

# MATH4280

Lecture Notes 1: Singular value decomposition (SVD)

# Singular value decomposition (SVD)

- One of the most important matrix factorizations
- Foundation of many data analytical tools
- Provide a systematic way to determine a low-dimensional approximation to high-dimensional data
  - e.g. image, audio, video, fluid flow
- It is a **data-driven** approach, patterns are discovered purely from data

# Definition of SVD

- We will analyze a large data set

$$\mathbf{X} = \begin{bmatrix} | & | & \cdots & | \\ \mathbf{x}_1 & \mathbf{x}_2 & \cdots & \mathbf{x}_m \\ | & | & \cdots & | \end{bmatrix} \quad \mathbf{X} \in \mathbb{C}^{n \times m}$$

- Each column  $\mathbf{x}_k \in \mathbb{C}^n$  represents one data, called **snapshot**
  - e.g. different images, state of a physical system at different times
- The dimension  $n$  is usually very large
- $m$  is the number of **snapshots**

- The SVD is a unique matrix factorization given by

$$\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^*$$

- Here,  $\mathbf{U} \in \mathbb{C}^{n \times n}$  and  $\mathbf{V} \in \mathbb{C}^{m \times m}$  are unitary matrices
- $\mathbf{\Sigma} \in \mathbb{R}^{n \times m}$  is a real matrix, with real and nonnegative entries on diagonal, and zeros off the diagonal
- In most cases,  $n \geq m$ , so we can write  $\mathbf{\Sigma} = \begin{bmatrix} \hat{\mathbf{\Sigma}} \\ \mathbf{0} \end{bmatrix}$
- The columns of  $\mathbf{U}$  and  $\mathbf{V}$  are called left singular vectors and right singular vectors respectively

- We can write

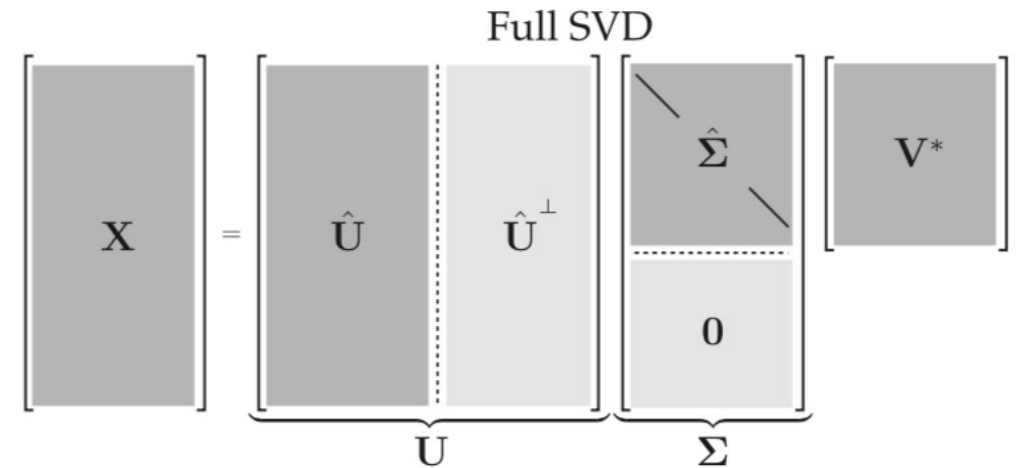
$$\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^* = \begin{bmatrix} \hat{\mathbf{U}} & \hat{\mathbf{U}}^\perp \end{bmatrix} \begin{bmatrix} \hat{\mathbf{\Sigma}} \\ \mathbf{0} \end{bmatrix} \mathbf{V}^* = \hat{\mathbf{U}}\hat{\mathbf{\Sigma}}\mathbf{V}^*$$

Schematic of SVD

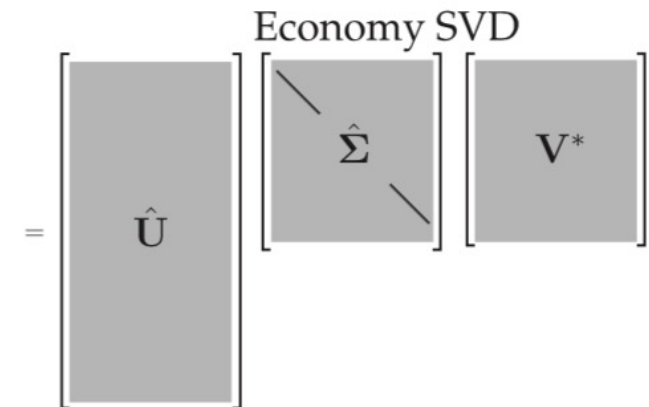
- The diagonal elements of

$$\hat{\mathbf{\Sigma}} \in \mathbb{C}^{m \times m}$$

are called the **singular values**, arranged from large to small



- The **rank** of  $\mathbf{X}$  is equal to the number of nonzero singular values



