



Principle Component Analysis (PCA)

Lecture 13

Rishabh Iyer

University of Texas at Dallas

Eigenvalues

- λ is an **eigenvalue** of a matrix $A \in \mathbb{R}^{n \times n}$ if the linear system $\underbrace{Ax = \lambda x}_{\text{has at least one non-zero solution}}$ has at least one non-zero solution
- If $Ax = \lambda x$ we say that λ is an eigenvalue of A with corresponding **eigenvector** x
- Could be multiple eigenvectors for the same λ

$$(A - I)\underline{x} = 0$$

Eigenvalues of Symmetric Matrices



- If $A \in \mathbb{R}^{n \times n}$ is symmetric, then it has n linearly independent eigenvectors v_1, \dots, v_n corresponding to n real eigenvalues
 - Moreover, it has n linearly independent **orthonormal** eigenvectors

$$\underbrace{v_i^T v_j = 0 \text{ for all } i \neq j}_{\text{orthogonal}}$$

$$\underbrace{v_i^T v_i = 1 \text{ for all } i}_{\text{normalized}}$$

$$A_{ij} = A_{ji}$$

$$\begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 2 & 3 \\ 2 & 3 & 1 \\ 3 & 1 & 2 \end{bmatrix} \leftarrow \text{symmetric}$$

Eigenvalues of Symmetric Matrices



- If $A \in \mathbb{R}^{n \times n}$ is symmetric, then it has n linearly independent eigenvectors v_1, \dots, v_n corresponding to n real eigenvalues
- A symmetric matrix is **positive definite** if and only if all of its eigenvalues are positive
 - The orthonormal eigenvectors form a **basis** of \mathbb{R}^n (similar to the standard coordinate axes)

Examples

- The 2×2 identity matrix has all of its eigenvalues equal to 1 (it is positive definite) with orthonormal eigenvectors $\begin{bmatrix} 1 \\ 0 \end{bmatrix}$ and $\begin{bmatrix} 0 \\ 1 \end{bmatrix}$
 - The matrix $\begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$ has eigenvalues $\overbrace{0}$ and $\overbrace{2}$ with orthonormal eigenvectors $\begin{bmatrix} -1 \\ \sqrt{2} \\ 1 \\ \sqrt{2} \end{bmatrix}$ and $\begin{bmatrix} 1 \\ \sqrt{2} \\ 1 \\ \sqrt{2} \end{bmatrix}$
 - The matrix $\begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}$ has eigenvalues 1 and 3 with orthonormal eigenvectors $\begin{bmatrix} -1 \\ \sqrt{2} \\ 1 \\ \sqrt{2} \end{bmatrix}$ and $\begin{bmatrix} 1 \\ \sqrt{2} \\ 1 \\ \sqrt{2} \end{bmatrix}$
- $\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} n = \overrightarrow{n}$
 $\begin{bmatrix} n_1 \\ n_2 \end{bmatrix} = \sum \begin{bmatrix} n_1 \\ n_2 \end{bmatrix}$

Eigenvalues

- Suppose $A \in \mathbb{R}^{n \times n}$ is symmetric
- Any $x \in \mathbb{R}^n$ can be written as $x = \sum_{i=1}^n c_i v_i$ where v_1, \dots, v_n are the eigenvectors of A
 - $Ax = \sum_{i=1}^n \lambda_i c_i v_i$
 - $A^2x = \sum_{i=1}^n \lambda_i^2 c_i v_i$
 - \vdots
 - $A^t x = \sum_{i=1}^n \lambda_i^t c_i v_i$

Eigenvalues

- Suppose $A \in \mathbb{R}^{n \times n}$ is symmetric
- Any $x \in \mathbb{R}^n$ can be written as $x = \sum_{i=1}^n c_i v_i$ where v_1, \dots, v_n are the eigenvectors of A
 - $c_i = \overbrace{v_i^T x}^{\text{projection}}$, this is the projection of x along the line given by v_i (assuming that v_i is a unit vector)

Eigenvalues of Symmetric Matrices



- Let $Q \in \mathbb{R}^{n \times n}$ be the matrix whose i^{th} column is v_i and $D \in \mathbb{R}^{n \times n}$ be the diagonal matrix such that $D_{ii} = \lambda_i$

- $Ax = QDQ^T x$

$$A = Q D Q^T$$

$$D = \begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ 0 & & & \lambda_n \end{bmatrix}$$

