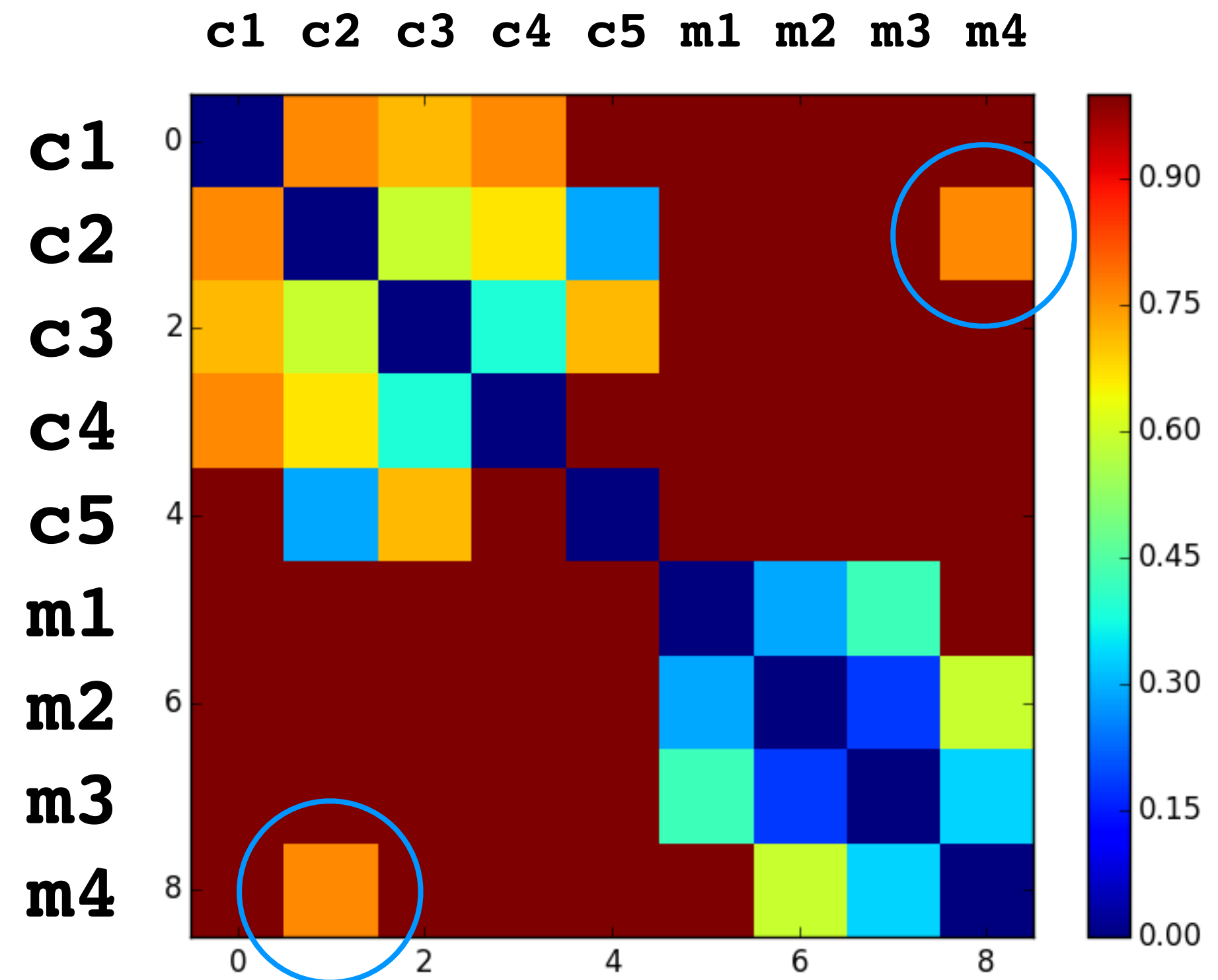
- ▶ Matrix size  $\sim 10^{11}$ , but only  $\sim 10^8$  in **sparse representation**

which is more similar?



# Doc dissimilarity matrix: example

- Define a **distance / dissimilarity** measure between docs  $D_{ij} = d(x_{i,.}, x_{j,.})$
- Clustering can be tricky
  - ▶ High dissimilarity **within topic**
  - ▶ Some **cross-topic** similarity



