

JPL-Caltech Virtual Summer School

Big Data Analytics



September 2 – 12, 2014

David R. Thompson

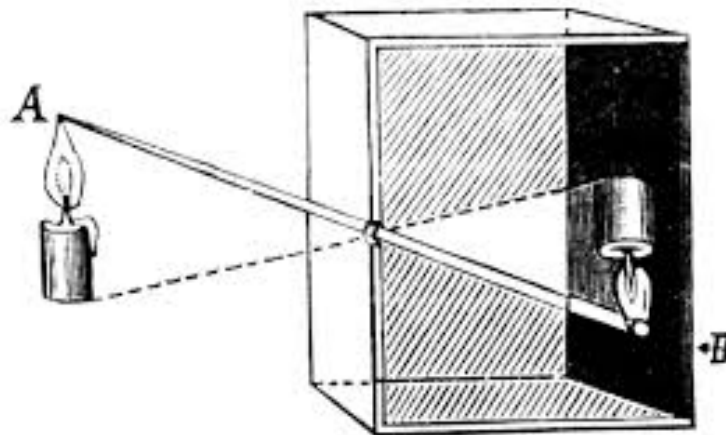
Jet Propulsion Laboratory, California Institute of Technology

Linear Dimensionality Reduction

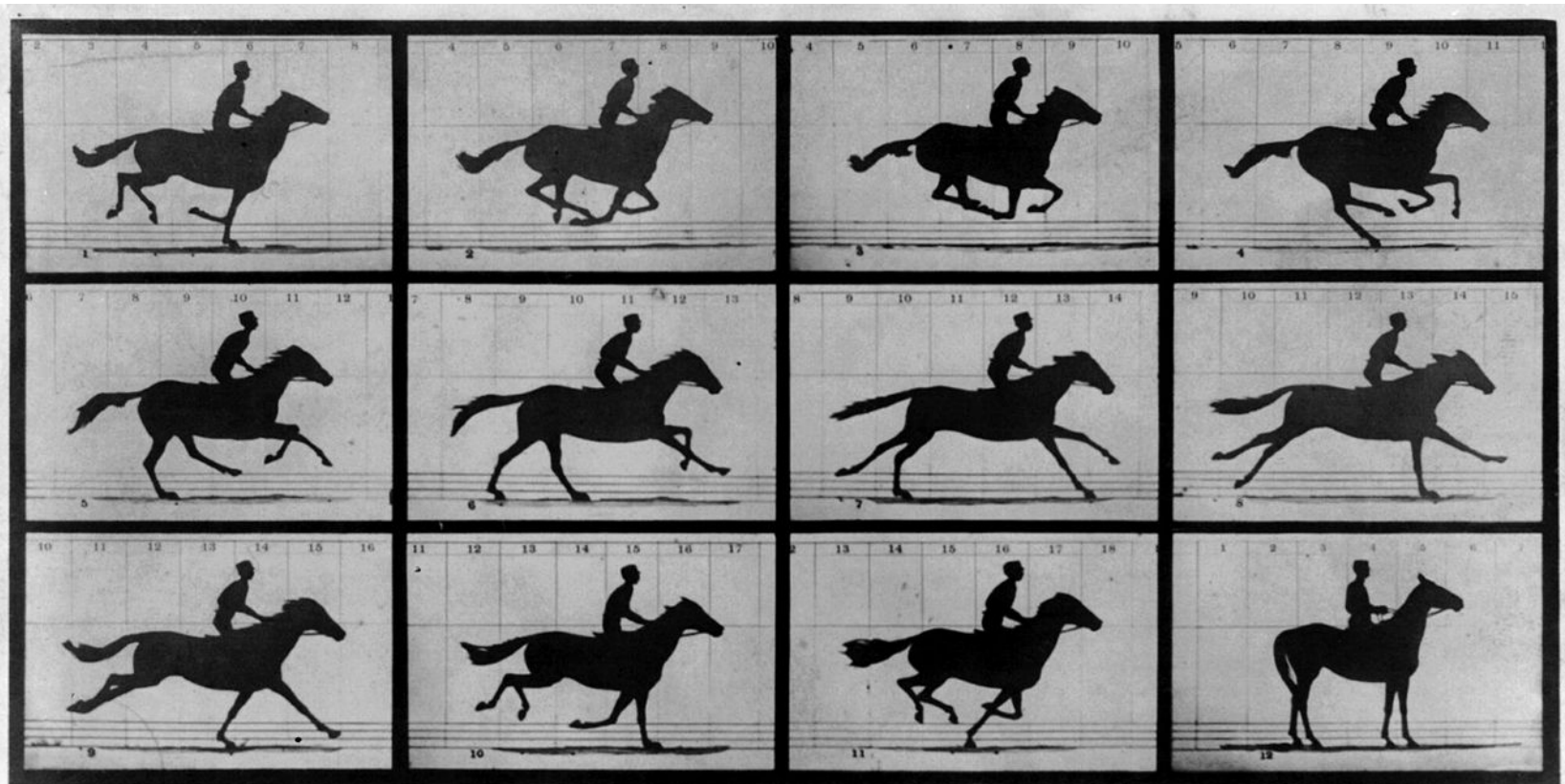
Copyright 2014 California Institute of Technology. All Rights Reserved. US Government Support Acknowledged.

