

Machine Learning, Machine Learning (extended)

9 – Unsupervised Learning: Dimensionality Reduction

Kashif Rajpoot

k.m.rajpoot@cs.bham.ac.uk

School of Computer Science

University of Birmingham

Outline

- High dimensionality
- Curse of dimensionality
- Dimensionality reduction
- Projection
- Preserving 'interesting' characteristics in data
- Principal component analysis (PCA)

High dimensionality

- Dimensionality
 - Number of attributes (i.e. features) in the data
- Technological advances lead to increasingly high dimensional data sets
 - Tens to hundreds to thousands...
 - Mass spectrometry, brain functional imaging, genomics, hyperspectral imaging, financial analysis
- High dimensional data => expected to give “more” information (features) about an object?
 - Not always, it could actually cause “curse of dimensionality”

Curse of dimensionality

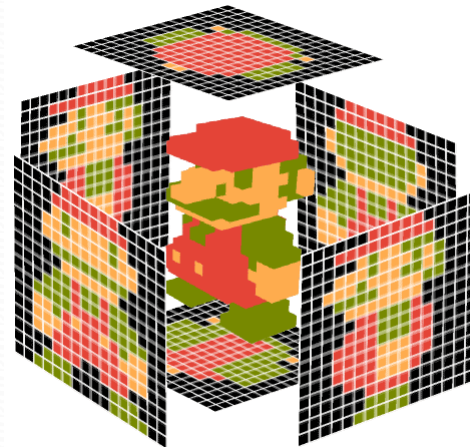
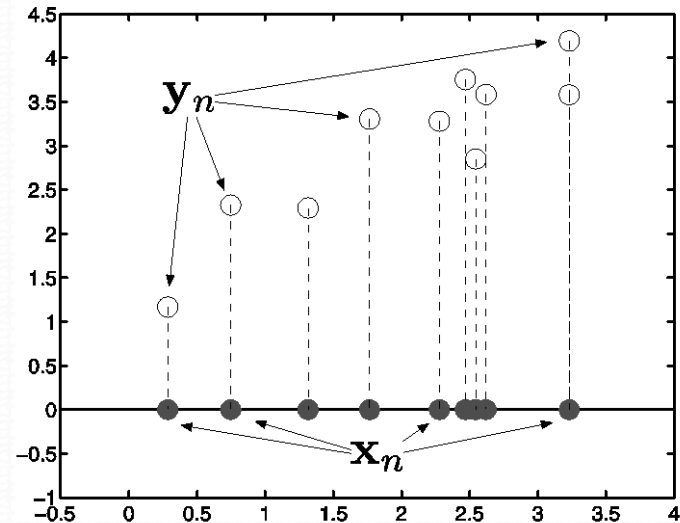
- Computational burden
 - For example: consider clustering 3000 samples of 2000-dimensional data with k-means algorithm in 10 classes?
- Visualization
 - Typically, we can visualize data only up to 3d/4d
 - What about higher-dimensional data?
- Parameter estimation
 - Higher dimensions need estimation of higher number of parameters (e.g. regression, classification)

Dimensionality reduction

- Dimensionality reduction aims to avoid the 'curse of dimensionality' by reducing the attributes/dimensions/features
1. Feature selection
 - *Sequential feature selection* (one of the simplest methods for feature selection)
 - Gradually add (remove) a feature to include (exclude)
 - Determine feature scoring/importance by cross-validation
 2. **Subspace projection**
 - Make new features by combining (i.e. linearly or non-linearly) the old features
 - Represent the data in fewer number of dimensions but 'preserving' the 'interesting' characteristics of data

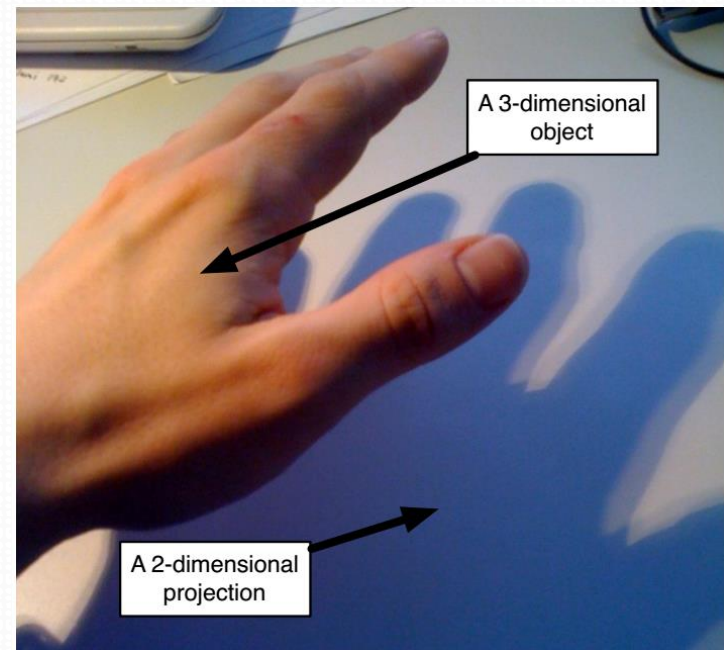
Projection

- Represent the data in fewer number of features but 'preserving' the 'interesting' characteristics of data?
- 3D world to 2D image projection



