

Announcements:

- 0.25 EC point for each lecture you attend for rest of the term.
- links.eecs16b.org/lecture-ec
- Lab Design Contest: See Ed post (EC opportunities!)

Last time:

- SVD: 3 forms

- "full" $A = U \Sigma V^T$
- "compact" $A = U_r \Sigma_r V_r^T$
- "outer-product" $A = \sum_{i=1}^r \sigma_i \vec{u}_i \vec{v}_i^T$

- Geometry of SVD: $A\vec{x} = U \Sigma V^T \vec{x}$

Pseudoinverse: $A \in \mathbb{R}^{m \times n}$

• $A = U_r \Sigma_r V_r^T \Rightarrow A^+ = V_r \Sigma_r^{-1} U_r^T$

$(m \times n)$ $(n \times n)$ $(n \times m)$ $(r \times r)$

• $A = \sum_{i=1}^r \sigma_i \vec{u}_i \vec{v}_i^T \Rightarrow A^+ = \sum_{i=1}^r \frac{1}{\sigma_i} \vec{v}_i \vec{u}_i^T$

$r \leq \min(m, n)$

• $Ax = y \Rightarrow x = A^+ y$, holds for $\begin{cases} m < n \\ m > n \\ m = n \end{cases}$ and any rank (r)

Last lecture: $A\vec{x} = \vec{y}$, $A \in \mathbb{R}^{m \times n}$

$$A = U_r \Sigma_r V_r^T$$
$$A^+ = V_r \Sigma_r^{-1} U_r^T$$

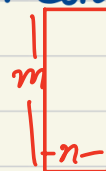
$\vec{x} = A^+ \vec{y}$, where A^+ was defined using SVD of A , is:

• unique solution of $A\vec{x} = \vec{y}$ when A is square and invertible ($m=n=r$) because $A^+ = A^{-1}$ in that case

• LS solution when A is tall ($m > n$), if full column rank also ($n=r$), then

$$A^+ = (A^T A)^{-1} A^T$$

which recovers LS solution studied before.

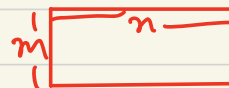


• Min. norm solution when A is wide ($n > m$) and infinitely many solutions exist. If full row rank ($m=r$),

$$A^+ = A^T (A A^T)^{-1}$$

Thus,

$$\vec{x}_{MN} = A^T (A A^T)^{-1} \vec{y}$$



⊗ $A = U_r \Sigma_r V_r^T \Rightarrow A = U_r \Sigma_r V^T$ (since $r=n$)

$$A^T = V \Sigma_r^T U_r^T \Rightarrow A^T A = V \Sigma_r U_r^T U_r \Sigma_r V^T$$
$$= V \Sigma_r^2 V^T$$

$$\Rightarrow (A^T A)^{-1} = V \Sigma_r^{-2} V^T$$

$$\Rightarrow (A^T A)^{-1} A^T = (V \Sigma_r^{-2} V^T) (V \Sigma_r U_r^T) = V \Sigma_r^{-1} U_r^T$$

$$\Rightarrow \vec{x}_{LS} = (A^T A)^{-1} A^T \vec{b} = A^+ \vec{b} = A^+$$

REMINDER!

If columns of Q are orthonormal,

$$Q^T Q = I$$

• $Q Q^T$ is a projection onto $\text{Col}(Q)$

TODAY:

• PCA (Principal Component Analysis)

SVD/PCA has applications in many domains

Healthcare: • predicting patient's health & susceptibility to diseases based on health risk factors.

Biology: • predicting which gene mutations are likely to cause cancer

Retail: • predicting which user will buy which product based on historical data

•
•

SVD
PCA

- data reduction technique
- data-driven generalization of the "Fourier transform" (FFT) tailored to specific problem
- used universally by big-tech companies!
 - Google: PageRank
 - Face Book: Face Recognition
 - Netflix: Recommender systems
 - Amazon → (Netflix prize)
 - \$\$\$
- simple & interpretable
- scalable

Low-rank Approximation

Given a high rank matrix $A \in \mathbb{R}^{m \times n}$ with

$$r \approx \min\{m, n\},$$

find an approximation with rank $l \ll \min\{m, n\}$.

$$A = U \begin{array}{|c|c|} \hline \begin{array}{c} 100 \\ 90.5 \\ 55 \\ \dots \\ 1 \\ 0.5 \\ 0.01 \\ \dots \end{array} & \mathbf{0} \\ \hline \end{array} V^T$$

