

# IOT607U Data Mining

## Week 4: Data preprocessing

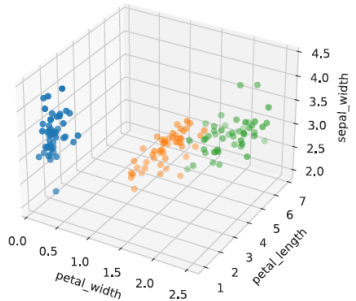
---

Dr Lin Wang

School of EECS, Queen Mary University of London

# Last week

1. Data Exploration
2. Data Summarization
3. Data Visualization



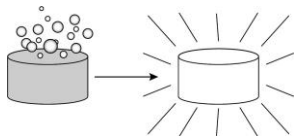
# This week's contents

1. Data Preprocessing: An Overview

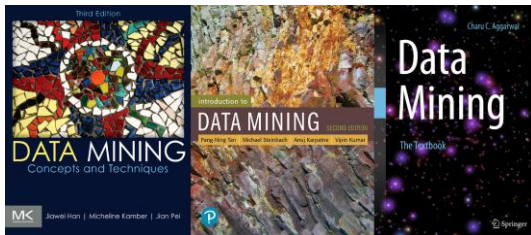
2. Data Cleaning

3. Data Normalization

4. Dimensionality Reduction



- Chapter 3 of J. Han, M. Kamber, J. Pei, “Data Mining: Concepts and Techniques”, 3rd edition, Elsevier/Morgan Kaufmann, 2012
- Section 2.3 of P.-N. Tan, M. Steinbach, A. Karpatne, V. Kumar, “Introduction to Data Mining”, 2nd edition, Pearson, 2019
- Chapter 2 of C. C. Aggarwal, “Data Mining: The Textbook”, Springer, 2015



# **Data Preprocessing: An Overview**

---

# Why Preprocess the Data?

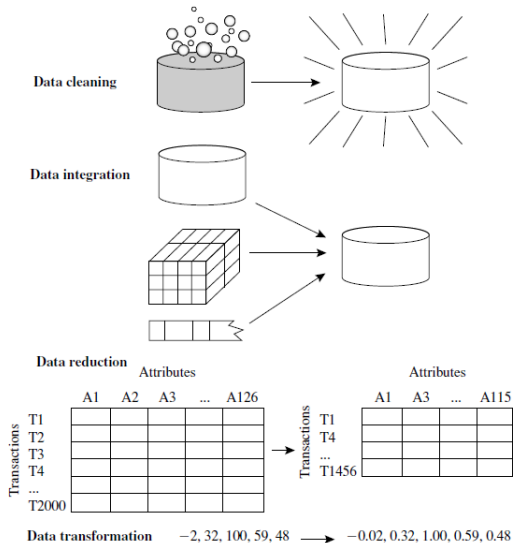
Data have **quality** if they satisfy the requirements of the intended use. There are many factors comprising data quality:

- **Accuracy**: incorrect attribute values
- **Completeness**: not recorded, unavailable
- **Consistency**: data not consistent with other recorded data
- **Timeliness**: data not updated in a timely fashion
- **Believability**: how much the data is trusted by users
- **Interpretability**: how easily the data can be understood

# What is Data Preprocessing? – Major Tasks

- **Data cleaning:** handle missing data, smooth noisy data, identify or remove outliers, and resolve inconsistencies
- **Data integration:** Integration of multiple databases
- **Data reduction:**
  - Subset selection
  - Dimensionality reduction
- **Data transformation:**
  - Normalization

# What is Data Preprocessing? – Major Tasks



Forms of data preprocessing.



# Data Cleaning

---

# DATA CLEANING

” *80 percent of a data scientist's valuable time is spent simply finding, cleansing, and organizing data, leaving only 20 percent to actually perform analysis...*

IBM Data Analytics

- Filling in missing values
- Smoothing noisy data
- Removing outlier
- Correcting inconsistencies

Data in the real world is incomplete, noisy, and inconsistent.

- **Incomplete**: lacking attribute values, lacking certain attributes of interest, or containing only aggregate data  
e.g. *Occupation* = "" (missing data)
- **Noisy**: containing noise, errors, or outliers  
e.g. *Salary* = "-10" (an error)
- **Inconsistent**: containing discrepancies in records  
e.g. *Age* = "42", *Birthday* = "03/07/2010"  
e.g. Rating was "1, 2, 3", now rating is "A, B, C"
- **Intentional** (disguised missing data)  
e.g. *Jan. 1 as everyone's birthday*

# MISSING DATA

Data is not always available, and many entries might have no recorded value for several attributes.

## SOURCES OF MISSING DATA

- User forgot to fill in a field.
- Data was lost while transferring manually from a legacy database.
- There was a programming error.
- Users chose not to fill out a field tied to their beliefs about how the results would be used or interpreted.

