

Principal Component Analysis

Daifeng Wang

`dai feng . wang @ wisc . edu`

University of Wisconsin, Madison

Based on slides from Xiaojin Zhu, Yingyu Liang and
UCL Linear Algebra & Matrices, MfD 2009, modified
by Daifeng Wang

Outline

- Basic review on linear algebra
- Introduction to dimensionality reduction
- Principal component analysis:
formulation and computation
- Applications

REVIEW ON LINEAR ALGEBRA

Vector

- Not a physics vector (magnitude, direction)
- Column of numbers e.g. intensity of same voxel at different time points

$$\begin{bmatrix} x1 \\ x2 \\ x3 \end{bmatrix}$$

Matrices

- Rectangular display of vectors in rows and columns
- Can inform about the same vector intensity at different times or different voxels at the same time
- Vector is just a $n \times 1$ matrix

$$\mathbf{A} = \begin{bmatrix} 1 & 2 & 3 \\ 5 & 4 & 1 \\ 6 & 7 & 4 \end{bmatrix} \quad \mathbf{C} = \begin{bmatrix} 1 & 4 \\ 2 & 7 \\ 3 & 8 \end{bmatrix} \quad \mathbf{D} = \begin{bmatrix} d_{11} & d_{12} & d_{13} \\ d_{21} & d_{22} & d_{23} \\ d_{31} & d_{32} & d_{33} \end{bmatrix}$$

Square (3 x 3)

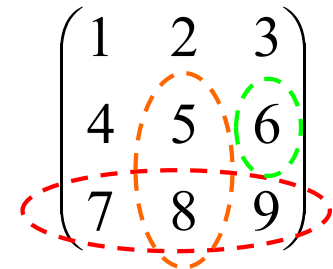
Rectangular (3 x 2)

d_{ij} : i^{th} row, j^{th} column

Defined as rows x columns (R x C)

Matrices in Python

- $X = [[1, 2, 3], [4, 5, 6], [7, 8, 9]]$
- Index from 0 to nrow/ncol-1
- := all row or column



Subscripting – each element of a matrix can be addressed with a pair of numbers;
[row first. column second]

e.g. $X[1, 2] = 6$

$$X[2, :] = (7 \ 8 \ 9)$$

$$X[1:2, 1] = \begin{pmatrix} 5 \\ 8 \end{pmatrix}$$

Transposition

$$\mathbf{b} = \begin{bmatrix} 1 \\ 1 \\ 2 \end{bmatrix}$$

$$\mathbf{b}^T = [1 \quad 1 \quad 2]$$

$$\mathbf{d} = [3 \quad 4 \quad 9]$$

$$\mathbf{d}^T = \begin{bmatrix} 3 \\ 4 \\ 9 \end{bmatrix}$$

column



row

row



column

$$\mathbf{A} = \begin{bmatrix} 1 & 2 & 3 \\ 5 & 4 & 1 \\ 6 & 7 & 4 \end{bmatrix}$$

$$\mathbf{A}^T = \begin{bmatrix} 1 & 5 & 6 \\ 2 & 4 & 7 \\ 3 & 1 & 4 \end{bmatrix}$$

Scalar multiplication

- Scalar \times matrix = scalar multiplication

$$\lambda \begin{pmatrix} a & b & c \\ d & e & f \end{pmatrix} = \begin{pmatrix} \lambda a & \lambda b & \lambda c \\ \lambda d & \lambda e & \lambda f \end{pmatrix}$$

Matrix Calculations

Addition

- Commutative: $\mathbf{A} + \mathbf{B} = \mathbf{B} + \mathbf{A}$
- Associative: $(\mathbf{A} + \mathbf{B}) + \mathbf{C} = \mathbf{A} + (\mathbf{B} + \mathbf{C})$

$$\mathbf{A} + \mathbf{B} = \begin{bmatrix} 2 & 4 \\ 2 & 5 \end{bmatrix} + \begin{bmatrix} 1 & 0 \\ 3 & 1 \end{bmatrix} = \begin{bmatrix} 2+1 & 4+0 \\ 2+3 & 5+1 \end{bmatrix} = \begin{bmatrix} 3 & 4 \\ 5 & 6 \end{bmatrix}$$

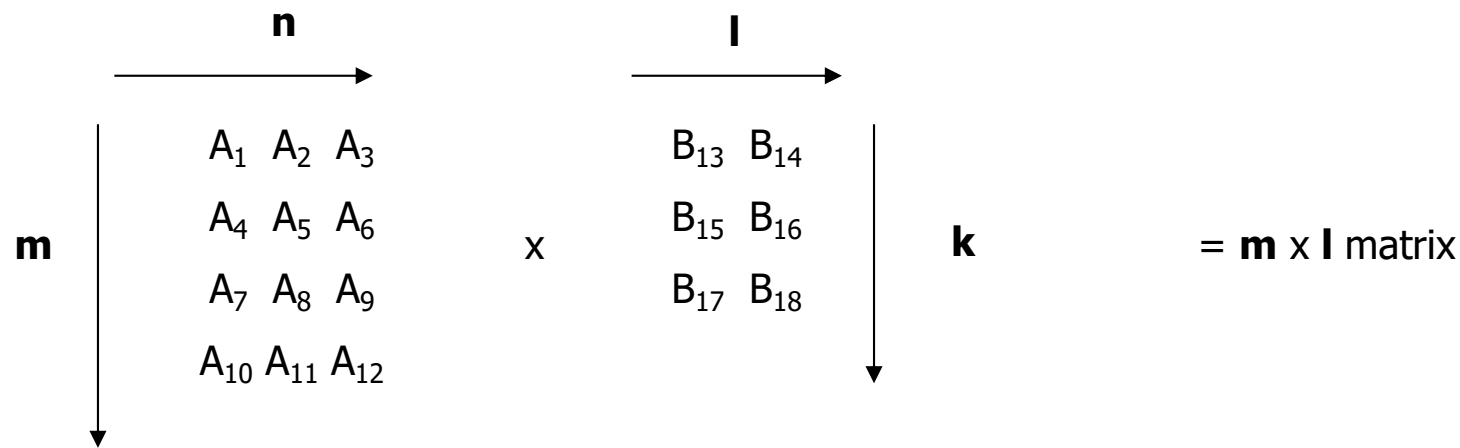
Subtraction

- By adding a negative matrix

$$\mathbf{A} - \mathbf{B} = \begin{bmatrix} 2 & 4 \\ 5 & 3 \end{bmatrix} - \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} = \begin{bmatrix} 2 & 4 \\ 5 & 3 \end{bmatrix} + \begin{bmatrix} -1 & -2 \\ -3 & -4 \end{bmatrix} = \begin{bmatrix} 1 & 2 \\ 2 & -1 \end{bmatrix}$$