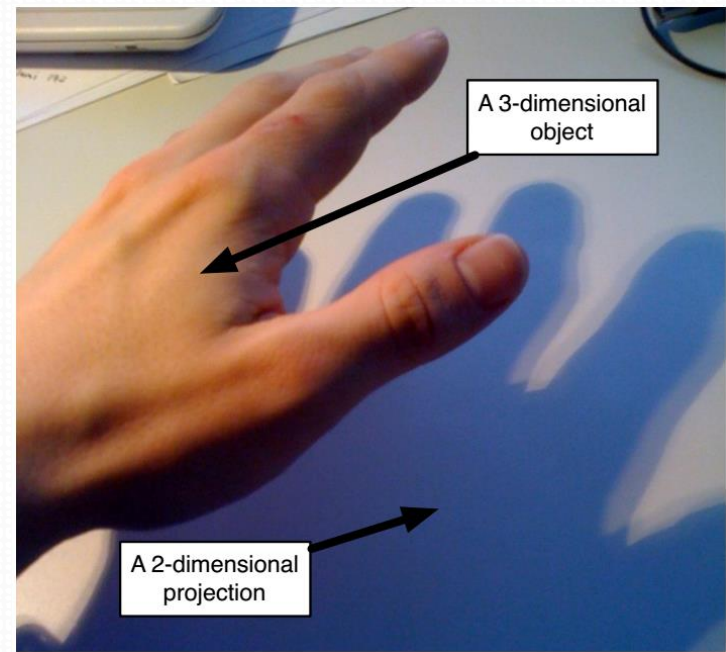
Projection

- Project M -dimensional data \mathbf{Y} containing N samples into a lower D -dimensional representation \mathbf{X} ($D \ll M$)
- $\mathbf{X} = \mathbf{Y}\mathbf{W}$
 - \mathbf{Y} is $N \times M$
 - \mathbf{W} is $M \times D$
 - \mathbf{X} is $N \times D$ (i.e. D -dimensional)
- What is \mathbf{W} ?
 - Defines the projection
 - Changing \mathbf{W} is like changing where the light is coming from or rotating the hand
- \mathbf{Y} is hand, \mathbf{X} is shadow



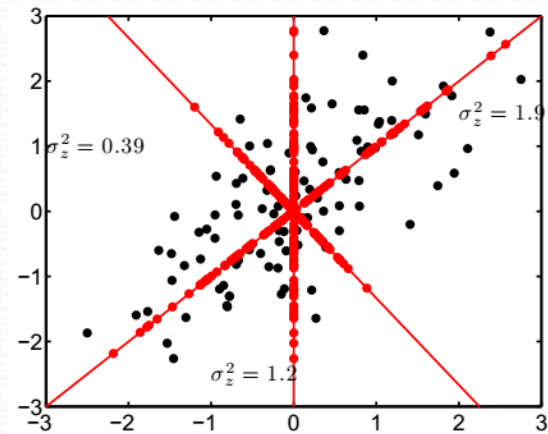
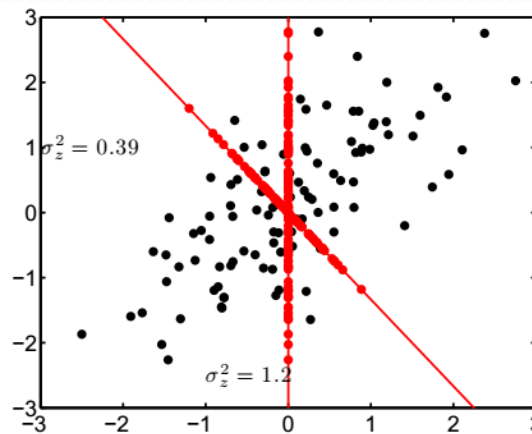
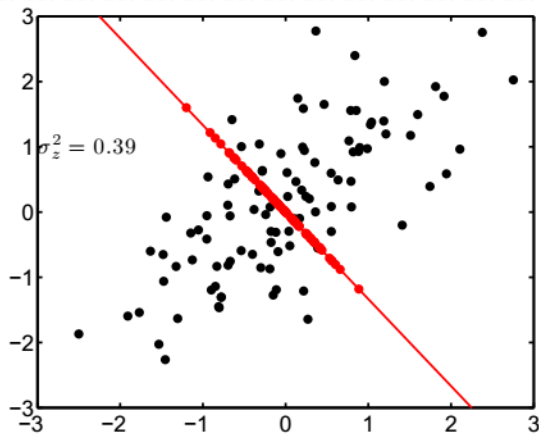
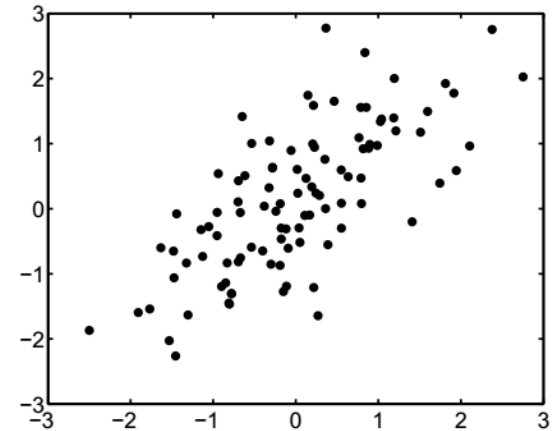
Projection

