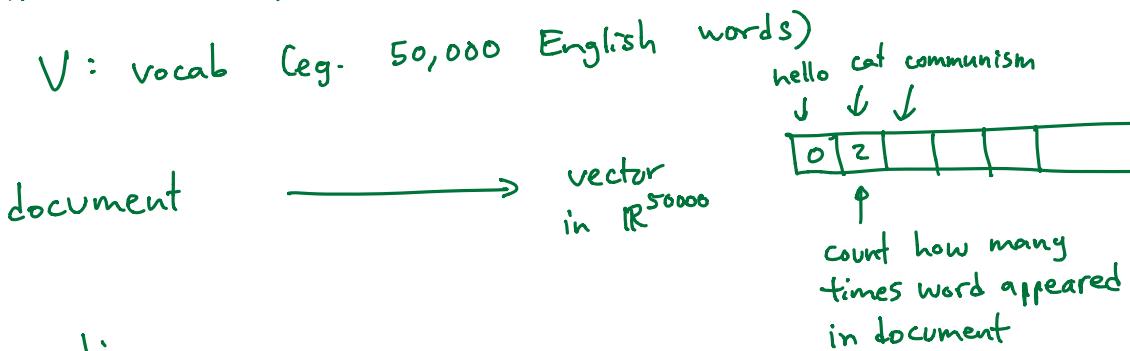


Singular value decomposition

DSE 220

Word embeddings

1 Traditional representation of documents: bag of words



The good:

- ① This allows us to use standard classification methods (e.g. logistic regression) for vector data.
- ② Has been reasonably successful in many applications.

The bad:

- ① Very high dimensional (∴ need lots of training data).

Possible saving grace: many of these words are very close in meaning:
cat — kitten — kitty — pussy cat — - -

and many are fairly close, e.g.:

- different dinner foods

→ We should be able to generalize across these words

This is what word embeddings exploit.

2 Embeddings

$\varphi: V \rightarrow \mathbb{R}^{300}$, i.e. map each word w to a 300-dim vector $\varphi(w)$

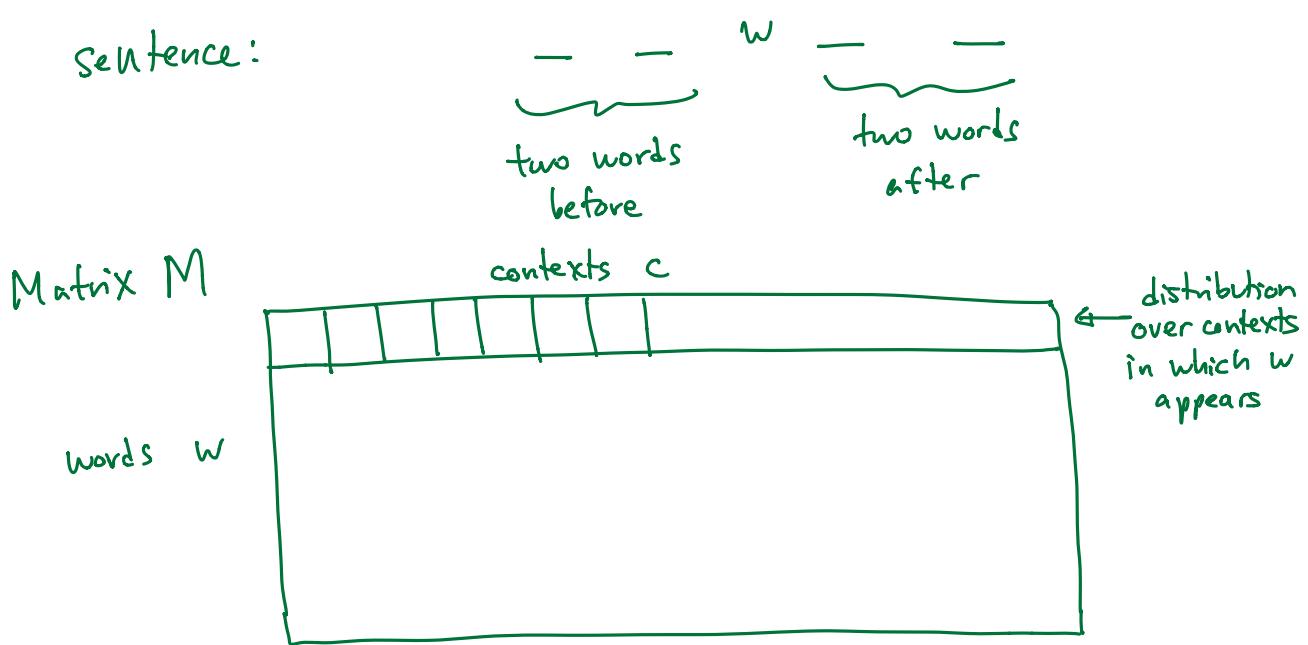
Goal: Words w, w' are "similar" $\hookrightarrow \varphi(w)$ is close to $\varphi(w')$
(more precisely, if $\varphi(w) \cdot \varphi(w')$ is large)

