

Dimensionality Reduction

PCA

Aarti Singh & Geoff Gordon

Machine Learning 10-701
May 3, 2021

Slides Courtesy: Tom Mitchell, Eric Xing, Lawrence Saul



MACHINE LEARNING DEPARTMENT

Carnegie Mellon.
School of Computer Science

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

High-Dimensional data

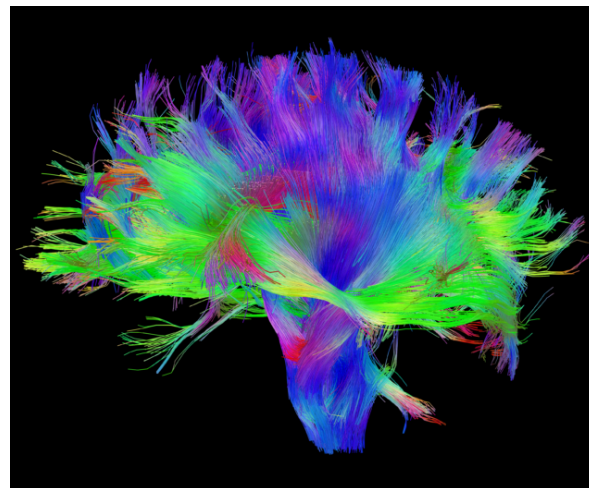
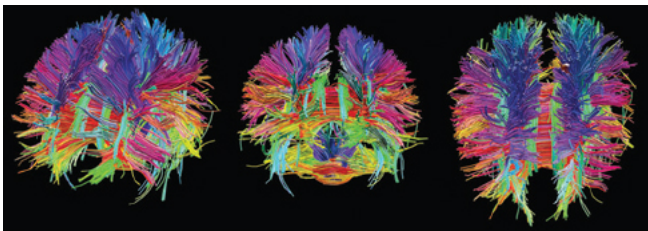
- High-Dimensions = Lot of Features

High resolution images

millions of pixels

Diffusion scans of Brain

300,000 brain fibers

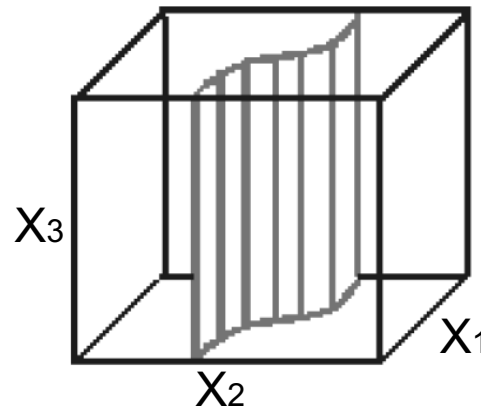


Curse of Dimensionality

- Why are more features bad?
 - Redundant features (not all words are useful to classify a document)
more noise added than signal
 - Hard to interpret and visualize
 - Hard to store and process data (computationally challenging)
 - Complexity of decision rule tends to grow with # features. Hard to learn complex rules as it needs more data (statistically challenging)

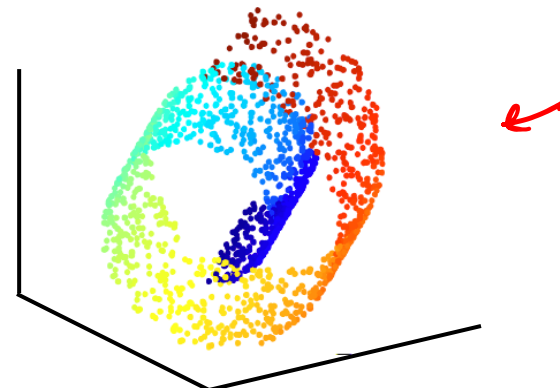
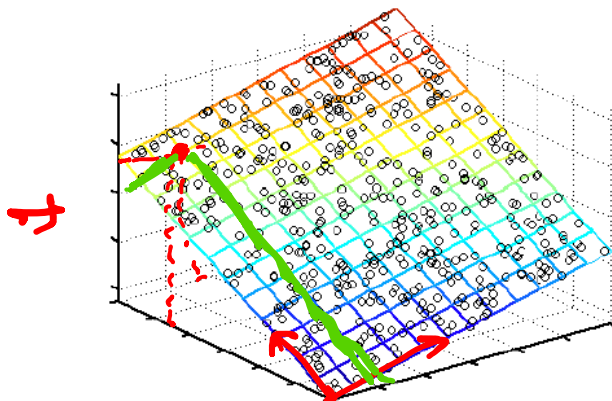
Dimensionality Reduction

- **Feature Selection** – Only a few features are relevant to the learning task



X_3 - Irrelevant

- **Latent features** – Some linear/nonlinear combination of features provides a more efficient representation than observed features



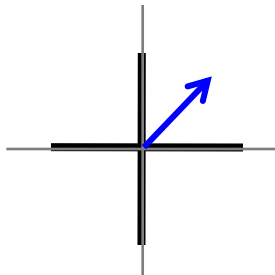
Feature Selection

- One Approach: **Regularization (MAP)**

Integrate feature selection into learning objective by penalizing number of features with non-zero weights

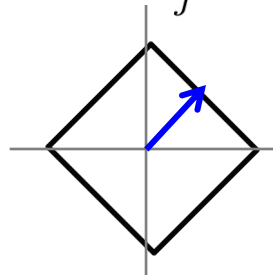
$$\widehat{W} = \arg \min_W \underbrace{\sum_{i=1}^n -\log P(Y_i|X_i; W)}_{\text{-ve log likelihood}} + \underbrace{\lambda \|W\|}_{\text{penalty}}$$

$$\|W\|_0 = \#\{W_j > 0\}$$



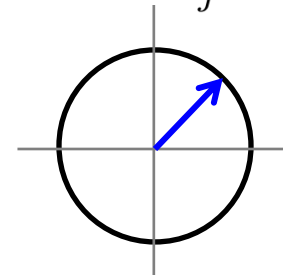
Minimizes # features
chosen

$$\|W\|_1 = \sum_j |W_j|$$



Convex
compromise

$$\|W\|_2 = \sum_j W_j^2$$



Small weights of
features chosen

Latent Features

Combinations of observed features provide more efficient representation, and capture underlying relations that govern the data