- Can throw away some eigenvectors to approximate this quantity

$$\underbrace{\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_k}_{\text{top } k \text{ eigenvalues}} \geq \dots \geq \lambda_n$$

- For example, let Q_k be the matrix formed by keeping only the top k eigenvectors and D_k be the diagonal matrix whose diagonal consists of the top k eigenvalues

$$\underbrace{Ax}_{} = Q D Q^T x$$

$$Q D Q^T \approx Q_k D_k Q_k^T$$

Frobenius Norm



- The Frobenius norm is a matrix norm given by

$$\|A\|_F = \sqrt{\sum_{i=1}^n \sum_{j=1}^n |A_{ij}|^2}$$

- $Q_k D_k Q_k^T$ is the best **rank** k approximation of the symmetric matrix A with respect to the Frobenius norm

$$\overbrace{Q_k D_k Q_k^T} = \underbrace{\operatorname{argmin}_{B \in \mathbb{R}^{n \times n} \text{ s.t. } \operatorname{rank}(B)=k} \|A - B\|_F}_{\text{ }}$$

Proof sketch:

Problem: $\min_{B \in \mathbb{R}^{n \times n}: \text{rank}(B) = k} \|A - B\|_F.$

know: $B^* = Q_R D_R Q_R^T$ is rank k .

$\therefore B^*$ is valid matrix.

$$\|A - B^*\|_F^2 \geq \|Q D Q^T - Q_R D_R Q_R^T\|^2$$

$\forall B: \text{rank}(B) \geq n$.

Principal Component Analysis



- Principle component analysis

- Can be used to reduce the dimensionality of the data while still maintaining a good approximation of the sample mean and variance

- Can also be used for selecting good features that are combinations of the input features
- Unsupervised – just finds a good representation of the data in terms of combinations of the input features

$x_i \in \mathbb{R}^n$

$$x_i \rightarrow \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}_k$$

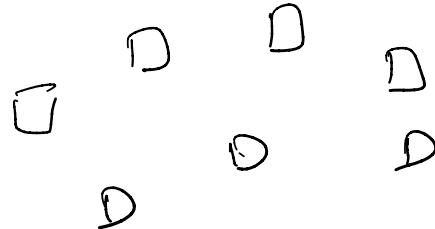
Basis: v_1, \dots, v_k

$$\begin{bmatrix} (v_1, v_1) \\ (v_1, v_2) \\ \vdots \end{bmatrix}$$

Motivation of PCA

$x \in \mathbb{R}^n$ [Very large]

$$1024 \times 1024 \rightarrow 10^6$$

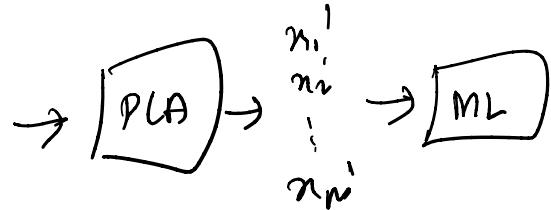


Reduce dim.: 100

Eigen-vector space

$X \in \mathbb{R}^{n \times p}$
 $\#$ examples
 $\#$ feats

n_1 Input
 n_2 images
 n_m



$$\begin{aligned} A &= \underbrace{X X^T}_{A_K} = Q D Q^T \\ &= \underbrace{Q_K D_K Q_K^T}_{\text{rank } K} \end{aligned}$$

Principal Component Analysis



- Input a collection of data points sampled from some distribution

$$x_1, \dots, x_p \in \mathbb{R}^n$$

x_i $\underbrace{\text{p.e.t samples}}$ n_2 # features

- Construct the matrix $W \in \mathbb{R}^{n \times p}$ whose i^{th} column is
- $$x_i - \frac{\sum_j x_j}{p} \leftarrow \text{mean}_n$$
- W $\underbrace{\text{size}}_{n \times p}$ n_1, n_2, \dots
- The matrix WW^T is the sample covariance matrix
 - WW^T is symmetric and positive semidefinite

$$WW^T \in \mathbb{R}^{n \times n}$$

Principal Component Analysis

$$\text{eig}(WW^T) \rightarrow \text{Sort} \rightarrow \underbrace{\lambda_1 > \lambda_2 > \dots}_{\leftarrow k} \geq \lambda_n \quad v_1, v_2, \dots, v_n$$

- PCA finds a set of orthogonal vectors that best explain the variance of the sample covariance matrix
 - From our previous discussion, these are exactly the eigenvectors of $\underbrace{WW^T}_{\leftarrow \text{Symmetric}}$
 - We can discard the eigenvectors corresponding to small magnitude eigenvalues to yield an approximation
 - Simple algorithm to describe, MATLAB and other programming languages have built in support for eigenvector computation