# Matrix approximation using SVD

- The optimal rank-r approximation to  $\mathbf{X}$  is given by the rank-r SVD truncation

$$\underset{\tilde{\mathbf{X}}, \text{ s.t. rank}(\tilde{\mathbf{X}})=r}{\operatorname{argmin}} \|\mathbf{X} - \tilde{\mathbf{X}}\|_F = \tilde{\mathbf{U}} \tilde{\Sigma} \tilde{\mathbf{V}}^*$$

Frobenius norm

- Here,  $\tilde{\mathbf{U}}$  and  $\tilde{\mathbf{V}}$  are the first  $r$  leading columns of  $\mathbf{U}$  and  $\mathbf{V}$
  - $\tilde{\Sigma}$  is the leading  $r \times r$  sub-block of  $\Sigma$
- 
- We can also use the following formula

$$\tilde{\mathbf{X}} = \sum_{k=1}^r \sigma_k \mathbf{u}_k \mathbf{v}_k^* = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^* + \sigma_2 \mathbf{u}_2 \mathbf{v}_2^* + \cdots + \sigma_r \mathbf{u}_r \mathbf{v}_r^*$$

Q: but how is it possible to compute the first  $r$  terms, without calculating all terms?

Full SVD

$$\begin{bmatrix} \mathbf{X} \end{bmatrix} = \underbrace{\begin{bmatrix} \tilde{\mathbf{U}} & \hat{\mathbf{U}}_{\text{rem}} & \hat{\mathbf{U}}^\perp \end{bmatrix}}_{\mathbf{U}} \begin{bmatrix} \tilde{\Sigma} & & \\ & \hat{\Sigma}_{\text{rem}} & \\ & & \mathbf{0} \end{bmatrix} \begin{bmatrix} \tilde{\mathbf{V}}^* \\ \mathbf{V}_{\text{rem}} \end{bmatrix}$$

Truncated SVD

$$\approx \begin{bmatrix} \tilde{\mathbf{U}} \end{bmatrix} \begin{bmatrix} \tilde{\Sigma} \end{bmatrix} \begin{bmatrix} \tilde{\mathbf{V}}^* \end{bmatrix}$$