Matrix Multiplication

"When A is a $m \times n$ matrix & B is a $k \times l$ matrix, AB is only possible if $n=k$. The result will be an $m \times l$ matrix"



Number of columns in A = Number of rows in B

Matrix multiplication

- Multiplication method:

Sum over product of respective rows and columns

$$\begin{pmatrix} 1 & 0 \\ 2 & 3 \end{pmatrix} \times \begin{pmatrix} 2 & 1 \\ 3 & 1 \end{pmatrix} = \begin{pmatrix} \mathbf{c_{11}} & \mathbf{c_{12}} \\ \mathbf{c_{21}} & \mathbf{c_{22}} \end{pmatrix} \quad \begin{array}{l} \text{Define output} \\ \text{matrix} \end{array}$$

A **B**

$$= \begin{bmatrix} (1 \times 2) + (0 \times 3) & (1 \times 1) + (0 \times 1) \\ (2 \times 2) + (3 \times 3) & (2 \times 1) + (3 \times 1) \end{bmatrix}$$

$$= \begin{pmatrix} \mathbf{2} & \mathbf{1} \\ \mathbf{13} & \mathbf{5} \end{pmatrix}$$

$$C_{i,j} = \sum_{k=1}^n A_{i,k} B_{k,j}$$

Matrix multiplication

- Matrix multiplication is NOT commutative
- $AB \neq BA$
- Matrix multiplication IS associative
- $A(BC) = (AB)C$
- Matrix multiplication IS distributive
- $A(B+C) = AB+AC$
- $(A+B)C = AC+BC$

Identity matrix

Is there a matrix which plays a similar role as the number 1 in number multiplication?

Consider the $n \times n$ matrix:

$$I_n = \begin{pmatrix} 1 & 0 & 0 & \cdots & 0 \\ 0 & 1 & 0 & \cdots & 0 \\ 0 & 0 & 1 & \cdots & 0 \\ \cdot & & & & \cdot \\ \cdot & & & & \cdot \\ \cdot & & & & \cdot \\ 0 & 0 & 0 & \cdots & 1 \end{pmatrix}.$$

For any $n \times n$ matrix A , we have $A I_n = I_n A = A$

For any $n \times m$ matrix A , we have $I_n A = A$, and $A I_m = A$ (so 2 possible matrices)

Matrix inverse

- **Definition.** A matrix A is called nonsingular or invertible if there exists a matrix B such that:

$$AB = BA = I_n$$

$$\begin{bmatrix} 1 & 1 \\ -1 & 2 \end{bmatrix} \times \begin{bmatrix} \frac{2}{3} & \frac{-1}{3} \\ \frac{1}{3} & \frac{1}{3} \end{bmatrix} = \begin{bmatrix} \frac{2+1}{3} & \frac{-1+1}{3} \\ \frac{-2+2}{3} & \frac{1+2}{3} \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

- **Notation.** A common notation for the inverse of a matrix A is A^{-1} . So:

$$AA^{-1} = A^{-1}A = I_n$$

- The inverse matrix is unique when it exists. So if A is invertible, then A^{-1} is also invertible and then $(A^T)^{-1} = (A^{-1})^T$

SciPy - Linear Algebra

Python For Data Science Cheat Sheet

SciPy - Linear Algebra

Learn More Python for Data Science Interactively at www.datacamp.com



SciPy

The SciPy library is one of the core packages for scientific computing that provides mathematical algorithms and convenience functions built on the NumPy extension of Python.



Interacting With NumPy

Also see NumPy

```
>>> import numpy as np
>>> a = np.array([1,2,3])
>>> b = np.array([(1+5j),2j,3j], (4j,5j,6j))
>>> c = np.array([(1.5,2,3), (4,5,6)], [(3,2,1), (4,5,6)])
```

Index Tricks

```
>>> np.mgrid[0:5,0:5]
>>> np.ogrid[0:2,0:2]
>>> np.r_[[3, [0]*5, -1:1:10]]
>>> np.c_[b,c]
```

Create a dense meshgrid
Create an open meshgrid
Stack arrays vertically (row-wise)
Create stacked column-wise arrays

Shape Manipulation

```
>>> np.transpose(b)
>>> b.flatten()
>>> np.hstack((b,c))
>>> np.vstack((a,b))
>>> np.hsplit(c,2)
>>> np.vsplit(d,2)
```

Permute array dimensions
Flatten the array
Stack arrays horizontally (column-wise)
Stack arrays vertically (row-wise)
Split the array horizontally at the 2nd index
Split the array vertically at the 2nd index

Polynomials

```
>>> from numpy import poly1d
>>> p = poly1d([3,4,5])
```

Create a polynomial object

Vectorizing Functions

```
>>> def myfunc(a):
>>>     if a < 0:
>>>         return a*2
>>>     else:
>>>         return a/2
>>> np.vectorize(myfunc)
```

Vectorize functions

Type Handling

```
>>> np.real(c)
>>> np.imag(c)
>>> np.real_if_close(c, tol=1e-100)
>>> np.cast['*f'](np.pi)
```

Return the real part of the array elements
Return the imaginary part of the array elements
Return a real array if complex parts close to 0
Cast object to a data type

Other Useful Functions

```
>>> np.angle(b, deg=True)
>>> g = np.linspace(0, np.pi, num=5)
>>> g[3:] += np.pi
>>> np.unwrap(g)
>>> np.logspace(0, 10, 3)
>>> np.select([c<4], [c*2])

>>> misc.factorial(a)
>>> misc.comb(10, 3, exact=True)
>>> misc.central_diff_weights(3)
>>> misc.derivative(myfunc, 1.0)
```

