

Randomized SVD

CSCI 4360/6360 Data Science II

Singular Value Decomposition (SVD)

- Given a (any!) matrix M , which is $n \times m$, it can be represented as

$$M = U\Sigma V^T$$

- U : $n \times n$, unitary matrix (orthogonal)
- Σ : $n \times m$, diagonal matrix of singular values
- V^T : $m \times m$, unitary matrix (orthogonal)

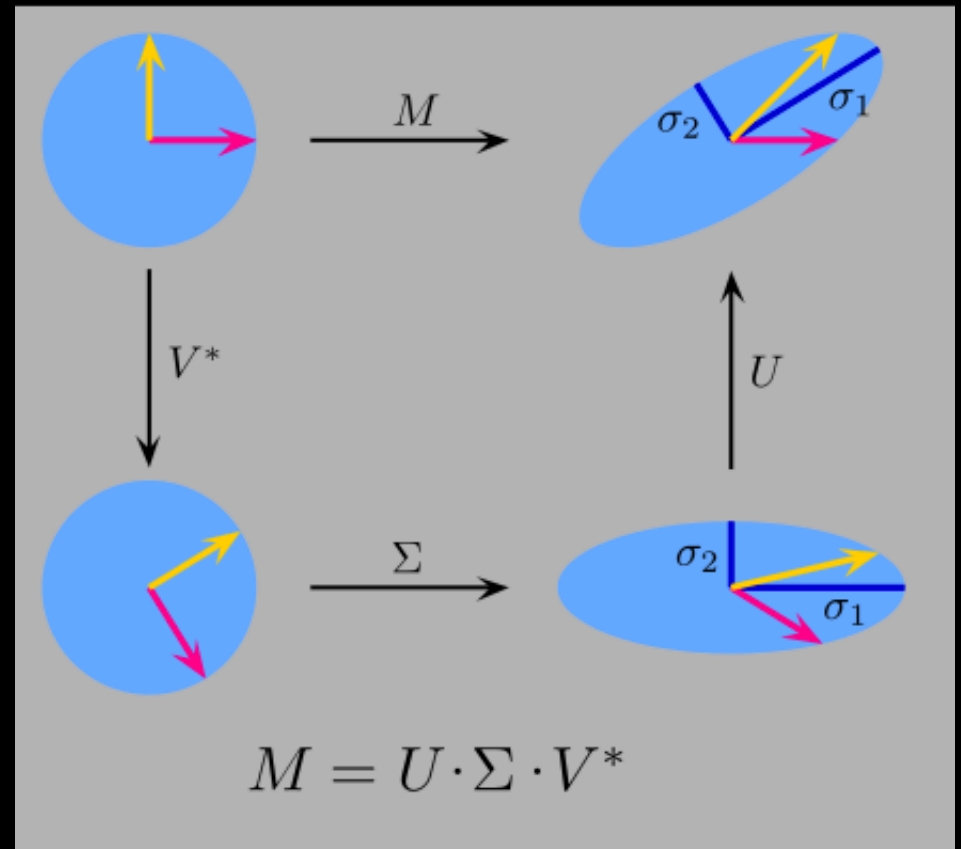
SVD

$$\mathbf{A}_{[m \times n]} = \mathbf{U}_{[m \times r]} \mathbf{\Sigma}_{[r \times r]} (\mathbf{V}_{[n \times r]})^T$$

- **A: Input data matrix**
 - $m \times n$ matrix (e.g., m documents, n terms)
- **U: Left singular vectors**
 - $m \times r$ matrix (m documents, r concepts)
- **Σ : Singular values**
 - $r \times r$ diagonal matrix (strength of each 'concept')
(r : rank of the matrix **A**)
- **V: Right singular vectors**
 - $n \times r$ matrix (n terms, r concepts)

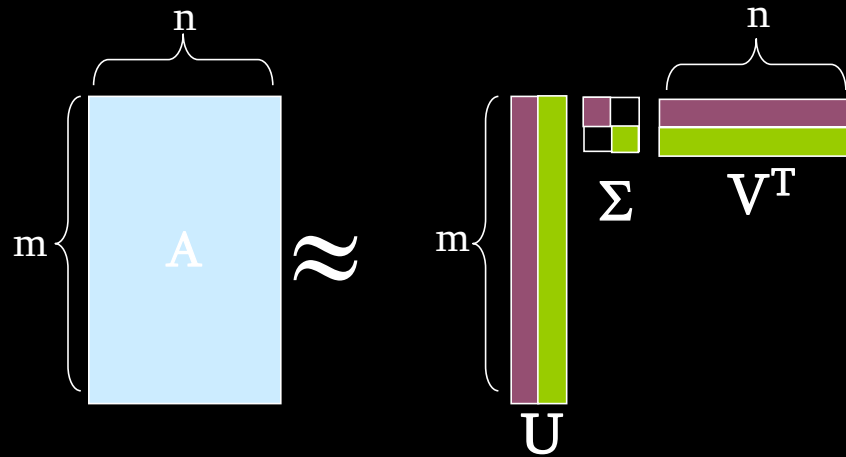
Singular Value Decomposition (SVD)

- Columns of U and V are orthonormal bases
- Singular values are the "strength" of each singular vector



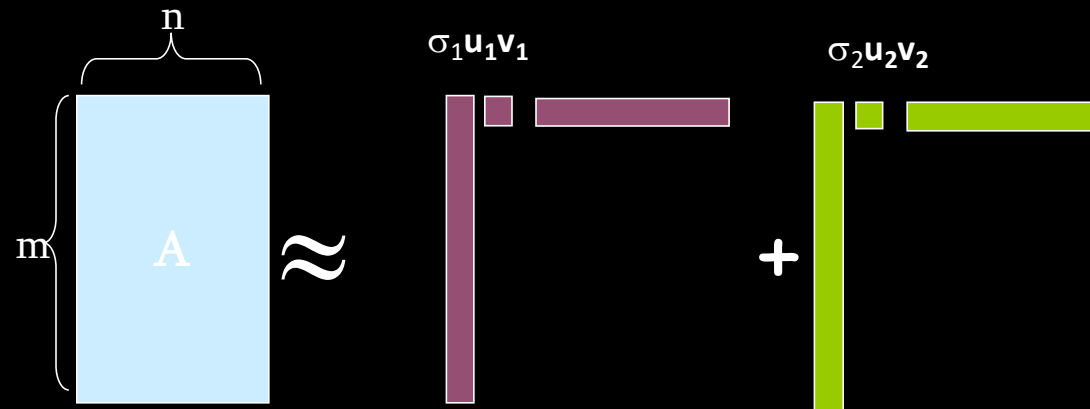
SVD

$$\mathbf{A} \approx \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T = \sum_i \sigma_i \mathbf{u}_i \circ \mathbf{v}_i^T$$



SVD

$$\mathbf{A} \approx \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T = \sum_i \sigma_i \mathbf{u}_i \circ \mathbf{v}_i^T$$



σ_i ... scalar
 \mathbf{u}_i ... vector
 \mathbf{v}_i ... vector

SVD - Properties

It is **always** possible to decompose a real matrix \mathbf{A} into $\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$, where

- $\mathbf{U}, \mathbf{\Sigma}, \mathbf{V}$: **unique**
- \mathbf{U}, \mathbf{V} : **column orthonormal**
 - $\mathbf{U}^T \mathbf{U} = \mathbf{I}; \mathbf{V}^T \mathbf{V} = \mathbf{I}$ (\mathbf{I} : identity matrix)
 - (Columns are orthogonal unit vectors)
- $\mathbf{\Sigma}$: **diagonal**
 - Entries (**singular values**) are **positive**, and sorted in decreasing order ($\sigma_1 \geq \sigma_2 \geq \dots \geq 0$)

Nice proof of uniqueness: <http://www.mpi-inf.mpg.de/~bast/ir-seminar-ws04/lecture2.pdf>