$\leftarrow m \qquad \leftarrow (n-m) \rightarrow$

Suppose $m = n = r = 10,000$
 $l = 10$

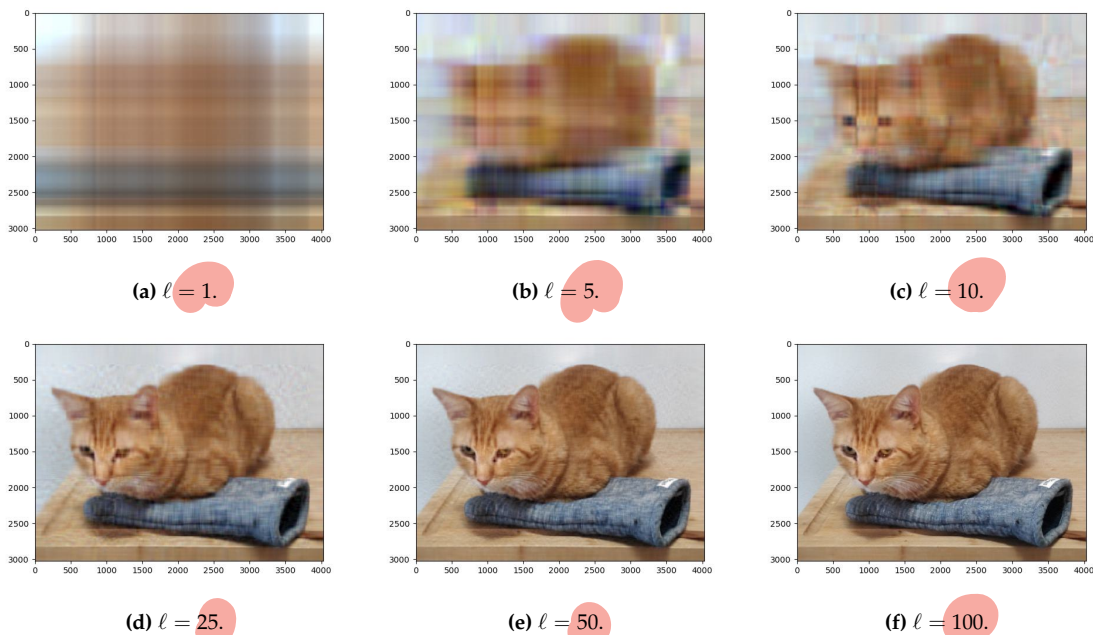
$$A \approx \sum_{i=1}^{10} \sigma_i \vec{u}_i \vec{v}_i^T$$

- A has 10^8 entries
- \vec{u}_i, \vec{v}_i of dimension 10,000 each, so 20,000 total per outer-product SVD; 10 terms \Rightarrow
200,000 versus 100,000,000

500x savings!

Figure 5: The author's friend's cat Snyder.

It can be represented as three matrices $A_R, A_G, A_B \in \mathbb{R}^{4032 \times 3024}$ corresponding to R, G, and B of the image. We perform a rank- ℓ approximation $A_R = U_{R;\ell} \Sigma_{R;\ell} V_{R;\ell}^T$, $A_G = U_{G;\ell} \Sigma_{G;\ell} V_{G;\ell}^T$, $A_B = U_{B;\ell} \Sigma_{B;\ell} V_{B;\ell}^T$, and then compose an image out of them, for different values of ℓ . The results are shown below.



By rank 100 approximation, the image is almost perfect. Now, the original image had $3 \times 4032 \times 3024 = 36578304$ entries; at rank 100, we have $3 \times 100 \times (4032 + 3024 + 1) = 2117100$ entries, so we need to store around 5% of the original image. Not bad!

See Appendix D for some code showing how these images were created.