- Different \mathbf{W} will result in different projections
- How to choose \mathbf{W} ?
 - Not all projections will represent our data 'well'
- We should choose a \mathbf{W} that preserves the 'interesting' data characteristics
 - such that $D \ll M$
 - M is the actual data's dimensionality
 - D is the projected data's dimensionality



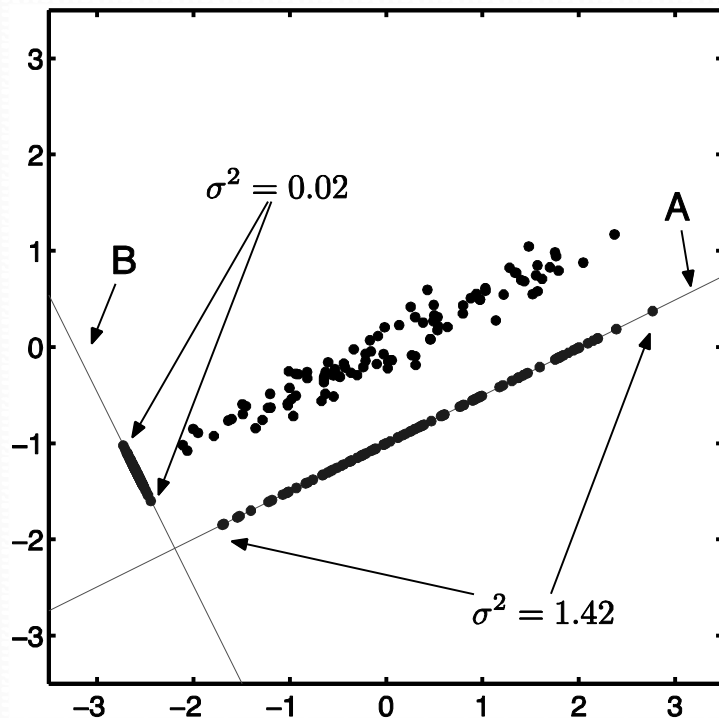
Preserve 'interesting' data characteristics

- Project 2-d data into 1-d?
- Pick some arbitrary \mathbf{w}
- Project the data onto it
- The position on the line is our 1-d representation
- Compute the variance on the line

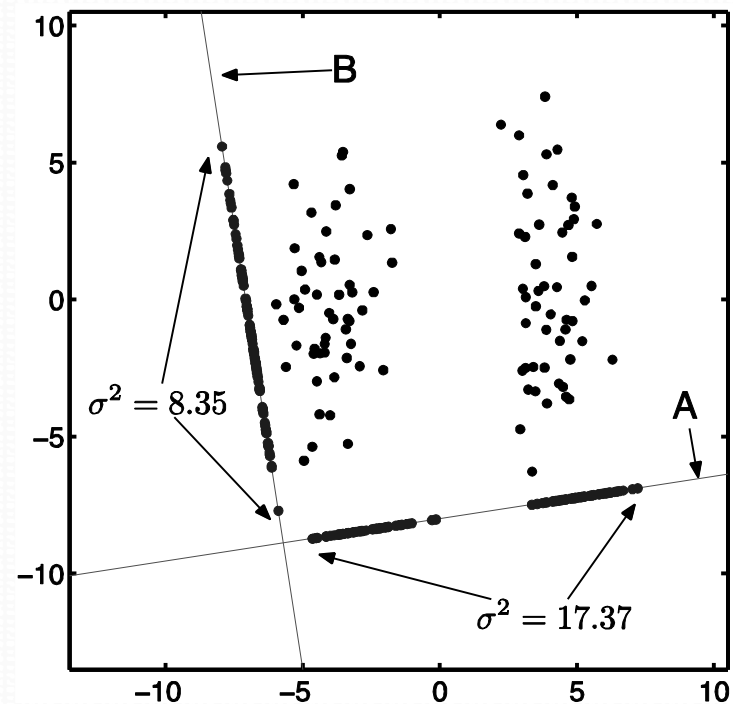


Preserve 'interesting' data characteristics

- Represent the data in fewer number of features but 'preserving' the 'interesting' characteristics of data?
 - High variance = highly "interesting" characteristics