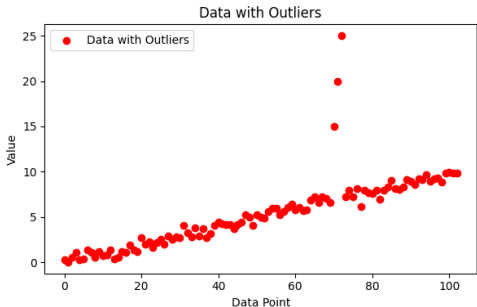
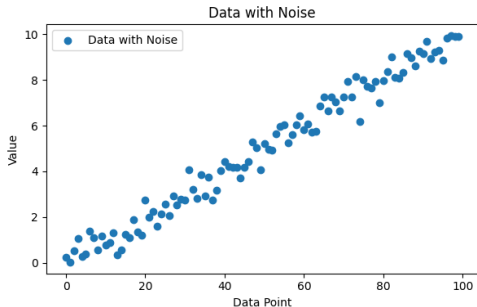
## HANDLING MISSING DATA

- **Ignore:** not effective when the % of missing values per attribute varies considerably
- **Fill in the missing value manually:** tedious + infeasible
- **Fill in the missing value automatically** with:
  - A global constant
  - The attribute mean/median
  - The most probable value: inference using interpolation/regression
- Useful Python functions
  - `isnull()`: detecting missing/null values
  - `dropna()`: removing rows/columns containing null values
  - `fillna()`: filling null values using a specified method

# Noisy Data

**Noise:** random error or variance in a measured variable

**Outlier:** data points that significantly differ from the rest of the data set

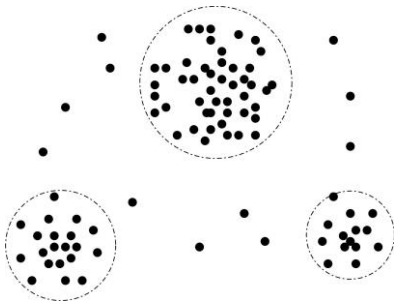


# How to Handle Noisy Data?

**Smoothing:** moving average, exponential smoothing

**Regression:** smooth by fitting the data into regression functions.

**Outlier Analysis:** detect and remove outliers.



**Semi-supervised:** combined computer and human inspection

# Data Transformation by Normalization

---



# Data Transformation

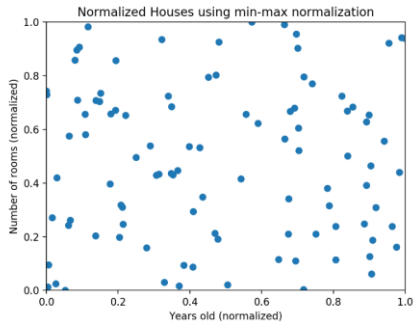
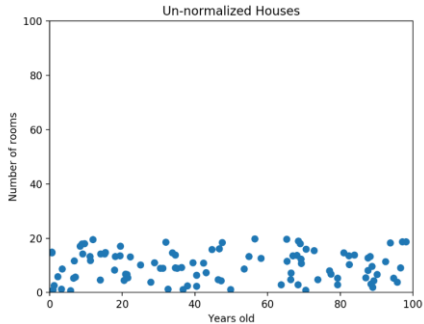
## Data Transformation

Data are transformed or consolidated so that the resulting mining process may be more efficient, and the patterns found may be easier to understand.

**Normalisation:** Scaled to fall within a smaller, specified range

- min-max normalisation
- z-score normalisation
- normalisation by decimal scaling

# Normalisation



# Normalisation

The measurement unit used in an attribute can affect data analysis. In **normalisation**, attributes are scaled to fall within a smaller, specified range, such that all attributes have an equal weight in the analysis.

**Min-max normalisation:** maps a value  $v$  of  $A$  to  $v'$  in the range  $[new\_min_A, new\_max_A]$ .

$$v' = \frac{v - min_A}{max_A - min_A}(new\_max_A - new\_min_A) + new\_min_A$$

**Example:** Normalise the income range 212,000 - 298,000 to [0.0, 1.0]. What is the normalised value for 273,600?

# Normalisation

**z-score normalisation:** the values for an attribute  $A$  are normalised based on the mean and standard deviation of  $A$ .

$$v' = \frac{v - \mu_A}{\sigma_A}$$

Useful when the actual minimum and maximum of attribute  $A$  are unknown, or when there are outliers that dominate the min-max normalisation.

**Example:** Suppose that the mean and standard deviation of the values for the attribute income are 254,000 and 216,000, respectively. What is the z-score normalised value for 273,600?

# Normalisation

**Normalisation by decimal scaling:** normalises by moving the decimal point of values of attribute  $A$ .

$$v' = \frac{v}{10^j}$$

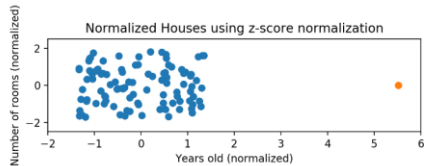
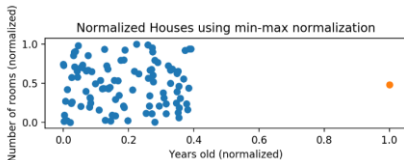
where  $j$  is the smallest integer such that  $\max(|v'|) < 1$ .

**Example:** Suppose that the recorded values of  $A$  range from -986 to 917. The maximum absolute value of  $A$  is 986. To normalise by decimal scaling, we therefore divide each value by 1000 ( $j = 3$ ) so that -986 normalizes to -0.986 and 917 normalises to 0.917.