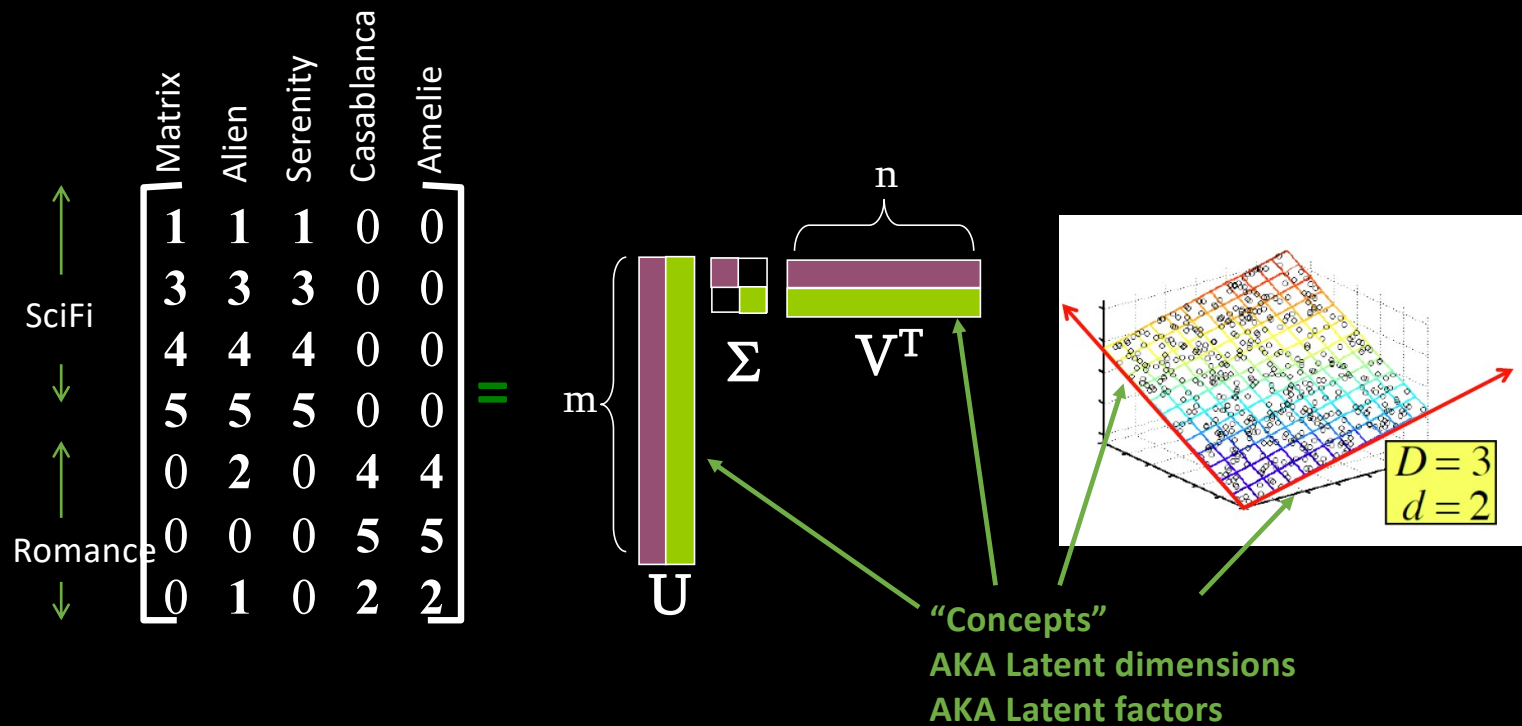
Why randomize SVD?

- **Runtime**

- We're good at generating [pseudo-]random numbers
- Can easily parallelize / distribute matrix algebra
- SVD, like PCA, runs $O(n^3)$, making anything beyond $\sim 10^3$ infeasible

SVD – Example: Users-to-Movies

- $A = U \Sigma V^T$ - example: Users to Movies



SVD – Example: Users-to-Movies

• **$A = U \Sigma V^T$ - example: Users to Movies**

$$\begin{array}{c}
 \begin{array}{c} \uparrow \\ \text{SciFi} \\ \downarrow \\ \uparrow \\ \text{Romnce} \\ \downarrow \end{array}
 \begin{bmatrix}
 \text{Matrix} & \text{Alien} & \text{Serenity} & \text{Casablanca} & \text{Amelie} \\
 1 & 1 & 1 & 0 & 0 \\
 3 & 3 & 3 & 0 & 0 \\
 4 & 4 & 4 & 0 & 0 \\
 5 & 5 & 5 & 0 & 0 \\
 0 & 2 & 0 & 4 & 4 \\
 0 & 0 & 0 & 5 & 5 \\
 0 & 1 & 0 & 2 & 2
 \end{bmatrix}
 =
 \begin{bmatrix}
 0.13 & 0.02 & -0.01 \\
 0.41 & 0.07 & -0.03 \\
 0.55 & 0.09 & -0.04 \\
 0.68 & 0.11 & -0.05 \\
 0.15 & -0.59 & 0.65 \\
 0.07 & -0.73 & -0.67 \\
 0.07 & -0.29 & 0.32
 \end{bmatrix}
 \times
 \begin{bmatrix}
 12.4 & 0 & 0 \\
 0 & 9.5 & 0 \\
 0 & 0 & 1.3
 \end{bmatrix}
 \times
 \begin{bmatrix}
 0.56 & 0.59 & 0.56 & 0.09 & 0.09 \\
 0.12 & -0.02 & 0.12 & -0.69 & -0.69 \\
 0.40 & -0.80 & 0.40 & 0.09 & 0.09
 \end{bmatrix}
 \end{array}$$

SVD – Example: Users-to-Movies

• **$A = U \Sigma V^T$ - example: Users to Movies**

Matrix A (User ratings):

	Matrix	Alien	Serenity	Casablanca	Amelie
SciFi	1	1	1	0	0
	3	3	3	0	0
	4	4	4	0	0
	5	5	5	0	0
Romnce	0	2	0	4	4
	0	0	0	5	5
	0	1	0	2	2

Matrix U (User latent factors):

0.13	0.02	-0.01
0.41	0.07	-0.03
0.55	0.09	-0.04
0.68	0.11	-0.05
0.15	-0.59	0.65
0.07	-0.73	-0.67
0.07	-0.29	0.32

Matrix Σ (Singular values):

12.4	0	0
0	9.5	0
0	0	1.3

Matrix V^T (Movie latent factors):

0.56	0.59	0.56	0.09	0.09
0.12	-0.02	0.12	-0.69	-0.69
0.40	-0.80	0.40	0.09	0.09

Equation: $A = U \Sigma V^T$

SVD – Example: Users-to-Movies

- $A = U \Sigma V^T$ - example: U is “user-to-concept” similarity matrix

$$\begin{array}{c} \uparrow \\ \text{SciFi} \\ \downarrow \\ \uparrow \\ \text{Romnce} \\ \downarrow \end{array} \begin{array}{c} \text{Matrix} \\ \text{Alien} \\ \text{Serenity} \\ \text{Casablanca} \\ \text{Amelie} \end{array} \begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 2 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 1 & 0 & 2 & 2 \end{bmatrix} = \begin{array}{c} \text{SciFi-concept} \\ \text{Romance-concept} \end{array} \begin{bmatrix} 0.13 & 0.02 & -0.01 \\ 0.41 & 0.07 & -0.03 \\ 0.55 & 0.09 & -0.04 \\ 0.68 & 0.11 & -0.05 \\ 0.15 & -0.59 & 0.65 \\ 0.07 & -0.73 & -0.67 \\ 0.07 & -0.29 & 0.32 \end{bmatrix} \times \begin{bmatrix} 12.4 & 0 & 0 \\ 0 & 9.5 & 0 \\ 0 & 0 & 1.3 \end{bmatrix} \times \begin{bmatrix} 0.56 & 0.59 & 0.56 & 0.09 & 0.09 \\ 0.12 & -0.02 & 0.12 & -0.69 & -0.69 \\ 0.40 & -0.80 & 0.40 & 0.09 & 0.09 \end{bmatrix}$$

SVD – Example: Users-to-Movies

• $A = U \Sigma V^T$ - example:

SciFi

Romnce

$$\begin{bmatrix}
 1 & 1 & 1 & 0 & 0 \\
 3 & 3 & 3 & 0 & 0 \\
 4 & 4 & 4 & 0 & 0 \\
 5 & 5 & 5 & 0 & 0 \\
 0 & 2 & 0 & 4 & 4 \\
 0 & 0 & 0 & 5 & 5 \\
 0 & 1 & 0 & 2 & 2
 \end{bmatrix}
 =
 \begin{bmatrix}
 0.13 & 0.02 & -0.01 \\
 0.41 & 0.07 & -0.03 \\
 0.55 & 0.09 & -0.04 \\
 0.68 & 0.11 & -0.05 \\
 0.15 & -0.59 & 0.65 \\
 0.07 & -0.73 & -0.67 \\
 0.07 & -0.29 & 0.32
 \end{bmatrix}
 \begin{bmatrix}
 12.4 & 0 & 0 \\
 0 & 9.5 & 0 \\
 0 & 0 & 1.3
 \end{bmatrix}
 \begin{bmatrix}
 0.56 & 0.59 & 0.56 & 0.09 & 0.09 \\
 0.12 & -0.02 & 0.12 & -0.69 & -0.69 \\
 0.40 & -0.80 & 0.40 & 0.09 & 0.09
 \end{bmatrix}$$

