Machine Learning CS 4641



Dimension Reduction

Nakul Gopalan Georgia Tech

These slides are adopted based on slides from Le Song, Chao Zhang, Mahdi Roozbahani and Barnabás Póczos.

Outline

- Overview
- Principle Component Analysis: Main Idea
- The PCA Algorithm
- PCA and SVD
- Summary

Motivating Example: Data Visualization

53 blood and urine samples (features) from 65 people

		H-WBC	H-RBC	H-Hgb	H-Hct	H-MCV	H-MCH	H-MCHC
(A1	8.0000	4.8200	14.1000	41.0000	85.0000	29.0000	34.0000
	A2	7.3000	5.0200	14.7000	43.0000	86.0000	29.0000	34.0000
	A3	4.3000	4.4800	14.1000	41.0000	91.0000	32.0000	35.0000
	A4	7.5000	4.4700	14.9000	45.0000	101.0000	33.0000	33.0000
$\langle $	A5	7.3000	5.5200	15.4000	46.0000	84.0000	28.0000	33.0000
	A6	6.9000	4.8600	16.0000	47.0000	97.0000	33.0000	34.0000
	A7	7.8000	4.6800	14.7000	43.0000	92.0000	31.0000	34.0000
	A8	8.6000	4.8200	15.8000	42.0000	88.0000	33.0000	37.0000
	A9	5.1000	4.7100	14.0000	43.0000	92.0000	30.0000	32.0000

• Matrix format (65x53)

Instances

Features

Difficult to see the correlations of different features

Motivating Example: Data Visualization

Is there a representation better than the coordinate axes?

Is it really necessary to show all the 53 dimensions?

... what if there are strong correlations between the features?

How could we find the *smallest* subspace of the 53-D space that keeps the *most information* about the original data?

A Solution: Dimension Reduction

Another Example: Dimension Reduction for Text



What are the relations between data points?





0.5

Bag-of-Words Representations



Term-Document Data Matrix – Bag-of-words

	database	SQL	index	regression	likelihood	linear
d1	24	21	9	0	0	3
d2	32	10	5	0	3	0
d3	12	16	5	0	0	0
d4	6	7	2	0	0	0
d5	43	31	20	0	3	0
d6	2	0	0	18	7	16
d7	0	0	1	32	12	0
d8	3	0	0	22	4	2
d9	1	0	0	34	27	25
d10	6	0	0	17	4	23

••• Many more features

Solution: Dimension Reduction

What is Dimension Reduction?

- The process of reducing the number of random variables under consideration
 - One can combine, transform or select variables
 - One can use linear or nonlinear operations



vector in R^d

Applications of Dimension Reduction

- The dimension-reduced data can be used for
 - Visualizing, exploring and understanding the data
 - Aggregating weak signals in the data
 - Cleaning the data
 - Speeding up subsequent learning task
 - Building simpler model later
- Key questions of a dimensionality reduction algorithm
 - What is the criterion for carrying out the reduction process?
 - What are the algorithm steps?

Outline

- Overview
- Principle Component Analysis: Main Idea
- The PCA Algorithm
- PCA and SVD
- Summary

PCA: Dimension Reduction by Capturing Variation

- There are many criteria (geometric based, information theory based, etc.)
- One criterion: want to capture variation in data
 - variations are "signals" or information in the data
 - need to normalize each variables first
- In the process, also discover variables or dimensions highly correlated
 - represent highly related phenomena
 - combine them to form a stronger signal
 - lead to simpler presentation

Capturing Variation in Data



Two Equivalent Perspectives of PCA



PCA:

- Orthogonal projection of the data onto a lower-dimension linear space that...
 - Imaximizes variance of projected data (purple line)
 - Image: mean squared distance between
 - data point and
 - projections (sum of blue lines)



Outline

- Overview
- Principle Component Analysis: Main Idea
- The PCA Algorithm
- PCA and SVD
- Summary

What is variance equation?

$$Var(x) = \frac{1}{n} \sum_{i=1}^{n} (x_i - \mu)^2$$

Formulating the Problem

- Given *n* data points, $\{x_1, x_2, ..., x_n\} \in \mathbb{R}^d$ with their mean $\mu = \frac{1}{n} \sum_{i=1}^n x_i$
- Find a direction $w \in \mathbb{R}^d$ where

$$\|w\| = \sqrt{\sum_{j \in d} w_j^2} = 1$$

We constrain the norm of w to be equal to one to avoid having very large variance in each new dimension. • Given *n* data points, $\{x_1, x_2, ..., x_n\} \in \mathbb{R}^d$ with their mean μ

$$\|w\| = \sqrt{\sum_{j \in d} w_j^2} = 1$$
 $\mu = \frac{1}{n} \sum_{i=1}^n x_i$

 Such that the variance (or variation) of the data along direction w is maximized

$$\max_{||w||=1} \frac{1}{n} \sum_{i=1}^{n} (x_i w - \mu w)^2$$
variance in new feature space