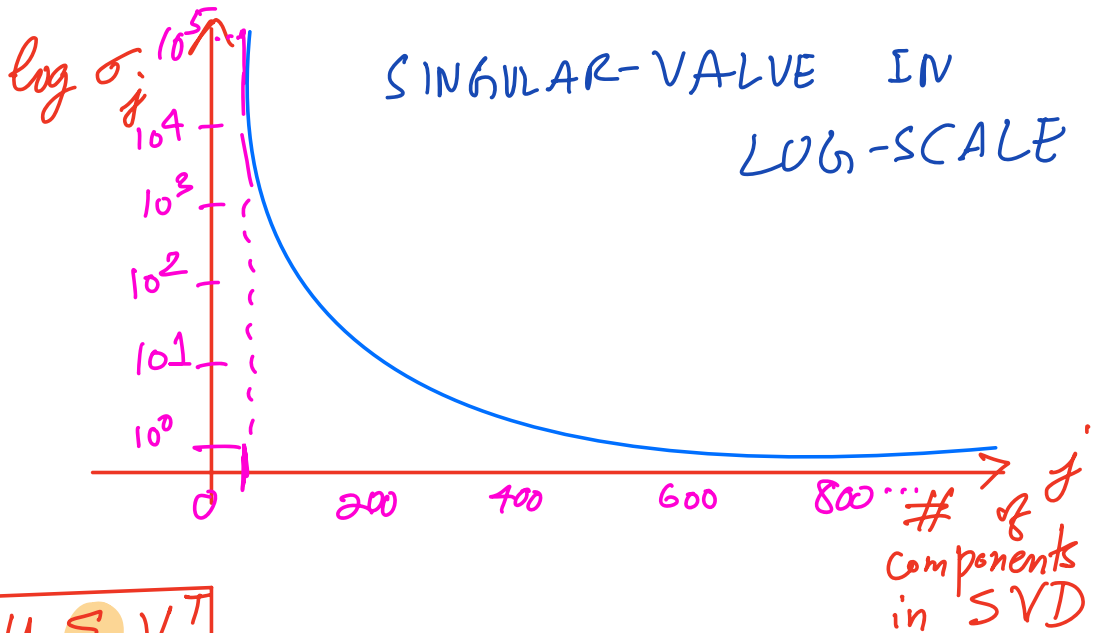
6.2 PCA

Sadly, no cats for this example.

Suppose we, as course staff, have m students in our class, and n assignments. Let $A \in \mathbb{R}^{m \times n}$ be a matrix, such that the i^{th} student's grade in the j^{th} assignment is A_{ij} .

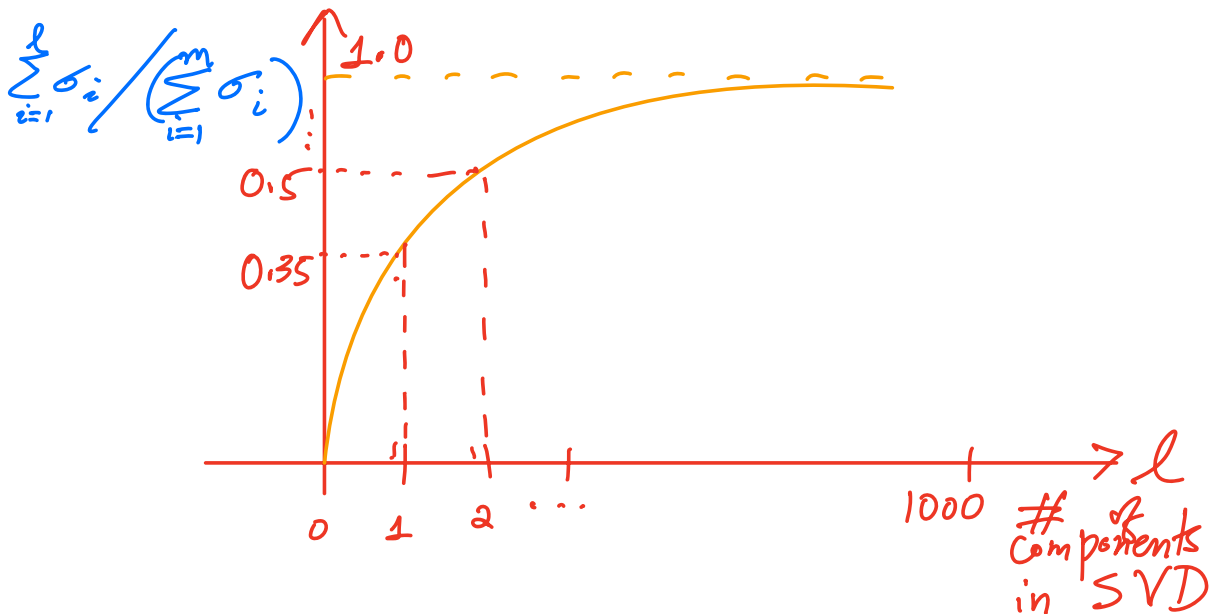
- If we consider the assignments to be the data points, then A is a data matrix with column data.