**Note:** it is necessary to save the normalisation parameters (e.g. the mean and standard deviation if using z-score normalisation) so that future data can be normalized in a uniform manner.

# MIN-MAX vs Z-SCORE

- Min-max normalization: Guarantees all features will have the exact same scale but does not handle outliers well.
- Z-score normalization: Handles outliers, but does not produce Normalized data with the exact same scale.



# Dimensionality Reduction

---

## Data Reduction

Obtain a reduced representation of the data set.

- Much smaller in volume but yet produces almost the same analytical results
- **Why data reduction?** Complex analysis may take a very long time to run on a complete data set.
- **Methods**
  - Attribute subselect selection
  - Dimensionality reduction: principal component analysis (PCA)  
multidimensional scaling (MDS)  
t-SNE



# **Dimensionality Reduction**

---

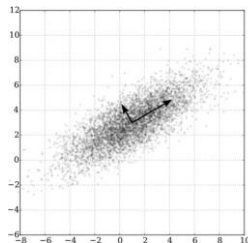
**Principal component analysis (PCA)**

# Principal Component Analysis

## Principal Component Analysis (PCA)

A statistical procedure that uses a transformation to convert a set of observations of possibly correlated variables into a set of values of linearly uncorrelated variables called **principal components**.

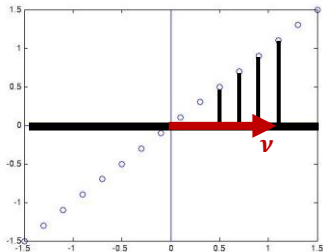
- The original data is **projected** onto a much smaller space, resulting in dimensionality reduction.
- **Method:** find the **eigenvectors** of the **covariance** matrix, and these eigenvectors define the new space.



# Projection

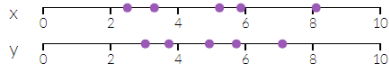
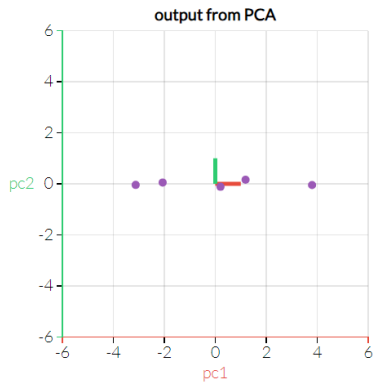
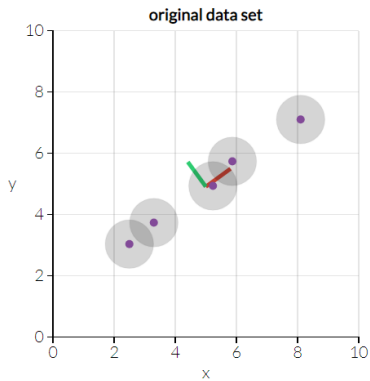
- Data points:  $x_i, i = 1, \dots, n$
- Direction defined by  $v$
- Projecting  $x_i$  towards the direction, the new value is

$$y_i = v^T x_i$$

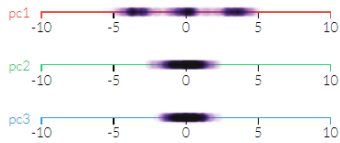
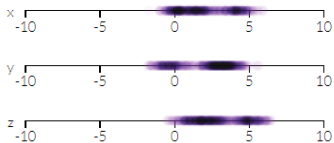
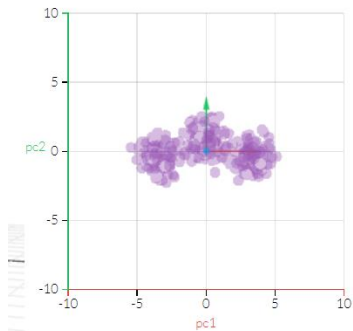
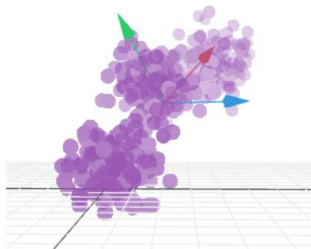


$$v = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

# 2D Example



# 3D Example



# Principal Component Analysis

- To find the principal components (projections)
  - The projection should be **orthogonal**
  - After projection, the data should be mostly spread out, i.e. with the highest **variance**
- PCA based on eigenvalue decomposition.