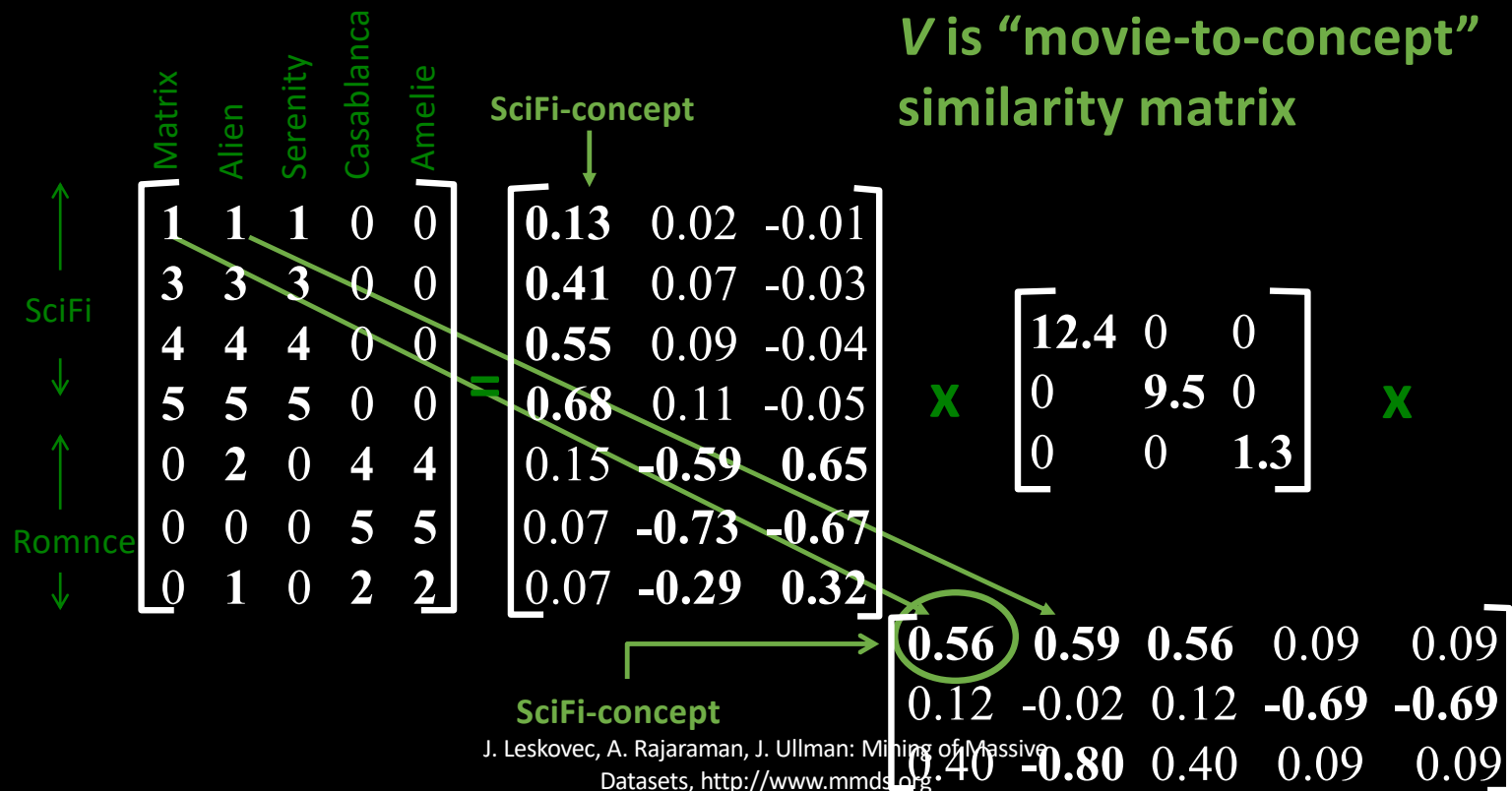
SciFi-concept

"strength" of the SciFi-concept

J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massive Datasets, <http://www.mmids.org>

SVD – Example: Users-to-Movies

• $A = U \Sigma V^T$ - example:



SVD - Interpretation #1

‘**movies**’, ‘**users**’ and ‘**concepts**’:

- U : user-to-concept similarity matrix
- V : movie-to-concept similarity matrix
- Σ : its diagonal elements:
‘strength’ of each concept

SVD - Interpretation #2

More details

- Q: How exactly is dim. reduction done?
- A: Set smallest singular values to zero

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 2 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 1 & 0 & 2 & 2 \end{bmatrix} = \begin{bmatrix} 0.13 & 0.02 & -0.01 \\ 0.41 & 0.07 & -0.03 \\ 0.55 & 0.09 & -0.04 \\ 0.68 & 0.11 & -0.05 \\ 0.15 & -0.59 & 0.65 \\ 0.07 & -0.73 & -0.67 \\ 0.07 & -0.29 & 0.32 \end{bmatrix} \times \begin{bmatrix} 12.4 & 0 & 0 \\ 0 & 9.5 & 0 \\ 0 & 0 & \cancel{1.3} \end{bmatrix} \times \begin{bmatrix} 0.56 & 0.59 & 0.56 & 0.09 & 0.09 \\ 0.12 & -0.02 & 0.12 & -0.69 & -0.69 \\ 0.40 & -0.80 & 0.40 & 0.09 & 0.09 \end{bmatrix}$$

SVD - Interpretation #2

More details

- Q: How exactly is dim. reduction done?
- A: Set smallest singular values to zero

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 2 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 1 & 0 & 2 & 2 \end{bmatrix} \approx \begin{bmatrix} 0.13 & 0.02 & -0.01 \\ 0.41 & 0.07 & -0.03 \\ 0.55 & 0.09 & -0.04 \\ 0.68 & 0.11 & -0.05 \\ 0.15 & -0.59 & 0.65 \\ 0.07 & -0.73 & -0.67 \\ 0.07 & -0.29 & 0.32 \end{bmatrix} \times \begin{bmatrix} 12.4 & 0 & 0 \\ 0 & 9.5 & 0 \\ 0 & 0 & \cancel{1.3} \end{bmatrix} \times \begin{bmatrix} 0.56 & 0.59 & 0.56 & 0.09 & 0.09 \\ 0.12 & -0.02 & 0.12 & -0.69 & -0.69 \\ 0.40 & -0.80 & 0.40 & 0.09 & 0.09 \end{bmatrix}$$

SVD - Interpretation #2

More details

- Q: How exactly is dim. reduction done?
- A: Set smallest singular values to zero

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 2 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 1 & 0 & 2 & 2 \end{bmatrix} \approx \begin{bmatrix} 0.13 & 0.02 & -0.01 \\ 0.41 & 0.07 & -0.03 \\ 0.55 & 0.09 & -0.04 \\ 0.68 & 0.11 & -0.05 \\ 0.15 & -0.59 & 0.65 \\ 0.07 & -0.73 & -0.67 \\ 0.07 & -0.29 & 0.32 \end{bmatrix} \times \begin{bmatrix} 12.4 & 0 & 0 \\ 0 & 9.5 & 0 \\ 0 & 0 & 1.3 \end{bmatrix} \times \begin{bmatrix} 0.56 & 0.59 & 0.56 & 0.09 & 0.09 \\ 0.12 & -0.02 & 0.12 & -0.69 & -0.69 \\ 0.40 & -0.80 & 0.40 & 0.09 & 0.09 \end{bmatrix}$$

J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massive
Datasets, <http://www.mmids.org>

SVD - Interpretation #2

More details

- Q: How exactly is dim. reduction done?
- A: Set smallest singular values to zero

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 2 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 1 & 0 & 2 & 2 \end{bmatrix} \approx \begin{bmatrix} 0.13 & 0.02 \\ 0.41 & 0.07 \\ 0.55 & 0.09 \\ 0.68 & 0.11 \\ 0.15 & -0.59 \\ 0.07 & -0.73 \\ 0.07 & -0.29 \end{bmatrix} \times \begin{bmatrix} 12.4 & 0 \\ 0 & 9.5 \end{bmatrix} \times \begin{bmatrix} 0.56 & 0.59 & 0.56 & 0.09 & 0.09 \\ 0.12 & -0.02 & 0.12 & -0.69 & -0.69 \end{bmatrix}$$

SVD - Interpretation #2

More details

- **Q: How exactly is dim. reduction done?**
- **A: Set smallest singular values to zero**

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 2 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 1 & 0 & 2 & 2 \end{bmatrix} \approx \begin{bmatrix} 0.92 & 0.95 & 0.92 & 0.01 & 0.01 \\ 2.91 & 3.01 & 2.91 & -0.01 & -0.01 \\ 3.90 & 4.04 & 3.90 & 0.01 & 0.01 \\ 4.82 & 5.00 & 4.82 & 0.03 & 0.03 \\ 0.70 & 0.53 & 0.70 & 4.11 & 4.11 \\ -0.69 & 1.34 & -0.69 & 4.78 & 4.78 \\ 0.32 & 0.23 & 0.32 & 2.01 & 2.01 \end{bmatrix}$$

$$\|A-B\|_F = \sqrt{\sum_{ij} (A_{ij}-B_{ij})^2}$$

is "small"

Frobenius norm:

$$\|M\|_F = \sqrt{\sum_{ij} M_{ij}^2}$$

SVD – Best Low Rank Approx.

$$A = U \Sigma V^T$$

The diagram shows three matrices: a square matrix A , a tall rectangular matrix U , a small square matrix Σ labeled "Sigma", and a wide rectangular matrix V^T . They are arranged in the equation $A = U \Sigma V^T$.

B is best approximation of A

$$B = U \Sigma_k V^T$$

The diagram shows the best low-rank approximation B of matrix A . It is represented as $B = U \Sigma_k V^T$, where U is a tall rectangular matrix, Σ_k is a small square matrix labeled "Sigma" with a blue top-left corner, and V^T is a wide rectangular matrix with a blue top section and a black bottom section. The entire equation is enclosed in a blue rounded rectangle.