Return the angle of the complex argument
Create an array of evenly spaced values (number of samples)
Unwrap
Create an array of evenly spaced values (log scale)
Return values from a list of arrays depending on conditions
Factorial
Combine N things taken at k time
Weights for Np-point central derivative
Find the n-th derivative of a function at a point

Linear Algebra

Also see NumPy

You'll use the `linalg` and `sparse` modules. Note that `scipy.linalg` contains and expands on `numpy.linalg`.

```
>>> from scipy import linalg, sparse
```

Creating Matrices

```
>>> A = np.matrix(np.random.random((2,2)))
>>> B = np.asmatrix(b)
>>> C = np.mat(np.random.random(10,5))
>>> D = np.mat([[3,4], [5,6]])
```

Basic Matrix Routines

Inverse

```
>>> A.I
>>> linalg.inv(A)
>>> A.T
>>> A.H
>>> np.trace(A)
```

Inverse
Inverse
Transpose matrix
Conjugate transposition
Trace

Norm

```
>>> linalg.norm(A)
>>> linalg.norm(A,1)
>>> linalg.norm(A,np.inf)
```

Frobenius norm
L1 norm (max column sum)
L inf norm (max row sum)

Rank

```
>>> np.linalg.matrix_rank(C)
```

Matrix rank

Determinant

```
>>> linalg.det(A)
```

Determinant

Solving linear problems

```
>>> linalg.solve(A,b)
>>> E = np.mat(a).T
>>> linalg.lstsq(D,E)
```

Solver for dense matrices
Solver for dense matrices
Least-squares solution to linear matrix equation

Generalized inverse

```
>>> linalg.pinv(C)
>>> linalg.pinv2(C)
```

Compute the pseudo-inverse of a matrix (least-squares solver)
Compute the pseudo-inverse of a matrix (SVD)

Creating Sparse Matrices

```
>>> F = np.eye(3, k=1)
>>> G = np.mat(np.identity(2))
>>> C[C > 0.5] = 0
>>> H = sparse.csr_matrix(C)
>>> I = sparse.csc_matrix(D)
>>> J = sparse.dok_matrix(A)
>>> E.todense()
>>> sparse.isspmatrix_csc(A)
```

Create a 2x2 identity matrix
Create a 2x2 identity matrix
Compressed Sparse Row matrix
Compressed Sparse Column matrix
Dictionary Of Keys matrix
Sparse matrix to full matrix
Identify sparse matrix

Sparse Matrix Routines

Inverse

```
>>> sparse.linalg.inv(I)
```

Inverse

Norm

```
>>> sparse.linalg.norm(I)
```

Norm

Solving linear problems

```
>>> sparse.linalg.spsolve(H, I)
```

Solver for sparse matrices

Sparse Matrix Functions

```
>>> sparse.linalg.expm(I)
```

Sparse matrix exponential

Asking For Help

```
>>> help(scipy.linalg.diagsvd)
>>> np.info(np.matrix)
```

Matrix Functions

Addition

```
>>> np.add(A,D)
```

Addition

Subtraction

```
>>> np.subtract(A,D)
```

Subtraction

Division

```
>>> np.divide(A,D)
```

Division

Multiplication

```
>>> np.multiply(D,A)
>>> np.dot(A,D)
>>> np.vdot(A,D)
>>> np.inner(A,D)
>>> np.outer(A,D)
>>> np.tensordot(A,D)
>>> np.kron(A,D)
```

Multiplication
Dot product
Vector dot product
Inner product
Outer product
Tensor dot product
Kronecker product

Exponential Functions

```
>>> linalg.expm(A)
>>> linalg.expm2(A)
>>> linalg.expm3(D)
```

Matrix exponential
Matrix exponential (Taylor Series)
Matrix exponential (eigenvalue decomposition)

Logarithm Function

```
>>> linalg.logm(A)
```

Matrix logarithm

Trigonometric Functions

```
>>> linalg.sinm(D)
>>> linalg.cosm(D)
>>> linalg.tanm(A)
```

Matrix sine
Matrix cosine
Matrix tangent

Hyperbolic Trigonometric Functions

```
>>> linalg.sinhm(D)
>>> linalg.coshm(D)
>>> linalg.tanhm(A)
```

Hyperbolic matrix sine
Hyperbolic matrix cosine
Hyperbolic matrix tangent

Matrix Sign Function

```
>>> np.sign(A)
```

Matrix sign function

Matrix Square Root

```
>>> linalg.sqrtm(A)
```

Matrix square root

Arbitrary Functions

```
>>> linalg.funm(A, lambda x: x*x)
```

Evaluate matrix function

Decompositions

Eigenvalues and Eigenvectors

```
>>> la, v = linalg.eig(A)
```

Solve ordinary or generalized eigenvalue problem for square matrix
Unpack eigenvalues
First eigenvector
Second eigenvector
Unpack eigenvalues

```
>>> l1, l2 = la
>>> v[:,0]
>>> v[:,1]
```

```
>>> linalg.eigvals(A)
```

Singular Value Decomposition

```
>>> U, s, Vh = linalg.svd(B)
```

Singular Value Decomposition (SVD)

```
>>> M, N = B.shape
```

```
>>> Sig = linalg.diagsvd(s, M, N)
```

Construct sigma matrix in SVD

LU Decomposition

```
>>> F, L, U = linalg.lu(C)
```

LU Decomposition

Sparse Matrix Decompositions

```
>>> la, v = sparse.linalg.eigs(F, 1)
```

Eigenvalues and eigenvectors

```
>>> sparse.linalg.svds(H, 2)
```

SVD

DataCamp
Learn Python for Data Science Interactively



INTRODUCTION TO DIMENSIONALITY REDUCTION

Big & High-Dimensional Data

- High-Dimensions = Lot of Features

Document classification

Features per document =
thousands of words/unigrams
millions of bigrams, contextual
information