• Given samples  $X = \{\mathbf{x}_1, \dots, \mathbf{x}_D\}$ ,  $\mathbf{x}_i \in R^{M \times 1}$

- compute sample mean:  $\boldsymbol{\mu} = \frac{1}{D} \sum_i \mathbf{x}_i$
- compute sample variance:  $\boldsymbol{\Sigma} = \frac{1}{D} \sum_i (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^T$
- compute eigenvalue and eigenvector of  $\boldsymbol{\Sigma} = \boldsymbol{\Phi} \boldsymbol{\Lambda} \boldsymbol{\Phi}^{-1}$ :
  - $\boldsymbol{\Phi} = [\mathbf{v}_1, \dots, \mathbf{v}_M]$  – projection
  - $\boldsymbol{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_M)$  - the variance after projection
- order eigen values  $\lambda_1 > \dots > \lambda_M$  and remain the first  $K$  eigenvalues
- the eigenvectors  $\mathbf{v}_1, \dots, \mathbf{v}_K$  would be the  $K$  principal components

$$\text{(projections)} \quad V = \begin{bmatrix} \mathbf{v}_1^T \\ \vdots \\ \mathbf{v}_K^T \end{bmatrix} \rightarrow Y = VX$$

# Eigenvalue decomposition - review

- Eigenvalue decomposition of an  $M \times M$  matrix  $\mathbf{A}$ :
  - $\mathbf{A} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^{-1}$ , where  $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_M)$ ,  $\mathbf{Q} = [\mathbf{v}_1, \dots, \mathbf{v}_M]$
- $\mathbf{A}\mathbf{v}_i = \lambda_i\mathbf{v}_i$ 
  - $\lambda_i$  is the eigenvalue (scalar),  $i = 1, \dots, M$
  - $\mathbf{v}_i$  is the eigenvector ( $\mathbf{v}_i \in \mathbb{R}^{M \times 1}$ ),  $i = 1, \dots, M$
  - $\mathbf{v}_i$  and  $\mathbf{v}_j$  are orthogonal to each other,  $i \neq j$

$$\begin{array}{ccccccc} \mathbf{A} & & \mathbf{Q} & & \mathbf{\Lambda} & & \mathbf{Q}^{-1} \\ \left[ \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} \right] & = & \left[ \begin{array}{|c|c|c|} \hline \mathbf{v}_1 & \mathbf{v}_2 & \mathbf{v}_3 \\ \hline \end{array} \right] & \left[ \begin{array}{|c|c|c|} \hline \lambda_1 & 0 & 0 \\ \hline 0 & \lambda_2 & 0 \\ \hline 0 & 0 & \lambda_3 \\ \hline \end{array} \right] & \left[ \begin{array}{|c|c|c|} \hline \mathbf{v}_1 & \mathbf{v}_2 & \mathbf{v}_3 \\ \hline \end{array} \right]^{-1} \\ & & \underbrace{\hspace{1.5cm}} & & \underbrace{\hspace{1.5cm}} & & \underbrace{\hspace{1.5cm}} \\ & & \text{Eigen vectors} & & \text{Eigen values} & & \text{Eigen vectors} \\ & & \text{of} & & \text{of} & & \text{of} \\ & & \mathbf{A} & & \mathbf{A} & & \mathbf{A} \end{array}$$

## Eigenvalue decomposition - review

- Generalized eigenvalue decomposition for a pair of matrix  $A$  and  $B$ :  $Av_i = \lambda_i Bv_i$ 
  - $\lambda_i$  is the generalized eigenvalue,  $i = 1, \dots, M$
  - $v_i$  is the generalized eigenvector,  $i = 1, \dots, M$ ;
  - $v_i$  and  $v_j$  are orthogonal to each other
- Rayleigh quotient

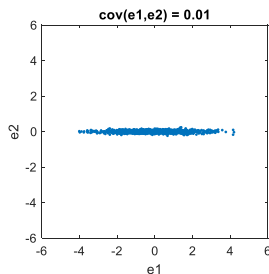
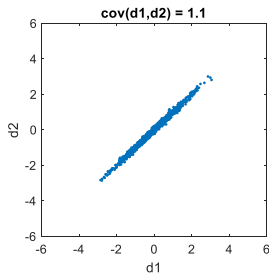
$$\lambda_i = \frac{v_i^T A v_i}{v_i^T B v_i}$$

When  $B = I$ ,  $\lambda_i = \frac{v_i^T A v_i}{v_i^T v_i}$



# Covariance

- $n$  data points:  $\mathbf{X} = \begin{bmatrix} \mathbf{x}^{d_1} \\ \mathbf{x}^{d_2} \end{bmatrix} = \begin{bmatrix} x_1^{d_1} & \dots & x_n^{d_1} \\ x_1^{d_2} & \dots & x_n^{d_2} \end{bmatrix}_{2 \times n}$
- The covariance between two directions ( $d_1, d_2$ ) is
$$\text{cov}(d_1, d_2) = \frac{\sum_{i=1}^n (x_i^{d_1} - \bar{x}^{d_1})(x_i^{d_2} - \bar{x}^{d_2})}{n}$$

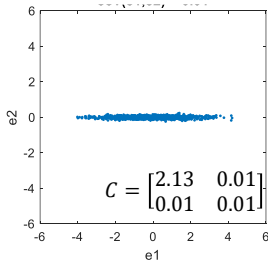
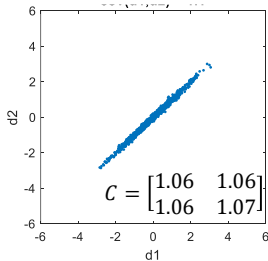