Relationship to PCA

- SVD can be applied to *any* matrix; PCA only works on symmetric covariance matrices

- However, there is a relationship

$$M^T M = V \Sigma^T U^T U \Sigma V^T = V (\Sigma^T \Sigma) V^T$$
$$M M^T = U \Sigma V^T V \Sigma^T U^T = U (\Sigma \Sigma^T) U^T$$

- Columns of V are eigenvectors of $M^T M$
- Columns of U are eigenvectors of $M M^T$
- Singular values are square roots of eigenvalues of $M^T M$ or $M M^T$

SVD: Drawbacks

+ Optimal low-rank approximation

in terms of Frobenius norm

- Interpretability problem:

- A singular vector specifies a linear combination of all input columns or rows

- Lack of sparsity:

- Singular vectors are **dense!**

$$\begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix} = \begin{bmatrix} \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \end{bmatrix} \begin{bmatrix} \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \end{bmatrix} \begin{bmatrix} \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \end{bmatrix} V^T$$

CUR Decomposition

Frobenius norm:

$$\|X\|_F = \sqrt{\sum_{ij} X_{ij}^2}$$

- **Goal: Express A as a product of matrices C,U,R**

Make $\|A - C \cdot U \cdot R\|_F$ small

- **“Constraints” on C and R:**

$$\begin{pmatrix} | & | & | \end{pmatrix} \approx \begin{pmatrix} | & | & | & | & | & | \end{pmatrix} \cdot \begin{pmatrix} U \end{pmatrix} \cdot \begin{pmatrix} R \end{pmatrix}$$

A C U R

CUR Decomposition

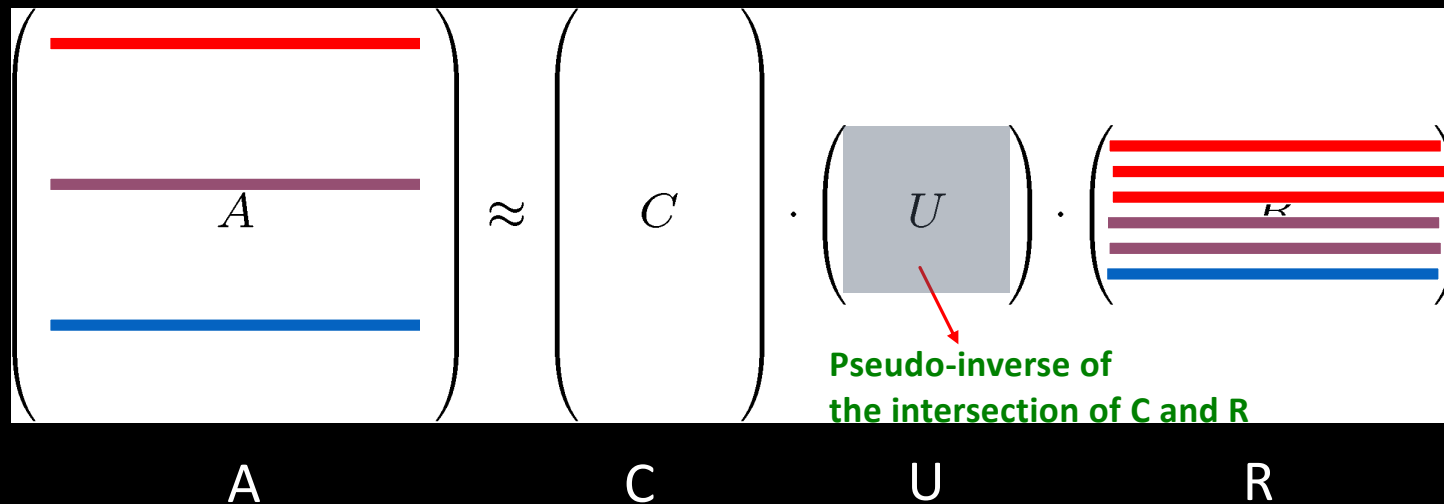
Frobenius norm:

$$\|X\|_F = \sqrt{\sum_{ij} X_{ij}^2}$$

- **Goal: Express A as a product of matrices C,U,R**

Make $\|A - C \cdot U \cdot R\|_F$ small

- “Constraints” on C and R:



CUR: How it Works

- Sampling columns (similarly for rows):

Note this is a randomized algorithm; the same column can be sampled more than once

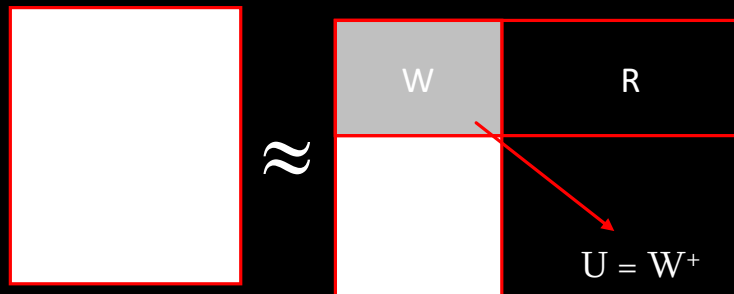
Input: matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$, sample size c

Output: $\mathbf{C}_d \in \mathbb{R}^{m \times c}$

1. for $x = 1 : n$ [column distribution]
2. $P(x) = \sum_i \mathbf{A}(i, x)^2 / \sum_{i,j} \mathbf{A}(i, j)^2$
3. for $i = 1 : c$ [sample columns]
4. Pick $j \in 1 : n$ based on distribution $P(x)$
5. Compute $\mathbf{C}_d(:, i) = \mathbf{A}(:, j) / \sqrt{cP(j)}$

Computing U

- Let \mathbf{W} be the “intersection” of sampled columns \mathbf{C} and rows \mathbf{R}
 - Let SVD of $\mathbf{W} = \mathbf{X} \mathbf{Z} \mathbf{Y}^T$
- Then: $\mathbf{U} = \mathbf{W}^+ = \mathbf{Y} \mathbf{Z}^+ \mathbf{X}^T$
 - \mathbf{Z}^+ : reciprocals of non-zero singular values: $Z_{ii}^+ = 1/Z_{ii}$
 - \mathbf{W}^+ is the “pseudoinverse”



Why pseudoinverse works?

$\mathbf{W} = \mathbf{X} \mathbf{Z} \mathbf{Y}$ then $\mathbf{W}^{-1} = \mathbf{X}^{-1} \mathbf{Z}^{-1} \mathbf{Y}^{-1}$

Due to orthonormality

$\mathbf{X}^{-1} = \mathbf{X}^T$ and $\mathbf{Y}^{-1} = \mathbf{Y}^T$

Since \mathbf{Z} is diagonal $\mathbf{Z}^{-1} = 1/Z_{ii}$

Thus, if \mathbf{W} is nonsingular, pseudoinverse is the true inverse

CUR: Pros & Cons

+ Easy interpretation

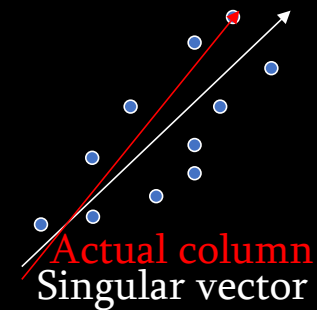
- Since the basis vectors are actual columns and rows

+ Sparse basis

- Since the basis vectors are actual columns and rows

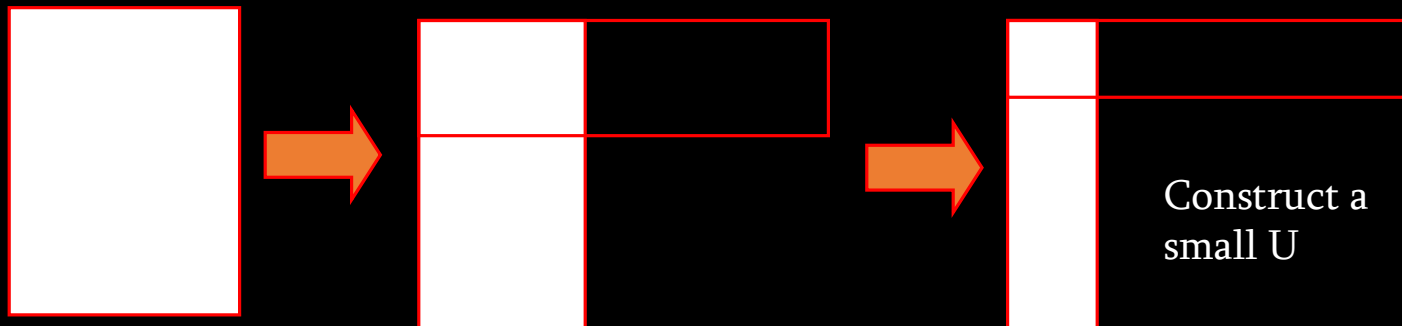
- Duplicate columns and rows

- Columns of large norms will be sampled many times



Solution

- **If we want to get rid of the duplicates:**
 - Throw them away
 - Scale (multiply) the columns/rows by the square root of the number of duplicates