E.g. $\varphi(\text{cat}) \approx \varphi(\text{kitten})$

- a) Words are similar if they appear in similar contexts

E.g. I had steak and potatoes for dinner
I had sausage and potatoes for dinner

⑤ Look at the contexts in which words appear



M is huge:

- 50,000 words (say)
 - Maybe 10^6 contexts ?
- } 50000×10^6

Each row M_w is a vector representation of the corresponding word that captures its meaning.

⑥ Problem: M_w is huge: $(10^6$ -dimensional.

Use PCA to reduce dimension to (say) 300.

This gives us an embedding of V .

Generalizing the spectral decomposition

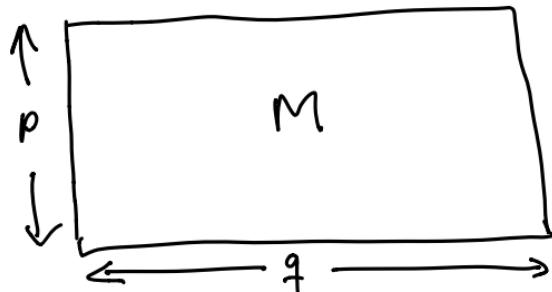
SVD: a popular way of
constructing embeddings.

$$\begin{array}{c} \uparrow \\ p \\ \downarrow \\ \xrightarrow{\quad p \quad} \end{array} M = \begin{array}{|c|c|} \hline & | \\ u_1 & \cdots & u_p \\ | & | & | \\ & \ddots & \\ & 0 & \\ \hline \end{array} \begin{array}{|c|c|} \hline \lambda_1 & \\ & \lambda_2 & \ddots & 0 \\ & 0 & & \lambda_p \\ \hline \end{array} \begin{array}{|c|c|} \hline & u_1 \\ & \vdots \\ & -u_p \\ \hline \end{array}$$

For **symmetric** matrices (e.g. covariance matrices), we have seen:

- Results about existence of eigenvalues and eigenvectors
- Eigenvectors form an alternative basis
- Resulting spectral decomposition, used in PCA

What about **arbitrary** matrices $M \in \mathbb{R}^{p \times q}$?



Data that we want to "embed"

Singular value decomposition (SVD)

M:
 $p \times q$



Any $p \times q$ matrix ($p \leq q$) has a **singular value decomposition**:

$$M = \underbrace{\begin{pmatrix} \uparrow & & \uparrow \\ u_1 & \cdots & u_p \\ \downarrow & & \downarrow \end{pmatrix}}_{p \times p \text{ matrix } U} \underbrace{\begin{pmatrix} \sigma_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \sigma_p \end{pmatrix}}_{p \times p \text{ matrix } \Lambda} \underbrace{\begin{pmatrix} \leftarrow & v_1 & \rightarrow \\ & \vdots & \\ \leftarrow & v_p & \rightarrow \end{pmatrix}}_{p \times q \text{ matrix } V^T}$$

- u_1, \dots, u_p are orthonormal vectors in \mathbb{R}^p (σ_1, u_1, v_1)
- v_1, \dots, v_p are orthonormal vectors in \mathbb{R}^q (σ_2, u_2, v_2)
- $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_p \geq 0$ are **singular values** (σ_3, u_3, v_3)

One line of Python ... ("thin SVD")