$$A = U_r \Sigma_r V_r^T$$



$$A = U_r \Sigma_r V_r^T$$

SINGULAR VALUE CUMULATIVE SUM



$$A_l = \sum_{i=1}^l \sigma_i \vec{u}_i \vec{v}_i^T$$

- Eckart-Young Theorem (Note 15) states that the SVD truncation above is more than a heuristic: A_l gives the least possible deviation from A that is possible with a rank- l matrix. More precisely,

A_l above solves: $\min_{B \in \mathbb{R}^{m \times n}} \|A - B\|_F$ Frobenius norm

such that $\text{rank}(B) = l$.

$$\|X\|_F = \left\| \begin{pmatrix} x_{11} & \dots & x_{1n} \\ \vdots & & \vdots \\ x_{m1} & \dots & x_{mn} \end{pmatrix} \right\|_F = \sqrt{x_{11}^2 + x_{12}^2 + \dots + x_{mn}^2}$$

- Suppose we want the best rank-1 approximation to $A \in \mathbb{R}^{m \times n}$, then over all rank-1 matrices $B \in \mathbb{R}^{m \times n}$, the winner is the rank-1 SVD decomposition $(\sigma_1 \vec{u}_1 \vec{v}_1^T)$, where σ_1 is the first (largest) singular value of A and \vec{u}_1, \vec{v}_1 are the first singular vectors of A .

"Best" is in the sense of minimum Frobenius norm of error.

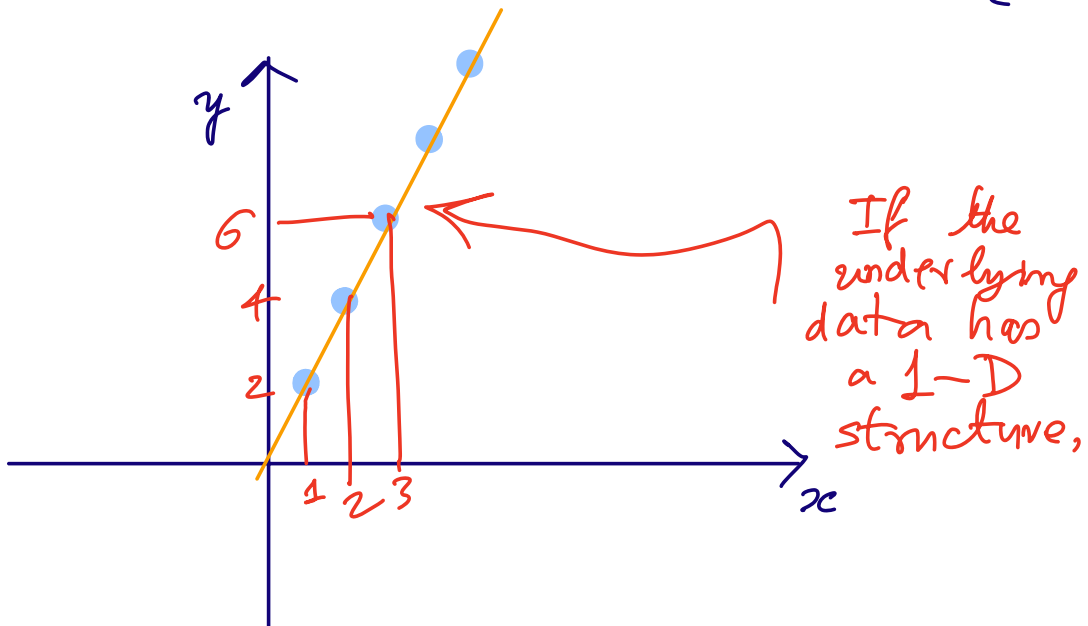
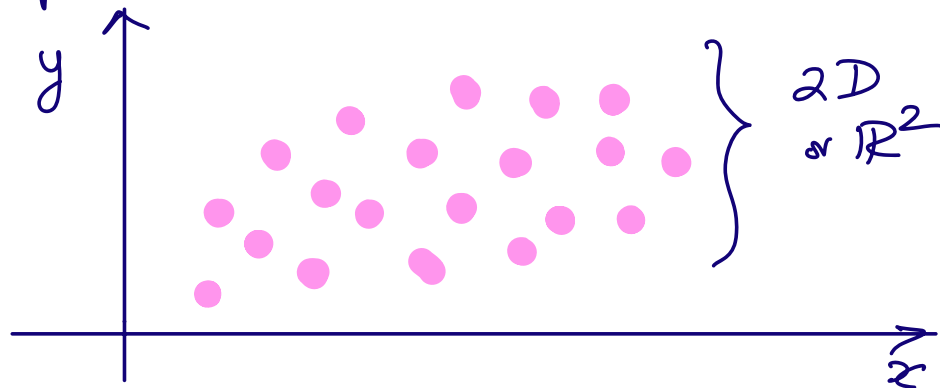
$$\|A - B\|_F^2 = \sum_{i=1}^m \sum_{j=1}^n |A_{ij} - B_{ij}|^2$$