(a) Data from a single, elongated Gaussian



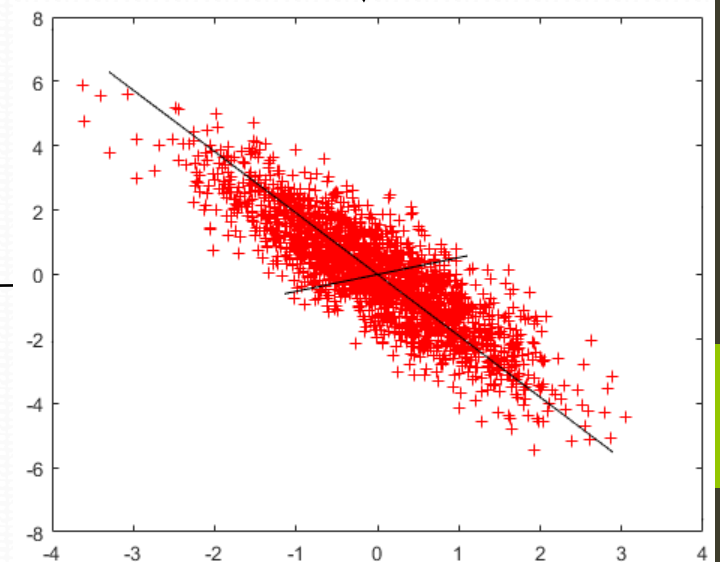
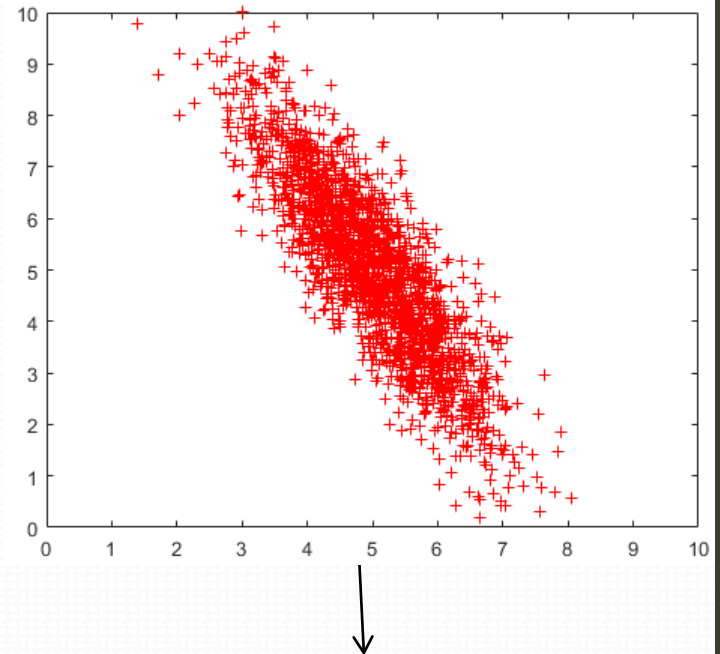
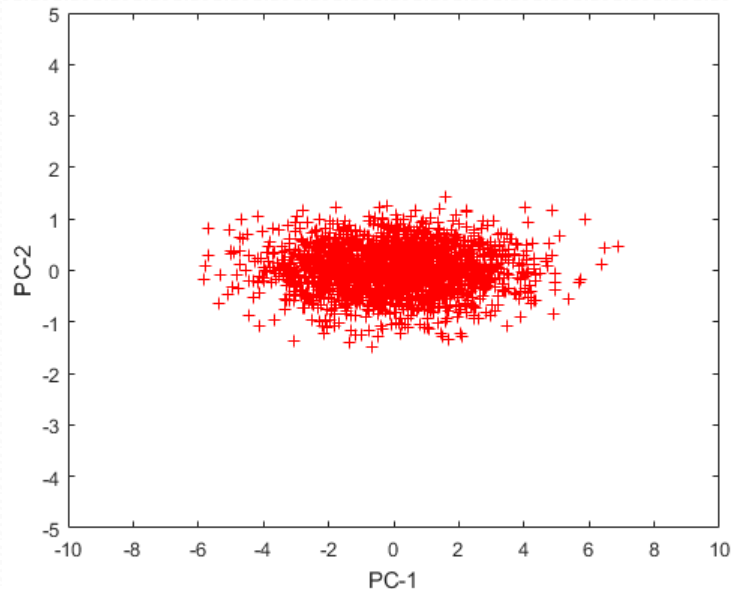
(b) Data from two Gaussians

Principal component analysis (PCA)