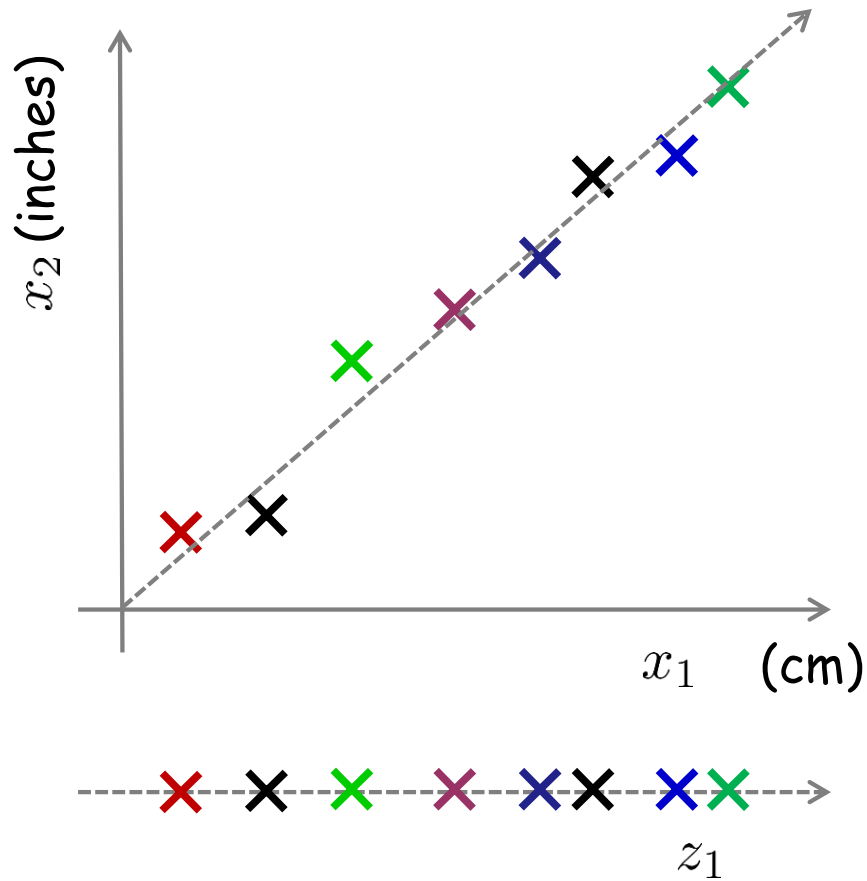
Surveys - Netflix

480189 users x 17770 movies

	movie 1	movie 2	movie 3	movie 4	movie 5	movie 6
Tom	5	?	?	1	3	?
George	?	?	3	1	2	5
Susan	4	3	1	?	5	1
Beth	4	3	?	2	4	2

- Big & High-Dimensional Data.
- Useful to learn lower dimensional representations of the data.
 - Given data points in D dimensions
 - Convert them to data points in $d < D$ dimensions
 - With minimal loss of information

Data Compression



Reduce data from
2D to 1D

$$(x_1 = \text{cm}, x_2 = \text{inch}) \rightarrow z_1$$

$$x^{(1)} \rightarrow z^{(1)}$$

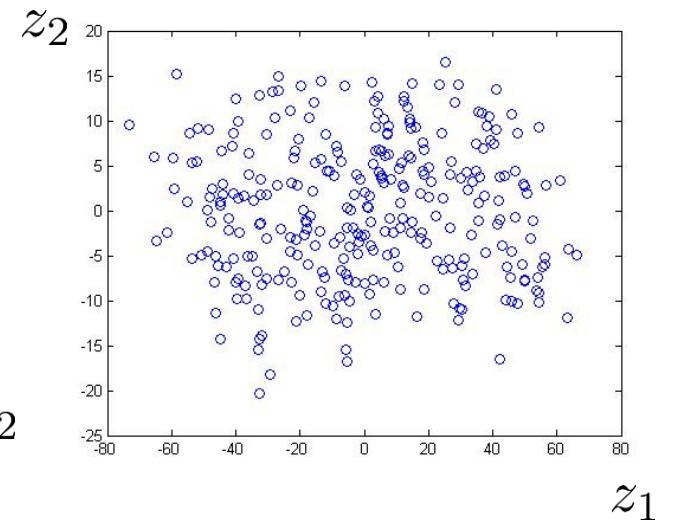
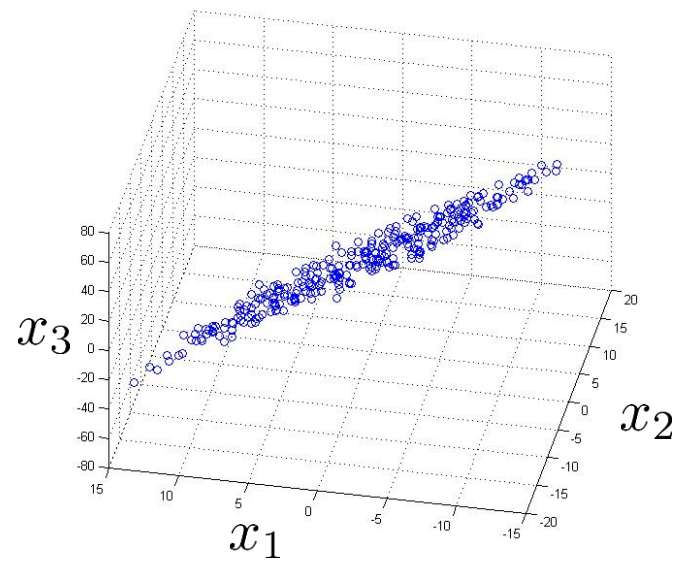
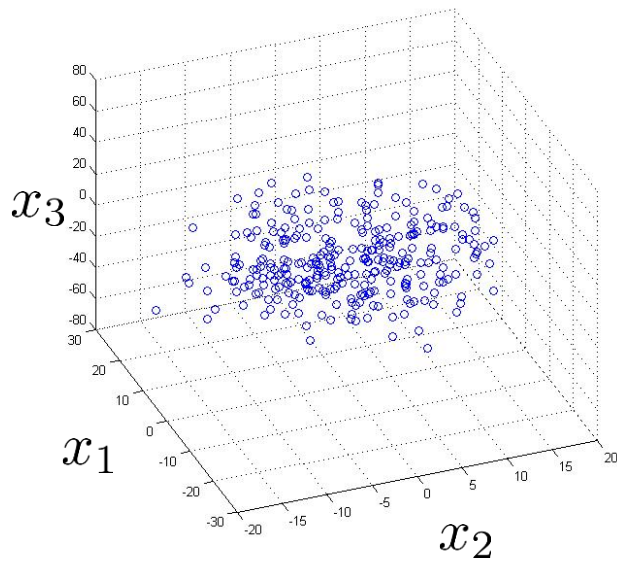
$$x^{(2)} \rightarrow z^{(2)}$$