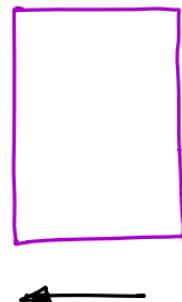


Low-rank approximation

Singular value decomposition of $p \times q$ matrix M (with $p \leq q$):

$$\underline{M} = \underbrace{\begin{pmatrix} u_1 & \cdots & u_p \end{pmatrix}}_{p \times p} \underbrace{\begin{pmatrix} \sigma_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \sigma_p \end{pmatrix}}_{p \times p} \underbrace{\begin{pmatrix} v_1 & \cdots & v_p \end{pmatrix}}_{p \times q}$$

\xrightarrow{k} \xrightarrow{k} \xrightarrow{k}



A concise approximation to M , for any $k \leq p$

$$\hat{M} = \underbrace{\begin{pmatrix} u_1 & \cdots & u_k \end{pmatrix}}_{p \times k} \underbrace{\begin{pmatrix} \sigma_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \sigma_k \end{pmatrix}}_{k \times k} \underbrace{\begin{pmatrix} v_1 & \cdots & v_k \end{pmatrix}}_{k \times q}$$

$\xrightarrow{p \times k}$ $\xrightarrow{k \times k}$ $\xrightarrow{k \times q}$

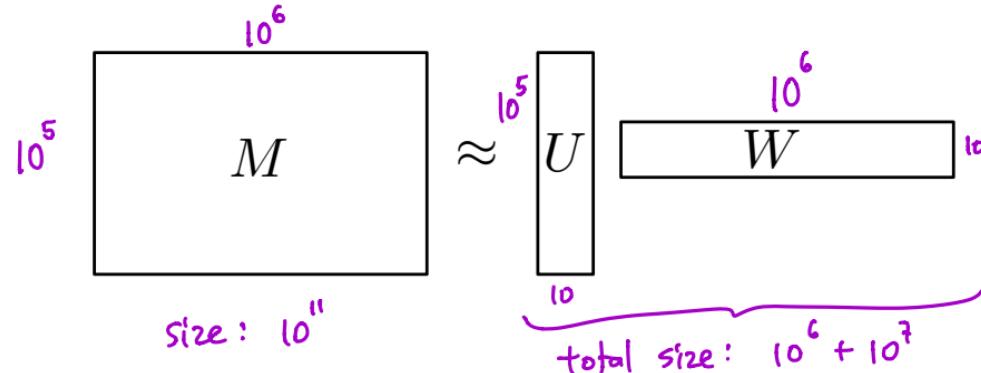
\hat{M} is the **best rank- k approximation** to M .

Optimality property

Let M be any $p \times q$ matrix.

Want to approximate M by a $p \times q$ matrix \hat{M} of the form UW :

- U is $p \times k$ and W is $k \times q$ $k = 10$
- $k \leq p, q$ is of our choosing

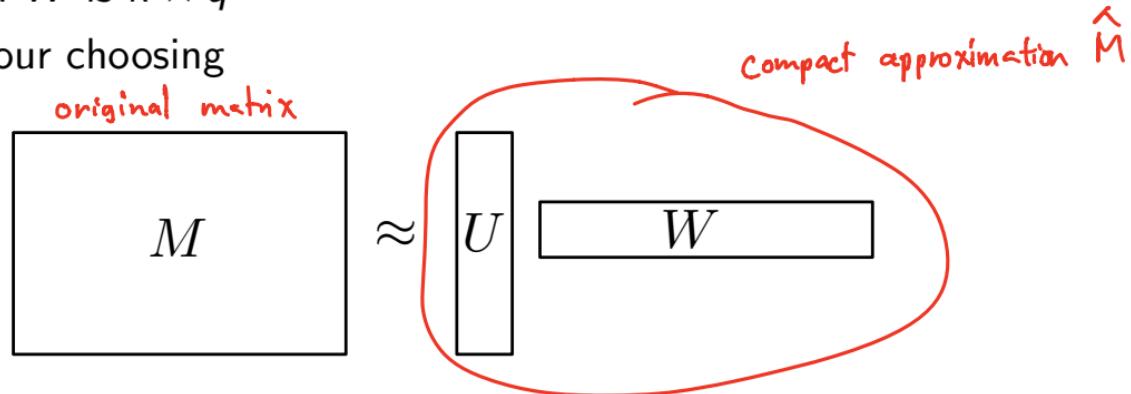


Optimality property

Let M be any $p \times q$ matrix.