- PCA chooses \mathbf{W} such that it transforms M -dimensional data \mathbf{Y} into D -dimensional representation \mathbf{X} with $D \ll M$
 - It preserves dimensions with high variance ('interesting' characteristics)
 - It removes dimensions which are redundant (i.e. not much 'interesting' – having low variance)
- Find the columns of \mathbf{W} one at a time
 - \mathbf{w}_d as the d^{th} column of \mathbf{W}
 - Each $M \times 1$ column defines a new dimension

PCA

- PCA determines the dominant modes of variation from within the data and then projects data onto this 'natural' coordinate system
- Matches the coordinate system to the shape of the data



PCA

- Consider \mathbf{w}_d as a new dimension, then the data projection in this dimension is computed as:

$$\mathbf{x}_d = \mathbf{Y}\mathbf{w}_d$$

- PCA chooses \mathbf{w}_d that maximizes the variance of \mathbf{x}_d

- $\sigma_d^2 = \frac{1}{N} \sum_{n=1}^N \left(x_d^{(n)} - \mu_d \right)^2$ where $\mu_d = \frac{1}{N} \sum_{n=1}^N x_d^{(n)}$

and $x_d^{(n)}$ denotes the n^{th} sample value for d^{th} attribute

- Each new column of \mathbf{W} is found such that it maximizes variance and is orthogonal (perpendicular) to the previous columns

PCA: process

- Search for $\mathbf{W} = [\mathbf{w}_1 \ \mathbf{w}_2 \ \dots \ \mathbf{w}_D]?$
- Fortunately, analytical solution exists
- \mathbf{W} are the *eigenvectors* of the *covariance matrix* $\mathbf{\Sigma}$ of \mathbf{Y}
 - The covariance matrix $\mathbf{\Sigma}$ describes the way different attributes (i.e. features or variables) co-vary

PCA: process

- Covariance matrix Σ describes the way attributes co-vary
 - $\Sigma_{i,j}$ denotes the covariance of i^{th} and j^{th} attributes of \mathbf{Y}

$$\Sigma_{i,j} = \frac{1}{N} \sum_{n=1}^N (y_i^{(n)} - \mu_i)(y_j^{(n)} - \mu_j)$$

where μ_i and μ_j denote the mean of i^{th} and j^{th} attributes, respectively, while $y_i^{(n)}$ and $y_j^{(n)}$ denote the n^{th} sample value for i^{th} and j^{th} attributes, respectively.

- The data \mathbf{Y} is mean subtracted (along each dimension) to translate the data to the centre of coordinate system
 - Each row is an object, each column is a dimension
 - $\tilde{\mathbf{Y}} = \mathbf{Y} - \bar{\mathbf{Y}}$, where each column of $\bar{\mathbf{Y}}$ is mean value μ across that particular attribute
- The covariance matrix can then be computed as:

$$\Sigma = \frac{1}{N} (\mathbf{Y} - \bar{\mathbf{Y}})^T (\mathbf{Y} - \bar{\mathbf{Y}}) = \frac{1}{N} \tilde{\mathbf{Y}}^T \tilde{\mathbf{Y}}$$

PCA: process

- PCA finds new coordinate vectors

$$\mathbf{W} = [\mathbf{w}_1 \ \mathbf{w}_2 \ \dots \ \mathbf{w}_D]$$

that align with data shape

- $\mathbf{w}_i^T \mathbf{w}_j = 0, \forall i \neq j$ new dimensions are orthogonal
- PCA looks to find the direction \mathbf{w} that maximizes variance in that direction

$$\mathcal{F} = \underset{\mathbf{w}}{\operatorname{argmax}} \operatorname{var}(\tilde{\mathbf{Y}}\mathbf{w})$$

- $\operatorname{var}(\tilde{\mathbf{Y}}\mathbf{w}) = (\tilde{\mathbf{Y}}\mathbf{w})^T \tilde{\mathbf{Y}}\mathbf{w} = \mathbf{w}^T \tilde{\mathbf{Y}}^T \tilde{\mathbf{Y}}\mathbf{w} = \mathbf{w}^T \boldsymbol{\Sigma} \mathbf{w}$

- So

$$\mathcal{F} = \underset{\mathbf{w}}{\operatorname{argmax}} \mathbf{w}^T \boldsymbol{\Sigma} \mathbf{w}$$

subject to $\mathbf{w}^T \mathbf{w} = \mathbf{1}$

(we could keep increasing \mathbf{w} to maximise \mathcal{F} , hence the need to constrain $\mathbf{w}^T \mathbf{w} = \mathbf{1}$)

PCA: process

$$\mathcal{F} = \underset{\mathbf{w}}{\operatorname{argmax}} \mathbf{w}^T \Sigma \mathbf{w}$$

subject to $\mathbf{w}^T \mathbf{w} = 1$

- Using the Lagrange multiplier “trick” (beyond our module scope):

$$\mathcal{F} = \underset{\mathbf{w}}{\operatorname{argmax}} (\mathbf{w}^T \Sigma \mathbf{w} - \lambda(\mathbf{w}^T \mathbf{w} - 1))$$