# Covariance matrix

- The covariance matrix of two directions  $d1$ ,  $d2$  is

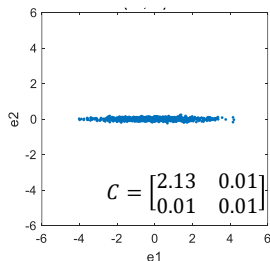
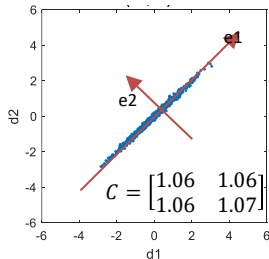
$$C(d1, d2) = \begin{bmatrix} \text{var}(d_1) & \text{cov}(d_1, d_2) \\ \text{cov}(d_2, d_1) & \text{var}(d_2) \end{bmatrix}$$

- The diagonal element (variance) indicates the spread along one direction
- The off-diagonal element (covariance) indicates the correlation (redundancy) between two directions
- Matrix form:  $C_X = XX^T$  (after removing the mean)



# Principal Component Analysis

- Principal component analysis, or PCA, is a dimensionality reduction method that is often used to reduce the redundancy in the data set
- To find some directions that preserve most information (variation/energy) of the data points after projection
- To diagonalize the covariance matrix
  - Minimize redundancy (covariance)
  - Maximize signal energy (variance)



# Principal Component Analysis

- $n$  data points:  $X = \begin{bmatrix} \mathbf{x}^{d_1} \\ \mathbf{x}^{d_2} \end{bmatrix} = \begin{bmatrix} x_1^{d_1} & \cdots & x_n^{d_1} \\ x_1^{d_2} & \cdots & x_n^{d_2} \end{bmatrix}_{2 \times n}$
- covariance matrix:  $\mathbf{C}_X = \mathbf{X}\mathbf{X}^T$
- one projection direction:  $\mathbf{v}_i = \begin{bmatrix} v_{i1} \\ v_{i2} \end{bmatrix}_{2 \times 1}$
- $n$  data points after projection:  $\mathbf{y}^{v_i} = \mathbf{v}_i^T \mathbf{X} = [y_1^{v_i}, \dots, y_n^{v_i}]$
- variance along  $\mathbf{v}_i$ :  $\text{var}(\mathbf{v}_i) = \mathbf{y}^{v_i}(\mathbf{y}^{v_i})^T = \mathbf{v}_i^T \mathbf{X}\mathbf{X}^T \mathbf{v}_i = \mathbf{v}_i^T \mathbf{C}_X \mathbf{v}_i$

$$\mathbf{C}_Y = \begin{bmatrix} \mathbf{v}_1^T \mathbf{C}_X \mathbf{v}_1 & \mathbf{v}_1^T \mathbf{C}_X \mathbf{v}_2 \\ \mathbf{v}_2^T \mathbf{C}_X \mathbf{v}_1 & \mathbf{v}_2^T \mathbf{C}_X \mathbf{v}_2 \end{bmatrix}$$

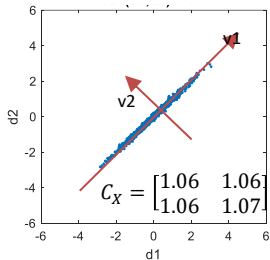
# Principal Component Analysis

$$\mathbf{C}_Y = \begin{bmatrix} \mathbf{v}_1^T \mathbf{C}_X \mathbf{v}_1 & \mathbf{v}_1^T \mathbf{C}_X \mathbf{v}_2 \\ \mathbf{v}_2^T \mathbf{C}_X \mathbf{v}_1 & \mathbf{v}_2^T \mathbf{C}_X \mathbf{v}_2 \end{bmatrix}$$


- Eigenvalue decomposition:  $\mathbf{C}_X \mathbf{v} = \lambda \mathbf{v}$
- eigenvalues:  $\{\lambda_1, \lambda_2\}$ , eigenvectors:  $\{\mathbf{v}_1, \mathbf{v}_2\}$
- Rayleigh quotient:  $\lambda = \frac{\mathbf{v}^T \mathbf{C}_X \mathbf{v}}{\mathbf{v}^T \mathbf{v}} \rightarrow \lambda = \mathbf{v}^T \mathbf{C}_X \mathbf{v} \text{ s.t. } \mathbf{v}^T \mathbf{v} = 1$
- To maximize the variance after projection (the diagonal item of the covariance matrix)
  - first direction:  $\mathbf{v}_1$  is the eigenvector corresponding to the largest eigenvalue
  - second direction:  $\mathbf{v}_2$  is the eigenvector corresponding to the second largest eigenvalue
  - ...
- The off-diagonal item of the covariance matrix
  - $\mathbf{y}^{v_1} = \mathbf{v}_1^T \mathbf{X}, \mathbf{y}^{v_2} = \mathbf{v}_2^T \mathbf{X}$
  - $\text{cov}(\mathbf{v}_1, \mathbf{v}_2) = \mathbf{y}^{v_1} (\mathbf{y}^{v_2})^T = \mathbf{v}_1^T \mathbf{C}_X \mathbf{v}_2 = \lambda_2 \mathbf{v}_1^T \mathbf{v}_2 = 0$
- $\mathbf{C}_Y = \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix}$

# Principal Component Analysis

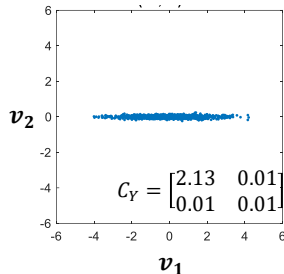
$$C_Y = \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix}$$



$$v_1 = \begin{bmatrix} 0.71 \\ 0.71 \end{bmatrix}, \lambda_1 = 2.13$$



$$v_2 = \begin{bmatrix} -0.71 \\ 0.71 \end{bmatrix}, \lambda_2 = 0.01$$



# Principal Component Analysis

- To find the principal components (projections)
  - The projection should be **orthogonal**
  - After projection, the data should be mostly spread out, i.e. with the highest **variance**
- PCA based on eigenvalue decomposition.