PCA in Practice



$$1024 \times 1024 \\ n \sim 1M$$

- Forming the matrix WW^T can require a lot of memory
(especially if $n \gg p$)
- Need a faster way to compute this without forming the matrix explicitly
- Typical approach: use the singular value decomposition (SVD)

Singular Value Decomposition (SVD)



- Every matrix $B \in \mathbb{R}^{n \times p}$ admits a decomposition of the form

$$\underline{\underline{B}} = U \Sigma V^T \quad U, V \rightarrow \text{Orthogonal matrices}$$

- where $U \in \mathbb{R}^{n \times n}$ is an orthogonal matrix, $\Sigma \in \mathbb{R}^{n \times p}$ is non-negative diagonal matrix, and $V \in \mathbb{R}^{p \times p}$ is an orthogonal matrix
- A matrix $C \in \mathbb{R}^{m \times m}$ is **orthogonal** if $C^T = C^{-1}$. Equivalently, the rows and columns of C are orthonormal vectors

$$U U^T = U^T U = \mathbb{I}$$
$$\Sigma = \begin{bmatrix} \sigma_1 & & & \\ & \ddots & & \\ & & \sigma_p & \\ & 0 & & \end{bmatrix}$$

Singular Value Decomposition (SVD)



- Every matrix $B \in \mathbb{R}^{n \times p}$ admits a decomposition of the form

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

Diagonal elements of Σ called
singular values

- where $U \in \mathbb{R}^{n \times n}$ is an orthogonal matrix, $\Sigma \in \mathbb{R}^{n \times p}$ is non-negative diagonal matrix, and $V \in \mathbb{R}^{p \times p}$ is an orthogonal matrix
- A matrix $C \in \mathbb{R}^{m \times m}$ is orthogonal if $C^T = C^{-1}$. Equivalently, the rows and columns of C are orthonormal vectors