Objectives

1. Introduction to dimensionality reduction and its relationship with feature selection
2. Understand Principal Component Analysis and its relation to SVD
3. Find your eigenface representation



Dimensionality reduction



Copyright, 1878, by MUYBRIDGE.

MORSE'S Gallery, 417 Montgomery St., San Francisco.

THE HORSE IN MOTION.

Illustrated by
MUYBRIDGE.

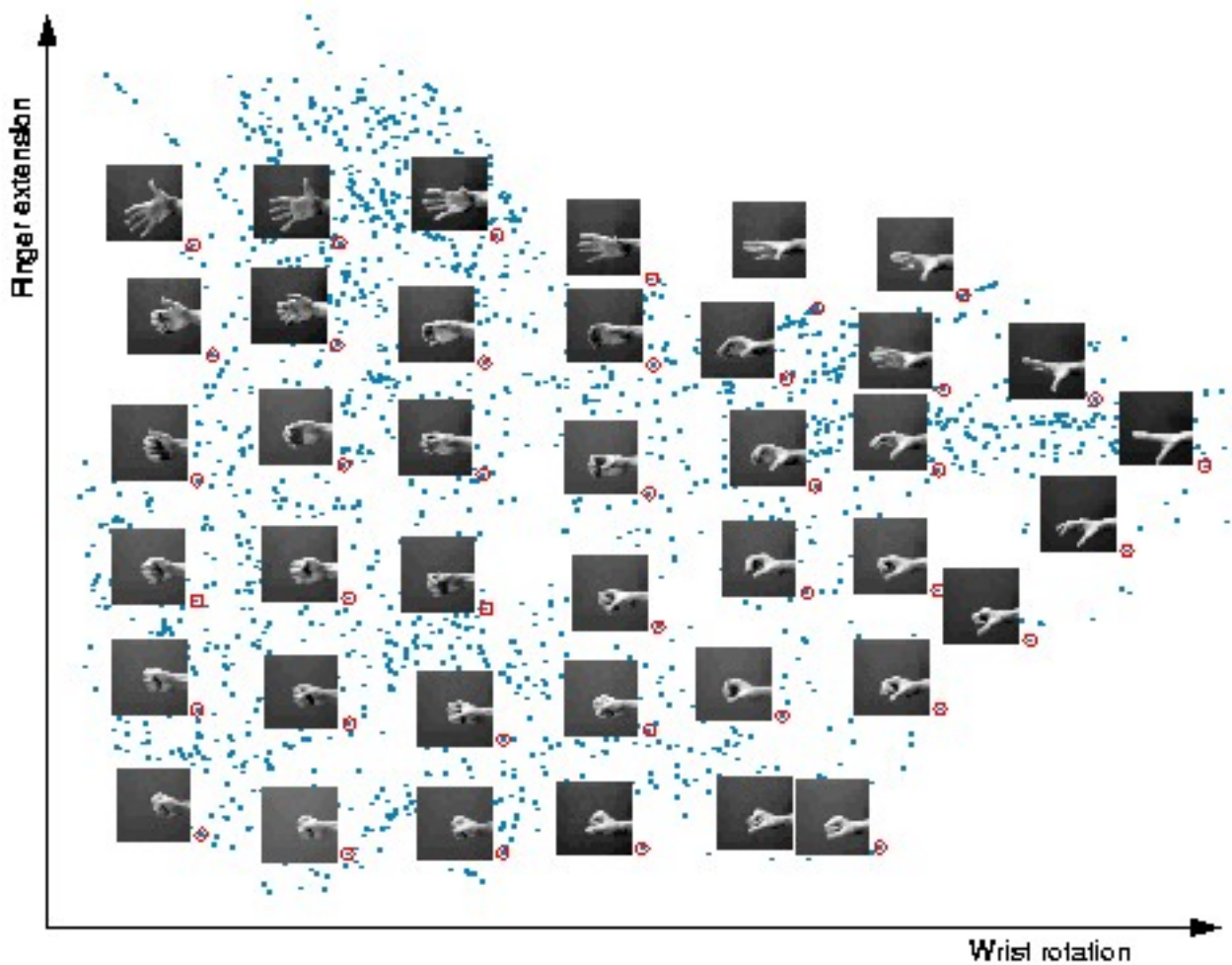
AUTOMATIC ELECTRO-PHOTOGRAPH.