E.g. Ego, personality and intelligence are hidden attributes that characterize human behavior instead of survey questions

Topics (sports, science, news, etc.) instead of documents

Often may not have physical meaning

- Linear

- Principal Component Analysis (PCA)**

- Factor Analysis

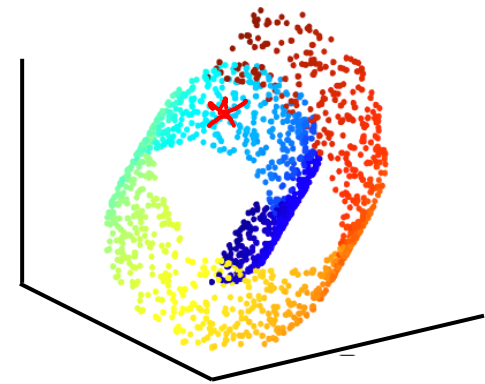
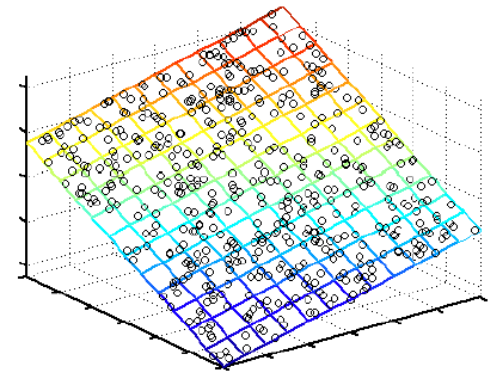
- Independent Component Analysis (ICA)

- Nonlinear

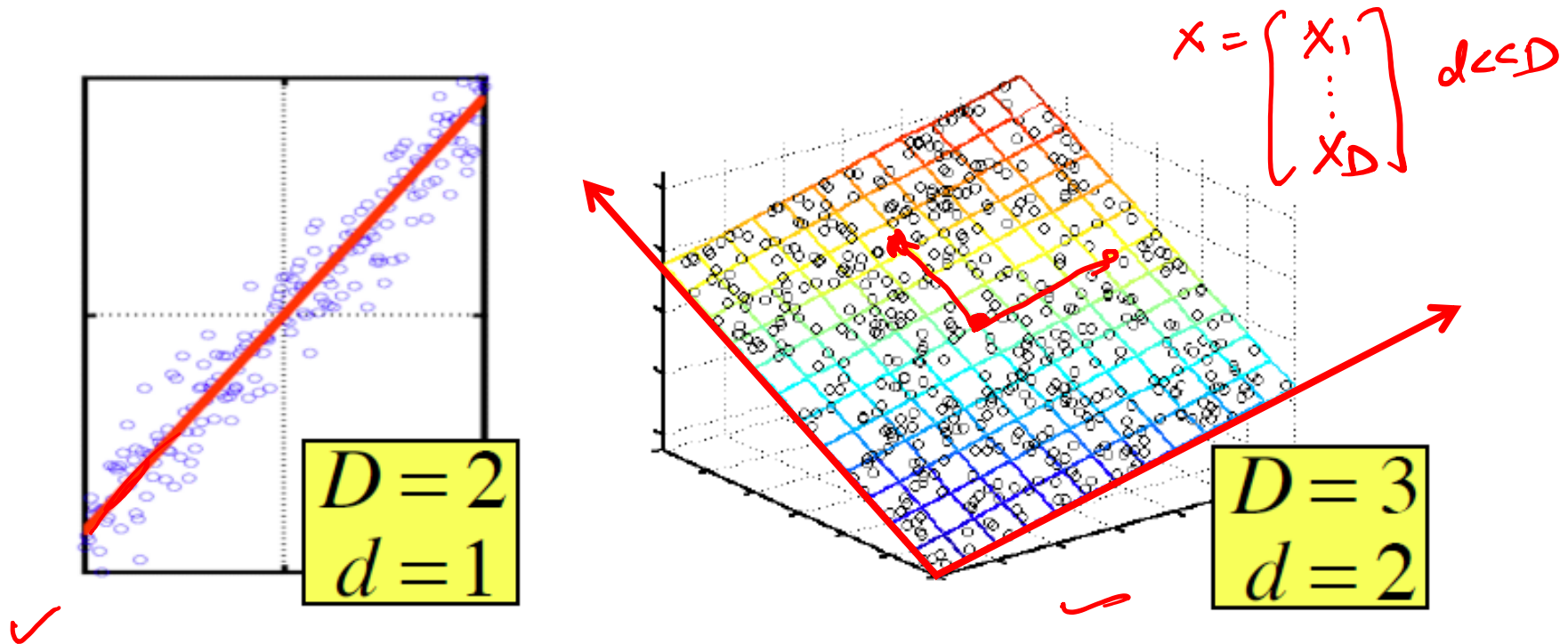
- Kernel PCA**

- Laplacian Eigenmaps, ISOMAP, LLE

- Autoencoders



Principal Component Analysis (PCA)



When 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 Eigen or Singular value decomposition

Data for PCA

Data $X = [x_1, x_2, \dots, x_n]$ where each data point x_i is D-dimensional vector

X is $D \times n$ matrix

Assume data are centered i.e. sample mean $\frac{1}{n} \sum_{i=1}^n x_i = 0$

What if data is not centered?