Want to approximate M by a $p \times q$ matrix \hat{M} of the form UW :

- U is $p \times k$ and W is $k \times q$
- $k \leq p, q$ is of our choosing



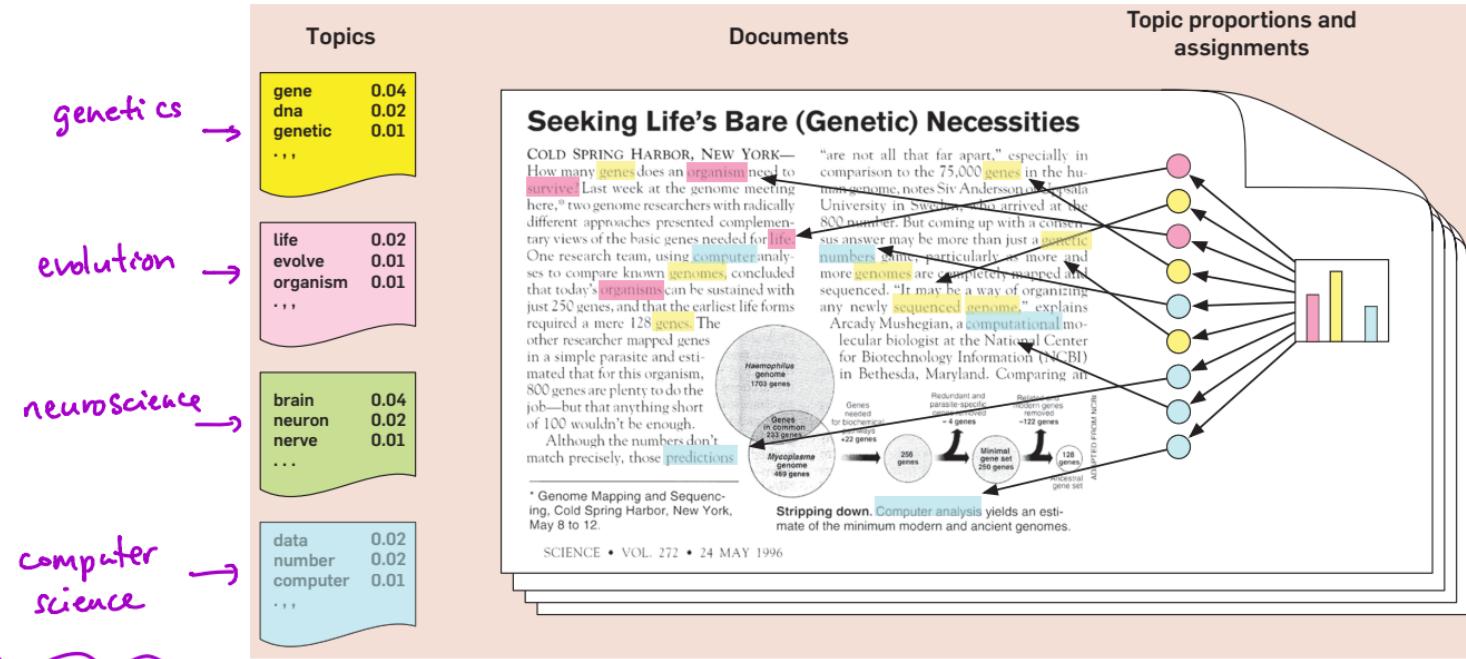
SVD yields the best such approximation \hat{M} , minimizing the squared error

$$\sum_{i,j} (M_{ij} - \hat{M}_{ij})^2.$$

Example: Topic modeling

Eg. Corpus of all Science articles over 100 years

Blei (2012):



Embedding of document: combination of topics \rightarrow vector in \mathbb{R}^k

Latent semantic indexing (LSI)

Given a large corpus of n documents:

- Fix a vocabulary, say of V words.
- Bag-of-words representation for documents: each document becomes a vector of length V , with one coordinate per word.
- The corpus is an $n \times V$ matrix, one row per document.

\uparrow \uparrow
#docs #words

Huge matrix

	cat	dog	house	boat	garden	...
Doc 1	4	1	1	0	2	
Doc 2	0	0	3	1	0	
Doc 3	0	1	3	0	0	
			:			

Let's find a concise approximation to this matrix M .