"SALLIE GARDNER," owned by LELAND STANFORD; running at a 1.40 gait over the Palo Alto track, 19th June, 1878.

The negatives of these photographs were made at intervals of twenty-seven inches of distance, and about the twenty-fifth part of a second of time; they illustrate consecutive positions assumed in each twenty-seven inches of progress during a single stride of the mare. The vertical lines were twenty-seven inches apart; the horizontal lines represent elevations of four inches each. The exposure of each negative was less than the two-thousandth part of a second.



Dimensionality reduction



[Tenenbaum et al., *Science* 2000]



Formal definitions

Data point as an n-dimensional column vector

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_n \end{bmatrix}$$

Data set as a [n x d] matrix

$$\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_d]$$



Formal definitions

Data point as an n-dimensional column vector

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_n \end{bmatrix}$$

Data set as a [n x d] matrix

$$\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_d]$$

Projected datapoint

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \dots \\ y_m \end{bmatrix}$$

New data set as a [m x d] matrix

$$\mathbf{Y} = [\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_d]$$

Dimensionality reduction is a mapping $\mathbf{x} \mapsto \mathbf{y}$



Linear Dimensionality reduction

Define $\mathbf{x} \mapsto \mathbf{y}$ to be a *linear* mapping:

$$\mathbf{Y} = \mathbf{A}\mathbf{X}$$

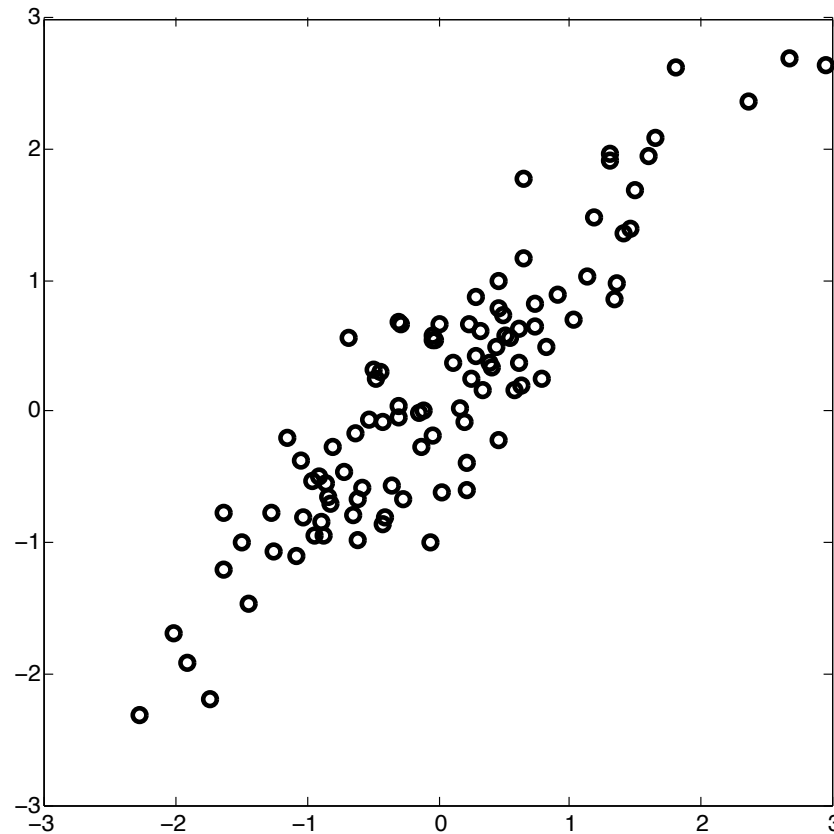
[m x d] projected data matrix [m x n] projection matrix [n x d] data matrix (zero-meanned)