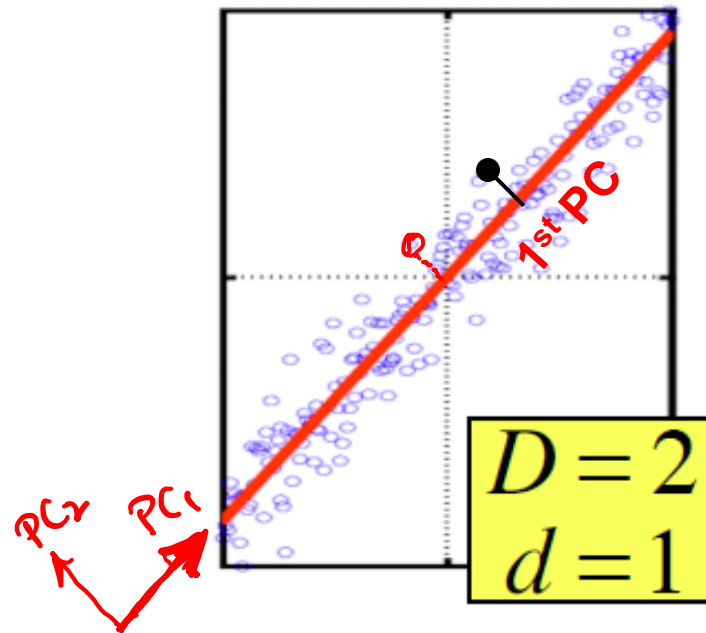
Subtract off sample mean from each data point

Since data matrix is centered, sample covariance matrix can be written as

$$S = \frac{1}{n} X X^T$$

$D \times n$ $n \times D$
 $D \times D$

Principal Component Analysis (PCA)



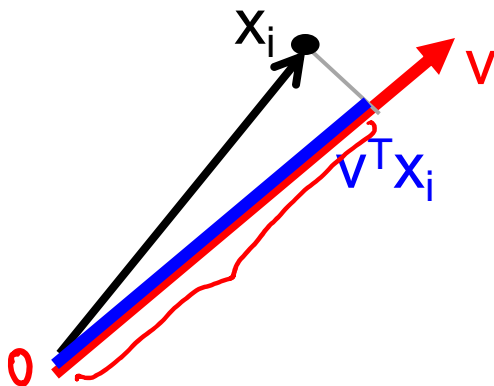
Principal Components (PC) are orthogonal directions that capture most of the variance in the data

1st PC – direction of greatest variability in data

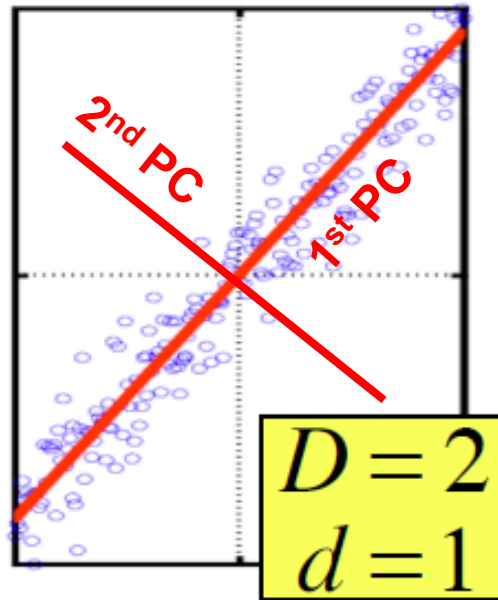
Projection of data points along 1st PC discriminate the data most along any one direction

Take a data point x_i (D-dimensional vector)

Projection of x_i onto the 1st PC v is $v^T x_i$



Principal Component Analysis (PCA)



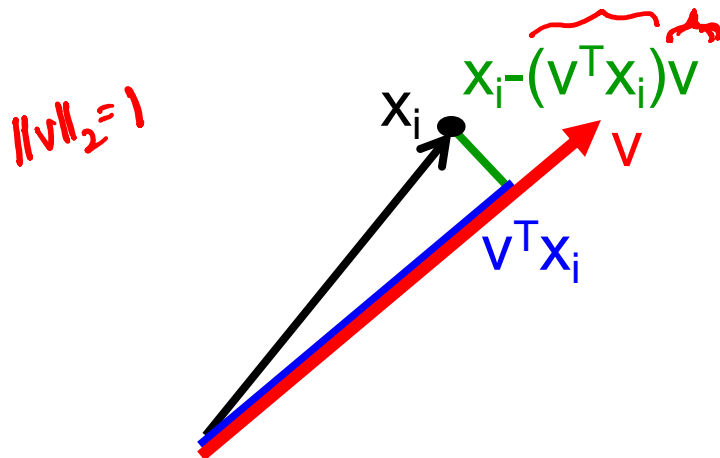
Principal Components (PC) are orthogonal unit norm directions that capture most of the variance in the data

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 ...



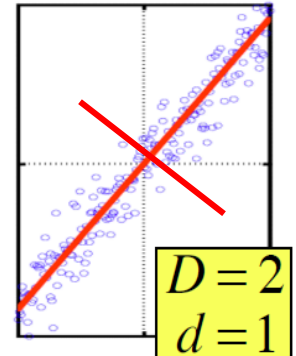
Principal Component Analysis (PCA)

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

Orthogonal and unit norm $v_i^T v_j = 0 \quad i \neq j$

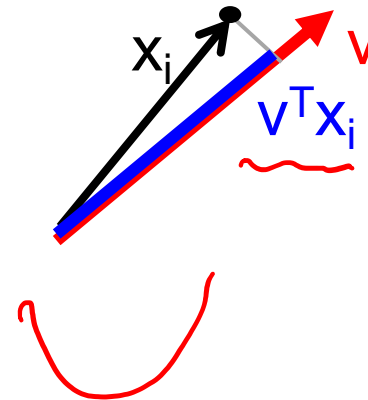
$$v_i^T v_i = 1 \quad \leftarrow \|v_i\|_2 = 1$$

Find vector that maximizes sample variance of projection



$$\frac{1}{n} \sum_{i=1}^n (\underbrace{v^T x_i}_{1 \times D \quad D \times n \quad n \times D \quad D \times 1})^2 = \frac{v^T X X^T v}{n} \quad \leftarrow$$

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



Poll:

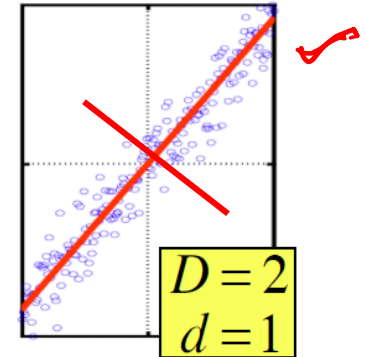
➤ Is this a convex optimization problem?