Schematic of truncated SVD

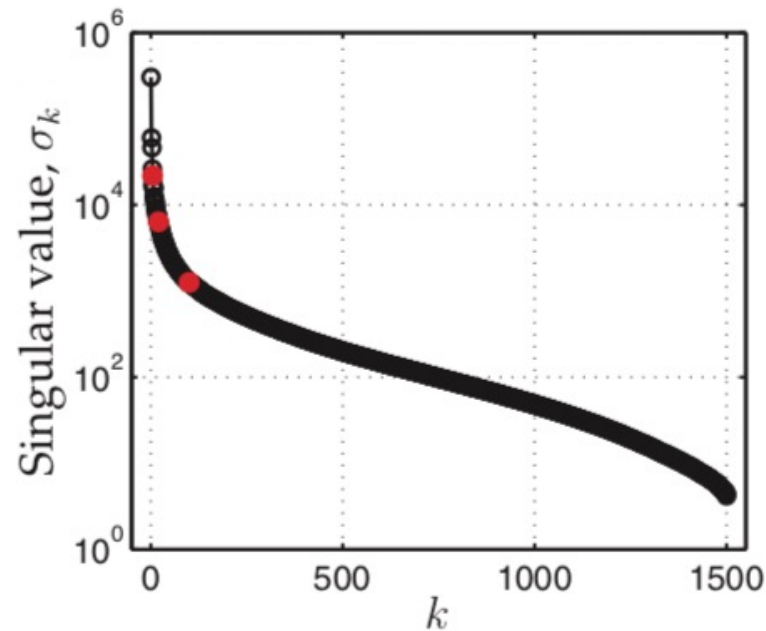
# Example: image compression

A digital image can be considered as a matrix

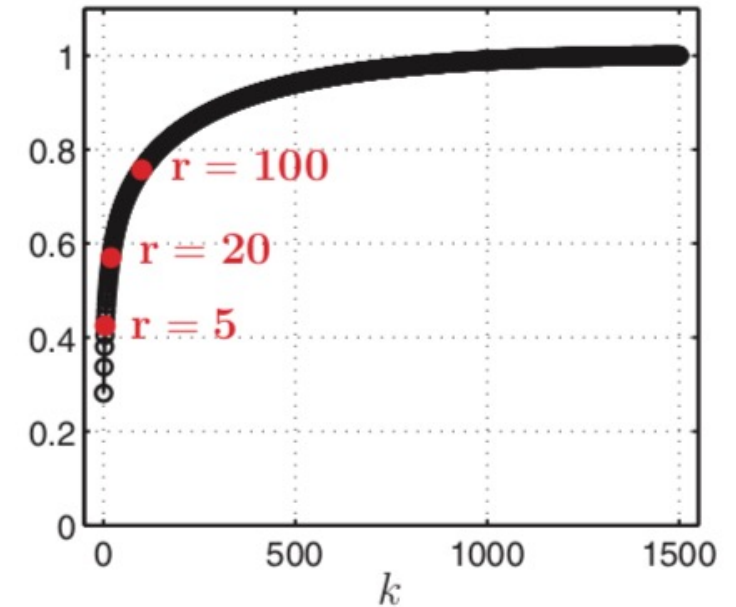
$X =$



An image of resolution 2000 x 1500  
The matrix  $X$  is 2000 x 1500



Singular values



Cumulative energy

it equals to the sum of  $\sigma_i^2$   
divided by the sum of  $\sigma_n^2$



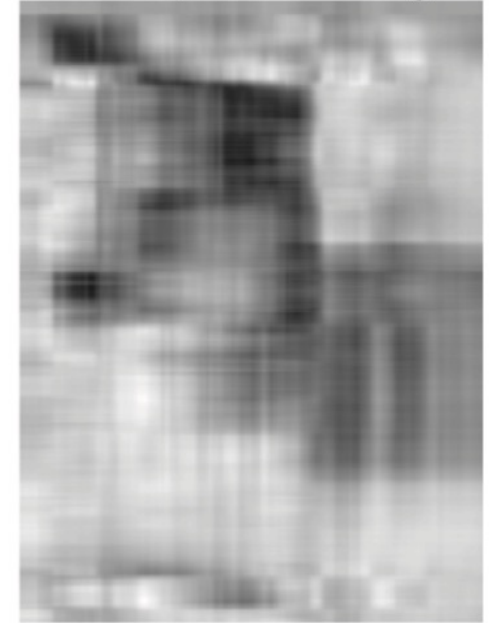
## Results using various values of $r$

- Good result when  $r=100$
- An example of compressibility

Original



$r = 5$ , 0.57% storage



$r = 20$ , 2.33% storage



$r = 100$ , 11.67% storage



# SVD and correlation matrix

A diagram illustrating the calculation of the row correlation matrix  $XX^*$ . On the left, a tall gray rectangle labeled  $X$  has a horizontal white oval at its bottom, representing a row vector. To its right is a wide gray rectangle labeled  $X^*$  with a vertical white oval at its left, representing a column vector. An equals sign follows, leading to a square gray rectangle labeled  $XX^*$  with a small white circle in its bottom-left corner, representing the inner product of the selected row and column.

Taking inner products of rows

A diagram illustrating the calculation of the column correlation matrix  $X^*X$ . On the left, a wide gray rectangle labeled  $X^*$  has a horizontal white oval at its top, representing a row vector. To its right is a tall gray rectangle labeled  $X$  with a vertical white oval at its right, representing a column vector. An equals sign follows, leading to a square gray rectangle labeled  $X^*X$  with a small white circle in its top-right corner, representing the inner product of the selected row and column. This final rectangle is enclosed in a red box, with a red arrow pointing from it towards the bottom right.

Taking inner products of columns

Compute this, because it is often cheaper (will discuss it in p.12)

Correlation matrix: This  $XX^*$  and  $X^*X$  provides an intuitive interpretation of the SVD, where the columns of  $U$  are eigenvectors of the correlation matrix  $XX^*$ , and columns of  $V$  are eigenvectors of  $X^*X$ .

Thus the columns of  $U$  are hierarchically ordered by how much correlation they capture in the columns of  $X$ .

- Using the definition of SVD

Correlation matrix

$$XX^* = U \begin{bmatrix} \hat{\Sigma} & 0 \\ 0 & 0 \end{bmatrix} V^* V \begin{bmatrix} \hat{\Sigma} & 0 \\ 0 & 0 \end{bmatrix} U^* = U \begin{bmatrix} \hat{\Sigma}^2 & 0 \\ 0 & 0 \end{bmatrix} U^*$$

$$X^*X = V \begin{bmatrix} \hat{\Sigma} & 0 \\ 0 & 0 \end{bmatrix} U^* U \begin{bmatrix} \hat{\Sigma} & 0 \\ 0 & 0 \end{bmatrix} V^* = V \hat{\Sigma}^2 V^*.$$

Stick with this option, Sigma is simpler

- Since  $U$  and  $V$  are unitary matrices

$$XX^*U = U \begin{bmatrix} \hat{\Sigma}^2 & 0 \\ 0 & 0 \end{bmatrix},$$

$$X^*XV = V \hat{\Sigma}^2.$$

why we need to put it into this way?

So we use columns of  $V$  to produce our principal axis

- Columns of  $U$  are eigenvectors of  $XX^*$ , and columns of  $V$  are eigenvectors of  $X^*X$  Try to demonstrate why this statement is true

- Note that we order the singular values in descending order
- Columns of  $U$  are hierarchically ordered by the amount of correlations captured in the columns of  $X$

# Method of snapshots

This is especially true when  $n \gg m$

- Computing  $U$  is expensive as the size of  $XX^*$  is large
- Computing  $V$  is relatively cheap as the size of  $X^*X$  is small
- We can compute the columns of  $U$  corresponding to nonzero singular values as follows

$$\tilde{U} = X\tilde{V}\tilde{\Sigma}^{-1}$$

(m x m)

# Pseudo-inverse

- Many physical systems may be represented as a linear system  $\mathbf{Ax} = \mathbf{b}$
- In the **overdetermined** case with **no solution**, we will find the **least-squares solution**  $\mathbf{x}$  that minimizes
- In the **underdetermined** case with **infinitely many solutions**, we will find the **minimum norm solution**  $\mathbf{x}$  that minimizes