An Optimization Problem

Manipulate the objective with linear algebra

$$\frac{1}{n}\sum_{i=1}^{n}(x_{i}w-\mu w)^{2} = \frac{1}{n}\sum_{i=1}^{n}((x_{i}-\mu)w)^{2} =$$

$$= \frac{1}{n} \sum_{i=1}^{n} \left((x_{i} - \mu)w \right)^{T} ((x_{i} - \mu)w) = \frac{1}{n} \sum_{i=1}^{n} w^{T} (x_{i} - \mu)^{T} (x_{i} - \mu)w$$

$$(AB)^{T} = B^{T}A^{T}$$

$$w^{T} \left(\frac{1}{n} \sum_{i=1}^{n} (x_{i} - \mu)^{T} (x_{i} - \mu) \right) w = w^{T}Cw$$

Covariance matrix

Equivalence to The Eigenvalue Problem

Claim:

$$\max_{|w||=1} w^T C w$$

Form lagrangian function of the optimization problem

$$L(w,\lambda) = w^{\mathsf{T}}Cw + \lambda(1 - w^{t}w)$$

- If w is a maximum of the original optimization problem, then there exists a λ, where (w, λ) is a stationary point of L(w, λ)
- This implies that

$$\frac{\partial L}{\partial w} = 0 = 2Cw - 2\lambda w \quad \Rightarrow \ Cw = \lambda w$$

Eigen-Value Problem

Eigen-value problem

d: dimension

• Given a symmetric matrix $C \in \mathbb{R}^{d \times d}$

C is also a positive semidefinite matrix

- Find a vector $w \in \mathbb{R}^d$ and ||w|| = 1
- Such that

$$Cw = \lambda w$$

- There will be multiple solution of $w_1, w_2, ..., w_d$ for its corresponding $\lambda_1, \lambda_2, ..., \lambda_d$
 - They are ortho-normal: $w_i^T w_i = 1$ $w_i^T w_j = 0$

Eigenvalues and Eigenvectors

• Given a square matrix $A \in \mathbb{R}^{d \times d}$ we say that $\lambda \in \mathbb{C}$ is an eigenvalue of A and $x \in \mathbb{C}^{d}$ is an eigenvector if

$$Ax = \lambda x, \qquad x \neq 0$$

- Intuitively this means that upon multiplying the matrix A with a vector x, we get the same vector, but scaled by a parameter λ
- Geometrically, we are transforming the matrix A from its original orthonormal basis/co-ordinates to a new set of orthonormal basis x with magnitude as λ

Computing Eigenvalues and Eigenvectors

We can rewrite the original equation in the following manner

$$Ax = \lambda x, \quad x \neq 0$$

$$\Rightarrow (A - \lambda I) x = 0, \quad x \neq 0$$

- This is only possible if (A λI) is singular, that is | (A λI) | =
 0.
- Thus, eigenvalues and eigenvectors can be computed.
 - Compute the determinant of $A \lambda I$.
 - This results in a polynomial of degree d.
 - Find the roots of the polynomial by equating it to zero.
 - The d roots are the d eigenvalues of A. They make $A \lambda I$ singular.
 - For each eigenvalue λ , solve $(A \lambda I) x$ to find an eigenvector x

Matrix Eigen Decomposition

- All the eigenvectors can be written together as AX = XΛ where the columns of X are the eigenvectors of A, and Λ is a diagonal matrix whose elements are eigenvalues of A
- If the eigenvectors of A are invertible, then $A = X\Lambda X^{-1}$
- There are several properties of eigenvalues and eigenvectors
 - $Tr(A) = \sum_{i=1}^{d} \lambda_i$
 - $|A| = \prod_{i=1}^{d} \lambda_i$
 - Rank of A is the number of non-zero eigenvalues of A
 - If A is non-singular then $1/\lambda_i$ are the eigenvalues of A^{-1}
 - The eigenvalues of a diagonal matrix are the diagonal elements of the matrix itself!

Principal Direction of the Data



Variance in the Principal Direction

Principal direction w satisfies

$$Cw = \lambda w = w\lambda$$

Variance in principal direction is

 $w^T C w$

$$= w^T w \lambda$$

$$= \lambda$$
 eigen-value

Multiple Principal Directions

- Directions w₁, w₂, ... which has
 - the largest variances
 - but are orthogonal to each other
- Take the eigenvectors w₁, w₂, ... of C corresponding to
 - the largest eigenvalue λ_1 ,
 - the second largest eigenvalue λ_2



Extra Principal Directions



Relations Between Principal Components

Principal component #1: points in the direction of the **largest variance**.