Principal Component Analysis (PCA)

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

Orthogonal and unit norm $v_i^T v_j = 0 \quad i \neq j$ \leftarrow
 $v_i^T v_i = 1$

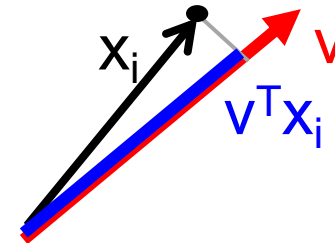
Find vector that maximizes sample variance of projection



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

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

Lagrangian: $\max_v v^T X X^T v - \lambda v^T v$



Wrap constraints into the objective function

$$\partial/\partial v = 0 \quad (X X^T - \lambda I) v = 0$$

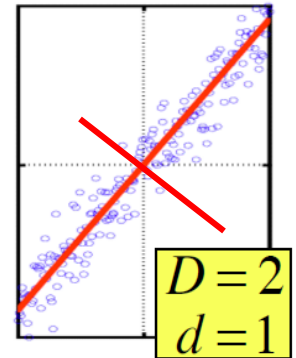
$$\Rightarrow \underline{(X X^T) v = \lambda v}$$

$$v^T \lambda v = \lambda$$

Principal Component Analysis (PCA)

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

Therefore, \mathbf{v} is the eigenvector of sample covariance matrix $\mathbf{X}\mathbf{X}^T$



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

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

Eigenvalues $\lambda_1 > \lambda_2 > \lambda_3 > \dots$

The 1st Principal component \mathbf{v}_1 is the eigenvector of the sample covariance matrix $\mathbf{X}\mathbf{X}^T$ associated with the largest eigenvalue λ_1

The 2nd Principal component \mathbf{v}_2 is the eigenvector of the sample covariance matrix $\mathbf{X}\mathbf{X}^T$ associated with the second largest eigenvalue λ_2

And so on ...

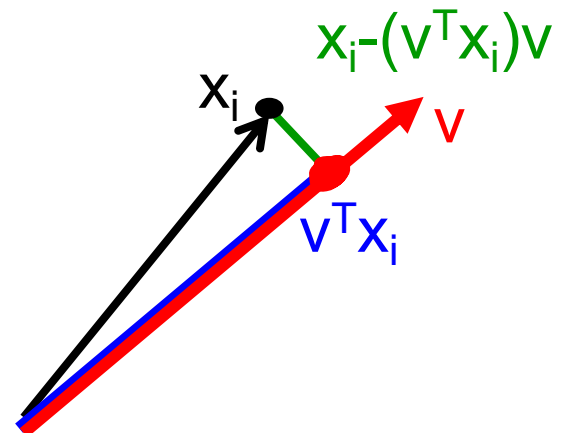
Another interpretation

Maximum Variance Subspace: PCA finds vectors v such that projections on to the vectors capture maximum variance in the data

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

Minimum Reconstruction Error: PCA finds vectors v such that projection on to the vectors yields minimum MSE reconstruction

$$\frac{1}{n} \sum_{i=1}^n \|x_i - \underbrace{(v^T x_i)v}_{\text{red arrow}}\|^2$$



$$\begin{aligned} \min_v \\ \text{st. } v^T v = 1 \end{aligned}$$

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

$$\cancel{x_i^T x_i} + \cancel{(v^T x_i)^2} \cancel{v^T v} - 2(v^T x_i)^2$$

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

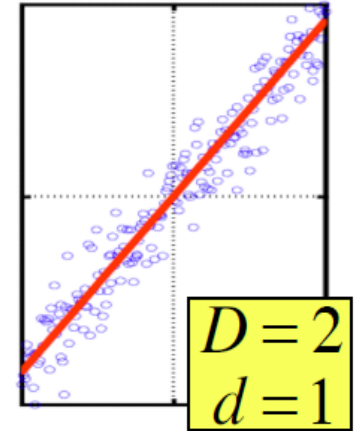
$$\equiv \max \sum_{i=1}^n (v^T x_i)^2$$

Dimensionality Reduction using PCA

The eigenvalue λ denotes the amount of variability captured along that dimension.

Zero eigenvalues indicate no variability along those directions => data lies exactly on a linear subspace

Only keep data projections onto principal components with non-zero eigenvalues, say v_1, \dots, v_d where $d = \text{rank}(XX^T)$



Original Representation
data point

$$x_i = [x_i^1, x_i^2, \dots, x_i^D]^T$$

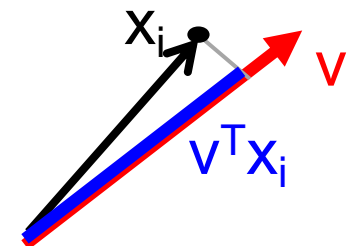
(D-dimensional vector)