- 
- We approximate the inverse of A by the inverse of the truncated SVD

Pseudo-inverse matrix of A

$$\mathbf{A}^\dagger \triangleq \tilde{\mathbf{V}} \tilde{\Sigma}^{-1} \tilde{\mathbf{U}}^* \implies \mathbf{A}^\dagger \mathbf{A} = \mathbf{I}_{m \times m}$$

- The above is called the **left pseudo-inverse** of A
- Applying to  $\mathbf{Ax} = \mathbf{b}$ ,

$$\mathbf{A}^\dagger \mathbf{A} \tilde{\mathbf{x}} = \mathbf{A}^\dagger \mathbf{b} \implies \tilde{\mathbf{x}} = \tilde{\mathbf{V}} \tilde{\Sigma}^{-1} \tilde{\mathbf{U}}^* \mathbf{b}$$

Our line

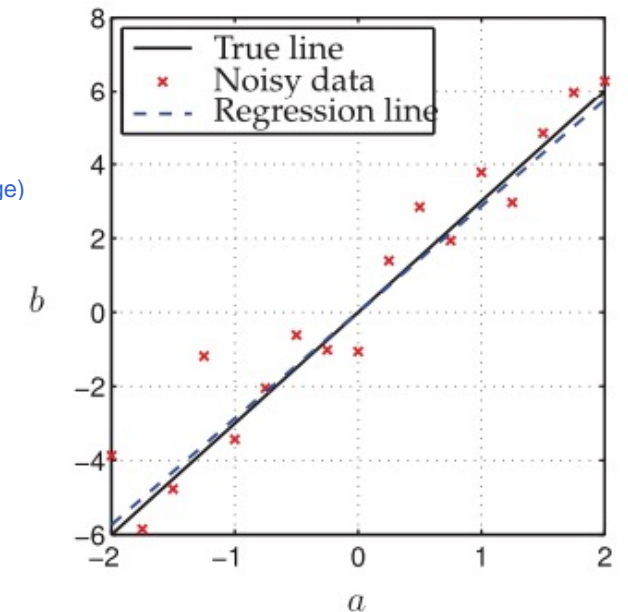
# Example: simple data fitting

- Given a set of data points  $(a_i, b_i)$ , fit a straight line centered at the origin with slope  $x$
- This results in the following problem

$$\begin{bmatrix} | \\ \mathbf{b} \\ | \end{bmatrix} = \begin{bmatrix} | \\ \mathbf{a} \\ | \end{bmatrix} x = \tilde{\mathbf{U}} \tilde{\Sigma} \tilde{\mathbf{V}}^* x.$$
$$\Rightarrow x = \tilde{\mathbf{V}} \tilde{\Sigma}^{-1} \tilde{\mathbf{U}}^* \mathbf{b}. \quad (\text{that is just applying the formulas in the previous page})$$

- We have  $\tilde{\Sigma} = \|\mathbf{a}\|_2$ ,  $\tilde{\mathbf{V}} = 1$ , and  $\tilde{\mathbf{U}} = \mathbf{a}/\|\mathbf{a}\|_2$ , so we obtain

$$x = \frac{\mathbf{a}^* \mathbf{b}}{\|\mathbf{a}\|_2^2}$$



# Principal Component Analysis (PCA)

- Provides a data-driven, hierarchical coordinate system to represent high-dimensional correlated data
- A number of measurements are collected, each measurement is a row of the large matrix  $X$  (where  $X$  is  $n \times m$ )
- We compute the row-wise mean given by

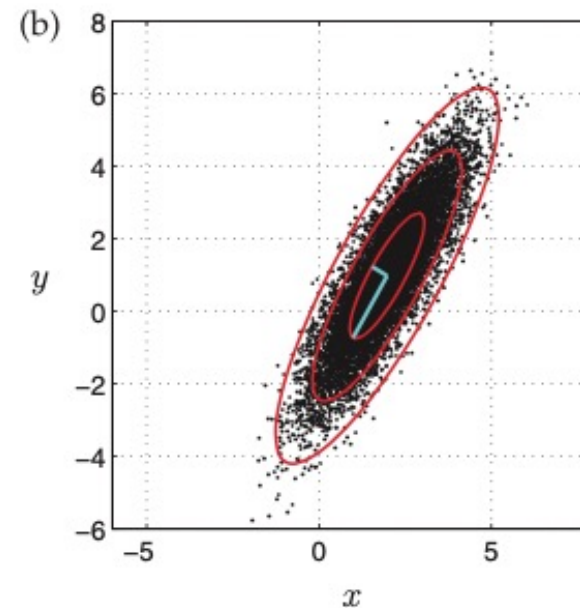
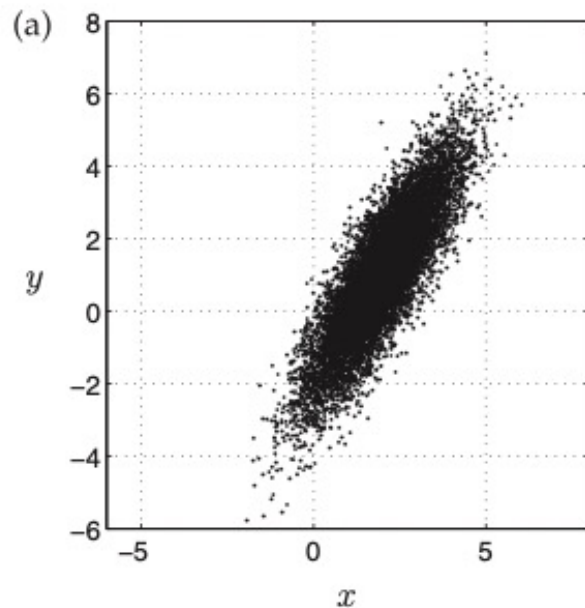
$$\bar{x}_j = \frac{1}{n} \sum_{i=1}^n X_{ij}.$$