⋮

$$x^{(m)} \rightarrow z^{(m)}$$

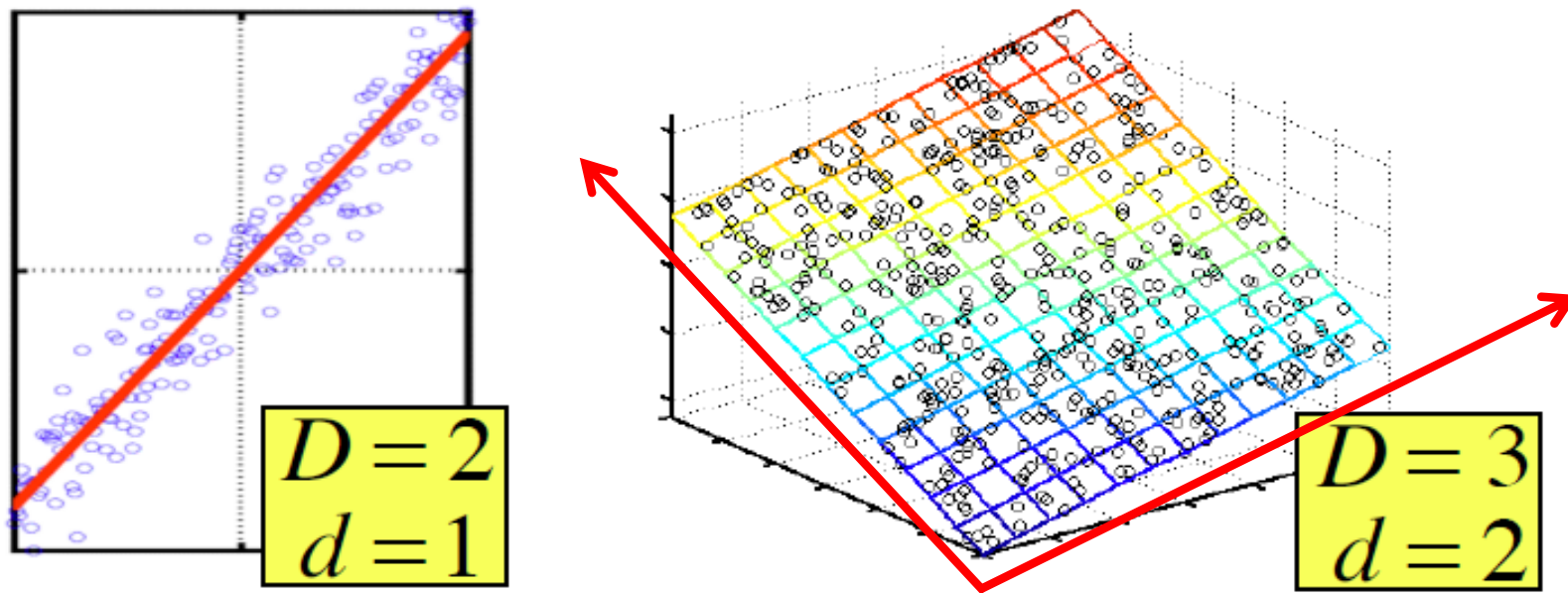
Data Compression

Reduce data from 3D to 2D



PRINCIPAL COMPONENT ANALYSIS (PCA)

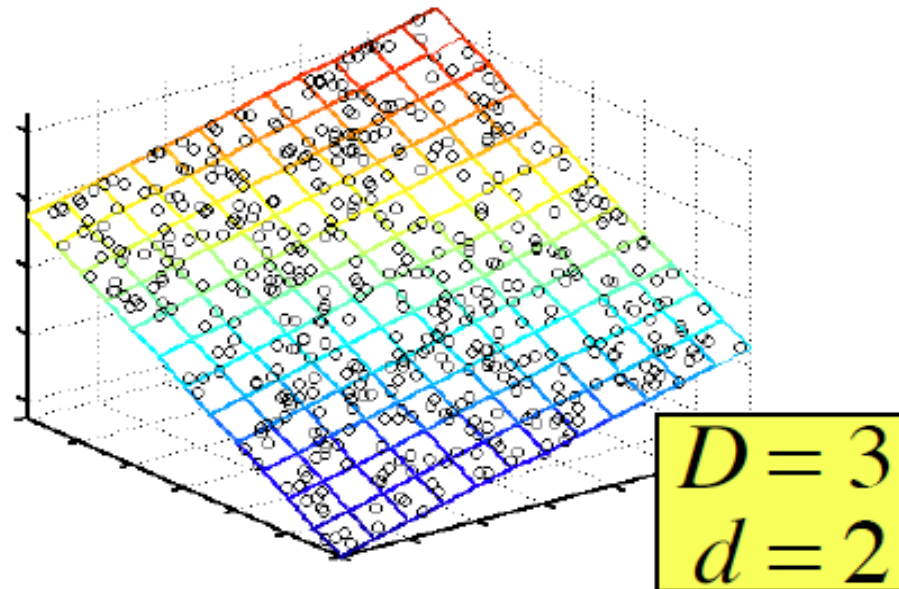
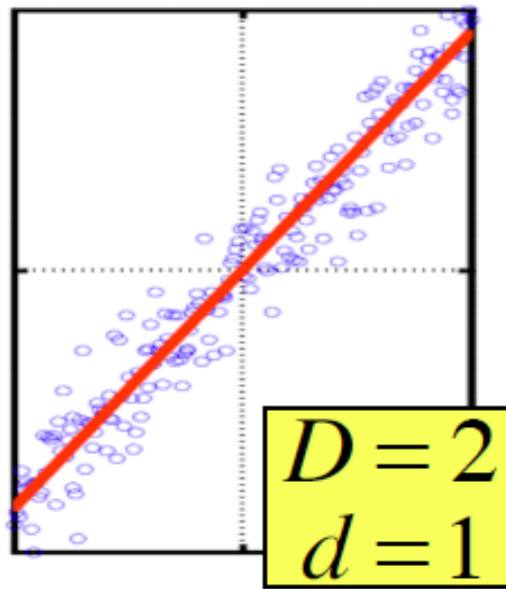
Principal Component Analysis (PCA)



In case where data lies on or near a low d -dimensional linear subspace, axes of this subspace are an effective representation of the data.