Transformed representation
projections

$$[v_1^T x_i, v_2^T x_i, \dots, v_d^T x_i]$$

(d-dimensional vector)

$$d \ll D$$

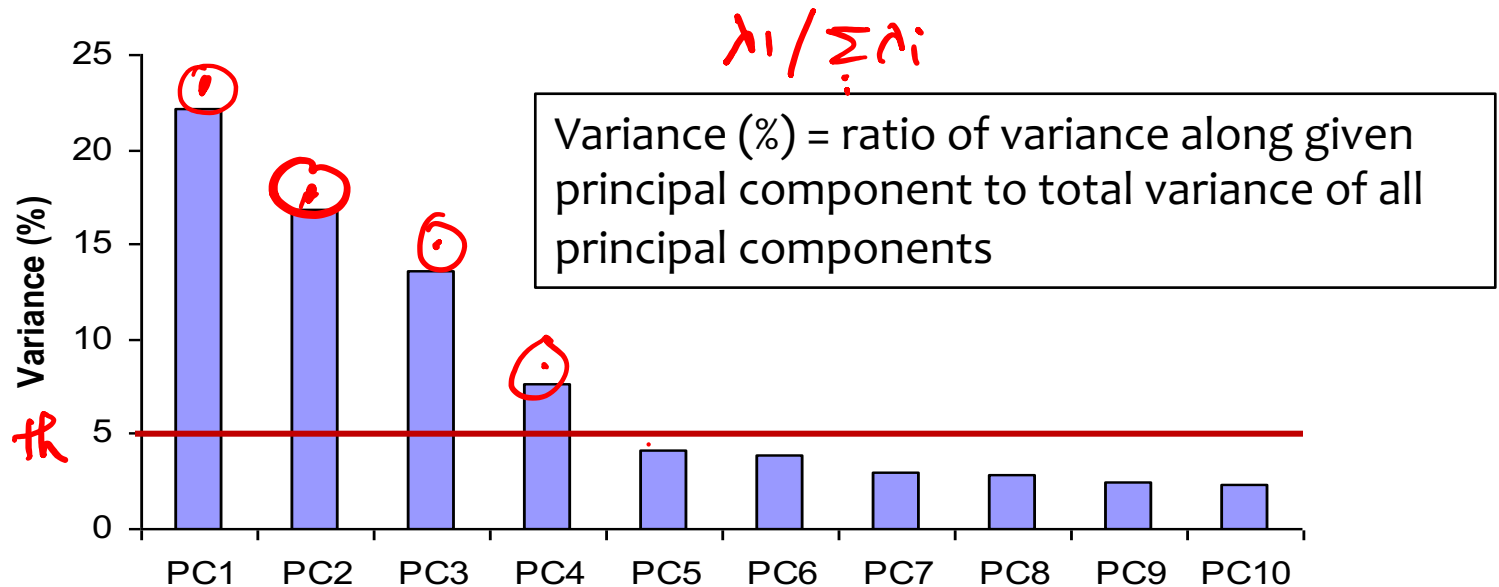


Dimensionality Reduction using PCA

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

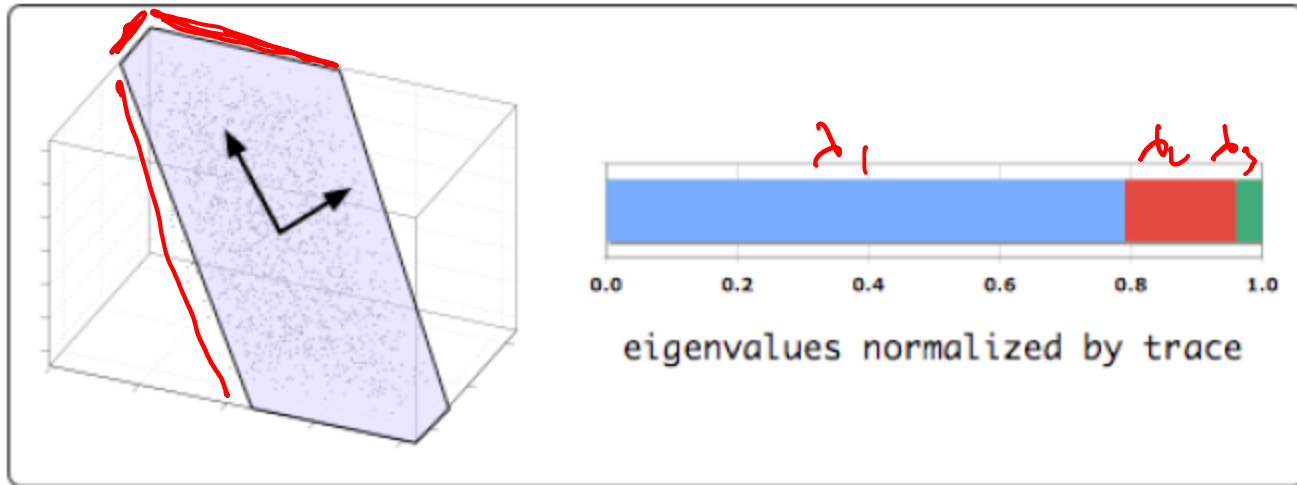
Only keep data projections onto principal components with **large** eigenvalues

Can *ignore* the components of lesser significance.