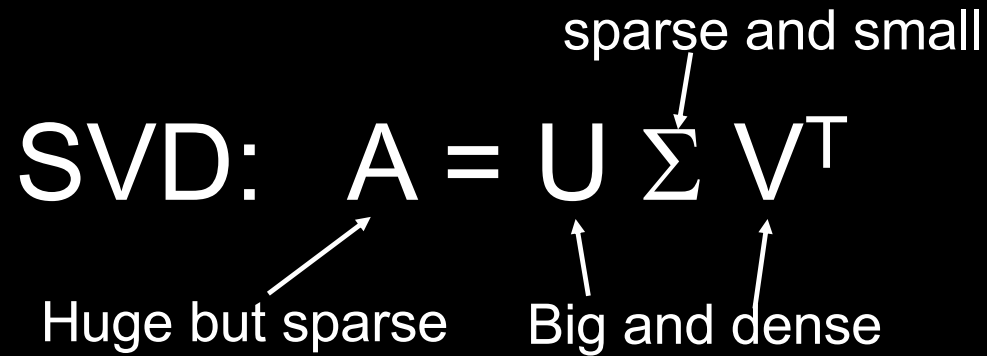


SVD vs. CUR

sparse and small

SVD: $A = U \Sigma V^T$

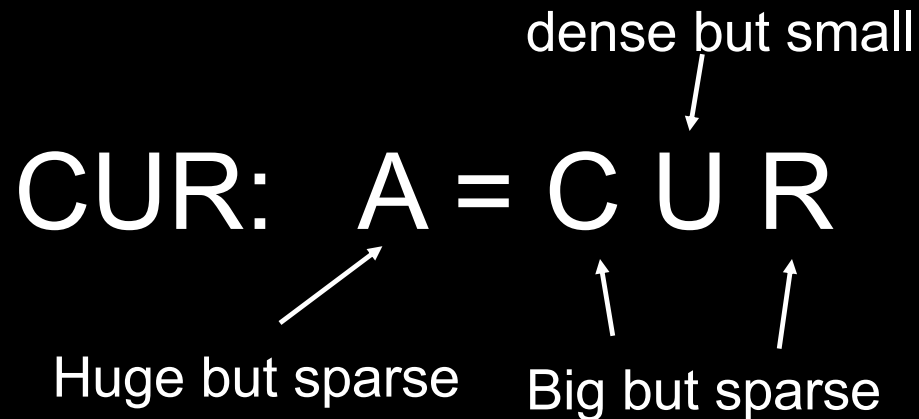
Huge but sparse Big and dense



dense but small

CUR: $A = C U R$

Huge but sparse Big but sparse



Stochastic SVD (SSVD)

- Uses **random projections** to find close approximation to SVD
- Combination of probabilistic strategies to maximize convergence likelihood
- Easily scalable to *massive* linear systems

Basic goal

- Matrix A
 - Find a low-rank approximation of A
 - Basic dimensionality reduction

$$\|A - QQ^*A\| < \epsilon$$



Preconditioning

Approximating range of A

- INPUT: A, k, p
- OUTPUT: Q

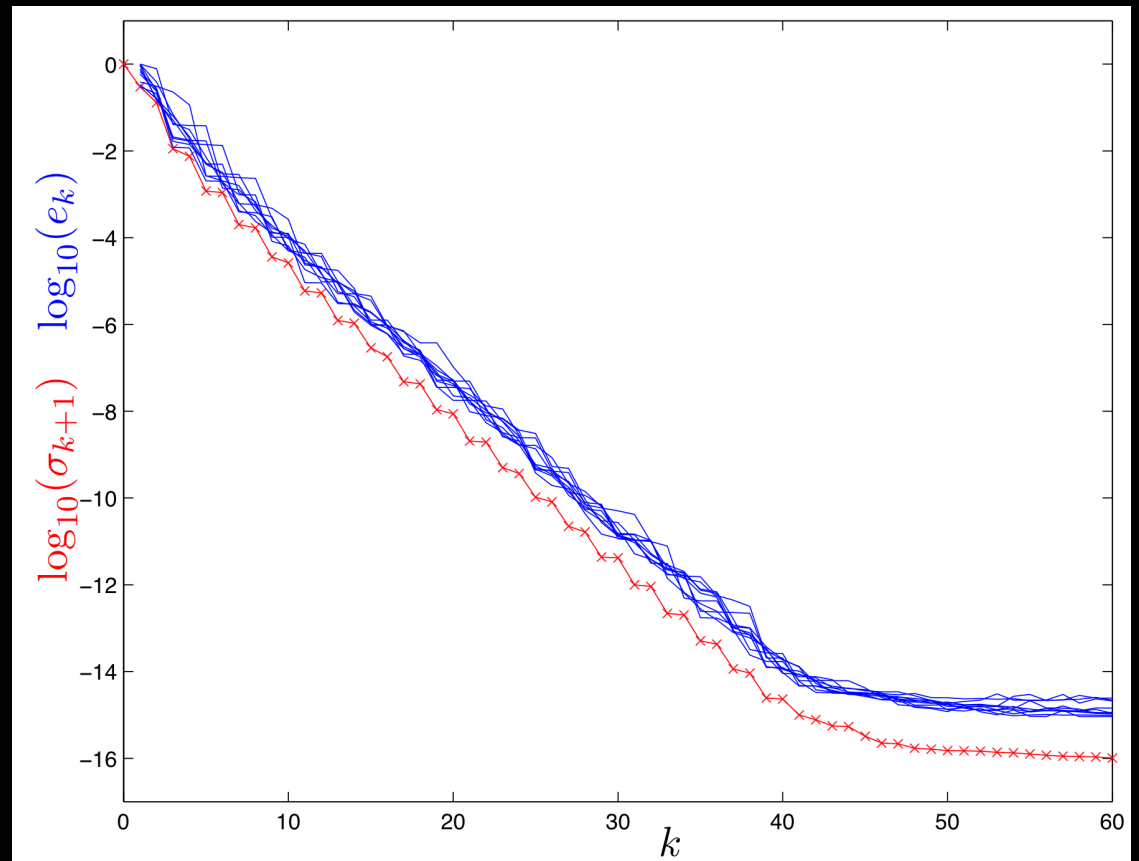
1. Draw Gaussian $n \times k$ test matrix Ω
2. Form product $Y = A\Omega$
3. Orthogonalize columns of $Y \rightarrow Q$

Approximating SVD of A

- INPUT: Q
 - OUTPUT: Singular vectors U
1. Form $k \times n$ matrix $B = Q^T A$
 2. Compute SVD of $B = \hat{U} \Sigma V^T$
 3. Compute singular vectors $U = Q \hat{U}$

Empirical Results

- 1000x1000 matrix
- Several runs of empirical results (blue) to theoretical lower bound (red)
- **Error seems to be systemic**



Power iterations

- Affects decay of eigenvalues / singular values

$$Y = \Omega.$$

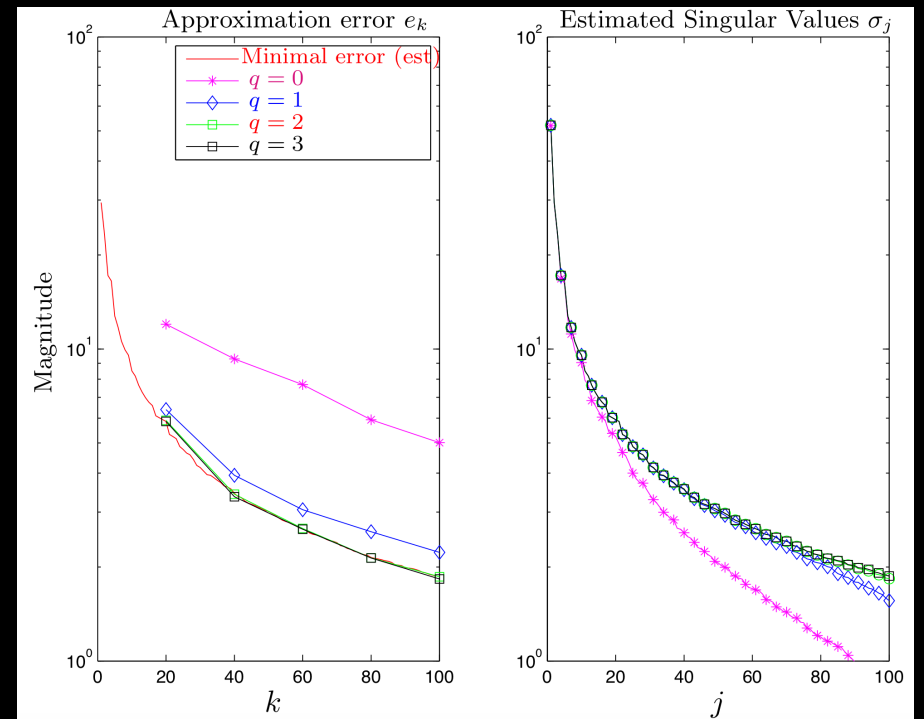
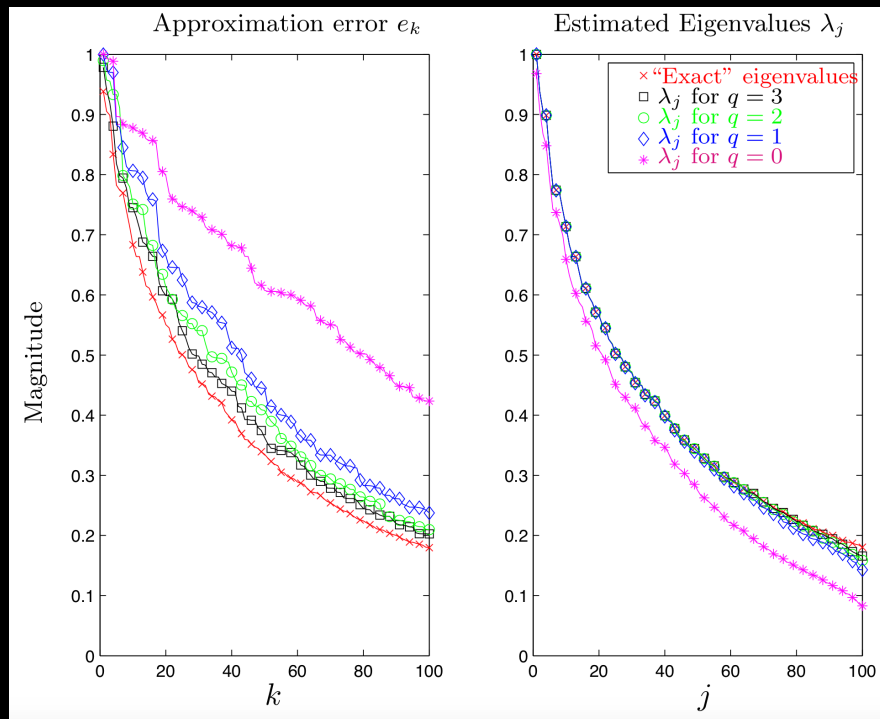
$$Y = (A A^*)^q A \Omega$$