Selects a *subspace* to best represent the data

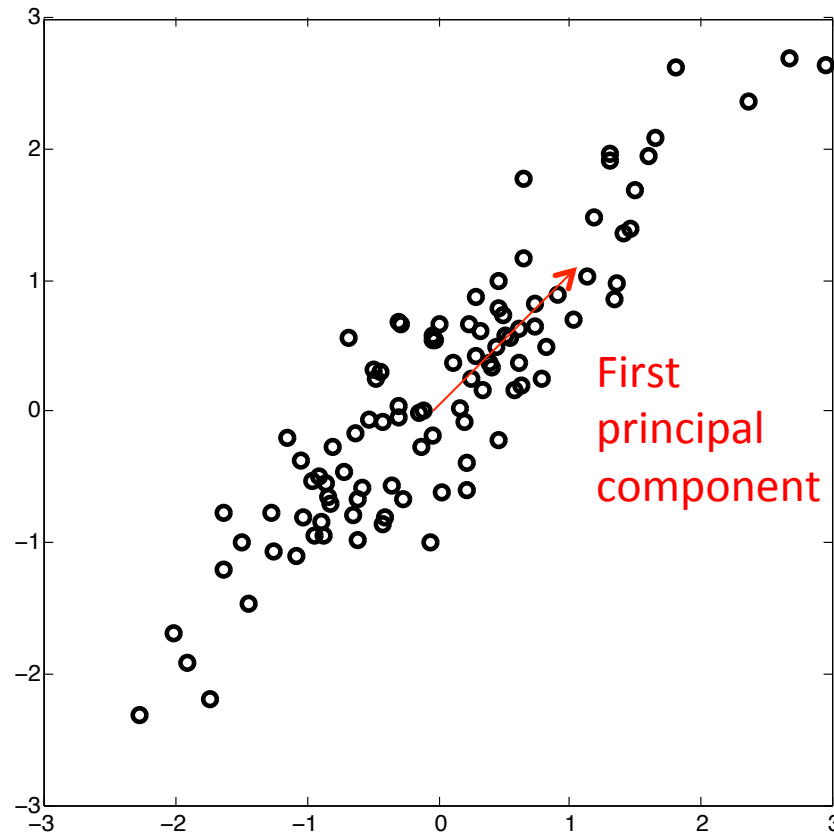
The most common method is **Principal Component Analysis (PCA)**



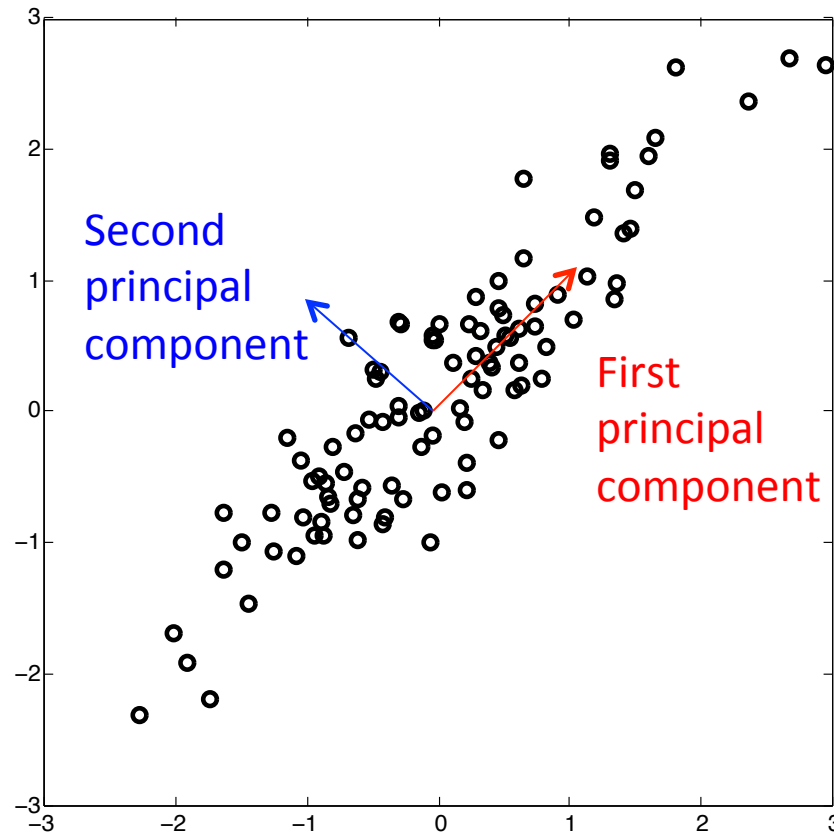
PCA selects orthogonal directions maximizing variance