Latent semantic indexing, cont'd

Use SVD to get an approximation to M : for small \underline{k} ,

e.g. $k=10$

$$\begin{array}{c} \uparrow \\ n \\ \downarrow \end{array} \left(\begin{array}{c} \xleftarrow{\hspace{1cm}} V \xrightarrow{\hspace{1cm}} \\ \left(\begin{array}{c} \xleftarrow{\hspace{1cm}} \text{doc 1} \xrightarrow{\hspace{1cm}} \\ \xleftarrow{\hspace{1cm}} \text{doc 2} \xrightarrow{\hspace{1cm}} \\ \xleftarrow{\hspace{1cm}} \text{doc 3} \xrightarrow{\hspace{1cm}} \\ \vdots \\ \xleftarrow{\hspace{1cm}} \text{doc } n \xrightarrow{\hspace{1cm}} \end{array} \right) \end{array} \approx \left(\begin{array}{c} \xleftarrow{\hspace{1cm}} \theta_1 \xrightarrow{\hspace{1cm}} \\ \xleftarrow{\hspace{1cm}} \theta_2 \xrightarrow{\hspace{1cm}} \\ \xleftarrow{\hspace{1cm}} \theta_3 \xrightarrow{\hspace{1cm}} \\ \vdots \\ \xleftarrow{\hspace{1cm}} \theta_n \xrightarrow{\hspace{1cm}} \end{array} \right) \underbrace{\left(\begin{array}{c} \xleftarrow{\hspace{1cm}} \Psi_1 \xrightarrow{\hspace{1cm}} \\ \vdots \\ \xleftarrow{\hspace{1cm}} \Psi_k \xrightarrow{\hspace{1cm}} \end{array} \right)}_{\substack{k \times V \text{ matrix } \Psi \\ 10 \times V}} \quad n \times k \text{ matrix } \Theta \quad n \times 10$$

Latent semantic indexing, cont'd

Use SVD to get an approximation to M : for small k ,

$$\underbrace{\begin{pmatrix} \leftarrow \text{doc 1} \rightarrow \\ \leftarrow \text{doc 2} \rightarrow \\ \leftarrow \text{doc 3} \rightarrow \\ \vdots \\ \leftarrow \text{doc } n \rightarrow \end{pmatrix}}_{n \times V \text{ matrix } M} \approx \underbrace{\begin{pmatrix} \leftarrow \theta_1 \rightarrow \\ \leftarrow \theta_2 \rightarrow \\ \leftarrow \theta_3 \rightarrow \\ \vdots \\ \leftarrow \theta_n \rightarrow \end{pmatrix}}_{n \times k \text{ matrix } \Theta} \underbrace{\begin{pmatrix} \leftarrow \Psi_1 \rightarrow \\ \vdots \\ \leftarrow \Psi_k \rightarrow \end{pmatrix}}_{k \times V \text{ matrix } \Psi}$$

Think of this as a *topic model* with k topics.

- Ψ_j is a vector of length V describing topic j : coefficient Ψ_{jw} is large if word w appears often in that topic.
- Each document is a combination of topics: $\theta_{ij} =$ weight of topic j in doc i .

Latent semantic indexing, cont'd

Use SVD to get an approximation to M : for small k ,
representation of document #1
in terms of the k topics

$$\underbrace{\begin{pmatrix} \leftarrow \text{doc 1} \rightarrow \\ \leftarrow \text{doc 2} \rightarrow \\ \leftarrow \text{doc 3} \rightarrow \\ \vdots \\ \leftarrow \text{doc } n \rightarrow \end{pmatrix}}_{n \times V \text{ matrix } M} \approx \underbrace{\begin{pmatrix} \leftarrow \theta_1 \rightarrow \\ \leftarrow \theta_2 \rightarrow \\ \leftarrow \theta_3 \rightarrow \\ \vdots \\ \leftarrow \theta_n \rightarrow \end{pmatrix}}_{n \times k \text{ matrix } \Theta} \underbrace{\begin{pmatrix} \leftarrow \Psi_1 \rightarrow \\ \vdots \\ \leftarrow \Psi_k \rightarrow \end{pmatrix}}_{k \times V \text{ matrix } \Psi} \xrightarrow{\text{in terms of the } k \text{ topics}}$$

← topic 1 ← topic k

Think of this as a *topic model* with k topics.