- Given samples  $X = \{\mathbf{x}_1, \dots, \mathbf{x}_D\}$ ,  $\mathbf{x}_i \in R^{M \times 1}$

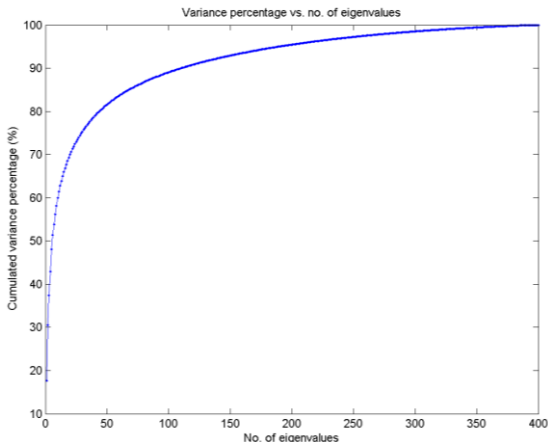
- compute sample mean:  $\boldsymbol{\mu} = \frac{1}{D} \sum_i \mathbf{x}_i$
- compute sample variance:  $\boldsymbol{\Sigma} = \frac{1}{D} \sum_i (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^T$
- compute eigenvalue and eigenvector of  $\boldsymbol{\Sigma} = \boldsymbol{\Phi} \boldsymbol{\Lambda} \boldsymbol{\Phi}^{-1}$ :
  - $\boldsymbol{\Phi} = [\mathbf{v}_1, \dots, \mathbf{v}_M]$  – projection
  - $\boldsymbol{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_M)$  - the variance after projection
- order eigen values  $\lambda_1 > \dots > \lambda_M$  and remain the first  $K$  eigenvalues
- the eigenvectors  $\mathbf{v}_1, \dots, \mathbf{v}_K$  would be the  $K$  principal components

$$\text{(projections)} \quad V = \begin{bmatrix} \mathbf{v}_1^T \\ \vdots \\ \mathbf{v}_K^T \end{bmatrix} \rightarrow Y = VX$$

# Principal Component Analysis

Determining K: cumulative variance proportion explained by the first  $n$  principal components

$$p(n) = \frac{\sum_{i=1}^n \lambda_i}{\sum_{i=1}^{N^2} \lambda_i}$$

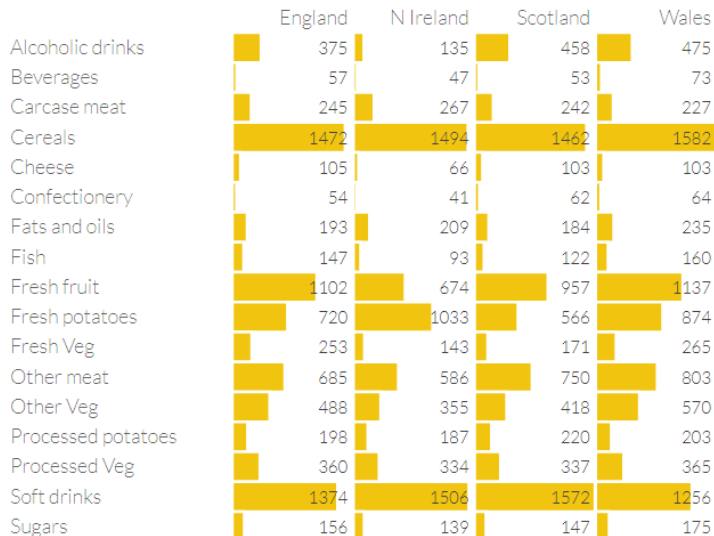




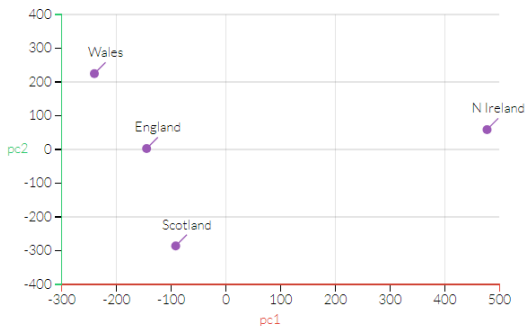
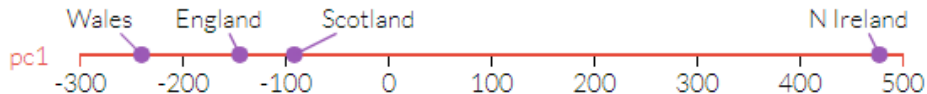
# Principal Component Analysis

- Application of PCA
  - Pattern extraction and visualization
  - Dimensionality reduction for classification
  - Data compression
- Advantages:
  - It helps in data compression and removes correlated features.
  - It helps in Speeding up other Data Mining Algorithms.
  - It converts high-dimensional data into low-dimensional data which improves and make visualization easy.
- Disadvantages:
  - It may lead to some amount of data loss.
  - It does not consider class separability and thus might not be optimal for classification problem.