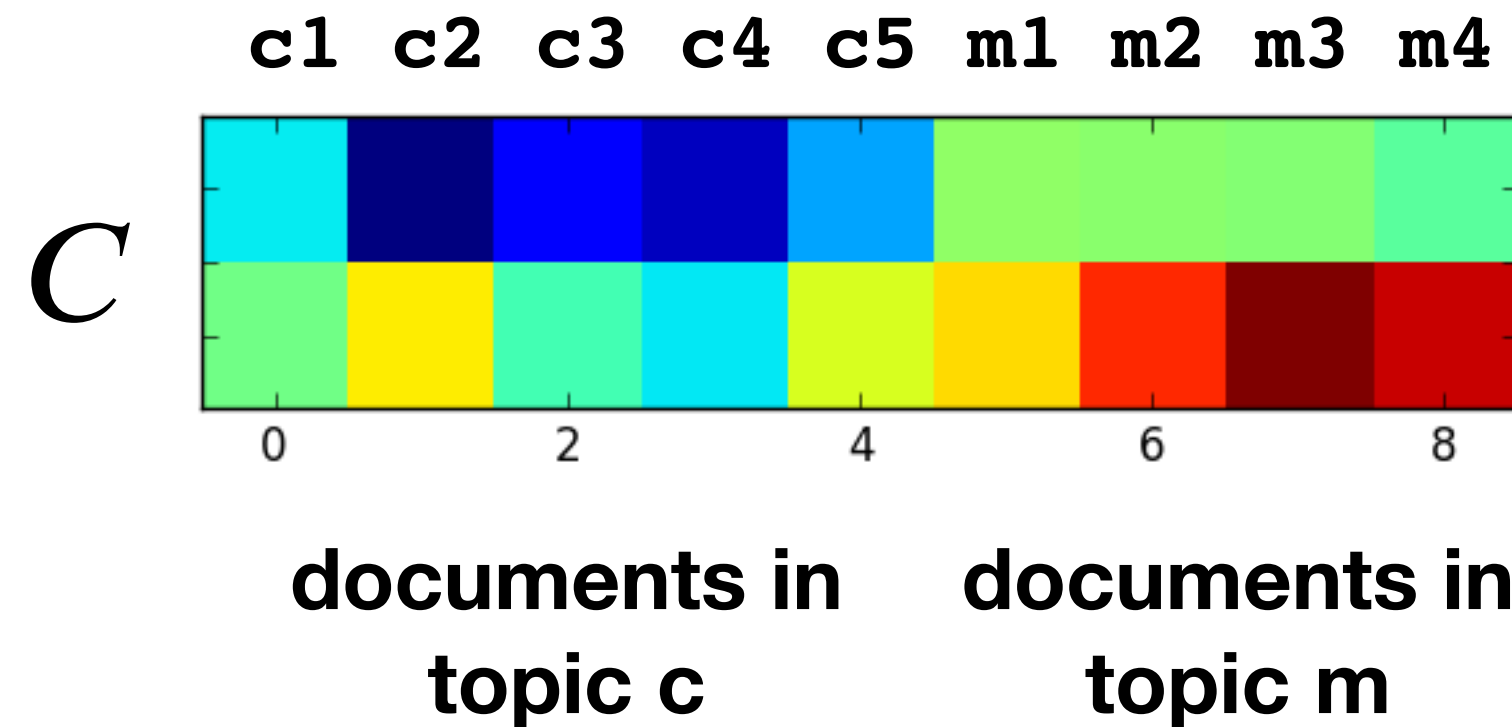
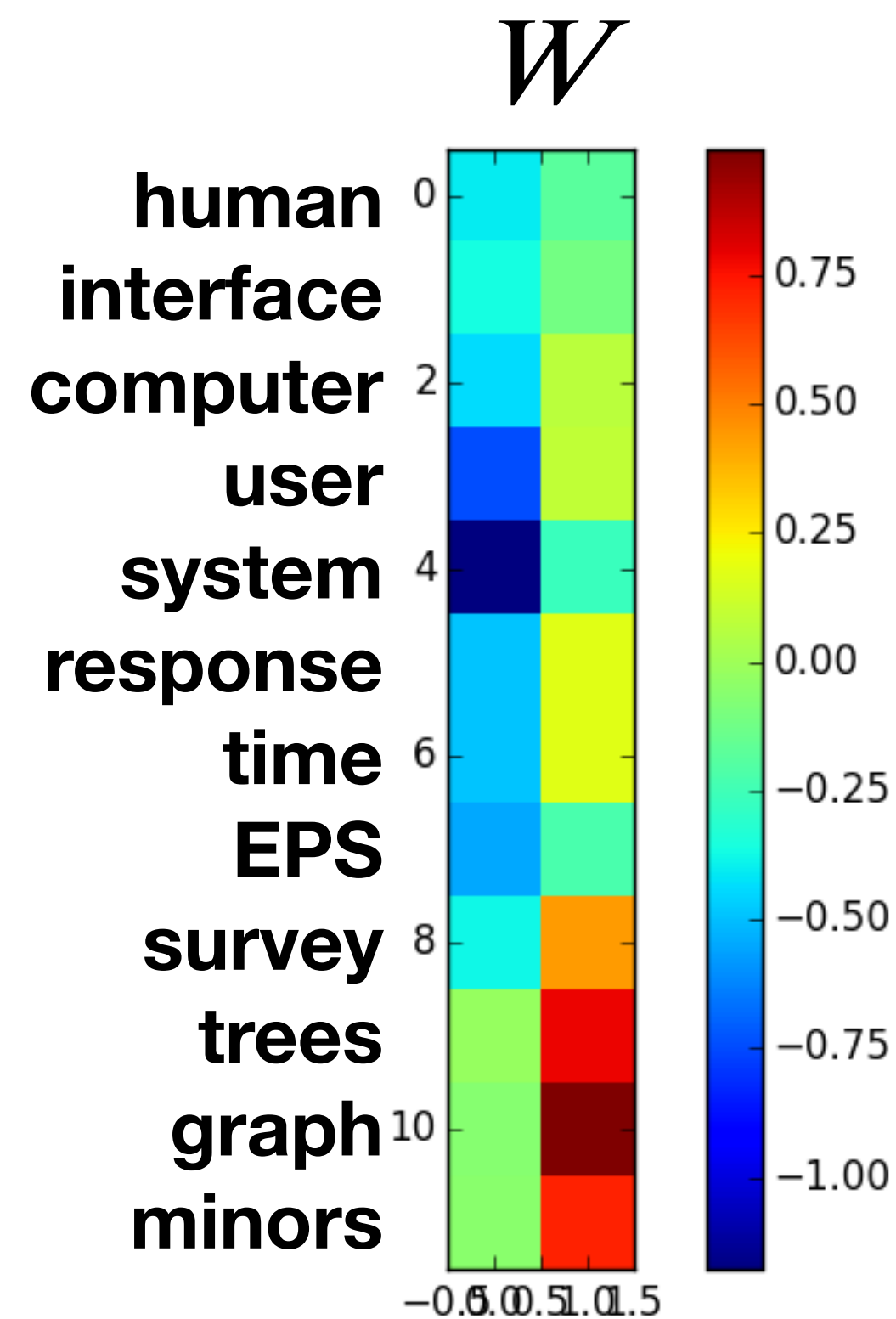
# Singular Value Decomposition (SVD)

- Useful to preserve **feature scale**:

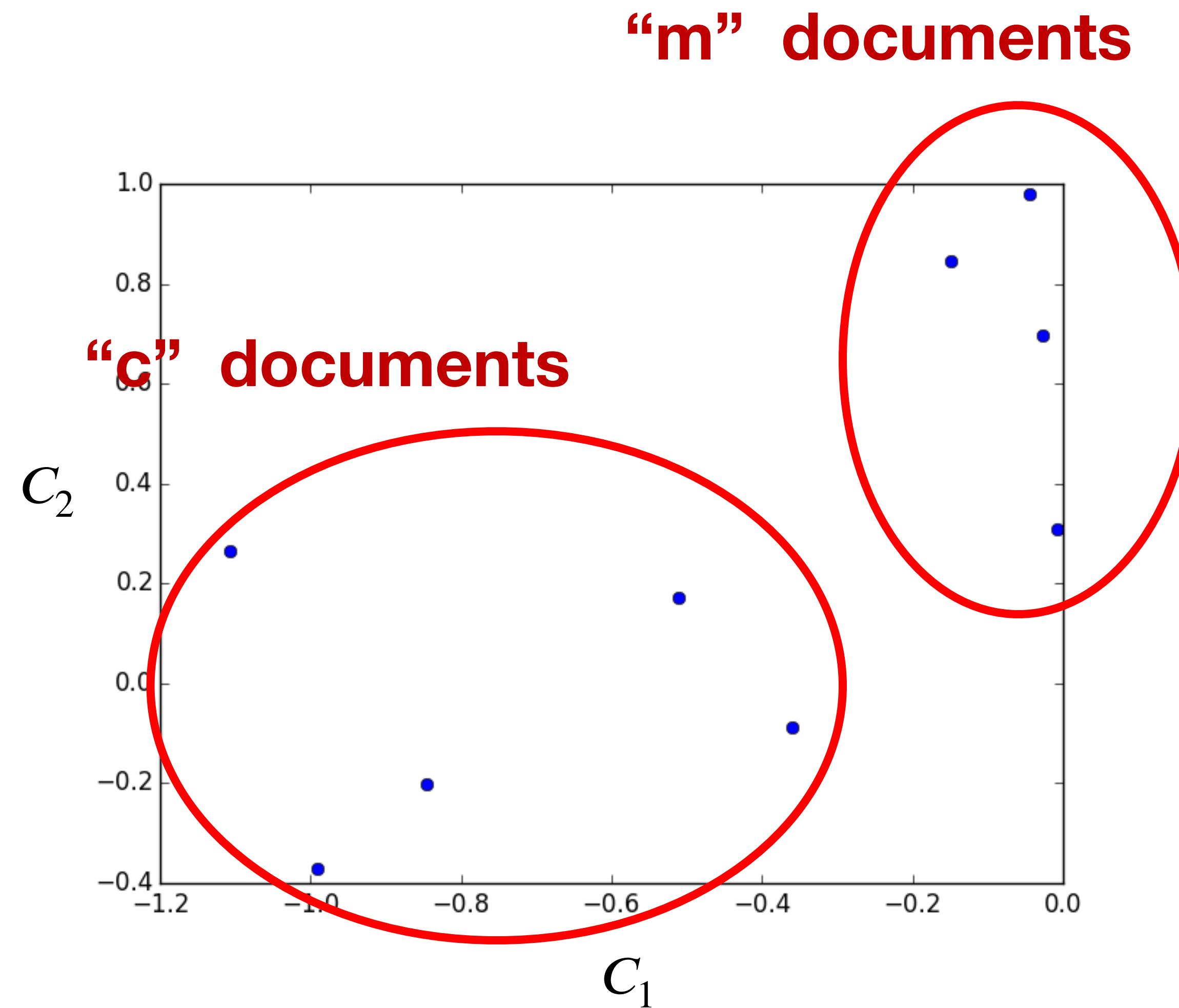
$$\begin{matrix} X \\ m \times n \end{matrix} \approx \begin{matrix} U \\ m \times k \end{matrix} \cdot \begin{matrix} D \\ k \times k \end{matrix} \cdot \begin{matrix} V^T \\ k \times n \end{matrix}$$

