# COMPSCI361: Machine Learning
## Data Preprocessing

Katerina Taskova and Jörg Simon Wicker

The University of Auckland
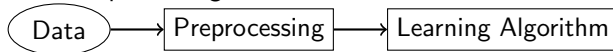
THE UNIVERSITY OF
AUCKLAND
NEW ZEALAND

SCIENCE
SCHOOL OF COMPUTER SCIENCE

# Week 5-8

- In weeks 5-8, we will cover:
  - Data Preprocessing

Data → Preprocessing → Learning Algorithm

# Week 5-8

- In weeks 5-8, we will cover:
  - Bayes Learning

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

# Week 5-8

- In weeks 5-8, we will cover:
  - Clustering



—

# Week 5-8

- In weeks 5-8, we will cover:
  - Association Rules

    If X buys *bread*, then X buys *milk* [support 50 %, confidence = 100 %]

| Bread | Eggs | Milk | Oranges |
|-------|------|------|---------|
| 1 | 1 | 1 | 0 |
| 0 | 0 | 1 | 0 |
| 1 | 0 | 1 | 0 |
| 0 | 1 | 0 | 1 |

—

# Data Preprocessing

# This week we will cover

Data Preprocessing
    Data Cleaning
    Missing Data
    Preprocessing and Evaluation
    Data Reduction
    Noisy Data
    Data Transformation and Data Discretization
    Imbalanced Data

# 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?

# This week 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
  - Numerosity reduction
  - Data compression
- Transformation and discretization
  - Normalization
  - Hierarchy generation

# Data Cleaning

- Basic assumption in machine learning?

- 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 |

# 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
  - Data history or changes of the data not recorded
- Missing data may need to be inferred
  - When, for example?

# What to Consider When Handling Missing Data?

- Missing completely at random (MCAR)
  - Completely unrelated to the data

| 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 |

  - —
  - Potential problem? Small sample size

# What to Consider When Handling Missing Data?

- **Missing at random (MAR)**
  - The fact the data are missing is related not to the missing attribute, but to some other data in the data set

| 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 |

—

  - Potential problem? Bias due to row-wise deletion

# What to Consider When Handling Missing Data?

- Missing not at random (MNAR)
  - There is a reason the data are missing and it is related to the attribute itself

| 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 |

—
  - Potential problem? Bias due to row-wise deletion

# How to Handle Missing Data – Imputation

■ Ignore the tuple



$$X$$

| 0 | 1 | 1 | 1 | ... |
|---|---|---|---|---|
| ? | ? | ? | 1 | ... |
| 1 | 0 | ? | ? | ... |
| ... | ... | ... | ... | ... |
| 1 | 0 | 1 | 0 | ... |

$$X'$$

| 0 | 1 | 1 | 1 | ... |
|---|---|---|---|---|
| ... | ... | ... | ... | ... |
| 1 | 0 | 1 | 0 | ... |

  ■ 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 infeasable

# How to Handle Missing Data – Imputation

- Fill in automatically
  - A global constant

|  |  | $X$ |  |  |
|---|---|---|---|---|
| sunny | warm | Mon | May | . . . |
| cloudy | ? | ? | July | . . . |
| sunny | cold | ? | ? | . . . |
| . . . | . . . | . . . | . . . | . . . |
| overcast | cold | Sat | June | . . . |

➜

|  |  | $X'$ |  |  |
|---|---|---|---|---|
| sunny | warm | Mon | May | . . . |
| cloudy | missing | missing | July | . . . |
| sunny | cold | missing | missing | . . . |
| . . . | . . . | . . . | . . . | . . . |
| overcast | cold | Sat | June | . . . |

  - 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 | 37 | ... |
| 11 | 12 | 38 | 90 | ... |
| 2 | 23 | 38 | 30 | ... |
| ... | ... | ... | ... | ... |
| 9 | 11 | 54 | 23 | ... |

  - Done automatically by many implementations
  - Changes relationship with other variables $\Rightarrow$ bias in data

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

|  | | $X \mid Y$ | | | |
|---|---|---|---|---|---|
| 12 | 2 | 22 | 38 | ... | 1 |
| 11 | ? | ? | 90 | ... | 0 |
| 2 | 23 | ? | ? | ... | 1 |
| ... | ... | ... | ... | ... | ... |
| 9 | 11 | 54 | 23 | ... | 0 |

➡

|  | | $X' \mid Y$ | | | |
|---|---|---|---|---|---|
| 12 | 2 | 22 | 38 | ... | 1 |
| 11 | 11 | 54 | 90 | ... | 0 |
| 2 | 23 | 22 | 38 | ... | 1 |
| ... | ... | ... | ... | ... | ... |
| 9 | 11 | 54 | 23 | ... | 0 |

  - Might change relationship with other variables other than class $\Rightarrow$ bias in data

- Fill in automatically
  - The most probable value



  - Inference-based such as Bayesian formula, decision tree, nearest neighbour,…

# More on Imputation

- Matrix decomposition approaches
  - Decompose matrix using, e.g, Singular Value Decomposition
    - Decompose the data matrix $X$ such that $X = U \Lambda V^T$
    - Create imputed matrix $X'$ by multiplying $U \times \Lambda \times V^T$

$$
\begin{bmatrix} x_{11} & \cdots & x_{1d} \\ \vdots & \ddots & \vdots \\ x_{n1} & \cdots & x_{nd} \end{bmatrix} \approx \begin{bmatrix} u_{11} & \cdots & u_{1k} \\ \vdots & \ddots & \vdots \\ u_{n1} & \cdots & u_{nk} \end{bmatrix} \begin{bmatrix} \lambda_{11} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \lambda_{nk} \end{bmatrix} \begin{bmatrix} v_{11} & \cdots & v_{1d} \\ \vdots & \ddots & \vdots \\ v_{k1} & \cdots & v_{kd} \end{bmatrix}
$$

  - Minimize the sum of squared errors

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

# 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?

# Conclusion

- Preprocessing is an important part in machine learning and data analysis
- Missing values can be caused by various reasons depending on what the reasons are, they must be addressed differently
- Various imputation approaches exist, they use the information of other instances and values to impute the missing values

# Literature

- Material in Chapter 3 in Han's *Data Mining*

Thank you for your attention!

https://ml.acukland.ac.nz