# Data Preprocessing

Kaiqi Zhao

The University of Auckland

# Content

- Data Cleaning
- Missing Data
    - Imputation
- Data Reduction
    - Dimensionality Reduction
    - PCA
    - Feature Selection
- Noisy Data
- Data Transformation and Data Discretization
    - Normalization
    - Discretization
- Imbalanced Data
    - Sampling
    - SMOTE

# Why preprocess?

- Preprocessing means to transform the data before we feed it to a learning algorithm

- Why would we do that?

- What would we for example do?

# In this topic we will...

- Talk about problems that can appear in data

- Introduce strategies to solve these problems

- Talk about feature selection, a very important technique in machine learning

# Major Tasks in Data Preprocessing

- Data cleaning
  - Missing values
  - Noisy data
  - Outliers

- Data reduction
  - Dimensionality reduction
  - Data compression

- Transformation and discretization
  - Normalization
  - Hierarchy generation

# Data Cleaning

# Data Cleaning

- Basic assumption in machine learning?
  - The distributions of the training and test data are the same

- But, real-world data are, in most cases, dirty

- This can lead to problems, if data are
  - Incomplete lacking attribute values, certain attributes, or containing only aggregate data
  - Noisy containing noise, errors, or outliers
  - Inconsistent containing discrepancies in codes or names
  - Intentially wrong for example, there are a lot of pictures with a GPS location just a bit west of Africa

# Missing Data

# Incomplete (Missing) Data

- Data are not always available
    - Many tuples have no recorded value for several attributes
    - E.g. customer income in sales data

- Missing data may be due to
    - Equipment malfunction
    - Inconsistent with other recorded data and thus deleted
    - Data not entered due to misunderstanding
    - Certain data may not be considered important at the time of entry
    - Not register history or changes of the data

- Missing data may need to be inferred
    - When, for example?

# What to Consider When Handling Missing Data?

- Why are data missing?

- Three types of missing data (Rubin, D. B., 1974)
  - **Missing completely at random (MCAR)**
    - Completely unrelated to the data
    - Potential problem?

| Name | Country | Income |
|------|---------|--------|
| Jane | NZ | $50k |
| Kate | NZ | $75k |
| Tom | US | $53k |
| George | UK | $64k |
| Mark | UK | $77k |
| Philippe | US | $80k |

MCAR →

| Name | Country | Income |
|------|---------|--------|
| Jane | NZ | |
| | NZ | $75k |
| Tom | US | |
| George | | $64k |
| | UK | $77k |
| Philippe | US | $80k |

Missing data not related to the data

# What to Consider When Handling Missing Data?

- Why are data missing?

- Three types of missing data (Rubin, D. B., 1974)

  - Missing at random (MAR)

    - The fact the data are missing is not related to the missing attribute itself, but to some other attributes in the data set

    - Potential problem?

| Name | Country | Income |
|------|---------|--------|
| Jane | NZ | $50k |
| Kate | NZ | $75k |
| Tom | US | $53k |
| George | UK | $64k |
| Mark | UK | $77k |
| Philippe | US | $80k |

MAR →

| Name | Country | Income |
|------|---------|--------|
| Jane | NZ | $50k |
| Kate | NZ | $75k |
| Tom | US | $53k |
| George | UK | |
| Mark | UK | |
| Philippe | US | $80k |

Missing income report from UK

# What to Consider When Handling Missing Data?

- Why are data missing?

- Three types of missing data (Rubin, D. B., 1974)
  - Missing not at random (MNAR)
    - The fact the data are missing is related to the missing attribute itself
    - Potential problem?

| Name | Country | Income |
|------|---------|--------|
| Jane | NZ | $50k |
| Kate | NZ | $75k |
| Tom | US | $53k |
| George | UK | $64k |
| Mark | UK | $77k |
| Philippe | US | $80k |

MNAR →

| Name | Country | Income |
|------|---------|--------|
| Jane | NZ | |
| Kate | NZ | $75k |
| Tom | US | |
| George | UK | |
| Mark | UK | $77k |
| Philippe | US | $80k |

People with income less than $70k might refuse to provide their income details

11

- Ignore the tuple



- Usually done when the class label is missing (classification)
- Not effective when the fraction of missing values varies considerably

# How to Handle Missing Data − Imputation

- Fill in the missing data manually



$X$

| 0 | 1 | 1 | 1 | … |
|---|---|---|---|---|
| ? | ? | ? | 1 | … |
| 1 | 0 | ? | ? | … |
| … | … | … | … | … |
| 1 | 0 | 1 | 0 | … |