- Ψ_j is a vector of length V describing topic j : coefficient Ψ_{jw} is large if word w appears often in that topic.
 - Each document is a combination of topics: $\theta_{ij} = \text{weight of topic } j \text{ in doc } i$.

Document i originally represented by i th row of M , a vector in \mathbb{R}^V .

Can instead use $\theta_i \in \mathbb{R}^k$, a more concise “semantic” representation.

Example: Collaborative filtering

Details and images from Koren, Bell, Volinksy (2009).

Recommender systems: matching customers with products.

- Given: data on prior purchases/interests of users
- Recommend: further products of interest

Prototypical example: Netflix.

Example: Collaborative filtering

Details and images from Koren, Bell, Volinksy (2009).

Recommender systems: matching customers with products.

- Given: data on prior purchases/interests of users
- Recommend: further products of interest

Prototypical example: Netflix.

A successful approach: **collaborative filtering**.

- Model dependencies between different products, and between different users.
- Can give reasonable recommendations to a relatively new user.

Example: Collaborative filtering

Details and images from Koren, Bell, Volinksy (2009).

Recommender systems: matching customers with products.

- Given: data on prior purchases/interests of users
- Recommend: further products of interest

Prototypical example: Netflix.

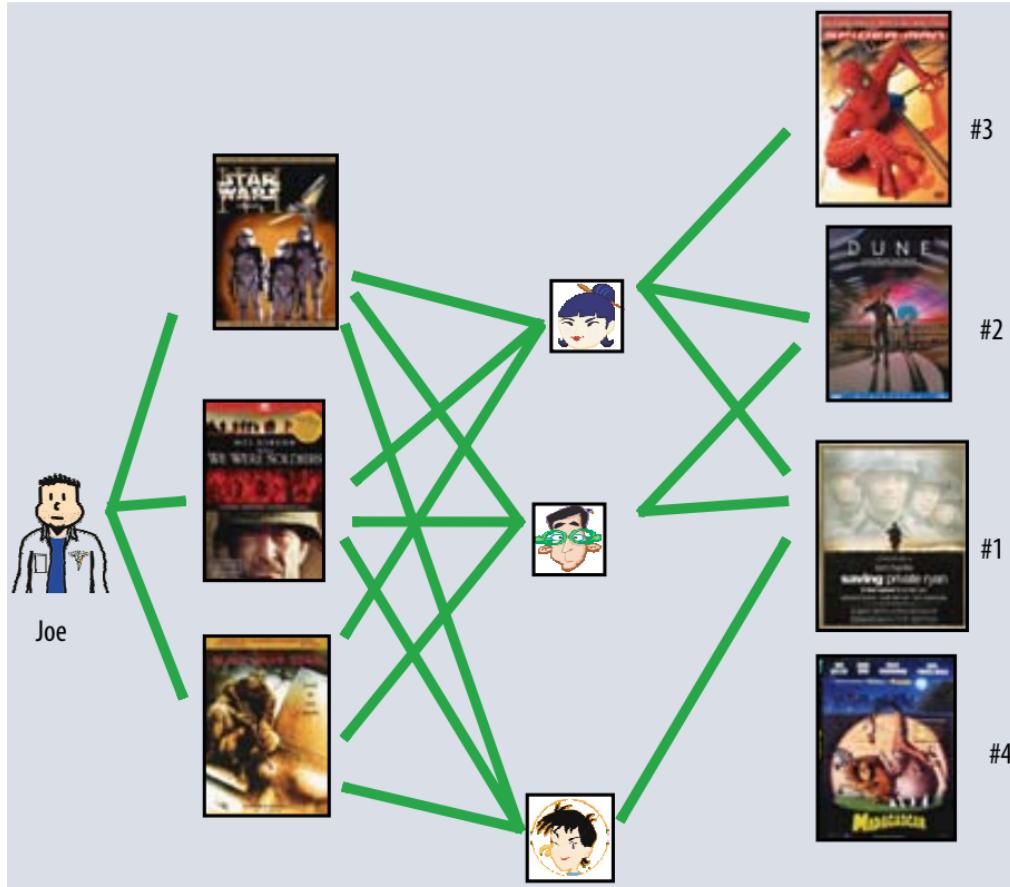
A successful approach: **collaborative filtering**.