Identifying the axes is known as **Principal Components Analysis**, and can be obtained by using classic matrix computation tools (Eigen or Singular Value Decomposition).

PCA formulation



- Reduce from 2-dimension to 1-dimension: Find a direction (a red vector $v_1 \in R^D$) onto which to project the data so as to minimize the projection error.
- Reduce from D -dimension to d -dimension: Find d vectors $v_i \in R^D, i = 1, 2, \dots, d$ onto which to project the data, so as to minimize the projection error.
- v_i is called a **principal component (PC)**

Principal Component Analysis

Input: N data points (D -dim vectors)

$$\mathbf{x} \in \mathbb{R}^D: \mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$$

Output:

- d principal components (PCs)

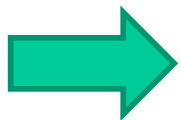
$$\mathbf{v}_j \in \mathbb{R}^D, j = 1, 2, \dots, d$$

s.t., Euclidean norm $\|\mathbf{v}_j\|_2 = (\sum_{k=1}^D v_j^2[k])^{\frac{1}{2}} = 1$

- For each \mathbf{x}_i , it's project coordinates on $\{\mathbf{v}_j\}$:

$$w_{i,j} = \mathbf{v}_j^T * \mathbf{x}_i, j = 1, 2, \dots, d$$

- Now \mathbf{x}_i , a D -dim vector can be represented by a d -dim vector ($d < D$)



$$[w_{i,1}, w_{i,2}, \dots, w_{i,d}]$$

Learning Representations

PCA, Kernel PCA, ICA, CCA: Powerful unsupervised learning techniques for extracting hidden (potentially lower dimensional) structure from high dimensional datasets.

Useful for:

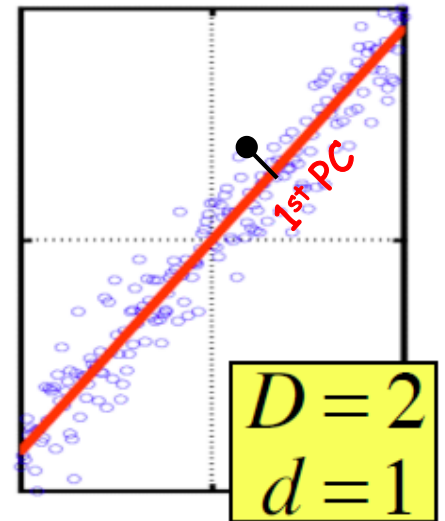
- Visualization
- More efficient use of resources (e.g., time, memory, communication)
- Statistical: fewer dimensions → better generalization
- Noise removal (improving data quality)
- Further processing by machine learning algorithms

PCA COMPUTATION

Principal Component Analysis (PCA)

Principal Components (PC) are orthogonal directions that capture most of the variance in the data.

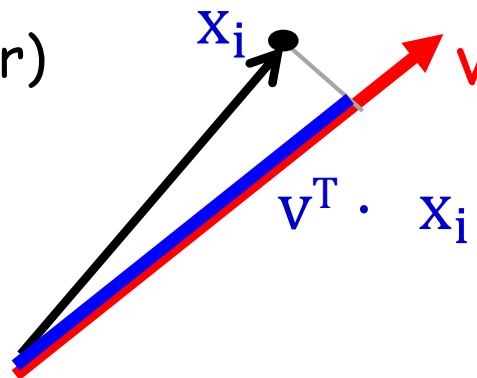
- First PC - direction of greatest variability in data.
- Projection of data points along first PC discriminates data most along any one direction (pts are the most spread out when we project the data on that direction compared to any other directions).



Quick reminder:

$\|v\|=1$, Point x_i (D-dimensional vector)

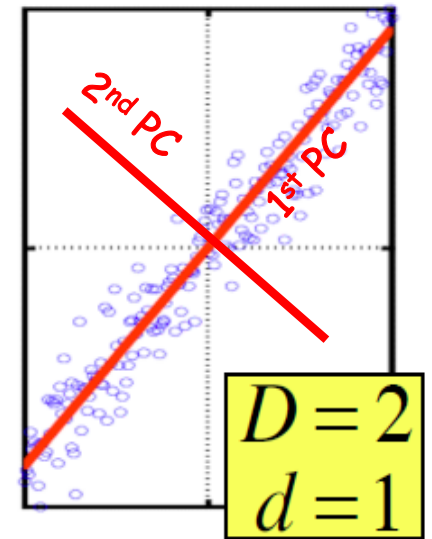
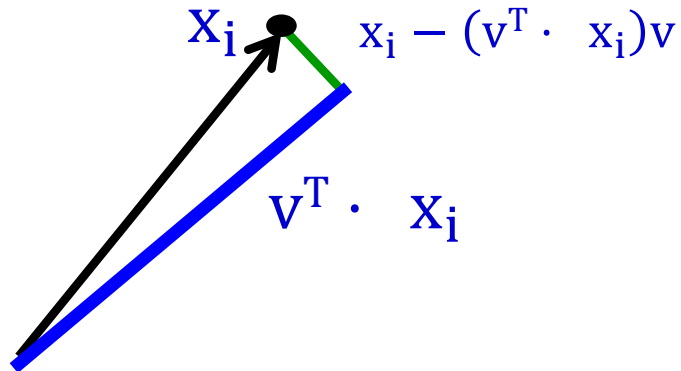
Projection of x_i onto v is $v^T \cdot x_i$



Principal Component Analysis (PCA)

Principal Components (PC) are orthogonal directions that capture most of the variance in the data.

- 1st PC - direction of greatest variability in data.