Power iterations

$$\begin{aligned}\mathbb{E}\|A - QQ^T A\|_2 &\leq \left(1 + \sqrt{\frac{k}{p-1}}\right) \sigma_{k+1} + \frac{e\sqrt{k+p}}{p} \cdot \left(\sum_{j>k} \sigma_j^2\right)^{1/2} \\ &\leq \left[1 + \frac{4\sqrt{k+p}}{p-1} \cdot \sqrt{\min\{m, n\}}\right] \sigma_{k+1} \\ &= C \cdot \sigma_{k+1}.\end{aligned}$$

Upshot: after only a single power iteration, the error is proportional to the *next* [uncomputed] singular value (times a constant C).

Empirical Results



Why does this work?

- Three primary reasons:

1. Johnson-Lindenstrauss Lemma

- Low-dimensional embeddings preserve pairwise distances

$$(1 - \varepsilon)\|u - v\|^2 \leq \|f(u) - f(v)\|^2 \leq (1 + \varepsilon)\|u - v\|^2$$

2. Concentration of measure

- Geometric interpretation of classical idea: regular functions of independent random variables rarely deviate far from their means

3. Preconditioning

- Condition number: how much change in output is produced from change in input (relation to #1)
- Q matrix lowers condition number while preserving overall system

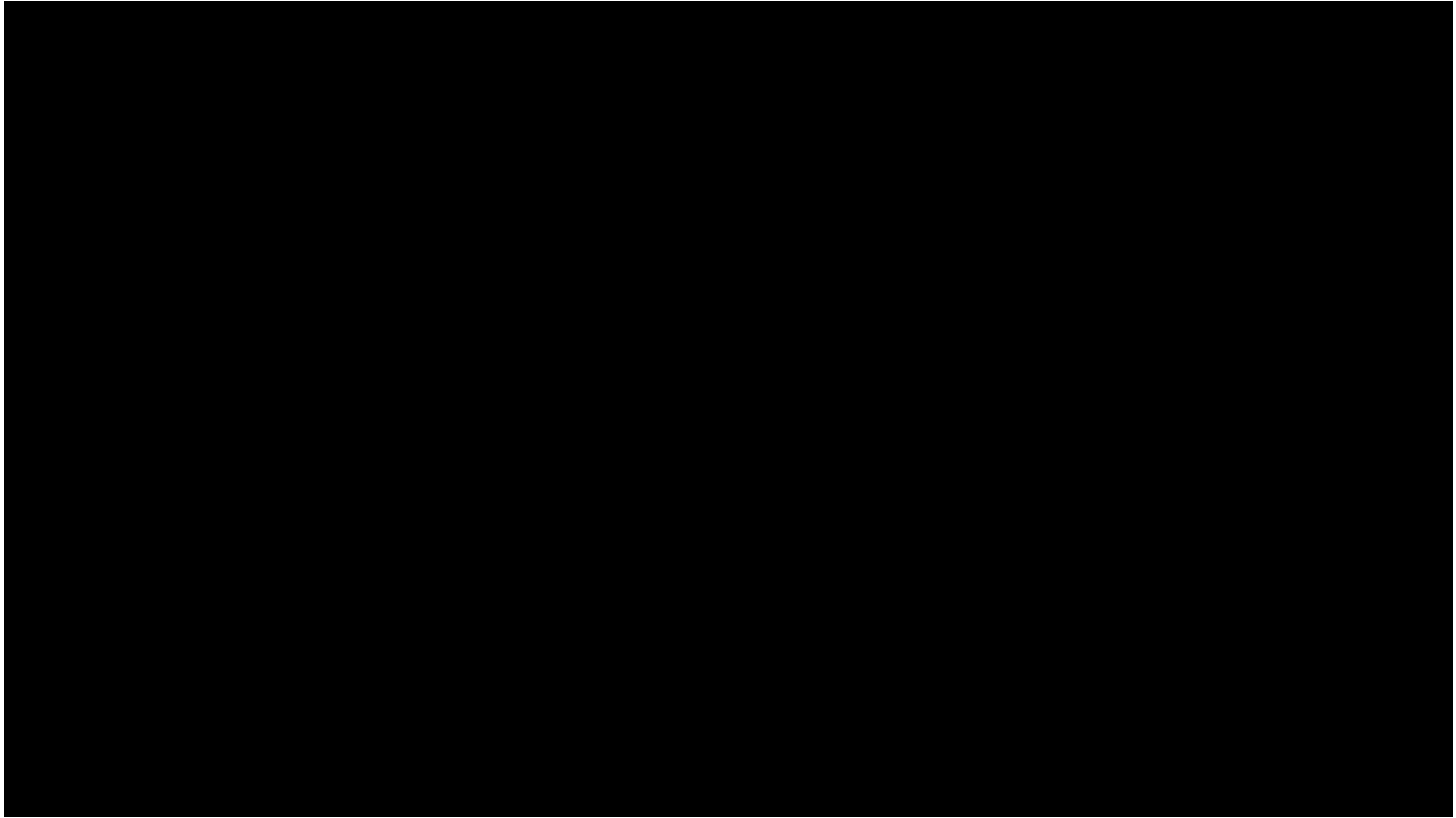
$$\kappa = \frac{|\lambda_{\max}|}{|\lambda_{\min}|}$$

Summary

- Relationship of SVD and PCA
 - PCA: eigenvectors and eigenvalues of the covariance matrix (or kernel matrix, for Kernel PCA)
 - SVD: Low-rank approximation for *any* matrix
- CUR
 - Randomly sample columns of data matrix A to use as basis
 - Interpretable and sparse, but potentially oversample high-magnitude columns
- SVD via SGD
 - Reframe SVD as a matrix completion problem
 - Use SGD in alternating least-squares to infer "missing" components
- SSVD
 - Full-blown Johnson-Lindenstrauss exploitation
 - Use random projections to approximate SVD to high accuracy
 - Requires some empirical tweaks (oversampling, power iterations)

References

- “Randomized methods for computing low-rank approximations of matrices”,
https://amath.colorado.edu/faculty/martinss/Pubs/2012_halko_dissertation.pdf
- “CUR decomposition for compression and compressed sensing of large-scale traffic data”, <https://dspace.mit.edu/openaccess-disseminate/1721.1/86879>