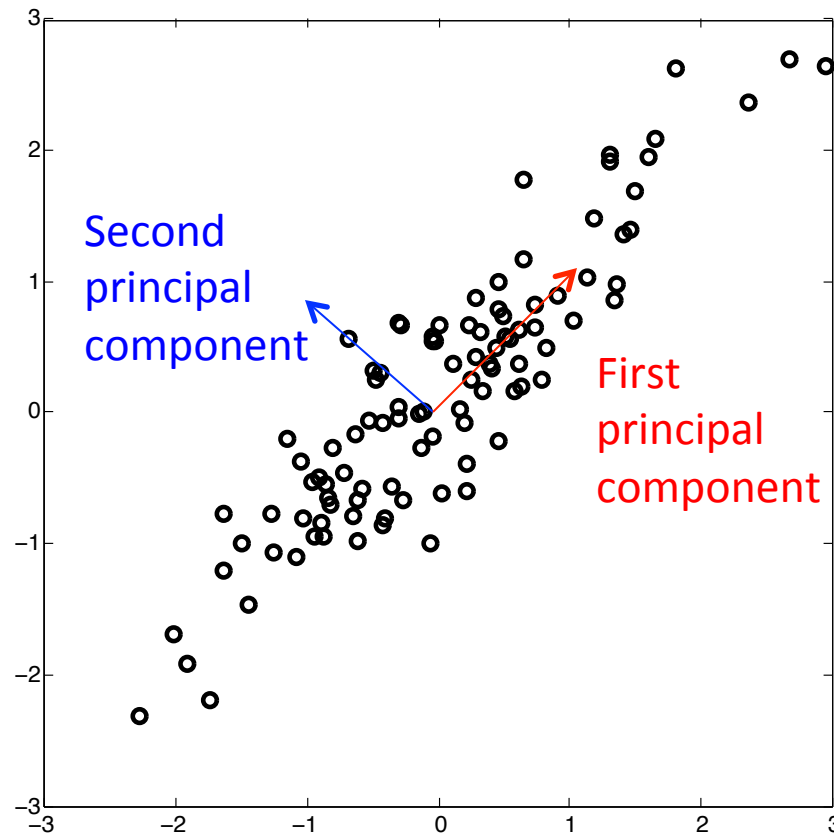
PCA selects orthogonal directions maximizing variance



PCA selects orthogonal directions maximizing variance



PCA selects orthogonal directions maximizing variance



Vectors have unit length - PCA provides an *orthonormal basis*



Covariance matrix

- An $m \times m$ matrix, with:

$$\Sigma_{\mathbf{X}}(i, j) = E[(\mathbf{x}_i - \mu_i)(\mathbf{x}_j - \mu_j)]$$

- Estimate using:

$$\Sigma_{\mathbf{X}} \equiv \frac{1}{d-1} \mathbf{X} \mathbf{X}^T$$

↑
Mean of j th
attribute

- Properties:

- Diagonal of $\Sigma_{\mathbf{X}}$ has the *variance* of \mathbf{x}
- Off-diagonal terms of $\Sigma_{\mathbf{X}}$ represent *covariances* of \mathbf{x}
- It's square, symmetric, positive semi-definite



Maximizing variance of a projection (one dimension)

$$\begin{aligned} E[\mathbf{v}^T \mathbf{x} - E[\mathbf{v}^T \mathbf{x}]]^2 &= E[(\mathbf{v}^T [\mathbf{x} - E\mathbf{x}])^2] \\ &= \mathbf{v}^T E[(\mathbf{x} - E\mathbf{x})(\mathbf{x} - E\mathbf{x})^T] \mathbf{v} \\ &= \mathbf{v}^T \Sigma_{\mathbf{x}} \mathbf{v} \end{aligned}$$