- With  $k = 2$ :

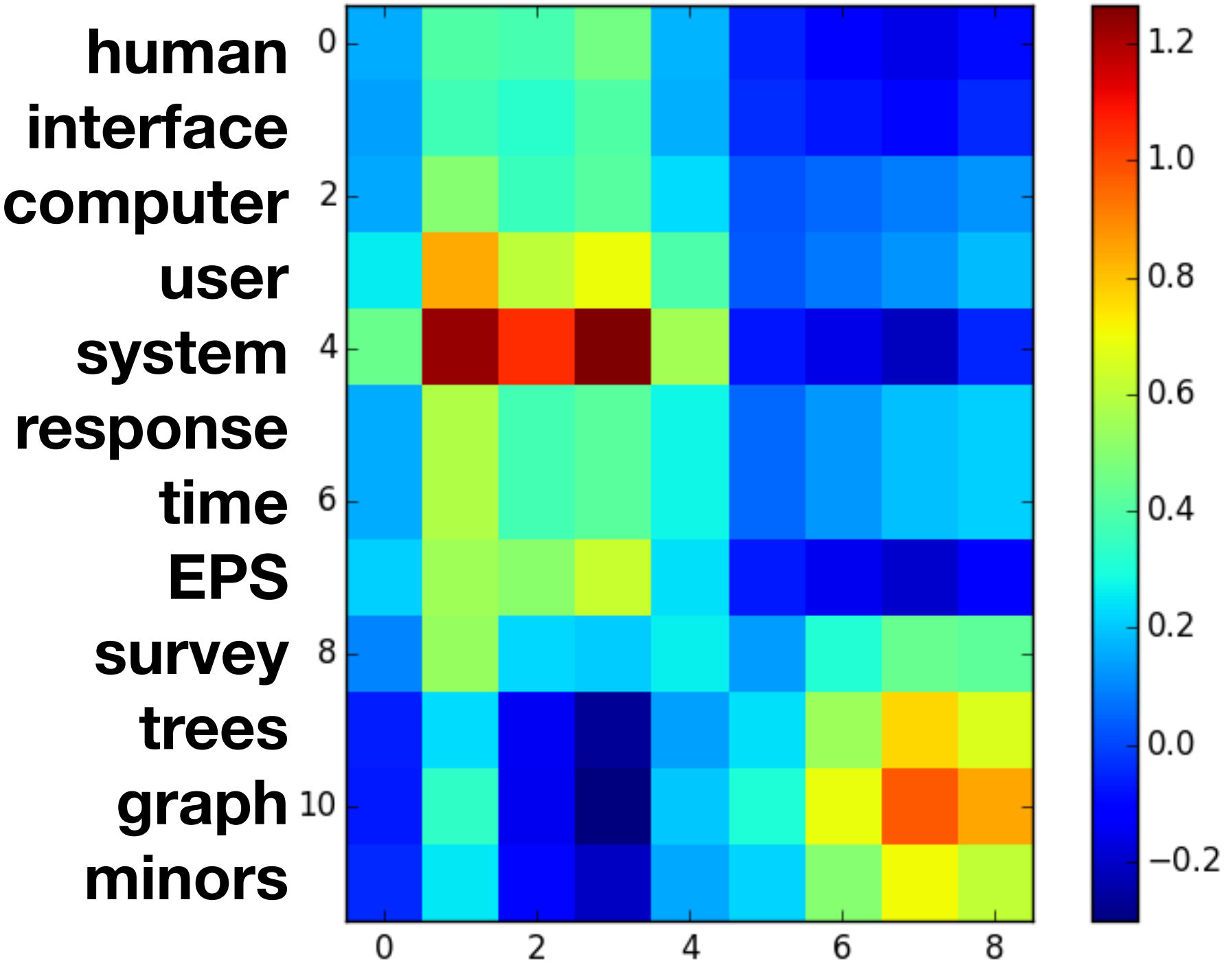
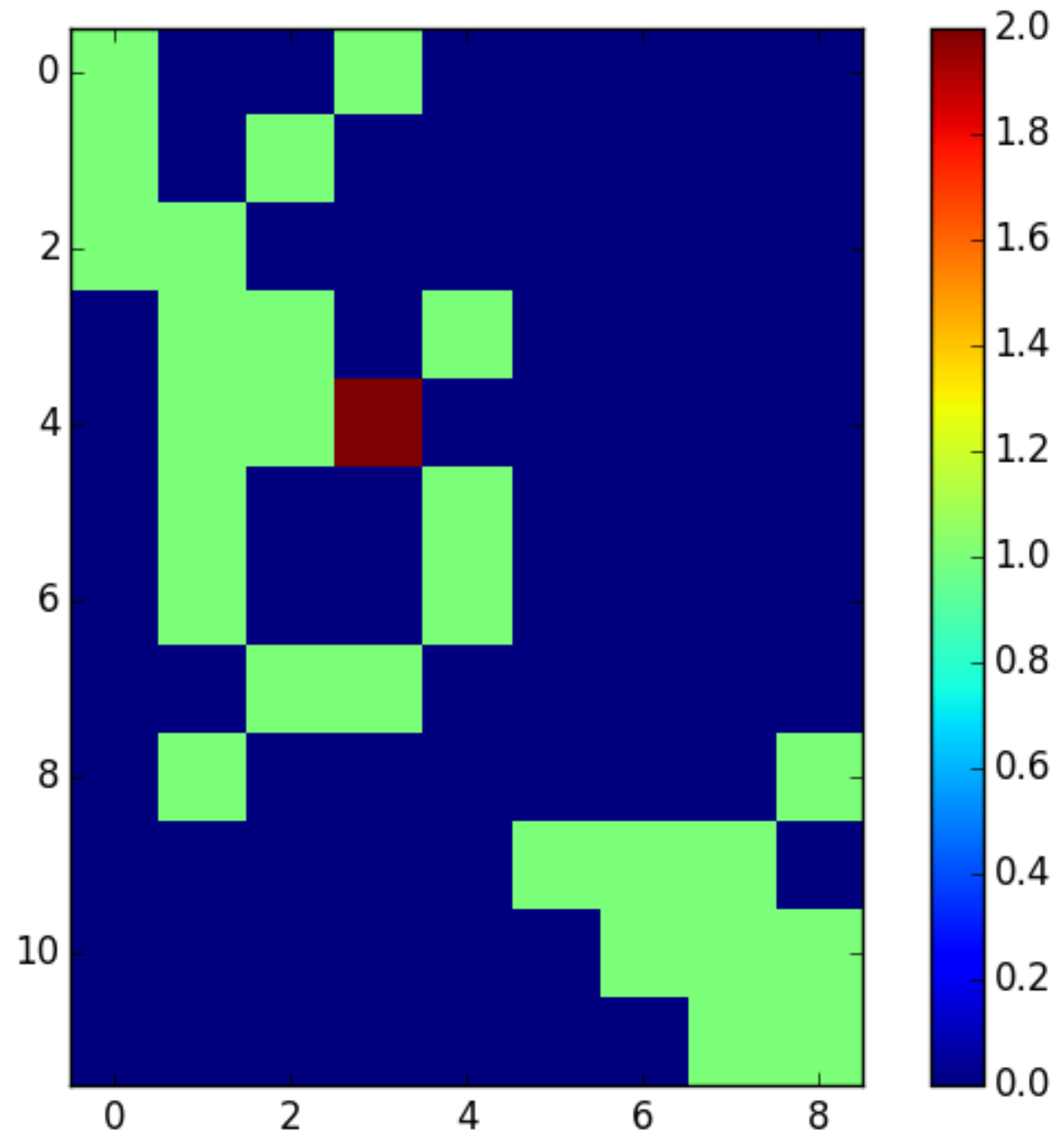
$$= \begin{matrix} W \\ m \times k \end{matrix} \cdot \begin{matrix} C \\ k \times n \end{matrix} \quad W = UD^{\frac{1}{2}} \quad C = D^{\frac{1}{2}}V^T$$