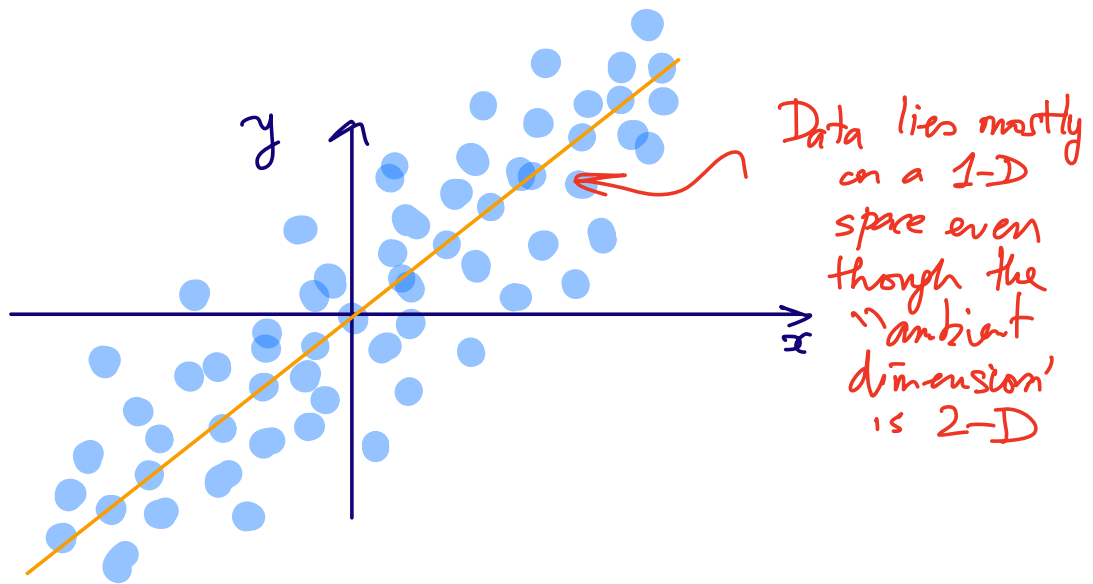
- Best Rank-2 approx. to A is $\sum_{i=1}^2 \sigma_i \vec{u}_i \vec{v}_i^T$
- Best Rank- l approx. to A is $\sum_{i=1}^l \sigma_i \vec{u}_i \vec{v}_i^T$

ordered from largest to smallest

PRINCIPAL COMPONENT ANALYSIS or PCA

"Principal Components", i.e. lower-dimensional structure in simple 2D data:





Suppose we have 2-D points (x_i, y_i) as follows: $\{(1, 2), (2, 4), (3, 6), (4, 8), (5, 10)\}$

$$A = \begin{bmatrix} 1 & 2 & 3 & 4 & 5 \\ 2 & 4 & 6 & 8 & 10 \end{bmatrix}$$

(This is the same example matrix we have seen before in our study of the SVD!)

Ex.:

$$A = \begin{bmatrix} 1 & 2 & 4 & 5 \\ 2 & 4 & 8 & 10 \end{bmatrix}_{2 \times 4}$$

$$A = \underbrace{\begin{bmatrix} \frac{1}{\sqrt{5}} & -\frac{2}{\sqrt{5}} \\ \frac{2}{\sqrt{5}} & \frac{1}{\sqrt{5}} \end{bmatrix}}_{U_{2 \times 2}} \underbrace{\begin{bmatrix} \sigma_1 & & & \\ 230 & 0 & 0 & 0 \\ & \sigma_2 & & \\ & 0 & 0 & 0 \end{bmatrix}}_{\Sigma_{2 \times 4}} \underbrace{\begin{bmatrix} \frac{1}{\sqrt{46}} & \frac{2}{23} & 2\sqrt{\frac{2}{23}} & \frac{5}{\sqrt{46}} \\ -\frac{5}{\sqrt{26}} & 0 & 0 & \frac{1}{\sqrt{26}} \\ -\frac{2}{\sqrt{273}} & 0 & \sqrt{\frac{13}{21}} & \frac{-10}{\sqrt{273}} \\ -\frac{1}{\sqrt{483}} & \sqrt{\frac{2}{23}} & \frac{-4}{\sqrt{483}} & \frac{-5}{\sqrt{483}} \end{bmatrix}}_{V^T_{4 \times 4}}$$