- And construct the mean matrix given by  $\bar{X} = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} \bar{\mathbf{x}}$
- The **principal components** are the singular vectors  $V$  of  $\mathbf{B} = \mathbf{X} - \bar{X}$  (or eigenvector of  $B^*B$ )

# An illustration

- A set of  $n$  data points in the  $m = 2$  dimensional space
- The mean is  $(2,1)$
- The **PCA modes** are obtained by the eigenvectors of the  $2 \times 2$  matrix  $B^*B$

PCA modes = the principal components / eigenvectors of the covariance matrix of the data.





# Example: eigenfaces

- Aim: use a large library of facial images to extract the most dominant correlation between images
- The result is a set of **eigenfaces**, which is a new coordinate system to represent the images
- The library contains images of 38 individuals, each of them has 64 images with various poses and lighting conditions
- The images of these 36 individuals will be used to construct the dominant correlations (PCA modes)
- The images of the other 2 individuals will be used to test the PCA modes



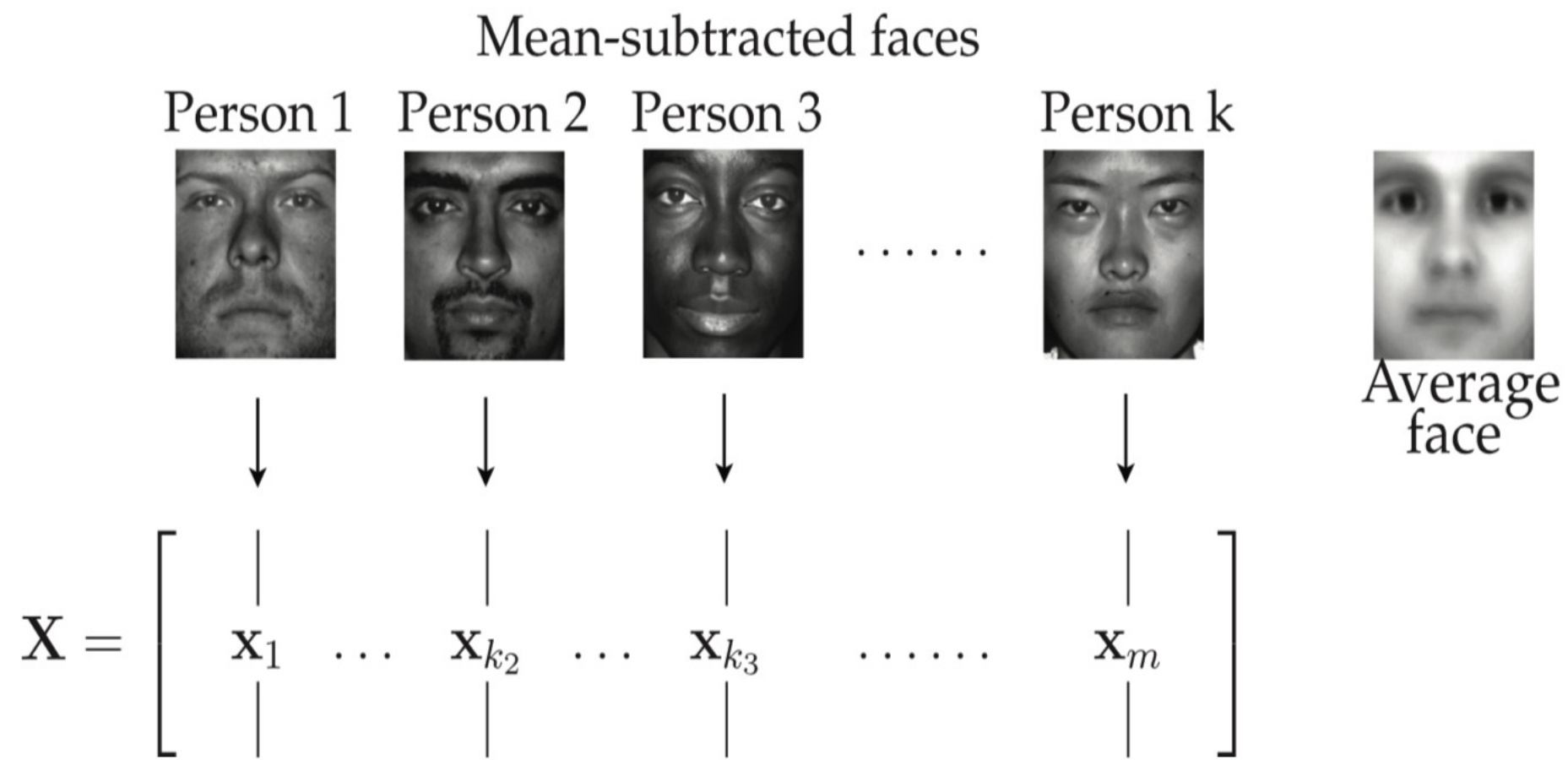
A single image for each individual



All images of a specific individual

Each image is 192 pixels tall and 168 pixels wide

Each image is arranged as a column vector (32256 x 1),  
and is subtracted by the column mean

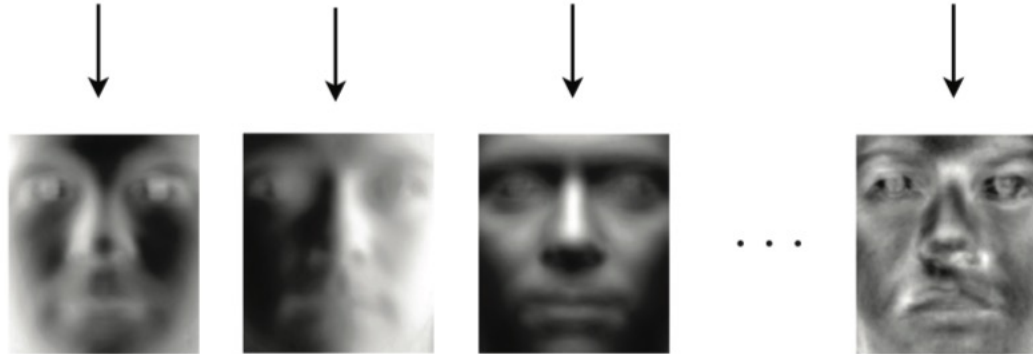


The matrix X has totally 2304 columns

Perform the SVD

$$\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^* \approx \tilde{\mathbf{U}}\tilde{\mathbf{\Sigma}}\tilde{\mathbf{V}}^*$$

$$\tilde{\mathbf{U}} = \left[ \begin{array}{c|c|c|c|c} | & | & | & \dots & | \\ \mathbf{u}_1 & \mathbf{u}_2 & \mathbf{u}_3 & & \mathbf{u}_r \\ | & | & | & & | \end{array} \right]$$



Eigenfaces

The first  $r$  columns of  $\mathbf{U}$  are the first  $r$  PCA modes

Take a test image and represent it using the first  $r$  PCA modes:

$$\tilde{\mathbf{x}}_{\text{test}} = \tilde{\mathbf{U}}\tilde{\mathbf{U}}^* \mathbf{x}_{\text{test}}$$

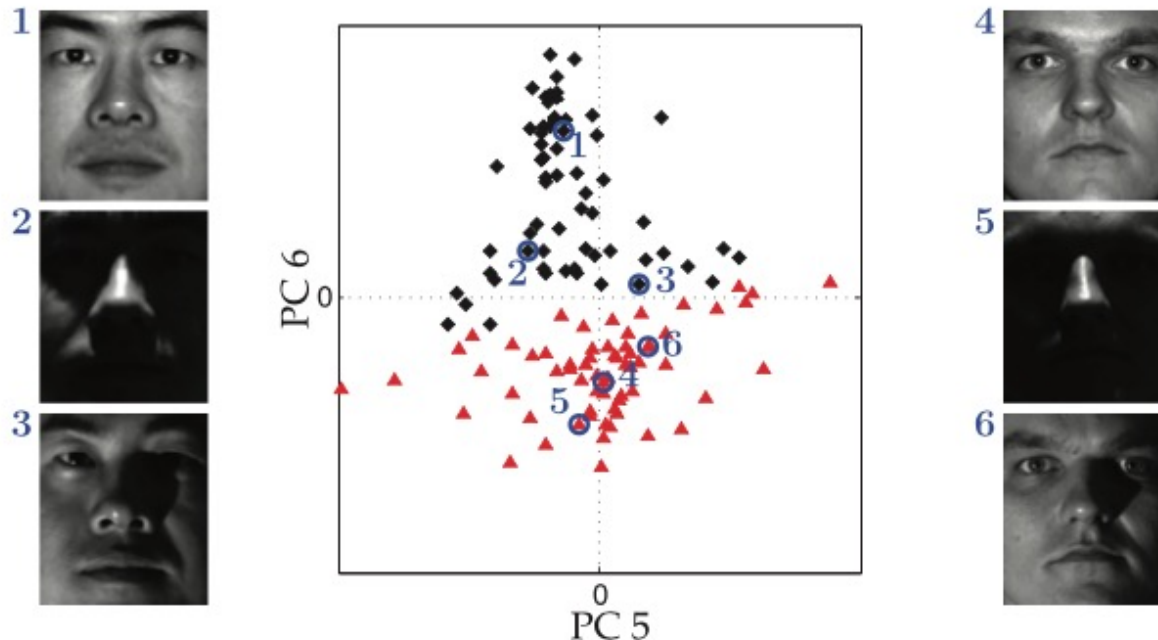
Note that this test image does not belong to the set of images for the PCA construction