- Model dependencies between different products, and between different users.
- Can give reasonable recommendations to a relatively new user.

Two strategies for collaborative filtering:

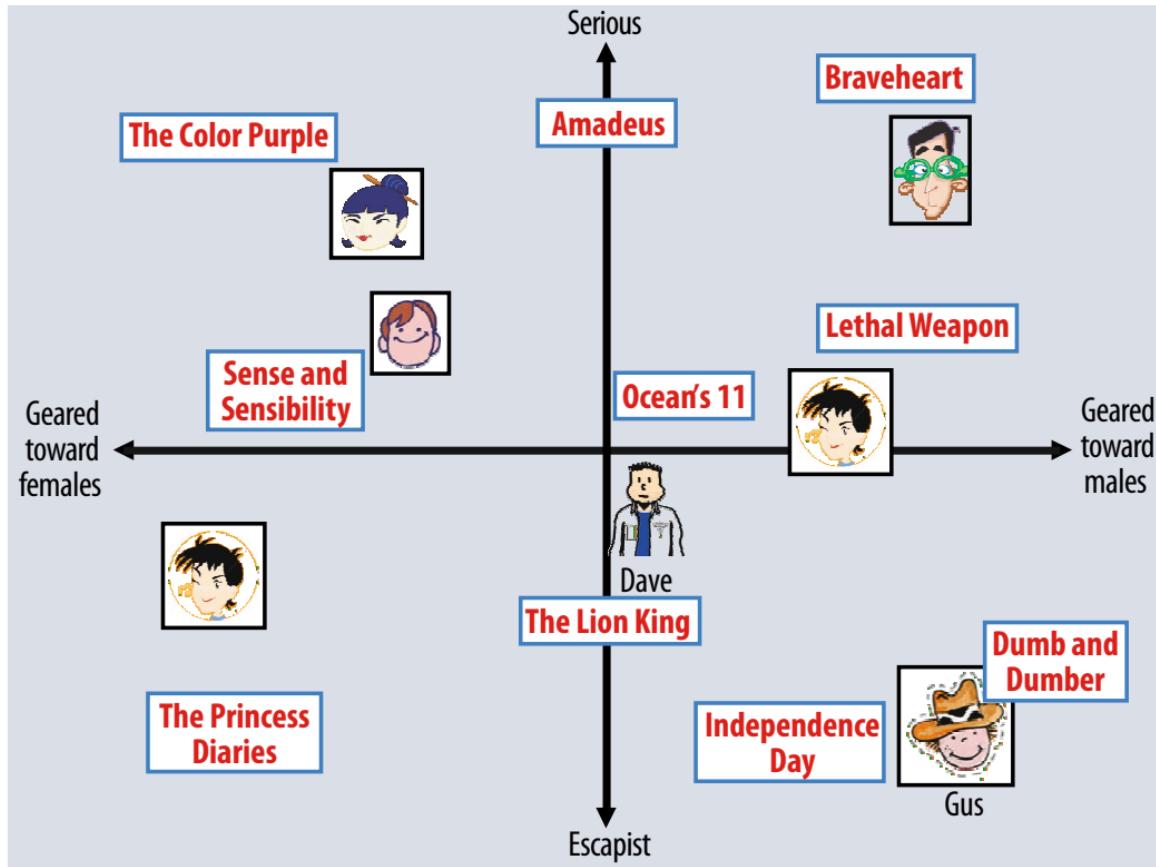
- Neighborhood methods
- Latent factor methods

Neighborhood methods



Latent factor methods

Embed all users and all movies into a common low-dimensional space.



The matrix factorization approach

User ratings are assembled in a large matrix M :

	Star Wars	Matrix	Casablanca	Camelot	Godfather	...
User 1	5	5	2	0	0	
User 2	0	0	3	4	5	
User 3	0	0	5	0	0	
	:					

- Not rated = 0, otherwise scores 1-5.
- For n users and p movies, this has size $n \times p$.
- Most of the entries are unavailable, and we'd like to predict these.

Idea: Get best low-rank approximation of M , and use to fill in the missing entries.