You might **lose some information**, but if the eigenvalues are small, you don't lose much

Example of PCA



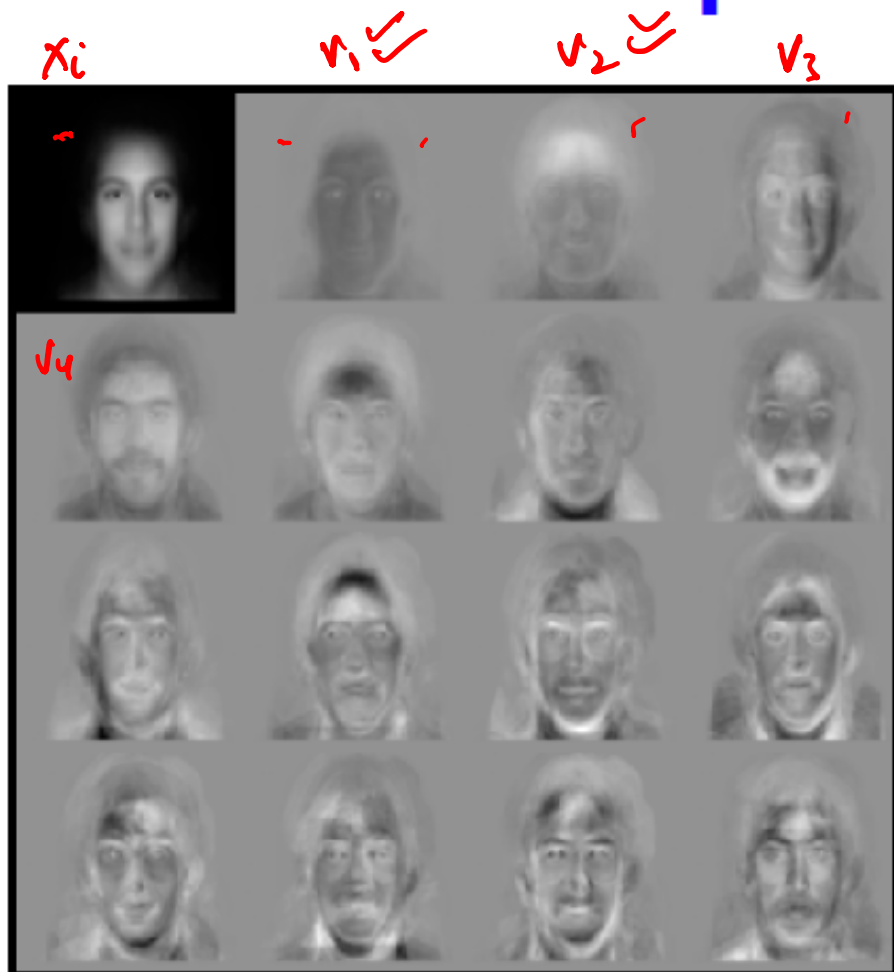
Eigenvectors and eigenvalues of covariance matrix for $n=1600$ inputs in $d=3$ dimensions.

$v \in \mathbb{R}^D$
 $D \times 1$

D pixels = features
 $(v_1^T x, v_2^T x)$

$v_i^T x_i$

Example: faces



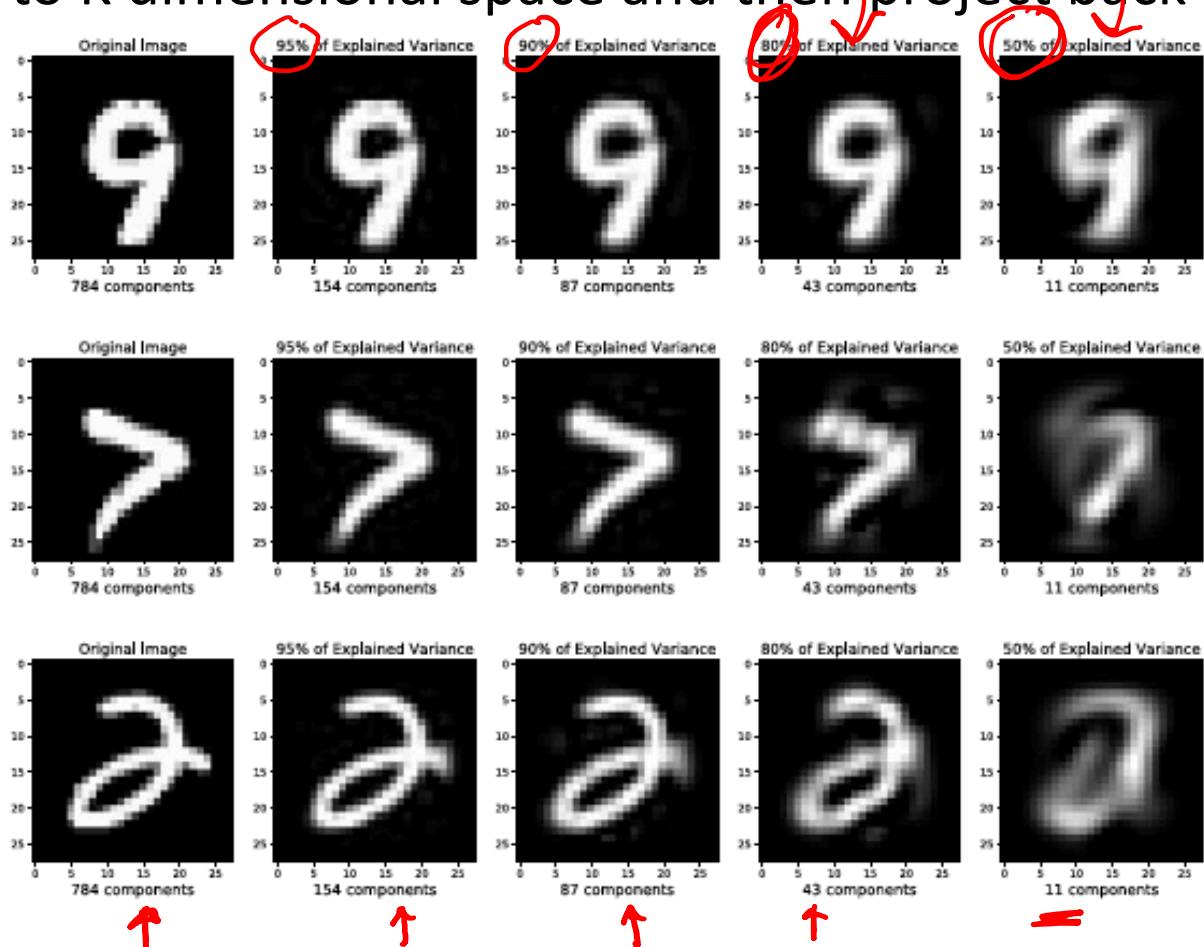
Eigenfaces
from 7562
images:

top left image
is linear
combination
of rest.