$$U U^T = I = V V^T$$

SVD and PCA

- Returning to PCA

- Let $W = U\Sigma V^T$ be the SVD of W

- $WW^T = U\Sigma V^T V \Sigma^T U^T = U \underbrace{\Sigma \Sigma^T}_{I} U^T$

- If we can compute the SVD of W , then we don't need to form the matrix WW^T

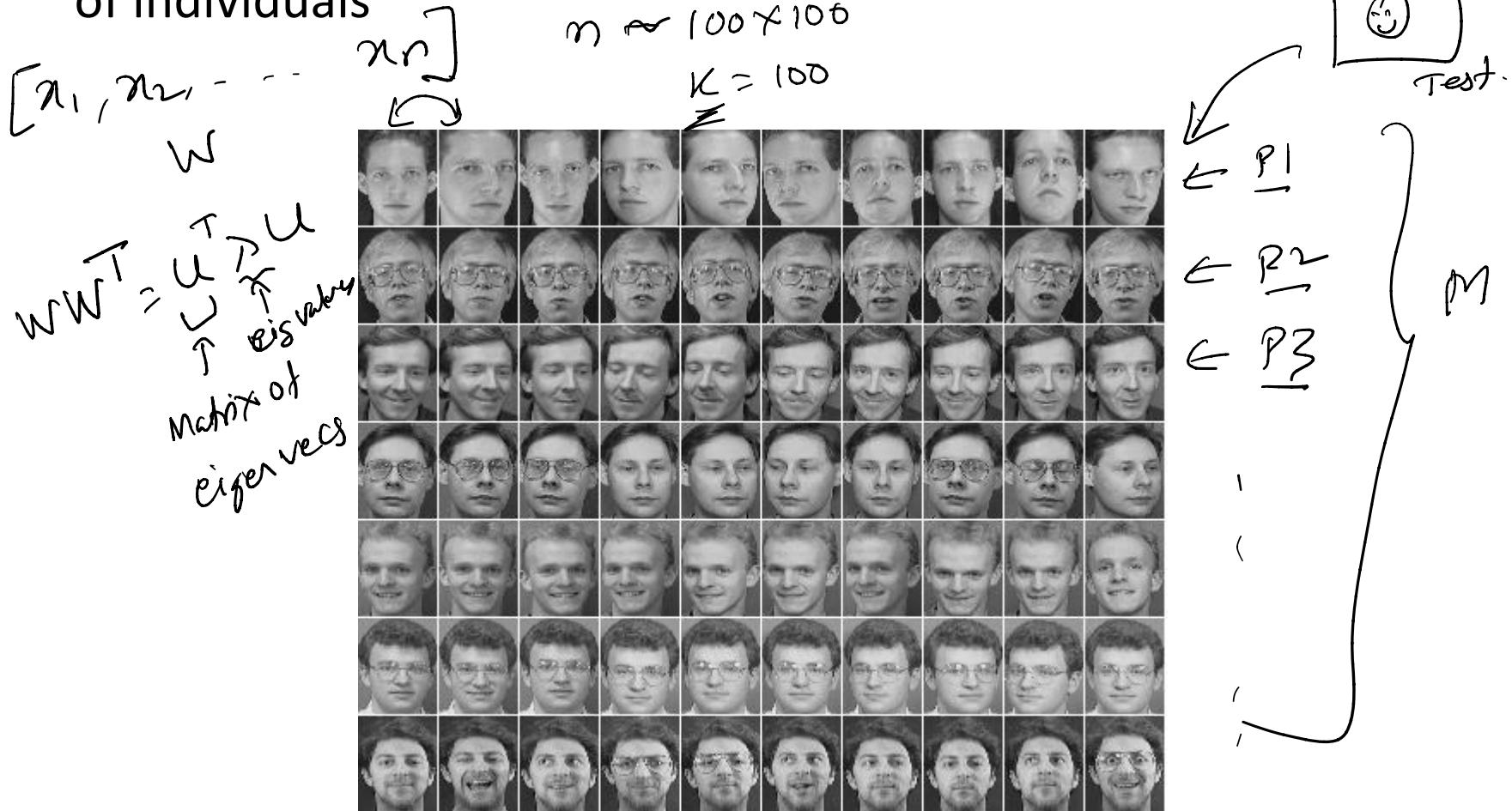
$$\begin{aligned} W &= U\Sigma V^T \\ WW^T &= U\Sigma \underbrace{V^T V}_{I} \Sigma^T U^T \\ &= U\Sigma \Sigma^T U^T \end{aligned}$$

SVD and PCA

- For any matrix A , AA^T is symmetric and positive semidefinite
 - Let $\overset{W}{A} = U\Sigma V^T$ be the SVD of $\overset{W}{A}$ $\text{eig}(WW^T) = \sum \Sigma^T$
 - $\overset{W}{AA^T} = U\Sigma V^T V\Sigma^T U^T = U\Sigma\Sigma^T U^T$
 - U must be a matrix of eigenvectors of $\overset{W}{AA^T}$ WW^T
 - The eigenvalues of $\overset{W}{AA^T}$ are all non-negative because $\Sigma\Sigma^T = \Sigma^2$ which are the square of the singular values of A

An Example: “Eigenfaces”

- Let's suppose that our data is a collection of images of the faces of individuals



$$W = [x_1, x_2, \dots, x_p]$$

$n = \# \text{ features}$
 $p = \# \text{ examples}$

$$\text{Cov}(w) = w w^T$$

Eigen-Decomposition:

$$v_1, v_2, \dots, v_n \in \text{eig-vec}(ww^T)$$

$$\lambda_1, \lambda_2, \dots, \lambda_n \in \text{eig-val}(ww^T)$$

$$x_i \in \mathbb{R}^n$$

$$\text{PCA}(x_i) = \begin{bmatrix} x_i^T v_1 \\ x_i^T v_2 \\ \vdots \\ x_i^T v_k \end{bmatrix} \in \mathbb{R}^k$$

$$v_1 = (0, 1)$$

$$\|\text{Cov}(w) - \text{Cov}_K(w)\|^2$$

$$\text{DFE}(m_i)$$

$v_i = \text{orthogonal}$

An Example: “Eigenfaces”

- Let's suppose that our data is a collection of images of the faces of individuals

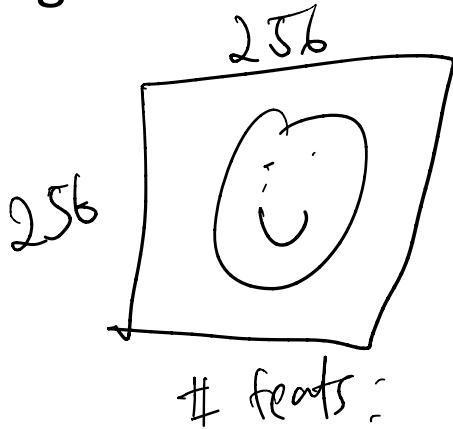
$$\underbrace{x_1, x_2, \dots, x_n}_{n_t}$$