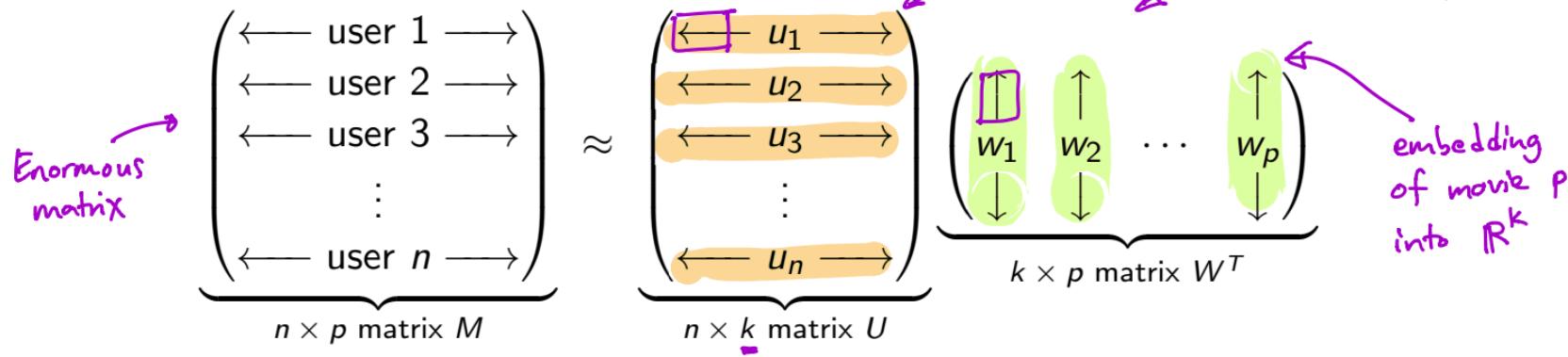
User and movie factors

Best rank- k approximation is of the form $M \approx UW^T$:

$$\underbrace{\begin{pmatrix} \leftarrow \text{user 1} \rightarrow \\ \leftarrow \text{user 2} \rightarrow \\ \leftarrow \text{user 3} \rightarrow \\ \vdots \\ \leftarrow \text{user } n \rightarrow \end{pmatrix}}_{n \times p \text{ matrix } M} \approx \underbrace{\begin{pmatrix} \leftarrow u_1 \rightarrow \\ \leftarrow u_2 \rightarrow \\ \leftarrow u_3 \rightarrow \\ \vdots \\ \leftarrow u_n \rightarrow \end{pmatrix}}_{n \times k \text{ matrix } U} \underbrace{\begin{pmatrix} w_1 & w_2 & \cdots & w_p \end{pmatrix}}_{k \times p \text{ matrix } W^T}$$

User and movie factors

Best rank- k approximation is of the form $M \approx UW^T$:



Thus user i 's rating of movie j is approximated as

$$M_{ij} \approx u_i \cdot w_j$$

to what extent does user i like movie j ?

embedding of user i in \mathbb{R}^k

embedding of movie j in \mathbb{R}^k

User and movie factors

Best rank- k approximation is of the form $M \approx UW^T$:

$$\underbrace{\begin{pmatrix} \leftarrow \text{user 1} \rightarrow \\ \leftarrow \text{user 2} \rightarrow \\ \leftarrow \text{user 3} \rightarrow \\ \vdots \\ \leftarrow \text{user } n \rightarrow \end{pmatrix}}_{n \times p \text{ matrix } M} \approx \underbrace{\begin{pmatrix} \leftarrow u_1 \rightarrow \\ \leftarrow u_2 \rightarrow \\ \leftarrow u_3 \rightarrow \\ \vdots \\ \leftarrow u_n \rightarrow \end{pmatrix}}_{n \times k \text{ matrix } U} \underbrace{\begin{pmatrix} w_1 & w_2 & \cdots & w_p \end{pmatrix}}_{k \times p \text{ matrix } W^T}$$

Thus user i 's rating of movie j is approximated as

$$M_{ij} \approx u_i \cdot w_j$$

"Latent" representation embeds users and movies in the same k -dimensional space:

- Represent i th user by $u_i \in \mathbb{R}^k$
- Represent j th movie by $w_j \in \mathbb{R}^k$

Top two Netflix factors