We see that the PCA modes can be used to represent images efficiently



# Classifying images

- Some PCA modes may capture the most common features
- Other PCA modes may be useful for distinguishing between images
- The following shows the 64 images of 2 individuals, using the 5th and the 6th PCA modes

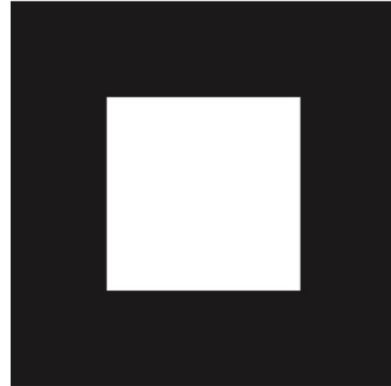


The 6th PCA modes can distinguish these two individuals

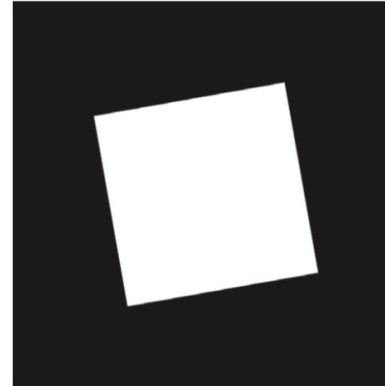
# Importance of data alignment

A 1000 x 1000 square matrix  
X with  
white = 1  
black = 0

(a) 0° Rotation

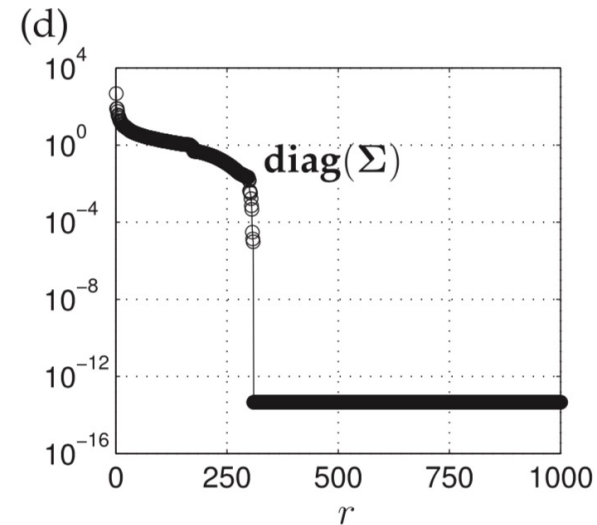
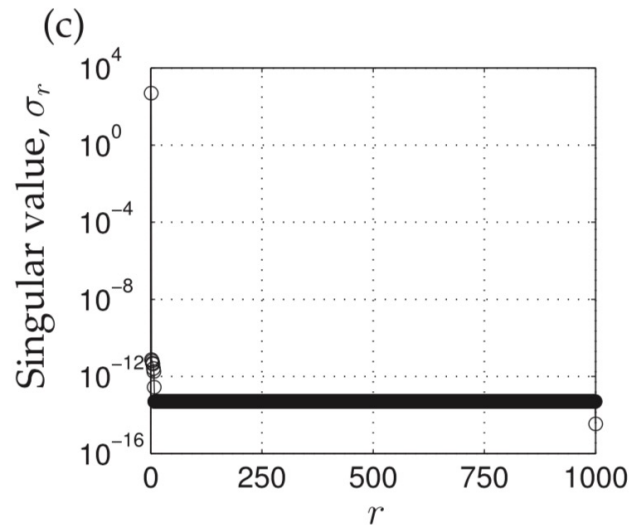


(b) 10° Rotation



A small modification of  
the matrix on the left

Very different behavior  
on the singular values





- A pitfall of the SVD/PCA is data misalignment
- It depends on the coordinate system in which the data is represented
- On the contrary, SVD is invariant under unitary transformations (inner product preserving)
- One should use SVD/PCA carefully given the above points



# Randomized SVD (rSVD)

- A more efficient algorithm for matrix decomposition focusing on extracting dominant low-rank structure in the matrix

Step 1: Construct a **random** projection matrix  $\mathbf{P} \in \mathbb{R}^{m \times r}$  to sample the column space of  $\mathbf{X} \in \mathbb{R}^{n \times m}$

We assume X is low rank.