# Principal Component Analysis



# Principal Component Analysis



# Dimensionality Reduction

---

Multidimensional scaling (MDS)

## Multidimensional scaling (MDS)

- Multidimensional scaling (MDS) is an unsupervised dimension reduction technique for data visualization
- Mapping data (observations) from a higher dimension down to a lower dimension, in such a way that the similarity (distance) between observations is preserved
- E.g. if observations are “similar” in 5D, they should be “similar” when mapped down to 2D with MDS
- MDS is commonly used to visualize complex, high-dimensional data, and to identify patterns and relationships that may not be apparent in the original space.

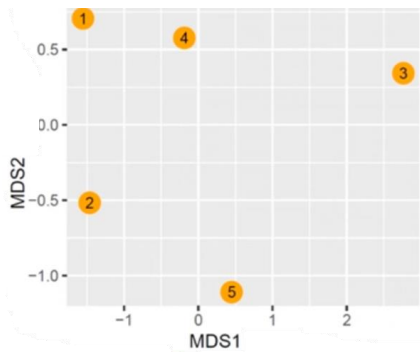
# Multidimensional scaling (MDS)



ID	DBP	SBP	BMI	Chol
1	82	132	18	192
2	86	133	20	240
3	98	145	35	140
4	85	139	22	170
5	87	145	28	218

Distance matrix

	1	2	3	4	5
1	0	1.426	4.356	1.463	2.741
2	1.426	0	4.327	2.051	2.314
3	4.356	4.327	0	3.093	2.866
4	1.463	2.051	3.093	0	1.809
5	2.741	2.314	2.866	1.809	0



# Multidimensional scaling (MDS)

- Classical multidimensional scaling (MDS) is also called Principal Coordinates Analysis (PCoA)
- General steps (very similar to PCA)
  - Compute distance matrix (e.g. Euclidean distance)
  - Eigenvalue decomposition
  - Retain eigenvectors with largest eigenvalues
  - Apply projection

# Dimensionality Reduction

---

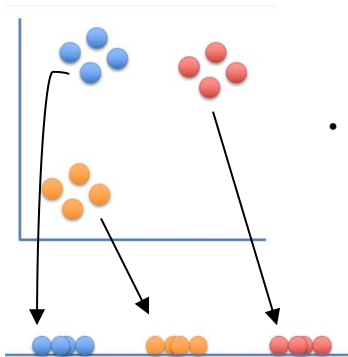
t-SNE



# t-SNE

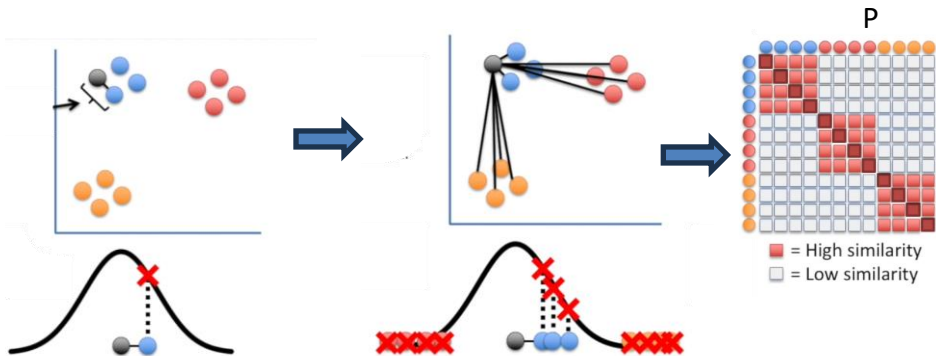
- t-SNE (t-distributed Stochastic Neighbor Embedding) is an unsupervised **non-linear** dimensionality reduction technique
- t-SNE is often used to visualize complex datasets into two and three dimensions, allowing us to understand more about underlying patterns and relationships in the data.
- t-SNE focuses on preserving the pairwise similarities (distance) between data points in a lower-dimensional space.
- It focuses on preserving small pairwise distances, e.g. the clustering of the data points.

