# Feature Engineering & Data Exploration

Dr Srećko Joksimović

University of
South Australia

# Feature Engineering

Applying **domain knowledge** to create new features that allow machine learning algorithms to **improve performance (or work at all)**
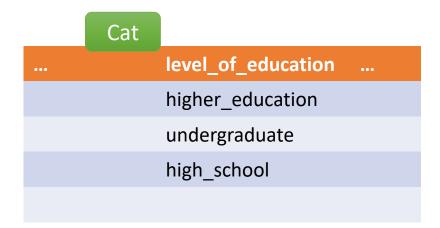
# Why so important?

- More important than algorithm selection

- Algorithms have no **_deeper_** understanding of data

- Intuition is good, but risky – remember to evaluate
  - For example, predicting health status or wellbeing:
    - Your weight in kilograms to the first decimal place (e.g., 86.4)
    - Your height in centimetres (e.g., 172).
    - Why not BMI, instead?

# Example: Date/Time Fields

| Date/Time |
|---|
| 2021-03-15 09:12 |

→

| ... | Year | Month | Day | Hour | Minute |
|---|---|---|---|---|---|
| | Num | Cat | Num | Num | Num |
| | 2021 | Mar | 15 | 9 | 12 |
| | | | | | |

University of
South Australia

# Example: Dummy variables



| Cat | | | | Num | Num | Num |
|---|---|---|---|---|---|---|
| ... | level_of_education | ... | | ... | higher_education | undergraduate | high_school |
| | higher_education | | | | 1 | 0 | 0 |
| | undergraduate | | | | 0 | 1 | 0 |
| | high_school | | | | 0 | 0 | 1 |

# Example: Text analysis

**Text** analysis, also known as **text** mining, is a machine learning technique used to automatically extract value from **text** data. With the help of natural language processing (NLP), **text** analysis tools are able to understand, analyze, and extract insights from your unstructured data.

**TEXT** appears 4 times

University of
South Australia

# Example: Text analysis

Text analysis, also known as **text** mining, is a machine learning technique used to automatically extract value from **text** data. With the help of natural language processing (NLP), **text** analysis tools are able to understand, analyze, and extract insights from your unstructured data.

| ... | text | learn | process |
|-----|------|-------|---------|
| | 4 | 1 | 1 |
| | | | |
| | | | |
| | | | |

University of
South Australia

# Example: Binning

| ... | value | bin |
|---|---|---|
| | 2 | Low |
| | 45 | Mid |
| | 3 | Low |
| | 85 | High |
| | 28 | Low |

University of
South Australia

# Example: Computed from Existing Features

| debt | income |
|------|--------|
| 10,134 | 100,000 |
| 85,234 | 134,000 |
| 8,112 | 21,500 |
| 0 | 45,900 |
| 17,534 | 52,000 |

| debt_income_ratio |
|-------------------|
| 0.10 |
| 0.64 |
| 0.38 |
| 0 |
| 0.34 |

# Outliers

# What are outliers?

An **outlier** is a data point that's significantly different from the remaining data

**University of South Australia**

# Detecting outliers

Using visualization plots like **boxplot** and **scatterplot**

**University of South Australia**

# Detecting outliers

Using a normal distribution
(mean and std)

**University of
South Australia**

# Handling outliers

- **Trimming:** Simply removing the outliers from our dataset.

- **Imputing:** We treat outliers as missing data, and we apply missing data imputation techniques.

- **Discretization:** We place outliers in edge bins with higher or lower values of the distribution.

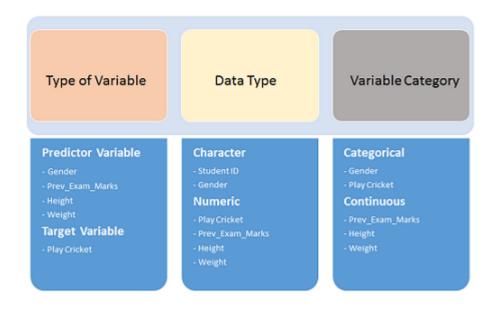- **Censoring:** Capping the variable distribution at the maximum and minimum values.