$X'$

| 0 | 1 | 1 | 1 | … |
|---|---|---|---|---|
| 1 | 0 | 0 | 1 | … |
| 1 | 0 | 1 | 1 | … |
| … | … | … | … | … |
| 1 | 0 | 1 | 0 | … |

- Tedious and sometimes infeasible

# How to Handle Missing Data − Imputation

- Fill in automatically
  - A global constant



$X$ → $X'$

- E.g. " missing"
- A new class

# How to Handle Missing Data − Imputation

- Fill in automatically

  - The attribute mean



$$X$$

| 12 | 2 | 22 | 38 | … |
| 11 | ? | ? | 90 | … |
| 2 | 23 | ? | ? | … |
| … | … | … | … | … |
| 9 | 11 | 54 | 23 | … |

$$X'$$

| 12 | 2 | 22 | 38 | … |
| 11 | 12 | 38 | 90 | … |
| 2 | 23 | 38 | 30 | … |
| … | … | … | … | … |
| 9 | 11 | 54 | 23 | … |

- A very commonly used method
- Changes relationship with other variables $\Rightarrow$ bias in data

# How to Handle Missing Data − Imputation

- Fill in automatically
  - The attribute mean of the samples belonging to the same class

$$X|Y \qquad\qquad X'|Y$$

| 12 | 2 | 22 | 38 | … | 1 |
|----|----|----|----|----|----|
| 11 | ? | ? | 90 | … | 0 |
| 2 | 23 | ? | ? | … | 1 |
| … | … | … | … | … | … |
| 9 | 11 | 54 | 23 | … | 0 |

$\longrightarrow$

| 12 | 2 | 22 | 38 | … | 1 |
|----|----|----|----|----|----|
| 11 | 12 | 38 | 90 | … | 0 |
| 2 | 23 | 38 | 30 | … | 1 |
| … | … | … | … | … | … |
| 9 | 11 | 54 | 23 | … | 0 |

- Might change relationship with other variables other than class ⇒ bias in data

# How to Handle Missing Data − Imputation

- Fill in automatically

  - The most probable value



- Inference-based such as decision tree, linear regression, Bayesian formula, nearest neighbour,...

# More on Imputation

- Matrix decomposition approaches
  - Decompose matrix using, for example, Singular Value Decomposition (SVD)
    - Decompose the data matrix $X$ such that $X = U\Lambda V$
    - Create imputed matrix $X'$ by multiplying $U \times \Lambda \times V$

$$
\underset{n \times d}{\overset{M}{\begin{bmatrix} x_{11} & \cdots & x_{1d} \\ \vdots & \ddots & \vdots \\ x_{n1} & \cdots & x_{nd} \end{bmatrix}}} \approx \underset{n \times k}{\overset{U}{\begin{bmatrix} u_{11} & \cdots & u_{1k} \\ \vdots & \ddots & \vdots \\ u_{n1} & \cdots & u_{nk} \end{bmatrix}}} \underset{k \times k}{\overset{\Lambda}{\begin{bmatrix} \lambda_{11} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \lambda_{kk} \end{bmatrix}}} \underset{k \times d}{\overset{V}{\begin{bmatrix} v_{11} & \cdots & v_{1d} \\ \vdots & \ddots & \vdots \\ v_{k1} & \cdots & v_{kd} \end{bmatrix}}}
$$

Minimize the Sum of Squared Errors

$$
\min_{U,\Lambda,\Sigma} \sum_{x_{ij} \in X} \left( x_{ij} - [U\Lambda V]_{ij} \right)^2
$$

# Even More on Imputation

- EM imputation
    - Expectation Maximization
    - Use other variables to impute the values (Expectation)
    - Check if value is most probable (Maximization)

- Multiple imputation (e.g., MICE)
    - 1. Impute missing values using appropriate model (for example using classifier / regression model to predict the missing value)
    - 2. Repeat the step multiple times (3-5)
    - 3. Carry out required full analysis of data (e.g. build classifier and evaluate)
    - 4. Average the results (predictions or evaluation)

- So what is the best approach?

# Preprocessing and Evaluation

- So now we know a preprocessing example

- Where would you put the preprocessing step in the evaluation?

- For example, for imputation:
  - Impute the values before splitting in train and test?
  - Impute the values in the training set – then how about the test set?