- 2nd PC - Next orthogonal (uncorrelated) direction of greatest variability

(remove all variability in first direction, then find next direction of greatest variability)

- And so on ...

Eigenvector and Eigenvalue

$$Ax = \lambda x$$

A: Square Matrix

x: Eigenvector

λ : Eigenvalue

Example

Show $x = \begin{bmatrix} 2 \\ 1 \end{bmatrix}$ is an eigenvector for $A = \begin{bmatrix} 2 & -4 \\ 3 & -6 \end{bmatrix}$

$$\text{Solution: } Ax = \begin{bmatrix} 2 & -4 \\ 3 & -6 \end{bmatrix} \begin{bmatrix} 2 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\text{But for } \lambda = 0, \lambda x = 0 \begin{bmatrix} 2 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

Thus, x is an eigenvector of A , and $\lambda = 0$ is an eigenvalue.

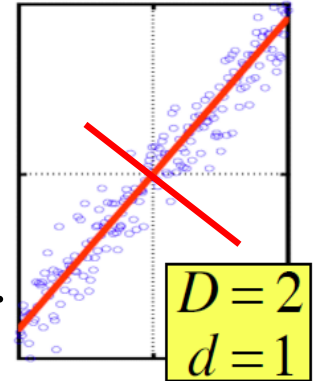
- ***The zero vector can not be an eigenvector***
- ***The value zero can be eigenvalue***

Principal Component Analysis (PCA)

Let v_1, v_2, \dots, v_d denote the d principal components.

$$v_i^T \cdot v_j = 0, i \neq j \text{ and } v_i^T \cdot v_i = 1, i = j$$

Assume data is centered (we extracted the sample mean).



Let $X = [x_1, x_2, \dots, x_n]$ (columns are the datapoints)

Find vector that maximizes sample variance of projected data

$$\sum_{i=1}^n (v^T x_i)^2 = v^T X X^T v$$

$$\max_v v^T X X^T v \quad \text{s.t.} \quad v^T v = 1$$

$$\text{Lagrangian: } \max_v v^T X X^T v - \lambda v^T v$$

Wrap constraints into the objective function

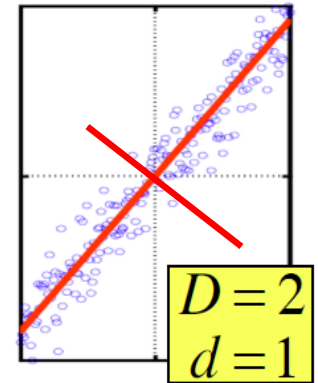
$$\frac{\partial}{\partial v} = 0 \quad (X X^T - \lambda I)v = 0 \quad \Rightarrow \quad (X X^T)v = \lambda v$$

Principal Component Analysis (PCA)

$(X X^T)v = \lambda v$, so v (the first PC) is the eigenvector of sample correlation/covariance matrix $X X^T$

Sample variance of projection $v^T X X^T v = \lambda v^T v = \lambda$

Thus, the eigenvalue λ denotes the amount of variability captured along that dimension (aka amount of energy along that dimension).



Eigenvalues $\lambda_1 \geq \lambda_2 \geq \lambda_3 \geq \dots$

- The 1st PC v_1 is the eigenvector of the sample covariance matrix $X X^T$ associated with the largest eigenvalue
- The 2nd PC v_2 is the eigenvector of the sample covariance matrix $X X^T$ associated with the second largest eigenvalue
- And so on ...

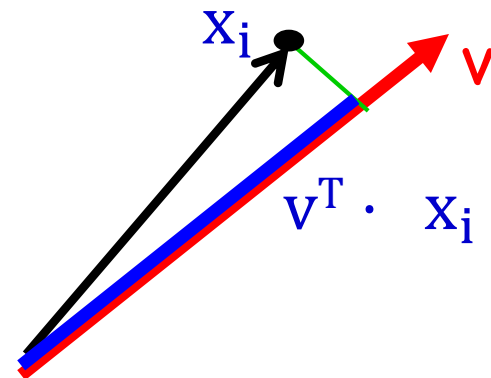
Two Interpretations

So far: **Maximum Variance Subspace**. PCA finds vectors \mathbf{v} such that projections on to the vectors capture maximum variance in the data

$$\sum_{i=1}^n (\mathbf{v}^T \mathbf{x}_i)^2 = \mathbf{v}^T \mathbf{X}\mathbf{X}^T \mathbf{v}$$

Alternative viewpoint: **Minimum Reconstruction Error**. PCA finds vectors \mathbf{v} such that projection on to the vectors yields minimum mean squared error (MSE) reconstruction