- To maximize, differentiate with respect to \mathbf{w} and set to 0:

$$\frac{\partial \mathcal{F}}{\partial \mathbf{w}} = 2\Sigma \mathbf{w} - 2\lambda \mathbf{w} = 0$$

$$\Sigma \mathbf{w} = \lambda \mathbf{w}$$

- This can be solved with eigenvalue/eigenvector method (beyond our module scope)

Aside: Eigenvector

- We obtained solution for \mathbf{w} :

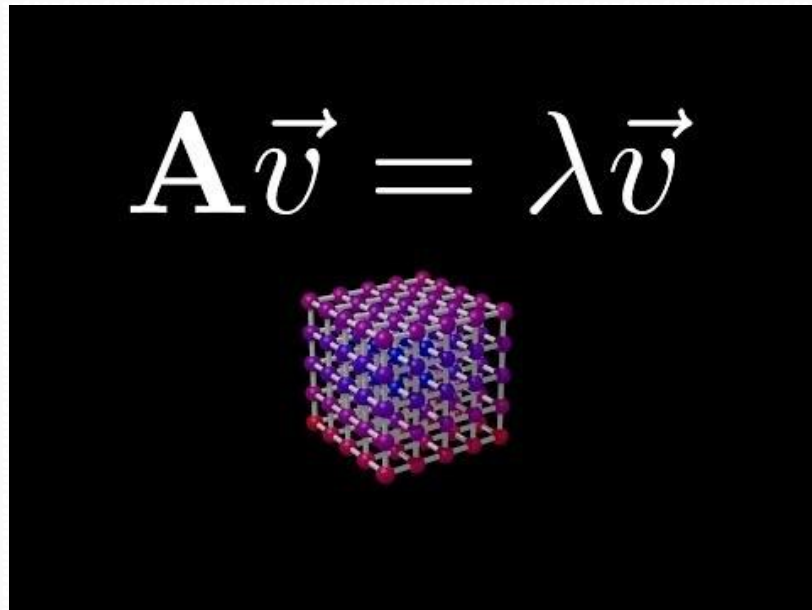
$$\mathbf{\Sigma}\mathbf{w} = \lambda\mathbf{w}$$

- Note that the multiplication of matrix $\mathbf{\Sigma}$ with vector \mathbf{w} changes it only by a scalar factor
- Such a vector \mathbf{w} is called eigenvector
 - λ is the eigenvalue, the scale by which eigenvector \mathbf{w} changes
- The eigenvector is a 'special' vector which has nice properties to become a transformation vector

Aside: Eigenvector

What is eigenvector?

<https://www.youtube.com/watch?v=ue3yoeZvt8E>



Introduction to eigenvalues and eigenvectors

<https://www.youtube.com/watch?v=PhfbEr2btGQ>



PCA: process

- The eigenvector and eigenvalue solution of covariance matrix $\mathbf{\Sigma}$ will provide M eigenvectors $[\mathbf{w}_1 \ \mathbf{w}_2 \ \dots \ \mathbf{w}_M]$ and M eigenvalues $[\lambda_1 \ \lambda_2 \ \dots \ \lambda_M]$:

$$\mathbf{\Sigma} \mathbf{w}_i = \lambda_i \mathbf{w}_i$$

- Which \mathbf{w} has highest variance?
- Multiplying both sides of above equation by \mathbf{w}_i^T :

$$\mathbf{w}_i^T \mathbf{\Sigma} \mathbf{w}_i = \lambda_i \mathbf{w}_i^T \mathbf{w}_i = \lambda_i$$

i.e. λ_i denotes variance along \mathbf{w}_i dimension since

$$\mathbf{w}_i^T \mathbf{\Sigma} \mathbf{w}_i = \text{var}(\tilde{\mathbf{Y}} \mathbf{w}_i)$$

- The principal components of data \mathbf{Y} are the eigenvectors $[\mathbf{w}_1 \ \mathbf{w}_2 \ \dots \ \mathbf{w}_M]$ of covariance matrix $\mathbf{\Sigma}$, ordered by eigenvalues $[\lambda_1 \ \lambda_2 \ \dots \ \lambda_M]$

PCA: process

- Having found the principal components (PCs)

$$\mathbf{W} = [\mathbf{w}_1 \ \mathbf{w}_2 \ \dots \ \mathbf{w}_M]$$

data \mathbf{Y} can now be transformed to these new dimensions:

$$\mathbf{X} = \tilde{\mathbf{Y}}\mathbf{W}$$

where $\tilde{\mathbf{Y}}$ is the mean-subtracted data

- \mathbf{W} is the projection (aka *loadings*)
- \mathbf{X} is the projected data (aka *scores*)

PCA: algorithmic workflow

1. Form the $N \times M$ zero-mean matrix, by subtracting mean
 - $\tilde{\mathbf{Y}} = \mathbf{Y} - \bar{\mathbf{Y}}$, where each column of $\bar{\mathbf{Y}}$ is mean value across that particular attribute