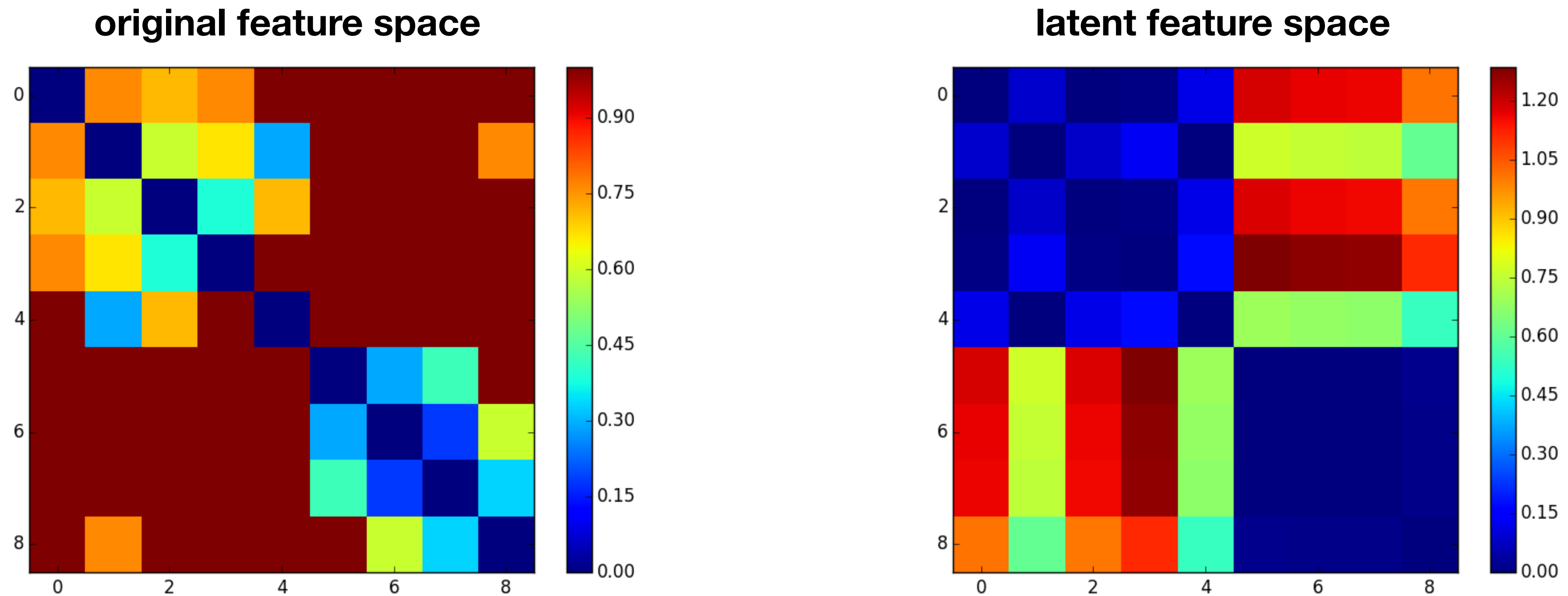
# Latent representation: example



# Reconstruction: example



# Distance matrix: example



- LSA results in more similarity **within topics**
  - Still some similarity across topics, but easily separable

# Today's lecture

---

Eigen-faces

Latent Semantic Analysis

**Collaborative Filtering**

# Recommendation Systems

---

- **Recommend** decisions / items the user may like
  - Need to understand both the **items** and the **users**
  - E.g. movies, books, groceries, restaurants (remember those?)
- Training data
  - **User information**  $\mapsto$  features: demographics, social network, past ratings
  - **Item metadata**  $\mapsto$  features: creator, genre, ingredients
- Learning output = decision function:
  - **Relevance** score, predicted **rating / ranking**



# Recommendation Systems: examples

**NETFLIX**  
Browse DVDs | Watch Instantly | Your Queue | Movies You'll Love | Friends & Community | DVD Rentals

**Movies You'll Love**  
Suggestions based on your ratings

To Get the Best Suggestions

1. Rate your genres.

New Suggestions for You  
Based on your recent ratings

**Amazon.com** | Your Amazon.com | Your Browsing History | Recommended For You | Amazon Betterizer | Improve Your Recommendations | Your