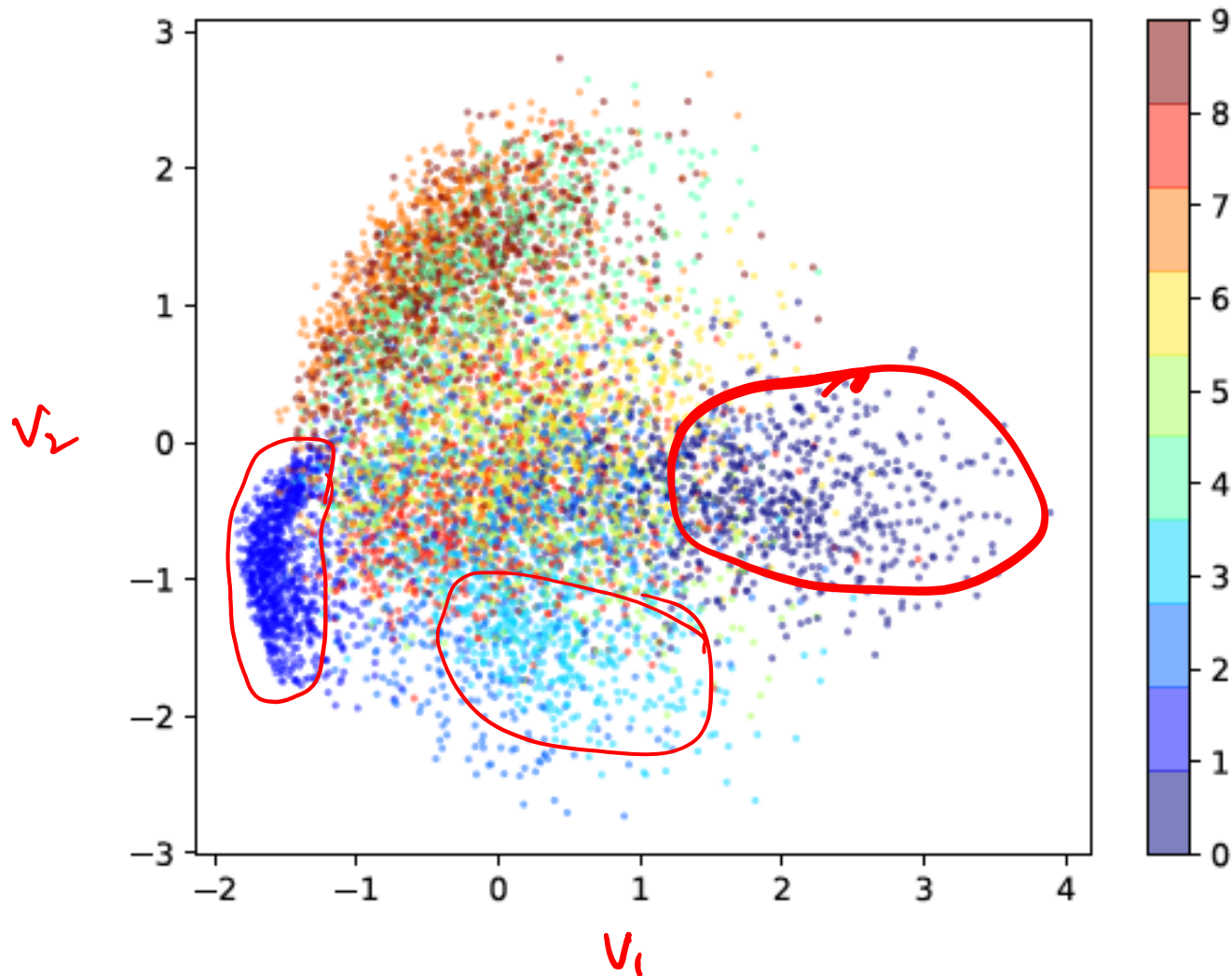
Sirovich & Kirby (1987)
Turk & Pentland (1991)