**University of
South Australia**

# Data Exploration

# Data Exploration

- Variable identification

- Univariate analysis

- Bi-variate analysis

# Variable identification

- Identify **Predictor** (Input) and **Target** (output) variables.

- Identify the data type and category of the variables.

| Student_ID | Gender | Prev_Exam_Marks | Height (cm) | Weight Caregory (kgs) | Play Cricket |
|------------|--------|-----------------|-------------|-----------------------|--------------|
| S001 | M | 65 | 178 | 61 | 1 |
| S002 | F | 75 | 174 | 56 | 0 |
| S003 | M | 45 | 163 | 62 | 1 |
| S004 | M | 57 | 175 | 70 | 0 |
| S005 | F | 59 | 162 | 67 | 0 |



Type of Variable

**Predictor Variable**
- Gender
- Prev_Exam_Marks
- Height
- Weight

**Target Variable**
- Play Cricket

Data Type

**Character**
- Student ID
- Gender

**Numeric**
- Play Cricket
- Prev_Exam_Marks
- Height
- Weight

Variable Category

**Categorical**
- Gender
- Play Cricket

**Continuous**
- Prev_Exam_Marks
- Height
- Weight

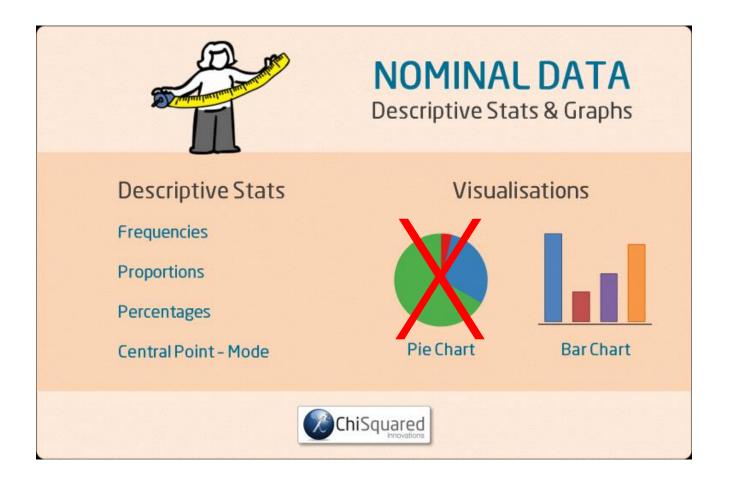**University of South Australia**

# Univariate analysis

- Exploring variables one by one

- Available methods – depending on the variable category

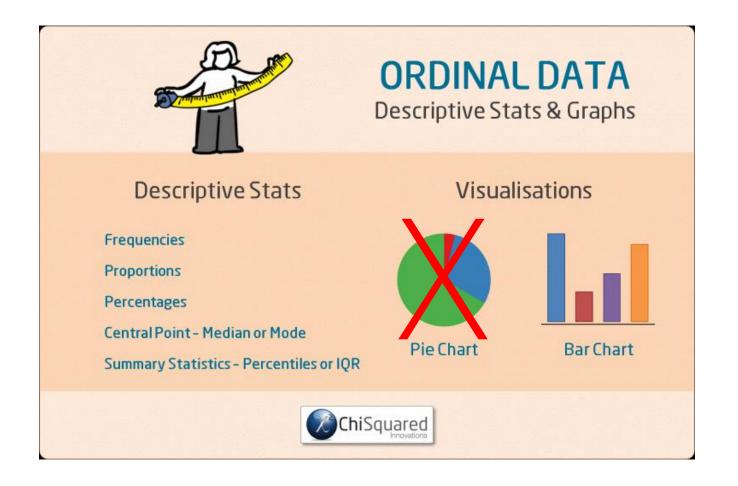Source: https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/

**University of
South Australia**

# Nominal data

Descriptive stats
and visualisations

University of
South Australia

# Ordinal data

Descriptive stats
and visualisations

University of
South Australia

# Interval data

Descriptive stats
and visualisations



INTERVAL DATA
Descriptive Stats & Graphs

**Descriptive Stats**

Central Point – Mean, Median or Mode

Range – Minimum & Maximum

Spread – Percentiles, IQR, SD

**Visualisations**

Boxplot

Histogram

ChiSquared
Innovations

**University of
South Australia**

# Ratio data

Descriptive stats
and visualisations

**University of
South Australia**

# Bi-variate analysis

- Exploring the association between two (or more) variables

- Common combinations:
  - Categorical & categorical
  - Continuous & continuous
  - Categorical & continuous