For a zero-mean covariance matrix, the unit vector minimizing this quantity is the top eigenvector of $\Sigma_{\mathbf{x}}$

$$(\mathbf{X}\mathbf{X}^T) \mathbf{v}_i = \lambda_i \mathbf{v}_i$$



Summary: PCA Recipe

1. Convert the data to have zero mean
2. Form \mathbf{A} using the top n eigenvectors of the sample covariance matrix \mathbf{XX}^T
 - Equivalently, use the left singular vectors of \mathbf{X} associated with the largest singular values
3. Project the data: $\mathbf{y} = \mathbf{A}(\mathbf{x} - \mu)$
4. To reconstruct: $\hat{\mathbf{x}} = \mathbf{A}^T \mathbf{y} + \mu$



Example: Eigenfaces

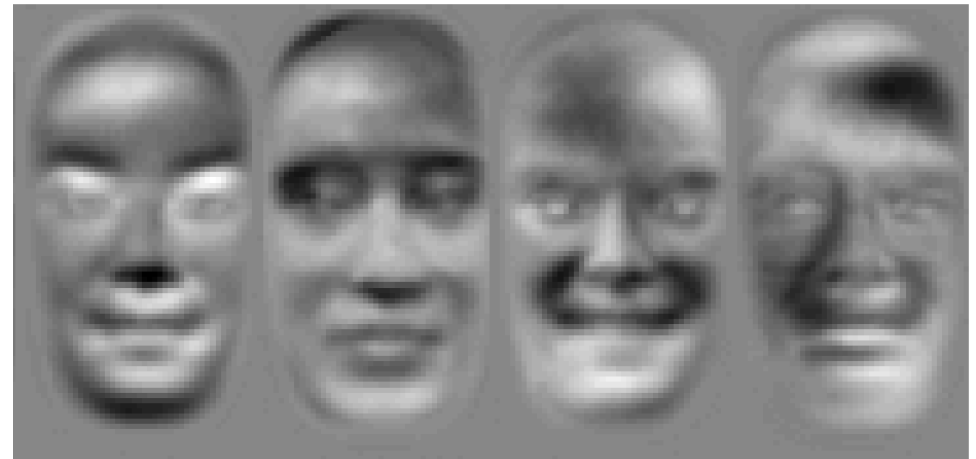
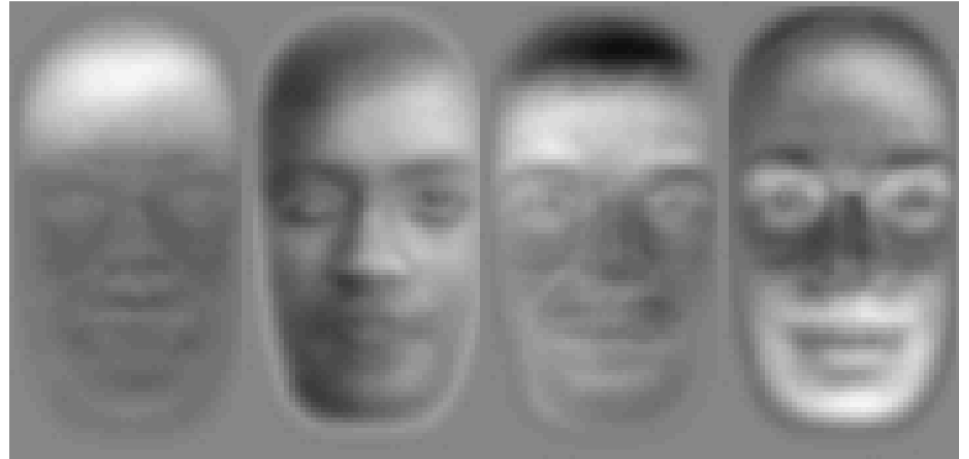


FERET frontal-view image database, from [Moghaddam, Wahid, and Pentland 1998]