There is a structure in the dataset

Many of columns are dependent, so there are certain structure in the dataset.

$$\mathbf{Z} = \mathbf{X}\mathbf{P}$$

lead to a big reduction in column space

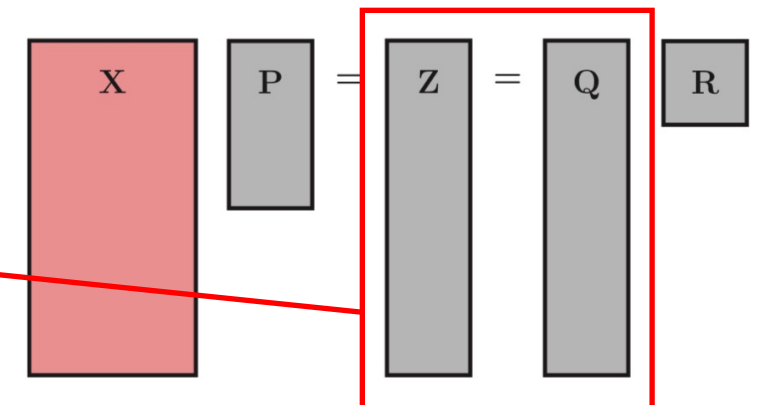
- $r$  is the target rank
- Likely that  $\mathbf{P}$  will project out important components of  $\mathbf{X}$
- $\mathbf{Z}$  approximates the column space of  $\mathbf{X}$  with high probability

Then we perform QR decomposition of  $\mathbf{Z}$  to obtain an **orthonormal basis for  $\mathbf{X}$**

why it is useful to obtain the orthonormal basis for  $\mathbf{X}$ ?

$$\mathbf{Z} = \mathbf{Q}\mathbf{R}$$

span the same column space of  $\mathbf{X}$



Step 2: Project X into a smaller space, and obtain a matrix Y

$$Y = Q^* X$$

what the heck is this?  
(Q: why it is a projection?)

Then we have

$$X \approx QY \quad \text{(better agreement when the singular values decay fast)}$$

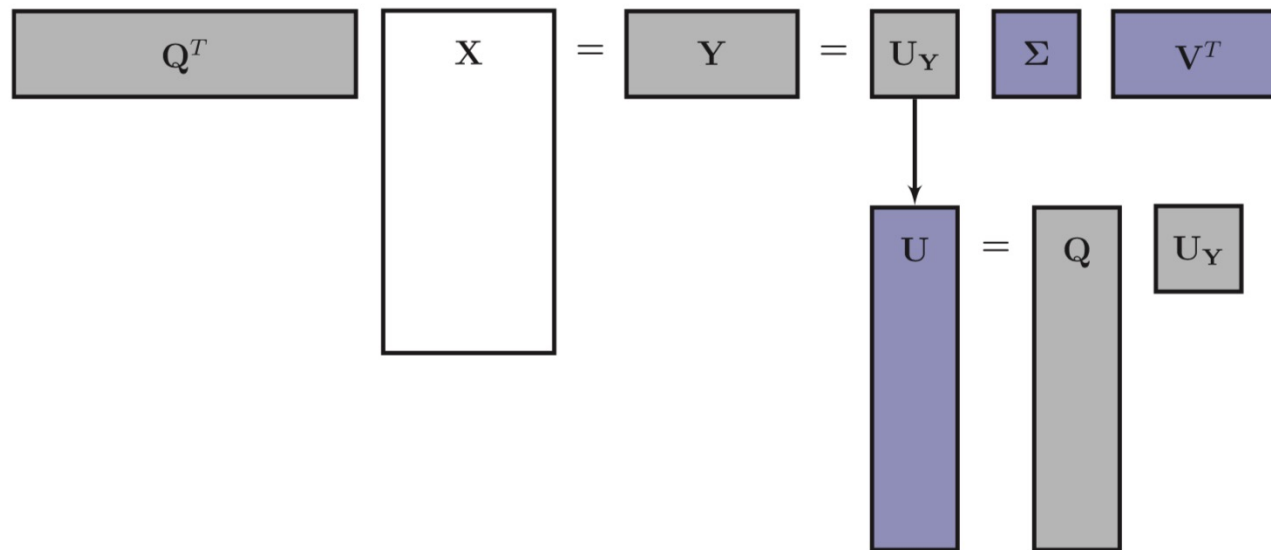
Perform SVD on the smaller matrix Y

$$Y = U_Y \Sigma V^*$$

Note that  $\Sigma$  and  $V$  are the same for Y and X

Step 3: Reconstruct the left singular vectors by

$$U = QU_Y$$



- Oversampling

- the matrix may not be of exactly rank  $r$
- increase the number of columns in the random matrix  $P$  from  $r$  to  $r + p$
- $p = 5$  or  $10$  works well

- Power iterations

- the matrix may have slowly decay singular values
- preprocess  $X$  by the power iterations

$$\mathbf{X}^{(q)} = (\mathbf{X}\mathbf{X}^*)^q \mathbf{X}$$

- the singular values decays more rapidly

$$\mathbf{X}^{(q)} = \mathbf{U}\mathbf{\Sigma}^{2q-1}\mathbf{V}^*$$

- Error bound

$$\mathbb{E}(\|\mathbf{X} - \mathbf{QY}\|_2) \leq \left(1 + \sqrt{\frac{r}{p-1}} + \frac{e\sqrt{r+p}}{p}\sqrt{m-r}\right)^{\frac{1}{2q+1}} \sigma_{k+1}(\mathbf{X})$$

L2-norm

Error would be small, if  $k$  is a good prediction of  $r$