Example: MNIST digits

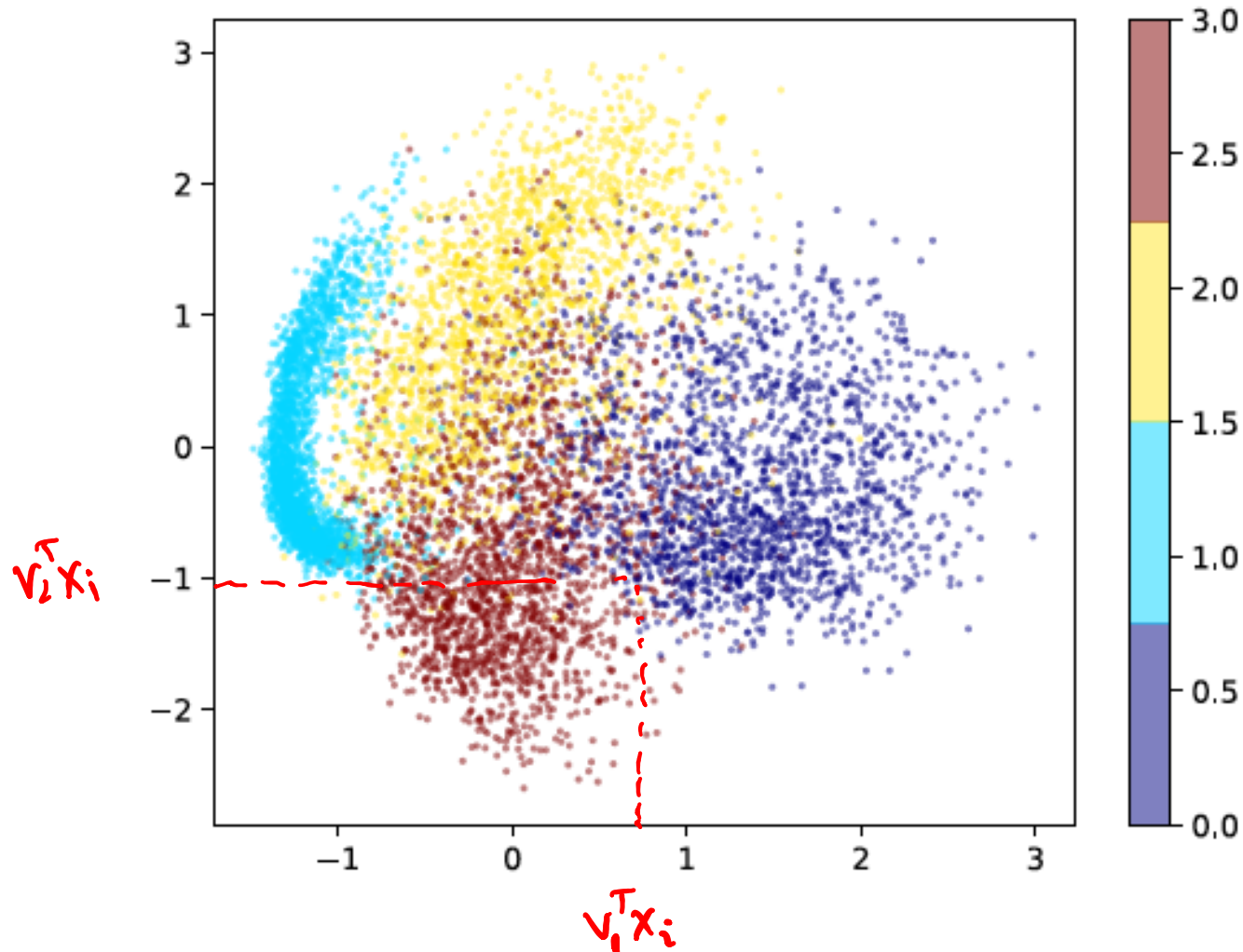
- 28x28 images = 784 PCA vectors
- Project to K dimensional space and then project back up



Projecting MNIST digits



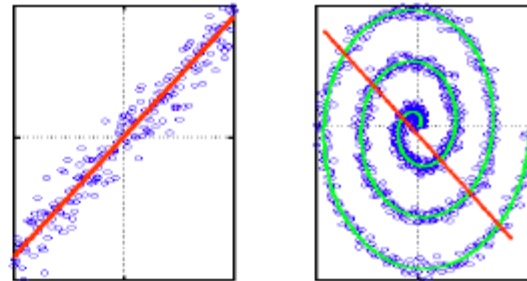
Projecting MNIST digits



Properties of PCA

- **Strengths**

- Eigenvector method
- No tuning parameters
- Non-iterative
- No local optima



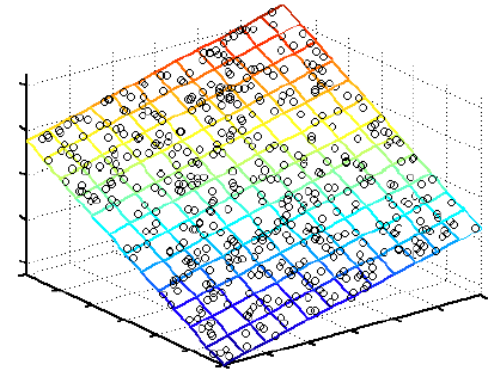
- **Weaknesses**

- Limited to second order statistics
- Limited to linear projections

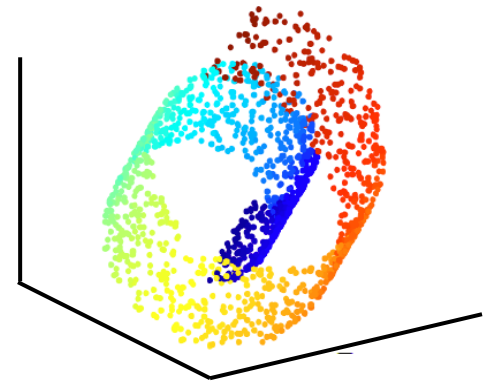
$X X^T$
 $x(j)x(k)$

Unsupervised Dimensionality Reduction

- Linear $X X^T \quad D \times D$
Principal Component Analysis (PCA)
Factor Analysis
Independent Component Analysis (ICA)



- Nonlinear $X^T X \quad n \times n$
Kernel PCA
Laplacian Eigenmaps, ISOMAP, LLE
Autoencoders



$$\max_V \quad \underline{V^T X X^T V} \quad \text{s.t.} \quad V^T V = I$$

$$\Rightarrow \quad \underline{X X^T V = \lambda V} \quad \leftarrow$$

$$\Rightarrow V = \frac{X X^T V}{\lambda}$$

$$V^T X X^T V = V^T X X^T \left(\frac{X X^T V}{\lambda} \right)$$

$$= \left(\frac{V^T X}{\sqrt{\lambda}} \right) \left(\underline{X^T X} \right) \left(\frac{X^T V}{\sqrt{\lambda}} \right)$$

$$= \underline{W^T X^T X W}$$

$$X X^T$$

$$D \times n \quad n \times D = D \times D$$

$$\downarrow$$

$$X^T X$$

$$n \times D \quad D \times n$$

$$X = U \Lambda V^T \leftarrow$$

$$W = \frac{X^T V}{\sqrt{\lambda}}$$

$$X X^T$$

$$X^T X$$

$$W^T W = \frac{V^T X X^T V}{\lambda} = 1$$

$$\max_W \quad \underline{W^T X^T X W} \quad \text{s.t.} \quad W^T W = I$$

$$\leftarrow K(x_i, x_j)$$