$$\sum_{i=1}^n \|\mathbf{x}_i - (\mathbf{v}^T \mathbf{x}_i)\mathbf{v}\|^2$$

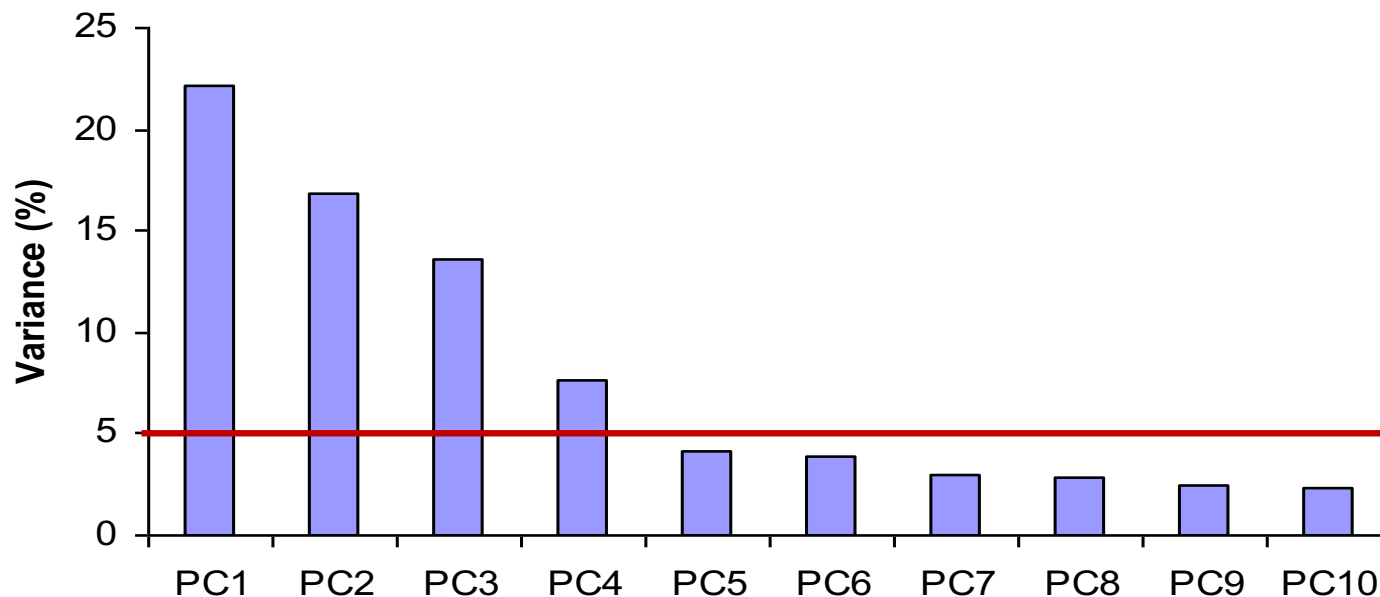


Dimensionality Reduction using PCA

In high-dimensional problems, data sometimes lies near a linear subspace, as noise introduces small variability

Only keep data projections onto principal components with **large** eigenvalues

Can **ignore** the components of smaller significance.



Might **lose some info**, but if eigenvalues are small, do not lose much

APPLICATION EXAMPLES

The space of all face images

- When viewed as vectors of pixel values, face images are extremely high-dimensional
 - 100x100 image = 10,000 dimensions
 - Slow and lots of storage
- But very few 10,000-dimensional vectors are valid face images
- We want to effectively model the subspace of face images

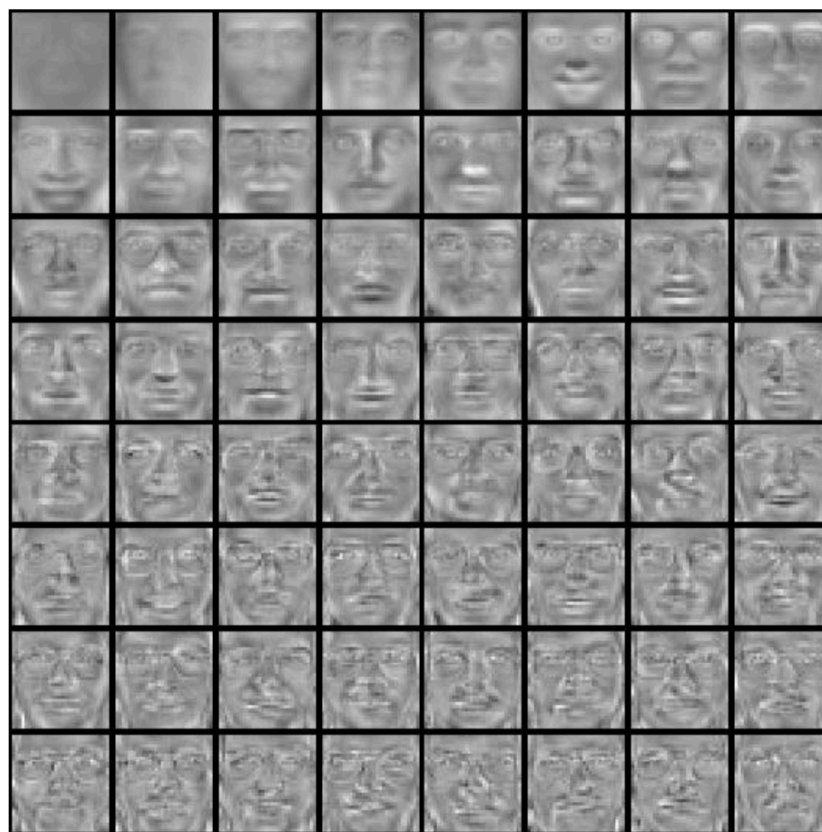
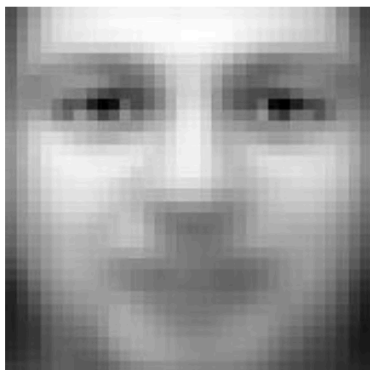


slide by Derek Hoiem

Eigenfaces example

Top eigenvectors: u_1, \dots, u_k

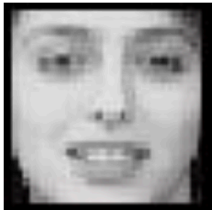
Mean: μ





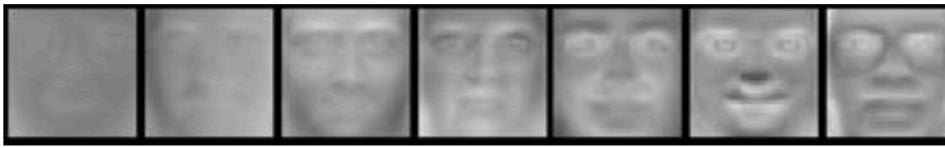
slide by Derek Hoiem

Representation and reconstruction

- Face \mathbf{x} in “face space” coordinates:


$$\mathbf{x} \rightarrow [\mathbf{u}_1^T (\mathbf{x} - \mu), \dots, \mathbf{u}_k^T (\mathbf{x} - \mu)]$$
$$= w_1, \dots, w_k$$

- Reconstruction:


$$=$$

$$+$$

$$\hat{\mathbf{x}} = \mu + w_1 \mathbf{u}_1 + w_2 \mathbf{u}_2 + w_3 \mathbf{u}_3 + w_4 \mathbf{u}_4 + \dots$$

slide by Derek Hoiem

Reconstruction



After computing eigenfaces using 400 face images from ORL face database

slide by Derek Hoiem