(FULL" SVD)

$$A = \underbrace{\begin{bmatrix} \vec{u}_1 \\ \frac{1}{\sqrt{5}} \\ \frac{2}{\sqrt{5}} \end{bmatrix}}_{\vec{u}_1} \underbrace{[\sigma_1]}_{\sqrt{230}} \underbrace{\begin{bmatrix} \vec{v}_1^T \\ \frac{1}{\sqrt{46}} & \sqrt{\frac{2}{23}} & 2\sqrt{\frac{2}{23}} & \frac{5}{\sqrt{46}} \end{bmatrix}}_{\vec{v}_1^T}$$

"COMPACT" SVD

- \vec{u}_1 is a "basis" for the Col. space (A)
- $(x, y) \in \text{Col}(A)$ if $y = 2x$

Generalizing PCA concept

Movie recommendation problem (Netflix app.)

	U_1 User 1	U_2 User 2	U_3 User 3	...	U_{1000} User 1000
V_1 Video 1	75	80	20		50
V_2 Video 2	30	20	85		70
...				\ddots	
V_{100} Video 100	90	95	30		20

Q-matrix

q_{ij} : rating of user j for video i .

• 100 videos (movies) / 1000 users: users who rate (score) the movies.

• Goal: learn the "low-dimensional" structure underlying this "big" chunk of data

(In practice, there are millions of users & thousands of movies)

$$\vec{s}_a = \begin{bmatrix} s_{a1} \\ s_{a2} \\ \vdots \\ s_{a1000} \end{bmatrix}$$

"action"
sensitivity vector for users $\in \mathbb{R}^{1000}$
Similarly, for $\vec{s}_b, \vec{s}_c, \vec{s}_d$

Consider

$$\begin{aligned} & \vec{a} \cdot \vec{s}_a \quad (100 \times 1) \quad (1 \times 1000) \quad \text{"outer product"} \\ & = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_{100} \end{bmatrix} \begin{bmatrix} s_{a1} & s_{a2} & \dots & s_{a1000} \end{bmatrix} \\ & = \begin{bmatrix} a_1 s_{a1} & a_1 s_{a2} & \dots & a_1 s_{a1000} \\ a_2 s_{a1} & & & \vdots \\ \vdots & & & \\ a_{100} s_{a1} & \dots & \dots & a_{100} s_{a1000} \end{bmatrix} \end{aligned}$$

$$Q = \vec{a} \cdot \vec{s}_a^T + \vec{b} \cdot \vec{s}_b^T + \vec{c} \cdot \vec{s}_c^T + \vec{d} \cdot \vec{s}_d^T$$

$$Q = U \Sigma V^T = \sum_{i=1}^r \sigma_i \vec{u}_i \vec{v}_i^T$$

r : rank of Q

The k principal components of matrix Q .

- along the columns: $\vec{u}_1, \vec{u}_2, \dots, \vec{u}_k$
- along the rows: $\vec{v}_1, \vec{v}_2, \dots, \vec{v}_k$

$$Q = \begin{matrix} & \begin{matrix} \vec{u}_1 & \vec{u}_2 & \dots & \vec{u}_{1000} \end{matrix} \\ \begin{matrix} \vec{v}_1 \\ \vec{v}_2 \\ \vdots \\ \vec{v}_{100} \end{matrix} & \begin{bmatrix} q_{11} & q_{12} & \dots & q_{1n} \\ q_{21} & q_{22} & \dots & q_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ q_{m1} & q_{m2} & \dots & q_{mn} \end{bmatrix} \end{matrix} \quad (m=100, n=1000)$$

- Data is organized by columns:

- i.e., each data point is a 100-dim. vector containing User j 's ratings for the 100 movies $\{q_{1j}, q_{2j}, \dots, q_{100j}\}$

Goal: Find the first principal component; that is, that vector (direction) that is "most informative" about the data.