# t-SNE



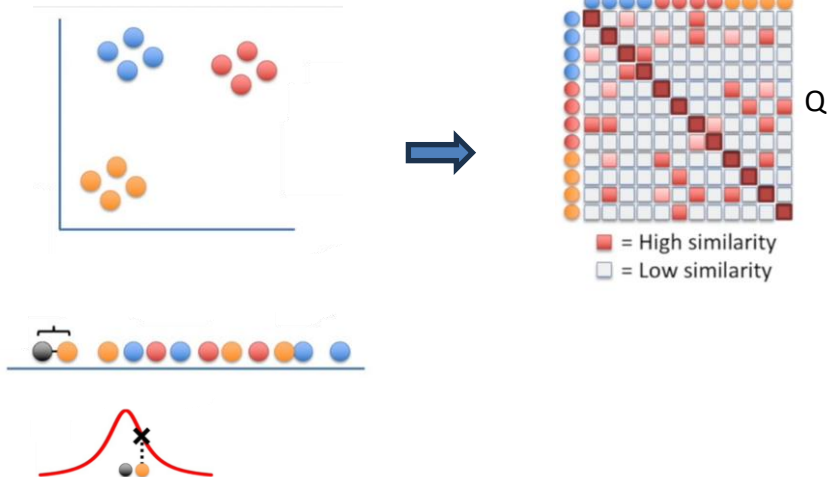
- What t-SNE does is to find way to project into a lower dimensional space so that the clustering in the high dimensional space is preserved.
- It focuses on preserving small pairwise distances, e.g. the data points within each cluster.

# t-SNE

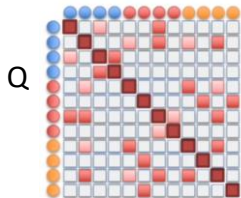


$$p_{ij} = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{d_{ij}^2}{\sigma^2}}$$

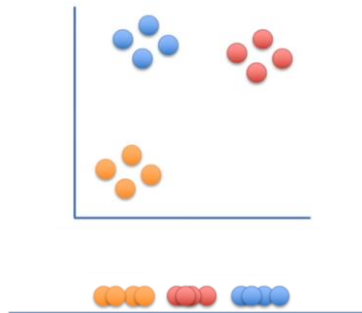
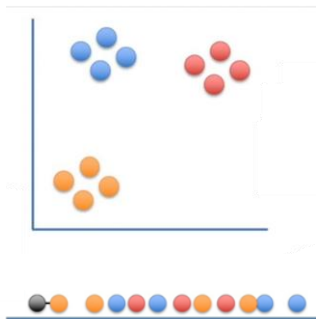
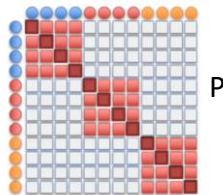
# t-SNE



# t-SNE



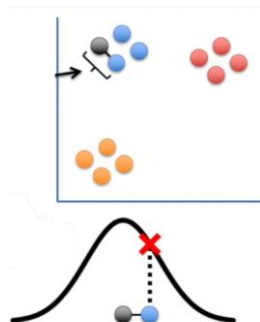
- Iterative optimization
- Minimize the difference between P (similarity matrix in high dimension) and Q (similarity matrix in low dimension), which is called KL divergence



1. Initialize points randomly / with PCA
2. Compute pairwise distances
3. Convert distances to similarity measures
4. Optimize KL divergence via gradient descent

# t-SNE

- Perplexity is a hyperparameter that measures how many neighbors we want to include for each point.
- It simply controls the variance of Gaussian kernel.
- A high perplexity means more data points are included in the neighborhood. It indicates a larger variance and a flat Gaussian distribution.
- A low perplexity means less points are included in the neighborhood. It indicates a peaked Gaussian distribution with small variance.
- Typical range [5, 50]



$$p_{ij} = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{d_{ij}^2}{\sigma^2}}$$

# t-SNE

- Original data: six 10-dimensional Gaussian balls ( $n=1000$  in each ball) far from each other
- Large perplexity tends to shrink into denser clusters

Perplexity 5



Perplexity 30



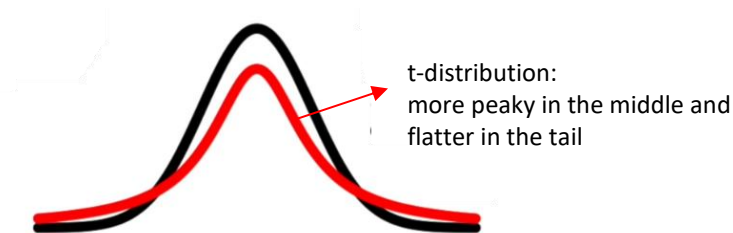
Perplexity 100





# t-SNE

- Similarity measure: t-distribution in stead of Gaussian distribution



The “t-distribution” is the “t” in t-SNE.

# PCA vs MDS vs t-SNE

- PCA: linear projection, preserve the variance in the data
  - Application: data compression, data visualization, feature extraction,
  - Method: eigen-value decomposition on the covariance matrix
- MDS: linear projection, preserve the pair-wise distance of all the data
  - Application: data visualization,
  - Method: eigen-value decomposition on the distance matrix
- t-SNE: nonlinear projection, preserve the pair-wise similarity in local neighborhoods,
  - Application: data visualization
  - Method: iterative optimization that minimize the KL divergency between similarity matrix

# Summary

**Data cleaning** routines attempt to fill in missing values, smooth out noise while identifying outliers, and correct inconsistencies in the data.

**Data transformation** routines convert the data into appropriate forms for mining.

**Data reduction** techniques obtain a reduced representation of the data while minimizing the loss of information content. Methods discussed include PCA, MDS and t-SNE.

**Questions?**

**also please use the forum on QM+**