2. Calculate the $M \times M$ covariance matrix Σ

$$\Sigma = \frac{1}{N} \tilde{\mathbf{Y}}^T \tilde{\mathbf{Y}}$$

3. Calculate the M eigenvalues ($\lambda_1, \lambda_2, \dots, \lambda_M$) and eigenvectors ($\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_M$) of Σ
4. Sort the eigenvalues and eigenvectors from largest to smallest by eigenvalue
5. Choose D eigenvectors corresponding to highest eigenvalues
6. Compute the scores \mathbf{X} (i.e. projections) by projecting data $\tilde{\mathbf{Y}}$ on to new coordinates \mathbf{W} (i.e. PCs or eigenvectors)

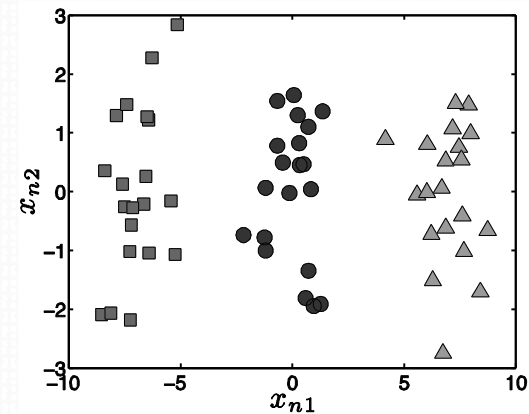
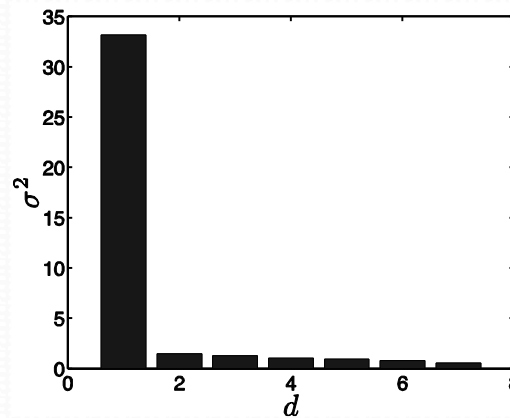
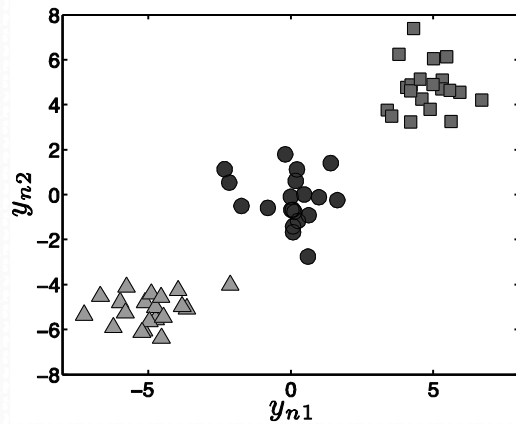
$$\mathbf{X} = \tilde{\mathbf{Y}}\mathbf{W}$$

How to choose D ?

- We get M eigenvectors from the $M \times M$ covariance matrix
 - i.e. M new dimensions
- How to choose D dimensions (such that $D \ll M$)?
 - Application domain knowledge
 - Visualization
 - Computational burden
 - 'interesting' structure (i.e. defined by variance)
 - Post-processing results
- Total variation ($\sum_{d=1}^M \lambda_d$)
 - Percentage variation preserved by D eigenvectors can be estimated as: $\left[\sum_{d=1}^D \lambda_d / \sum_{d=1}^M \lambda_d \right] * 100$

PCA: example

- Consider this 2D data for 3-classes
 - Five additional dimensions were added with random values $\mathcal{N}(0,1)$
 - Perform PCA and dimensionality reduction



(a) First two dimensions of the data objects \mathbf{y}_n

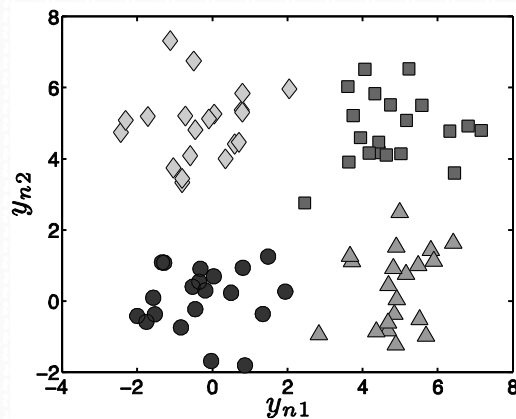
(b) The seven eigenvalues (variances of the projected dimensions)