Large-Scale Matrix Factorization with Distributed Stochastic Gradient Descent

Rainer Gemulla

talk pilfered from →



Peter J. Haas



Yannis Sismanis



Erik Nijkamp



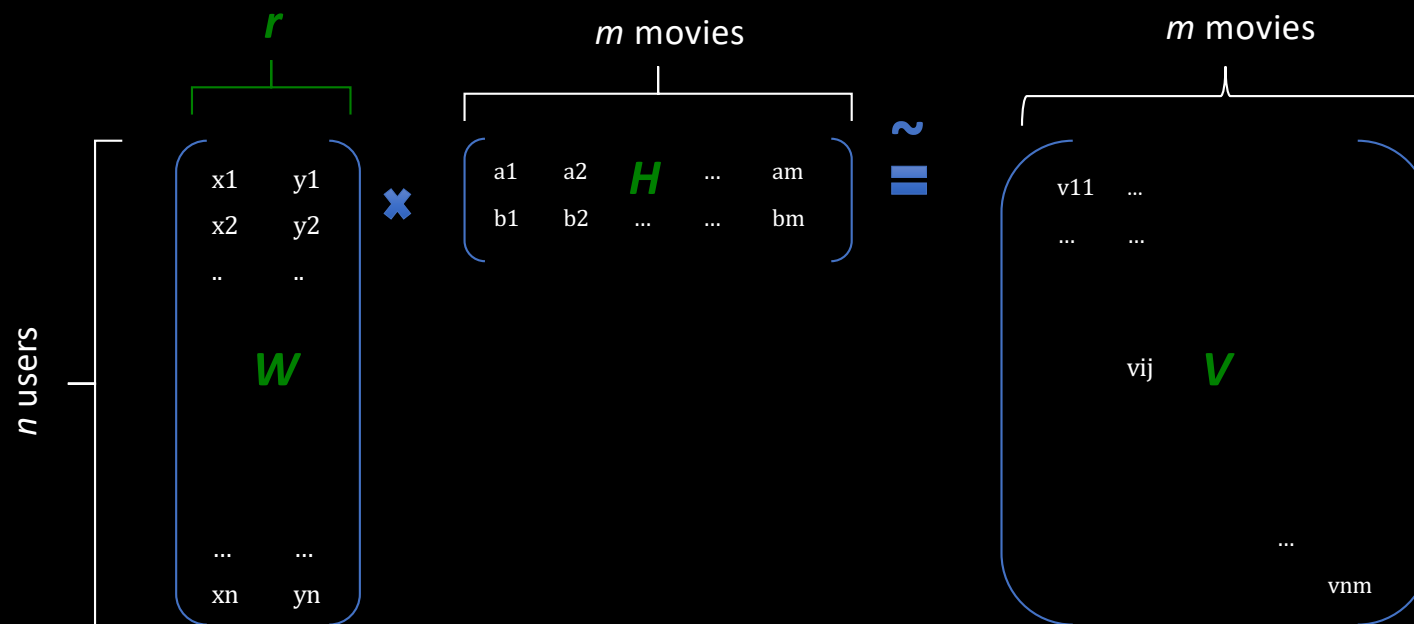
Collaborative Filtering

- ▶ Problem
 - ▶ Set of users
 - ▶ Set of items (movies, books, jokes, products, stories, ...)
 - ▶ Feedback (ratings, purchase, click-through, tags, ...)
- ▶ Predict additional items a user may like
 - ▶ Assumption: Similar feedback \implies Similar taste
- ▶ Example

	<i>Avatar</i>	<i>The Matrix</i>	<i>Up</i>
<i>Alice</i>	<i>(</i> ? <i>)</i>	4	2
<i>Bob</i>	3	2	<i>)</i> ?
<i>Charlie</i>	5	<i>(</i> ? <i>)</i>	3

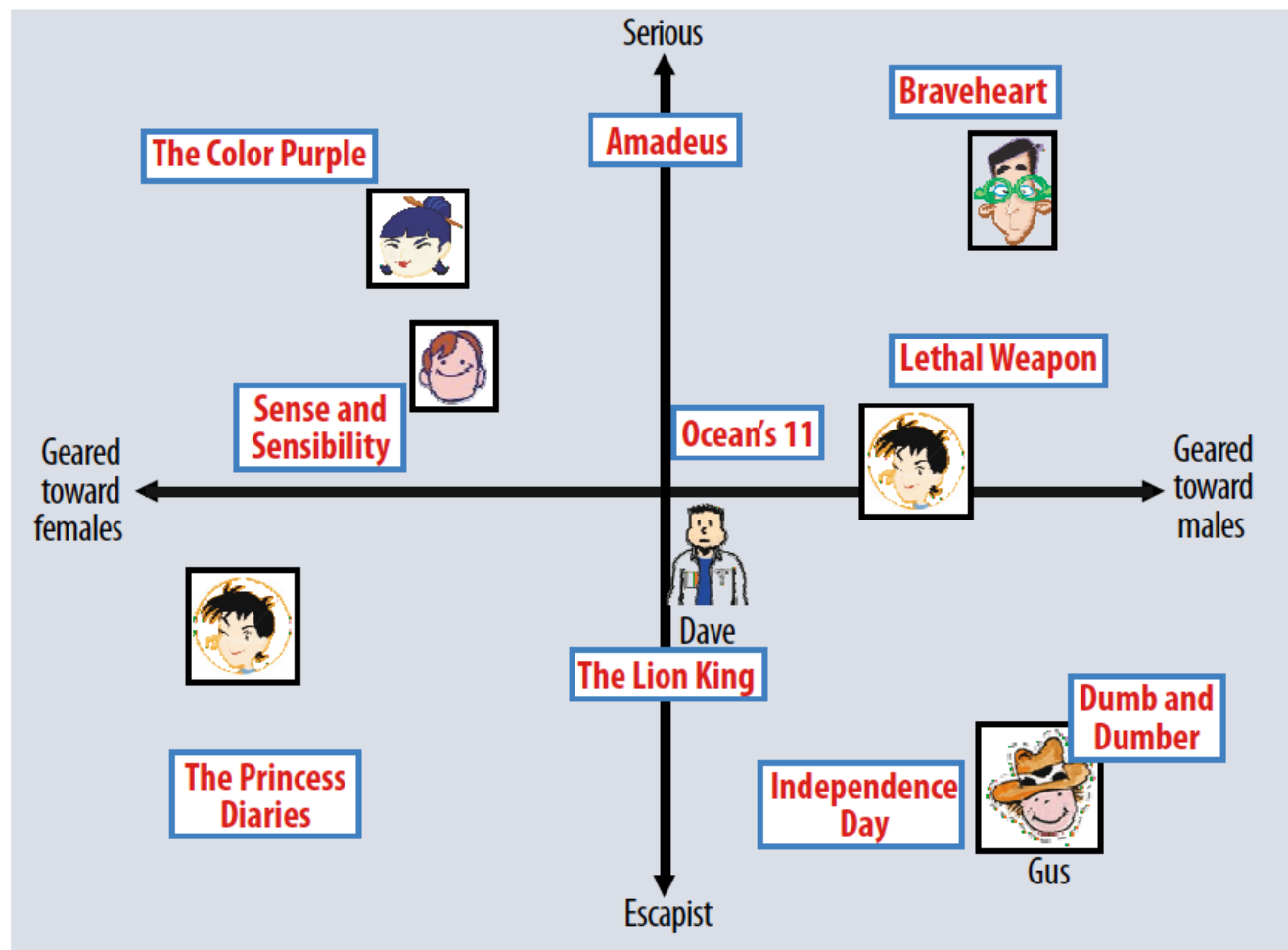
- ▶ Netflix competition: 500k users, 20k movies, 100M movie ratings, 3M question marks

Recovering latent factors in a matrix



$V[i,j]$ = user i 's rating of movie j

Semantic Factors (Koren et al., 2009)



Latent Factor Models

- ▶ Discover latent factors ($r = 1$)

	Avatar (2.24)	The Matrix (1.92)	Up (1.18)
Alice (1.98)		4 (3.8)	2 (2.3)
Bob (1.21)	3 (2.7)	2 (2.3)	
Charlie (2.30)	5 (5.2)		3 (2.7)

- ▶ Minimum loss

$$\min_{\mathbf{W}, \mathbf{H}} \sum_{(i,j) \in Z} (\mathbf{v}_{ij} - [\mathbf{WH}]_{ij})^2$$

Latent Factor Models

- Discover latent factors ($r = 1$)

	Avatar (2.24)	The Matrix (1.92)	Up (1.18)
Alice (1.98)	? (4.4)	4 (3.8)	2 (2.3)
Bob (1.21)	3 (2.7)	2 (2.3)	? (1.4)
Charlie (2.30)	5 (5.2)	? (4.4)	3 (2.7)

- Minimum loss

$$\min_{\mathbf{W}, \mathbf{H}, \mathbf{u}, \mathbf{m}} \sum_{(i,j) \in Z} (\mathbf{V}_{ij} - \mu - \mathbf{u}_i - \mathbf{m}_j - [\mathbf{WH}]_{ij})^2 + \lambda (\|\mathbf{W}\| + \|\mathbf{H}\| + \|\mathbf{u}\| + \|\mathbf{m}\|)$$

- Bias, regularization

Matrix factorization as SGD

require that the loss can be written as

$$L = \sum_{(i,j) \in Z} l(\mathbf{V}_{ij}, \mathbf{W}_{i*}, \mathbf{H}_{*j})$$

Algorithm 1 SGD for Matrix Factorization

Require: A training set Z , initial values \mathbf{W}_0 and \mathbf{H}_0

while not converged **do** {step}

 Select a training point $(i, j) \in Z$ uniformly at random.

$$\mathbf{W}'_{i*} \leftarrow \mathbf{W}_{i*} - \epsilon_n N \frac{\partial}{\partial \mathbf{W}_{i*}} l(\mathbf{V}_{ij}, \mathbf{W}_{i*}, \mathbf{H}_{*j})$$

$$\mathbf{H}_{*j} \leftarrow \mathbf{H}_{*j} - \epsilon_n N \frac{\partial}{\partial \mathbf{H}_{*j}} l(\mathbf{V}_{ij}, \mathbf{W}_{i*}, \mathbf{H}_{*j})$$