Each subsequent principal component

- is **orthogonal** to the previous ones, and
- points in the directions of the largest variance of the residual subspace

The PCA Algorithm

- Given *n* data points, $\{x_1, x_2, ..., x_n\} \in \mathbb{R}^d$ with mean
- Step 1: Estimate the mean and covariance matrix from data

$$\mu = \frac{1}{n} \sum_{i=1}^{n} x_i \quad and \quad C = \frac{1}{n} \sum_{i=1}^{n} (x_i - \mu)^T (x_i - \mu)$$
Principal directions

- Step 2: Take the eigenvectors w₁, w₂, ... of C corresponding to the largest eigenvalue λ₁, the second largest eigenvalue λ₂...
- Step 3: Compute reduced representation

$$z_{i} = \begin{pmatrix} (x_{i} - \mu_{1}) \\ \sigma_{1} \end{pmatrix} w_{1} \quad \frac{(x_{i} - \mu_{2})}{\sigma_{2}} w_{2} \dots \end{pmatrix} \qquad z \Rightarrow n \times k$$
Normalizing by
standard deviation

Outline

- Overview
- Principle Component Analysis: Main Idea
- The PCA Algorithm
- PCA and SVD
- Summary

Singular Value Decomposition

n: instances $X_{n \times d}$ d: dimensions X is a centered matrix – i.e., mean subtracted!!

$$U_{n \times n} \rightarrow unitary \ matrix \rightarrow U \times U^{T} = I$$
$$X = U\Sigma V^{T} \qquad \Sigma_{n \times d} \rightarrow diagonal \ matrix$$

Matrix compression: K dimensions out of d

According to PCA $\rightarrow Cw = \lambda w = w\lambda$

Covariance
$$C_{d \times d} = \frac{1}{n} \sum_{i=1}^{n} (x^{i} - \mu)^{T} (x^{i} - \mu) = \frac{X^{T} X}{n}$$

$$X = U\Sigma V^{T}$$

$$C = \frac{X^{T}X}{n}$$

$$C = \frac{V\Sigma^{T}U^{T}U\Sigma V^{T}}{n} = \frac{V\Sigma^{2}V^{T}}{n}$$

$$C = \frac{V\Sigma^2 V^T}{n} = V \frac{\Sigma^2}{n} V^T$$

$$CV = V\frac{\Sigma^2}{n}V^T V = V\frac{\Sigma^2}{n}$$

According to Eigen-decomposition definition $\rightarrow CV = V\Lambda$

V is the eigen vectors of covariance (Principal directions)

$$\lambda_i = \frac{\sigma_i^2}{n} \rightarrow$$
 The eigenvalues of covariance matrix

Let's project the data (X) on principal directions: $XV = U\Sigma V^T V = U\Sigma$

XV is independent linear combinations of the original data

Projection of one instance (x) on the first principal direction using k dimensions

$$\begin{aligned} \mathbf{p}_{1} &= \begin{bmatrix} u_{1\times 1} \Sigma_{1\times 1} , u_{1\times 2} \Sigma_{2\times 2} , \dots , u_{1\times k} \Sigma_{k\times k} \end{bmatrix} \\ \mathbf{p}_{2} &= \begin{bmatrix} u_{2\times 1} \Sigma_{1\times 1} , u_{2\times 2} \Sigma_{2\times 2} , \dots , u_{2\times k} \Sigma_{k\times k} \end{bmatrix} \\ \begin{aligned} & U \Rightarrow n \times k \\ & \Sigma \Rightarrow k \times k \\ & \text{Upper left corner} \end{aligned}$$



Principal components (Scores) or projections on principal directions

In fact, using the SVD to perform PCA makes much better sense numerically than forming the covariance matrix to begin with, since the formation of $X^T X$ can cause loss of precision.

Are Principal Components Good for Classification?



Why PCA potentially works in classification?

the dimension with the largest variance corresponds to the dimension with the largest entropy and thus encodes <u>the most information</u> (Information Theory). The smallest eigenvectors will often simply represent noise components, whereas the largest eigenvectors often correspond to the principal components that define the data.

Result with PCA – Algo. For Face detection

1. Treat each window in the image like a vector



2. Test whether x matches some y_j in the database



SSD:
$$(y_j - x)^2$$

Cross-correlation: $y_j \cdot x$
NCC, zero-mean NCC...

Slide by Derek Hoiem

Mean Face and Eigen Faces

Top eigenvectors: u₁,...u_k



Mean: µ



Slide by Derek Hoiem

Modern approaches

• Auto-encoders!

Example low dimensional embedding of MNIST dataset

From <u>Sebastian Pölsterl</u>'s blog



Outline

- Overview
- Principle Component Analysis: Main Idea
- The PCA Algorithm
- PCA and SVD
- Summary

Summary

PCA

- 。 Finds orthonormal basis for data
- 。 Sorts dimensions in order of "importance"
- 。 Discard low significance dimensions

Uses

- 。Get concise low-dimensional representations
- 。Remove noise
- Not magic
 - Doesn't know class labels
 - 。Can only capture linear variations

Image compression using PCA



PCs # 20



PCs # 50