(c) The data projected onto the first two principal components

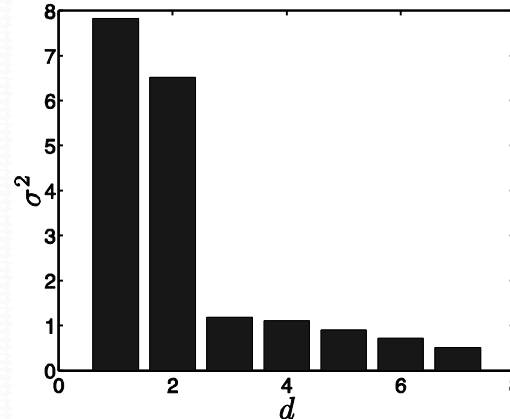
- It determines the dominant modes of variation from within the data and then projects data onto this 'natural' coordinate system

PCA: example

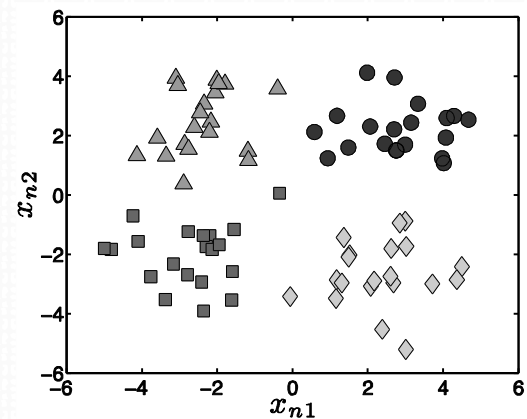
- Consider this 2D data for 4-classes
 - Five additional dimensions were added with random values $\mathcal{N}(0,1)$
 - Perform PCA and dimensionality reduction



(a) First two dimensions of the data objects y_n



(b) The seven eigenvalues (variances of the projected dimensions)



(c) The data projected onto the first two principal components

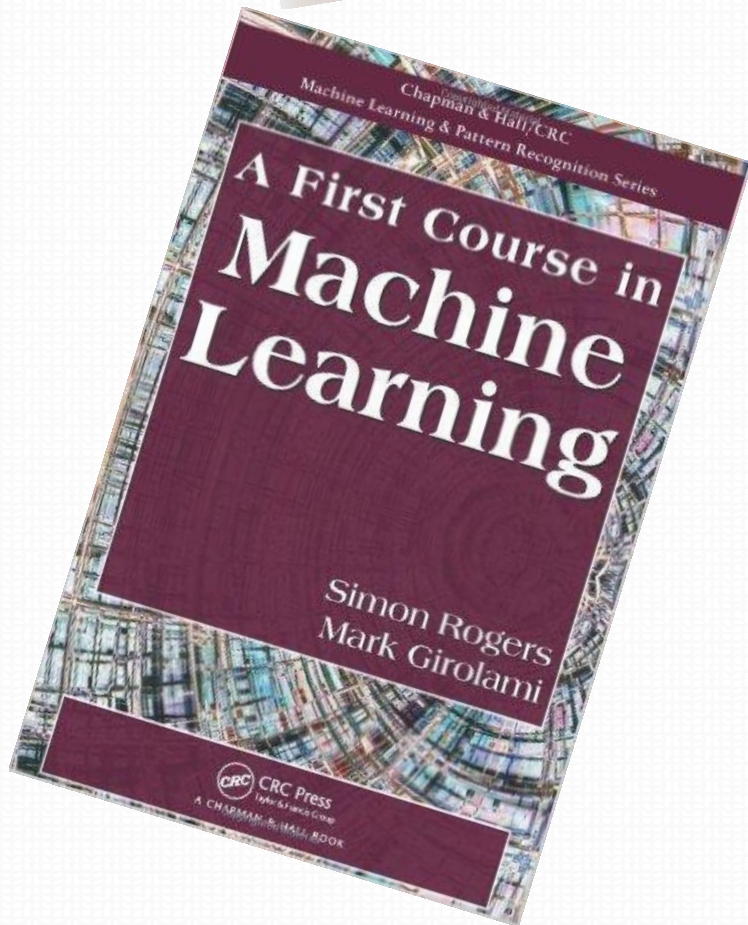
- It determines the dominant modes of variation from within the data and then projects data onto this 'natural' coordinate system

Summary

- “curse of dimensionality” with high-dimensional data can pose various problems
- Data projection to preserve ‘interesting’ characteristics within data
 - Avoid “curse of dimensionality”
- Dimensionality reduction is a form of unsupervised learning
- Dimensionality reduction is often used as a pre-processing step before classification or clustering
 - The ‘success’ of dimensionality reduction can be evaluated by subsequent processing operation

Exercise (ungraded)

- Experiment with MATLAB code – `pcaexample.m`
(from FCML book website)
- Experiment with MATLAB code – `pcaexample2.m`
(from FCML book website)
- Experiment with MATLAB code –
`rgbeyepcaexample.m` (from Canvas)



Author's material
(Simon Rogers)



Ata Kaban



Iain Styles



Thank You