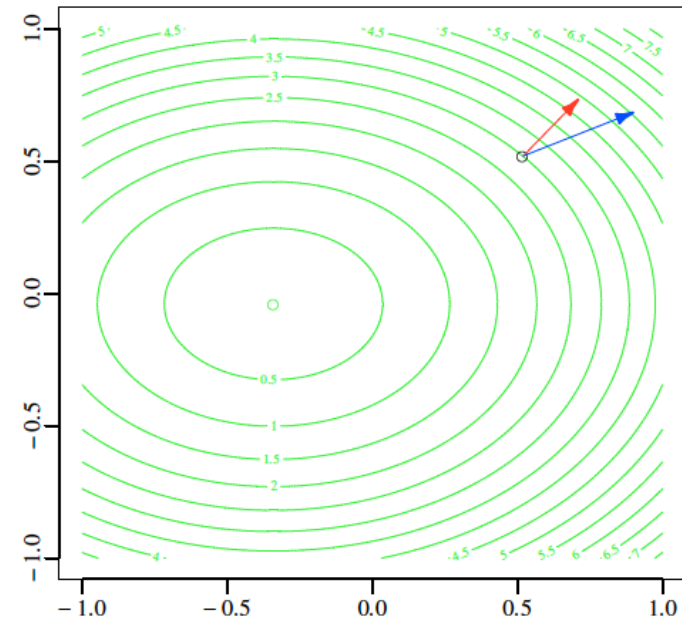
$$\mathbf{W}_{i*} \leftarrow \mathbf{W}'_{i*}$$

end while

why does this work?

Stochastic Gradient Descent

- ▶ Find minimum θ^* of function L
- ▶ Pick a starting point θ_0
- ▶ Approximate gradient $\hat{L}'(\theta_0)$

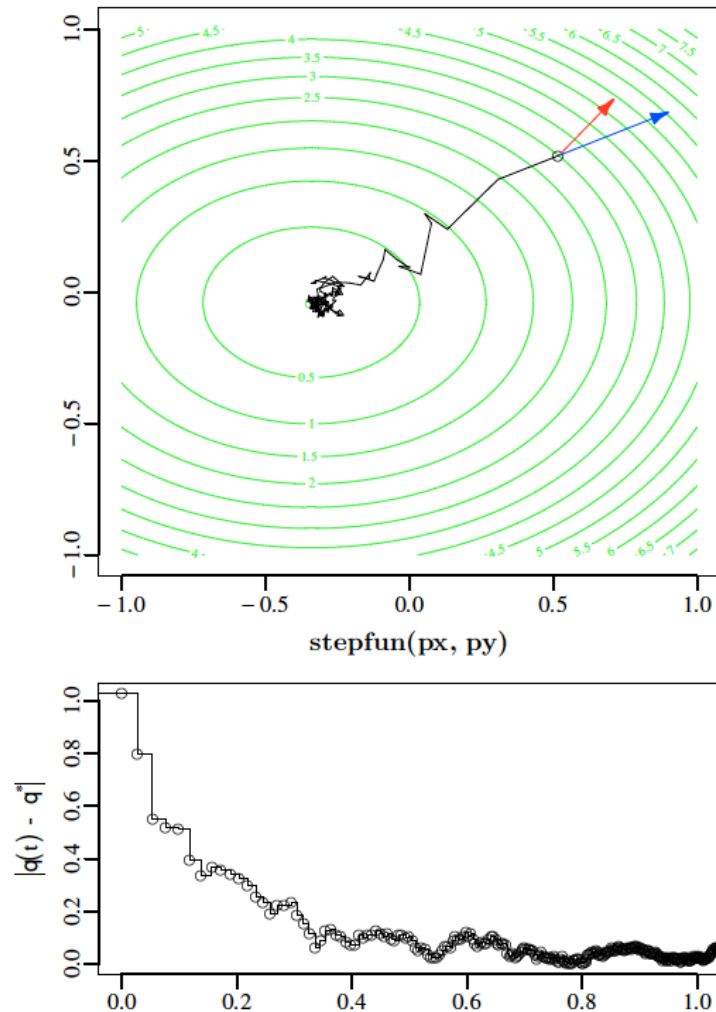


Stochastic Gradient Descent

- ▶ Find minimum θ^* of function L
- ▶ Pick a starting point θ_0
- ▶ Approximate gradient $\hat{L}'(\theta_0)$
- ▶ Jump “approximately” downhill
- ▶ Stochastic difference equation

$$\theta_{n+1} = \theta_n - \epsilon_n \hat{L}'(\theta_n)$$

- ▶ Under certain conditions, asymptotically approximates (continuous) gradient descent



Why does this work?

require that the loss can be written as

$$L = \sum_{(i,j) \in Z} l(\mathbf{V}_{ij}, \mathbf{W}_{i*}, \mathbf{H}_{*j})$$

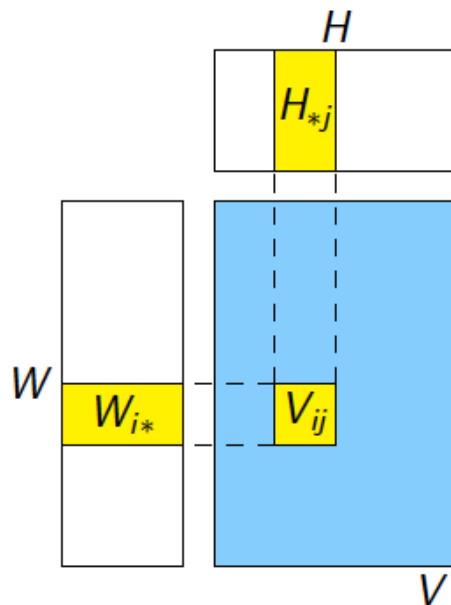
Algorithm 1 SGD for Matrix Factorization

Require: A training set Z , initial values \mathbf{W}_0 and \mathbf{H}_0
while not converged **do** {step}
 Select a training point $(i, j) \in Z$ uniformly at random.
 $\mathbf{W}'_{i*} \leftarrow \mathbf{W}_{i*} - \epsilon_n N \frac{\partial}{\partial \mathbf{W}_{i*}} l(\mathbf{V}_{ij}, \mathbf{W}_{i*}, \mathbf{H}_{*j})$
 $\mathbf{H}_{*j} \leftarrow \mathbf{H}_{*j} - \epsilon_n N \frac{\partial}{\partial \mathbf{H}_{*j}} l(\mathbf{V}_{ij}, \mathbf{W}_{i*}, \mathbf{H}_{*j})$
 $\mathbf{W}_{i*} \leftarrow \mathbf{W}'_{i*}$
end while

Key Claim

require that the loss can be written as

$$L = \sum_{(i,j) \in Z} l(\mathbf{V}_{ij}, \mathbf{W}_{i*}, \mathbf{H}_{*j})$$



$$\frac{\partial}{\partial W_{i'k}} L_{ij}(\mathbf{W}, \mathbf{H}) = \begin{cases} 0 & \text{if } i \neq i' \\ \frac{\partial}{\partial W_{ik}} l(\mathbf{V}_{ij}, \mathbf{W}_{i*}, \mathbf{H}_{*j}) & \text{otherwise} \end{cases}$$

$$\frac{\partial}{\partial H_{kj'}} L_{ij}(\mathbf{W}, \mathbf{H}) = \begin{cases} 0 & \text{if } j \neq j' \\ \frac{\partial}{\partial H_{kj}} l(\mathbf{V}_{ij}, \mathbf{W}_{i*}, \mathbf{H}_{*j}) & \text{otherwise} \end{cases}$$

Checking the claim

$$\frac{\partial}{\partial \mathbf{W}_{i*}} L(\mathbf{W}, \mathbf{H}) = \frac{\partial}{\partial \mathbf{W}_{i*}} \sum_{(i', j) \in Z} L_{i'j}(\mathbf{W}_{i'*}, \mathbf{H}_{*j}) = \sum_{j \in Z_{i*}} \frac{\partial}{\partial \mathbf{W}_{i*}} L_{ij}(\mathbf{W}_{i*}, \mathbf{H}_{*j}),$$

where $Z_{i*} = \{j : (i, j) \in Z\}$.

$$\frac{\partial}{\partial \mathbf{H}_{*j}} L(\mathbf{W}, \mathbf{H}) = \sum_{i \in Z_{*j}} \frac{\partial}{\partial \mathbf{W}_{*j}} L_{ij}(\mathbf{W}_{i*}, \mathbf{H}_{*j}),$$

where $Z_{*j} = \{i : (i, j) \in Z\}$.

Think for SGD for logistic regression

- LR loss = compare y and $\hat{y} = \text{dot}(\mathbf{w}, \mathbf{x})$
- similar but now update \mathbf{w} (user weights) and \mathbf{x} (movie weight)

Stochastic Gradient Descent on Netflix Data