**Amazon Betterizer**  
Take a minute to improve your shopping experience by telling us which things you like. This helps us provide  
[Learn more](#)

Search | Images | Maps | Play | YouTube | News | Gmail | Drive | Calendar | [Show different items](#) | [Show my new re](#)

Google |

Web | Images | Maps | Shopping | News | Nose | More

About 1,190,000,000 results (0.37 seconds)

**California Fish Grill Inc**  
[cafishgrill.com/](http://cafishgrill.com/)  
Score: **24** / 30 · 117 Google reviews

**Bistango**  
[www.bistango.com/](http://www.bistango.com/)  
Zagat: **25** / 30 · 247 Google reviews

**Ruth's Chris Steak House**  
[www.ruthschris.com/](http://www.ruthschris.com/)  
Zagat: **27** / 30 · 75 Google reviews

**Stonefire Grill**  
[www.stonefiregrill.com/](http://www.stonefiregrill.com/)  
Zagat: **23** / 30 · 47 Google reviews

**How the Grinch Stole Christmas!**  
Dr. Seuss

**JUMANJI**

**Winnie the Pooh**  
A Very Merry Pooh Year

Like | Like

# Paradigms

- **Recommendation Systems** help reduce information overload
  - Identify most **relevant** items to attend to

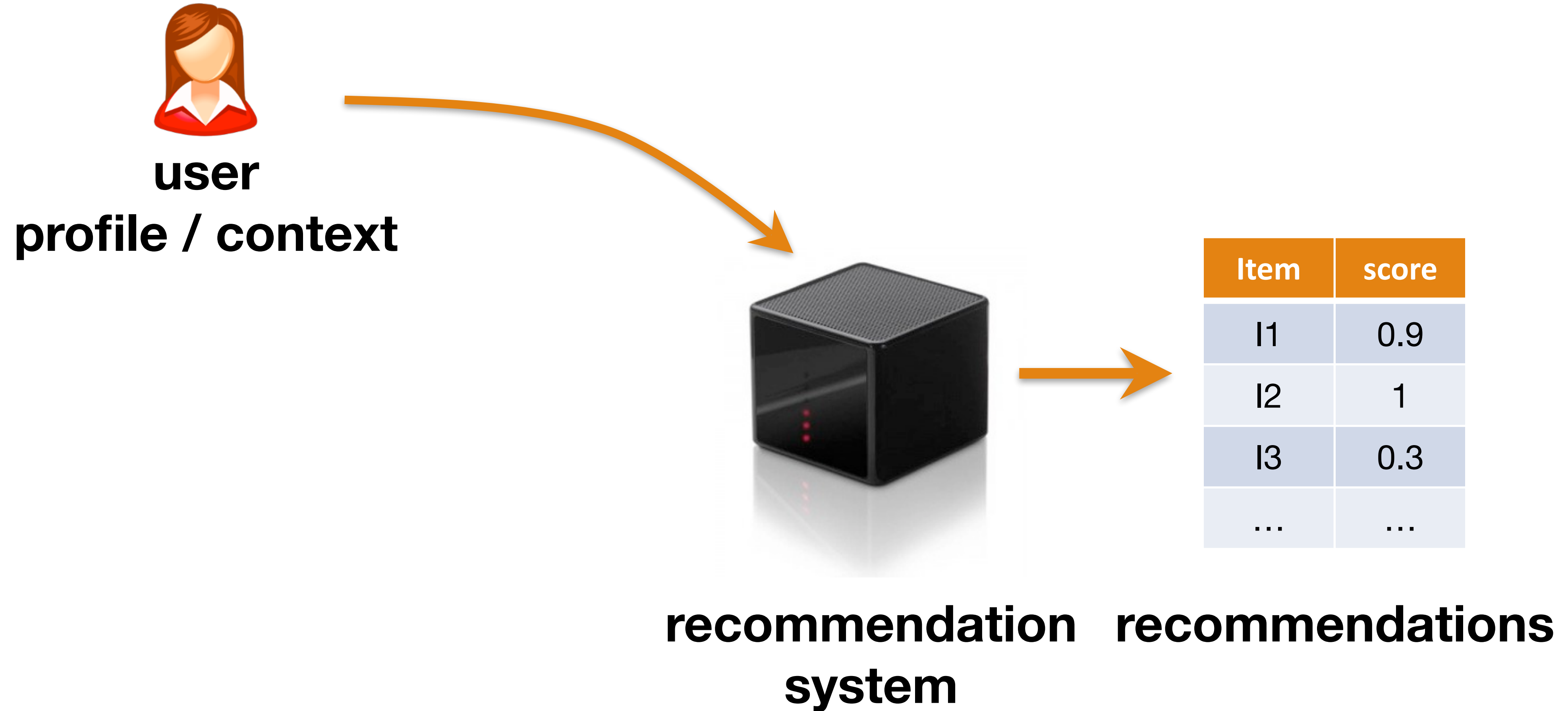


Item	score
I1	0.9
I2	1
I3	0.3
...	...