University of
South Australia

# Two categorical variables

*Table 6.1. Numerical Summary of Hometown Description*

| Hometown | Count | Proportion | Percent |
|---|---|---|---|
| Rural | 75 | $75/555 = 0.14$ | $0.14 \times 100$ |
| Suburb | 296 | $296/555 = 0.53$ | $0.53 \times 100$ |
| Small Town | 139 | $139/555 = 0.25$ | $0.25 \times 100$ |
| Big City | 45 | $45/555 = 0.08$ | $0.08 \times 100$ |
| Total | $n = 555$ | $555/555 = 1.0$ | $1.0 \times 100$ |

Numerical summary

| | Freshman | Sophomore | Junior | Senior | Total |
|---|---|---|---|---|---|
| Yes | 42 | 55 | 76 | 81 | 254 |
| No | 58 | 45 | 24 | 19 | 146 |
| Total | 100 | 100 | 100 | 100 | 400 |

Contingency table

| | Uneven Sidewalks | Even Sidewalks |
|---|---|---|
| High (over 20%) | 98 | 418 |
| Low (under 10%) | 9 | 301 |
| Total | 107 | 719 |

2x2 table

Statistical analysis:
» **Chi-square test**
» **Crammer's V**

**University of South Australia**

# Categorical and Continuous variable



Statistical analysis:
- » **T-test/Z-test**
- » **ANOVA**

**University of South Australia**

# Two continuous variables


strong, positive, linear


moderate, negative, linear


null / no relationship


strong, non-linear

Statistical analysis:
- » **Correlation**
- » **Regression**
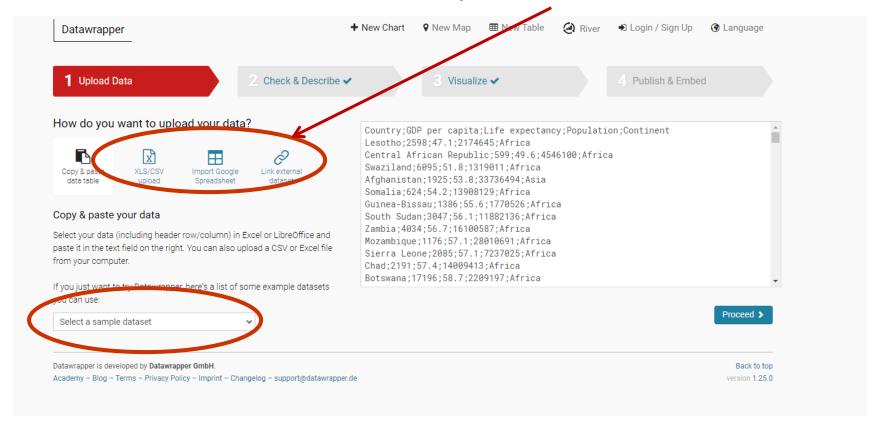
**University of South Australia**

# Hands-on

Explore visualisations (individually or in groups)

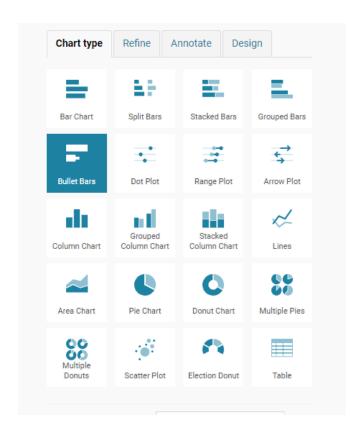- Go to https://www.datawrapper.de/

# Upload data or select a sample dataset

Choose one of the datasets provided in **Slides and Other Resources** folder

# Explore



☐ Choose at least **2 Chart types** to explore selected dataset.

☐ Explain why you selected these charts – what kind of information they provide?

☐ Discuss (compare and contrast) insights you obtained with each of the selected chart types. What those charts tell you about your dataset?