- The goal is, given the "training data", to correctly match new images to the training data
- Let's suppose that each image is an $s \times s$ array of pixels: $x_i \in \mathbb{R}^n$, $n = s^2$
- As before, construct the matrix $W \in \mathbb{R}^{n \times p}$ whose i^{th} column is $x_i - \sum_j \frac{x_j}{p}$

An Example: “Eigenfaces”

- Forming the matrix $\underbrace{WW^T}$ requires a lot of memory
 - s = 256 means $\underbrace{WW^T}$ is $\underbrace{65536 \times 65536}_{n=65536} = \text{very large}$ $\approx 26 \times 10^{10}$
 - Need a faster way to compute this without forming the matrix explicitly
 - Could use the singular value decomposition

256×265



An Example: “Eigenfaces”

- A different approach when $p \ll n$
 - Compute the eigenvectors of $A^T A$ (this is an $p \times p$ matrix)
 - Let v be an eigenvector of $\underbrace{A^T A}$ with eigenvalue λ
 - $AA^T A v = \lambda A v$
 - This means that $A v$ is an eigenvector of AA^T with eigenvalue λ (or 0)
 - Save the top k eigenvectors - called eigenfaces in this example

$$\underbrace{w^T w}$$

An Example: “Eigenfaces”

$$\begin{bmatrix} \mathbf{x}_1 & \dots & \mathbf{x}_P \end{bmatrix}$$

$\underbrace{\phantom{\mathbf{x}_1 \dots \mathbf{x}_P}}_{\mathbf{x}_t}$

- The data in the matrix is “training data”
 - Given a new image, we’d like to determine which, if any, member of the data set that it is most similar to
- Step 1: Compute the projection of the recentered, new image onto each of the k eigenvectors
 - This gives us a vector of weights c_1, \dots, c_k

An Example: “Eigenfaces”

- The data in the matrix is “training data”
 - Given a new image, we’d like to determine which, if any, member of the data set that it is most similar to
- Step 2: Determine if the input image is close to one of the faces in the data set
 - If the distance between the input and it's approximation is too large, then the input is likely not a face

$$\underline{z_i} = \text{PCA}(x_i) \quad z_i \in \mathbb{R}^n$$

An Example: “Eigenfaces”

- The data in the matrix is “training data”
 - Given a new image, we’d like to determine which, if any, member of the data set that it is most similar to
- Step 3: Find the person in the training data that is closest to the new input
 - Replace each group of training images by its average
 - Compute the distance to the i^{th} average $\|c - a^i\|$ where a^i are the coefficients of the average face for person i

An Example: “Eigenfaces”

- The data in the matrix is “training data”
 - Given a new image, we’d like to determine which, if any, member of the data set that it is most similar to
 - Step 3: Find the person in the training data that is closest to the new input
 - Replace each group of training images by its average
 - Compute the distance to the i^{th} average $\|c - a^i\|$ where a^i are the coefficients of the average face for person i

$$\begin{aligned} E[xx^T] &\notin \text{Co. } \\ E[(x - \bar{x})(x - \bar{x})^T] & \end{aligned}$$

Eigen Faces Example.

- ① Create mean-normalized Data set: (n_1, n_2, \dots, n_p)
each $n_i \in \mathbb{R}^n$
- ② Do SVD (\mathbf{W}), compute Eig Dec ($\mathbf{W}\mathbf{W}^T$)
get top k eig vec: v_1, v_2, \dots, v_n .
- ③ For each n_i , $z_i = \text{PCA}(n_i) = \begin{bmatrix} v_1^T n_i \\ v_2^T n_i \\ \vdots \\ v_n^T n_i \end{bmatrix}$
- ④ $D = (z_1, \dots, z_p)$
For test n_t , $z_t = \text{PCA}(n_t)$
use NN classifier to get ID of person!!

$$\text{Cov}(WW^T)$$

1, 5, 100, 500, ...
[]

$$\text{Cov}_K((x - \mu)(x - \mu)^T)$$

$\mu + \Delta t \alpha$

$$x_i \sim M_i + \underset{\substack{\uparrow \\ \text{Cov}}}{\text{Noise}}$$

$$x_i = M_i + z_i$$

$$\hat{x}_i = M_i + \text{Approx}(z_i)$$

\hat{x}_i better than $\text{Approx}(x_i)$