**recommendation system**      **recommendations**

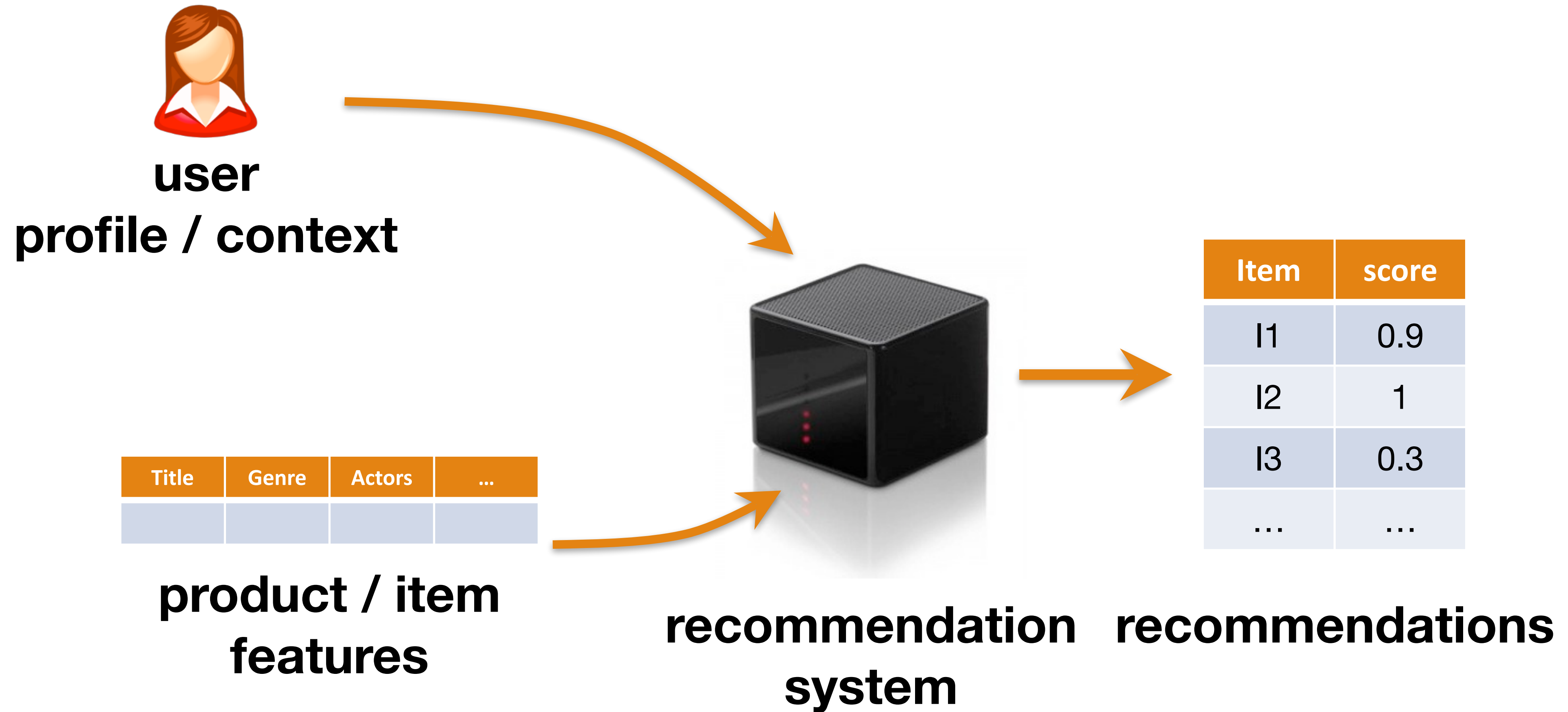
# Paradigms

- **Personalized** recommendations



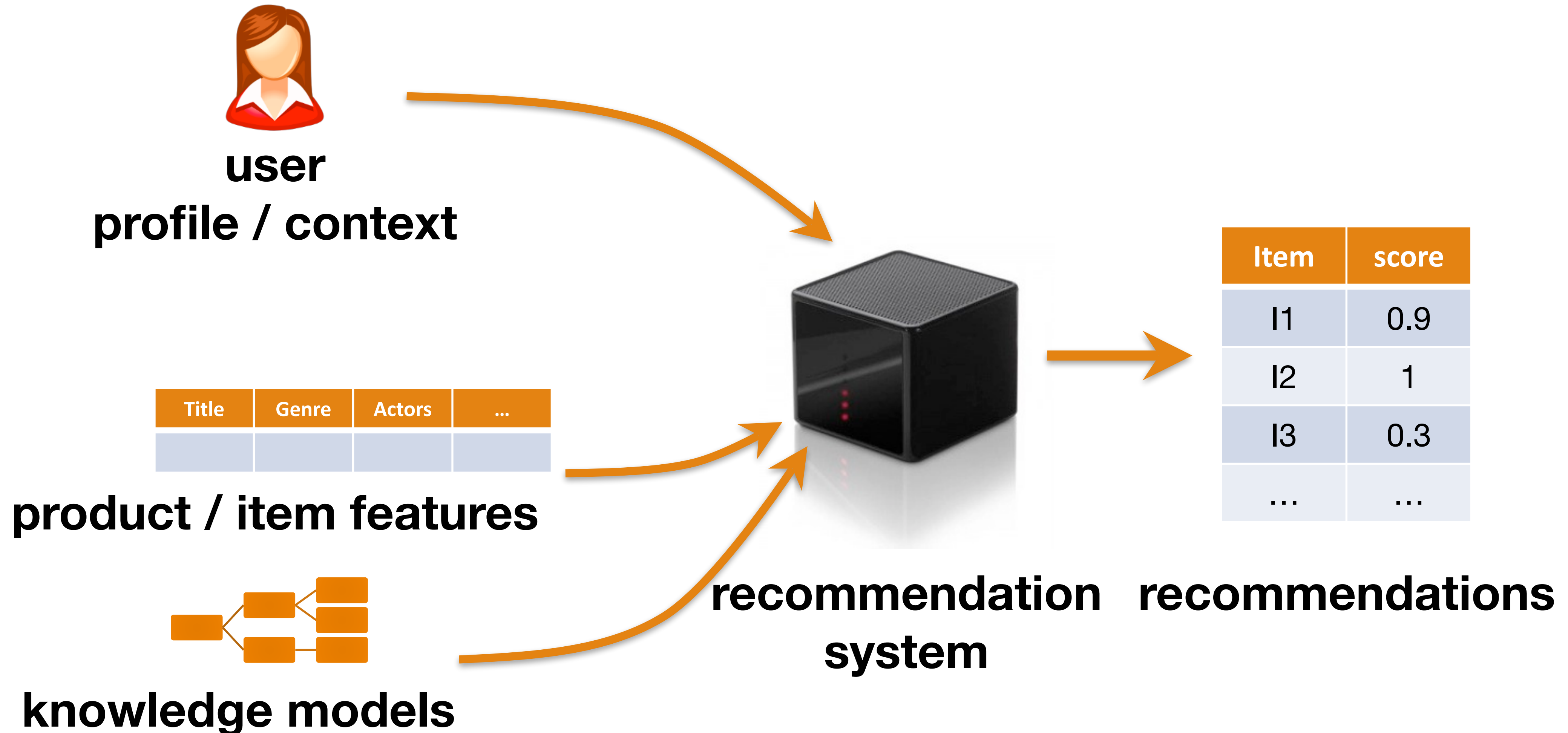
# Paradigms

- **Content-based**: “show me more of the things I liked before”