Interpreting the principal components

The vectors describe main axes of variation

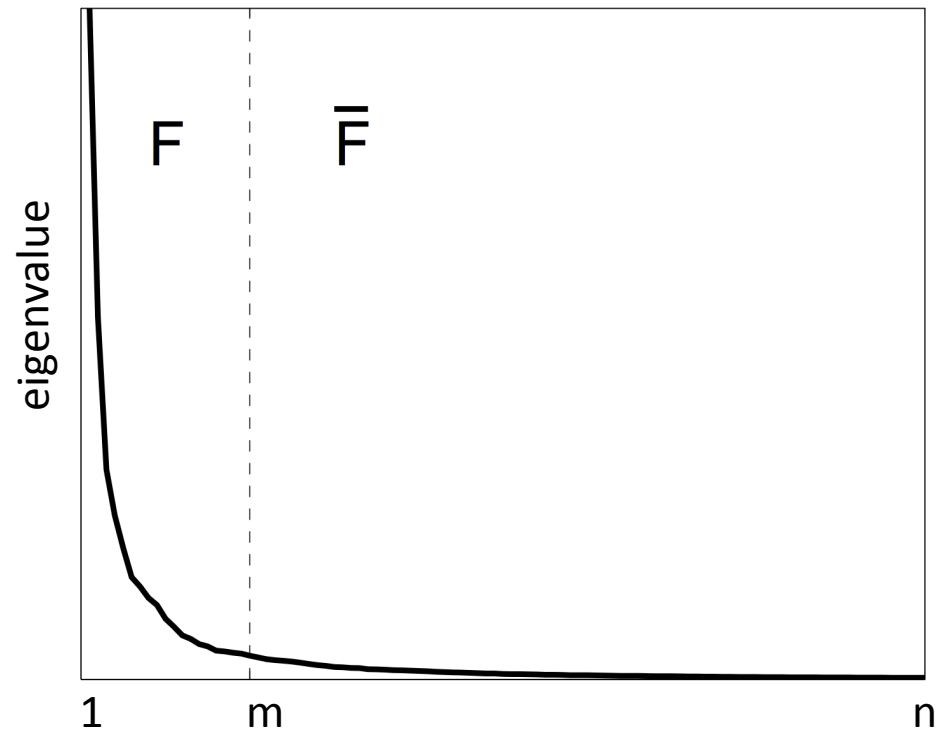


[Moghaddam, Wahid, and Pentland 1998]



Interpreting eigenvalues

Eigenvalue
“fall off”
suggests the
number of
non-noise
components



[Moghaddam, Wahid, and Pentland 1998]



Efficient implementations for large d or large n

- Use Singular Value Decomposition

$$\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^T$$

\mathbf{U} is a valid orthonormal basis

Diagonal matrix of singular values – use instead of covariance matrix eigenvalues

- Calculate one component at a time [e.g. Roweis, *NIPS 1998*]
- Use sequential estimation [e.g. Warmuth & Kuzmin, *JMLR 2008*]



Summary

- Principal Component Analysis (PCA) is a reliable, standard method for dimensionality reduction and visualization
- It finds an orthonormal basis to maximize the variance of the projected data
- The basis vectors and eigenvalues can provide insight about principal axes of variation in your data



Formal definitions

Data points $\mathbf{x} \in \mathbb{R}^n$, $\mathbf{y} \in \mathbb{R}^m$ $m \ll n$

$$\mathbf{x}_i = \begin{bmatrix} x_{i1} \\ x_{i2} \\ \dots \\ x_{in} \end{bmatrix}$$

Data point
having n
features

$$\mathbf{y}_i = \begin{bmatrix} y_{i1} \\ y_{i2} \\ \dots \\ y_{im} \end{bmatrix}$$

$$\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_d]$$

Data set as an $[n \times d]$ matrix

$$\mathbf{Y} = [\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_d]$$

Dimensionality reduction is a mapping $\mathbf{x} \mapsto \mathbf{y}$



Formal definitions

Data point as an n-dimensional row vector

$$\mathbf{x} = [x_1, x_2, \dots, x_n]$$

Projected datapoint

$$\mathbf{y} = [y_1, y_2, \dots, y_m]$$

Data set as a [d x n] matrix

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \dots \\ \mathbf{x}_d \end{bmatrix}$$

New data set as a [d x m] matrix

$$\mathbf{Y} = \begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \\ \dots \\ \mathbf{y}_d \end{bmatrix}$$

Dimensionality reduction is a mapping $\mathbf{x} \mapsto \mathbf{y}$