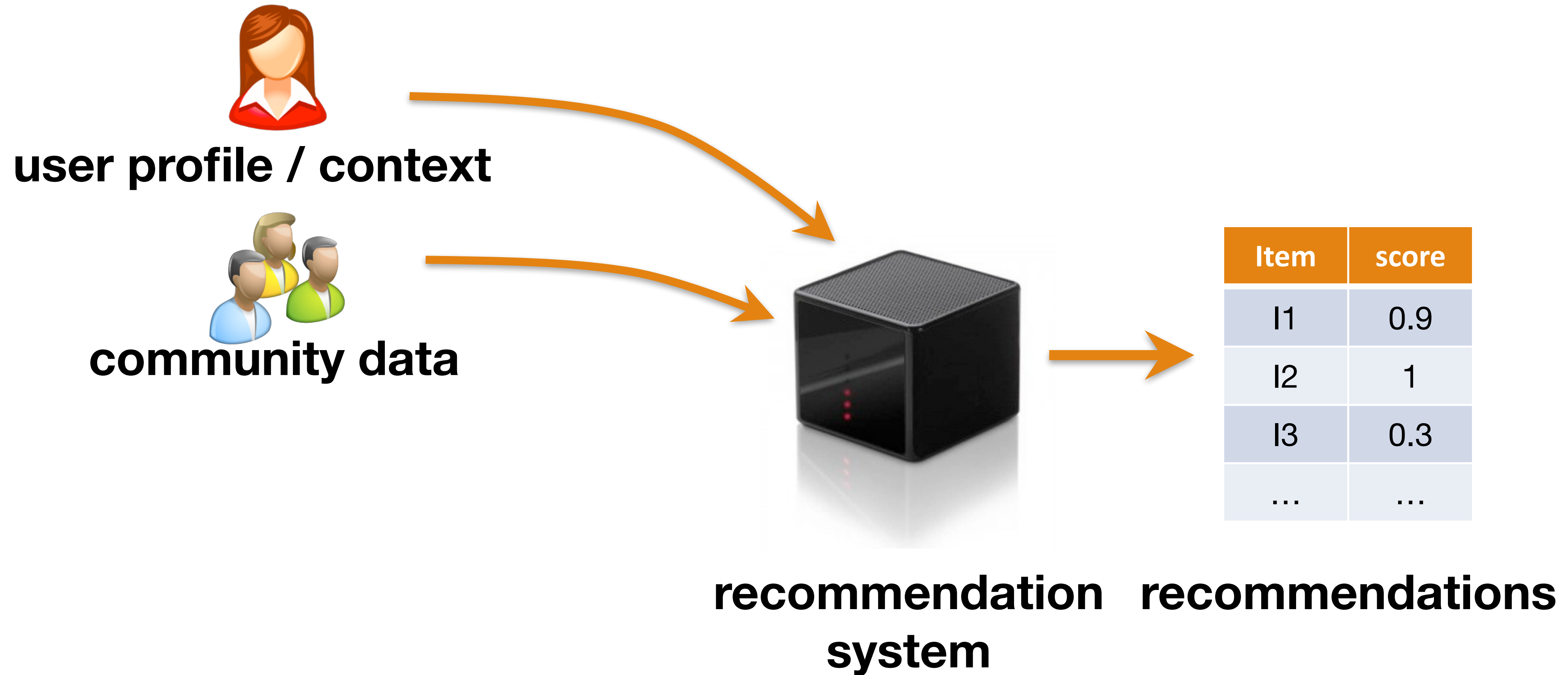
# Paradigms

- **Knowledge-based:** e.g., “tell me what fits my needs”



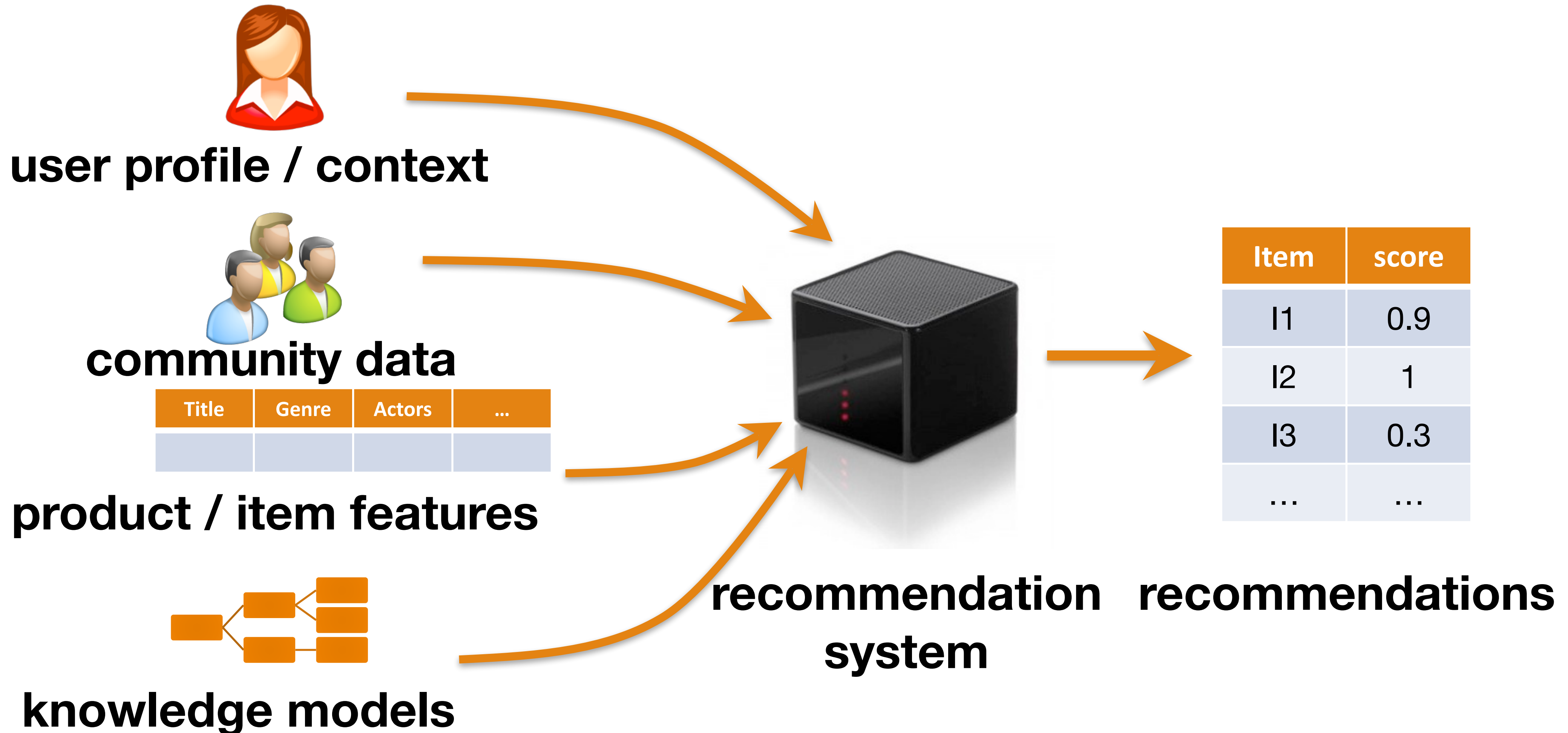
# Paradigms

- Collaborative Filtering: “tell me what’s popular among my peers”



# Paradigms

- **Hybrid**: combine information from many inputs and/or methods



# Measuring success

---

- **Prediction** perspective
  - Predict to what degree users like the item
  - Most common evaluation for research
  - Regression vs. “top-K” ranking
- **Interaction** perspective
  - Promote positive “feeling” in users (“satisfaction”, engagement)
  - Educate vs. persuade
- **Conversion** perspective
  - Measure commercial success (“hits”, click-through rates)
  - Optimize sales and profits



# Why are recommenders important?

- The **long tail** of item appeal
  - A few items are very popular
  - Most items are popular only with a few people
    - But everybody is interested in some obscure items
- **Goal:** recommend scarcely known items that the user might like



# Collaborative Filtering: example

**users**

1   2   3   4   5   6   7   8   9   10   11   12

**movies**

1	1		3		?	5			5		4	
2			5	4			4			2	1	3
3	2	4		1	2		3		4	3	5	
4		2	4		5			4			2	
5			4	3	4	2					2	5
6	1		3		3			2			4	

# Latent-space models

- **Model:** ratings matrix = user features · movie features

- ▶ Learn values from known ratings
- ▶ Extrapolate to unrated

**users**

	1		3		5		5		4			
			5	4			4		2	1	3	
<b>items</b>	2	4		1	2		3		4	3	5	
		2	4		5			4			2	
			4	3	4	2					2	5
	1		3		3			2			4	

$\approx$

**items**

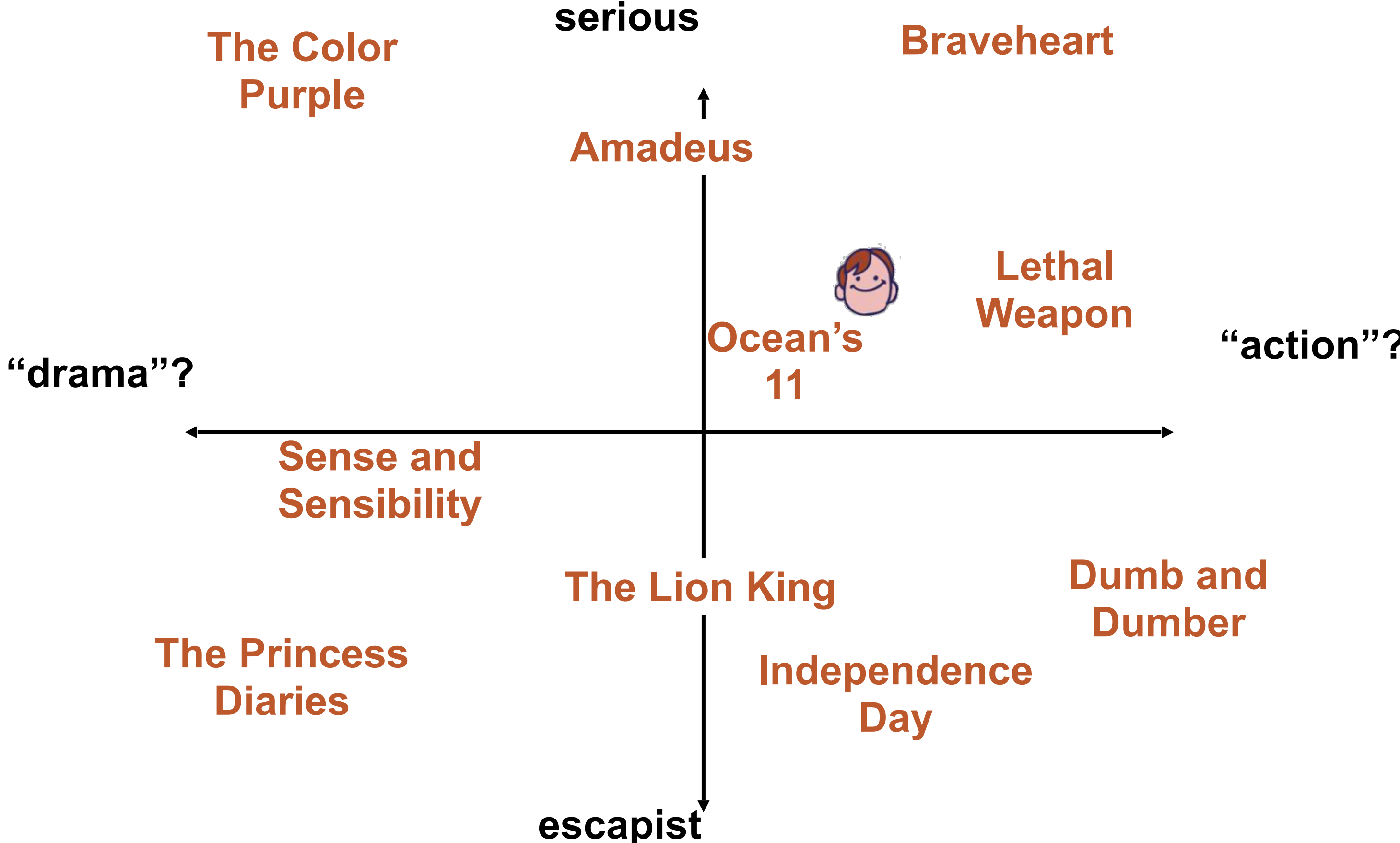
	1.	4.-	2.
	5.-	6.	5.
	2.-	3.	5.
	1.1	2.1	3.
	7.-	2.1	2.-
	1-	7.	3.



**users**

1.1	2.-	3.	5.	2.-	5.-	8.	4.-	3.	1.4	2.4	9.-
8.-	7.	5.	1.4	3.	1-	1.4	2.9	7.-	1.2	1.-	1.3
2.1	4.-	6.	1.7	2.4	9.	3.-	4.	8.	7.	6.-	1.

# Latent-space models



# Latent-space models

- Ideally, latent representation encodes some **meaning**
  - What kind of movie is this? Which movies is it similar to?
- Most data is **missing**  $\implies$  hard to perform SVD

▸  $\implies$  **gradient decent** on loss  $\mathcal{L}(U, V) = \sum_{u,m} \left( X_{mu} - \sum_k U_{mk} V_{ku} \right)^2$

```
# for user u, movie m, find the k'th eigenvector & coefficient by iterating:
predict_um = U[m,:].dot( V[:,u] ) # predict: vector-vector product
err = ( rating[u,m] - predict_um ) # find error residual
V_ku, U_mk = V[k,u], U[m,k] # make copies for update
U[m,k] += alpha * err * V_ku # Update our matrices
V[k,u] += alpha * err * U_mk # (compare to least-squares gradient)
```

# Latent-space models: extensions

- Add **lower-order** terms:  $r_{mu} \approx \mu + b_m + b_u + \sum_k U_{mk} V_{ku}$ 
  - $\mu$  = overall average rating (affected by user interface etc.)
  - $b_m + b_u$  = item and user biases
- Can add **non-linearity** (saturating?)
- Choose a **loss**, e.g. MSE or multilogistic
- Train using a gradient-based **optimizer**

# Ensembles for recommendations

---

- Many possible **models**:
  - Feature-based regression
  - (Weighted) kNN on items
  - (Weighted) kNN on users
  - Latent space representation
- Perhaps we should **combine** them?
- Use an **ensemble** average, or a stacked ensemble
  - **Stacked ensemble** = train a weighted combination of model predictions

# Demo

---

- <http://www.benfrederickson.com/matrix-factorization/>



# Logistics

---

assignments

- Assignment 5 **due Thursday**

project

- Final report **due next Thursday**

evaluations

- Evaluations **due end of next week**

final exam

- **Review:** next Thursday
- **Final:** Thursday, March 18, 1:30–3:30pm