# A Review of Software Tools for Data Analytics

Timothy Patterson B.Sc. Ph.D.

ulster.ac.uk

# What is data analytics?

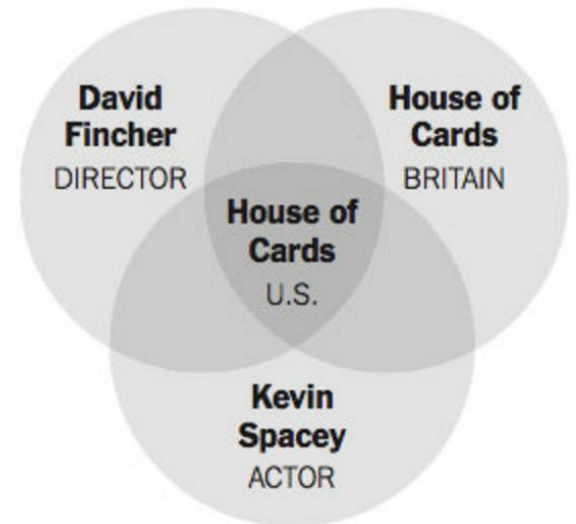"Data analytics is the science of *examining* and *drawing conclusions* from data"

# Where is data analytics used?

**NETFLIX**

- Netflix is a data driven company [2]

    - E.g. when you pause, fast forward, rewind; what series do you watch from start to finish

    - Customer retention

    - Content creation

**Circles of Proven Success**

Netflix determined that the overlap of these three areas would make "House of Cards" a successful entry into original programming.

A NETFLIX ORIGINAL SERIES
**HOUSE of CARDS**

**David Fincher** DIRECTOR

**House of Cards** BRITAIN

**House of Cards** U.S.

**Kevin Spacey** ACTOR

THE NEW YORK TIMES

Ulster University

[2]
https://blog.kissmetrics.com/how-netflix-uses-analytics/

# Where is data analytics used?

- *"We want to know what every product in the world is. We want to know who every person in the world is. And, we want to have the ability to connect them together in a transaction."* Neil Ashe, Walmart e-commerce CEO (2013)

[3] https://www.dezyre.com/article/how-big-data-analysis-helped-increase-walmarts-sales-turnover/109

# Common tasks in data analytics

- Data exploration
    - Missing values
    - Outlier detection and treatment
    - Visualization
- Feature selection / engineering
- Classification / prediction

Ulster
University

# Common tasks in data analytics
## Data exploration

- The quality of your input will determine the quality of your output and may take up to 70% of the project time [4]

- Data exploration tasks:

    - Variable identification

    - Uni/Bi-variate analysis

    - Treatment of missing values

    - Detection of outliers

    - Variable transformation / creation

Ulster
University

[4] http://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/

# Data exploration
## Variable identification

| 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**

**Data Type**

**Variable Category**

**Predictor Variable**

- Gender
- Prev_Exam_Marks
- Height
- Weight

**Target Variable**

- Play Cricket

**Character**

- Student ID
- Gender

**Numeric**

- Play Cricket
- Prev_Exam_Marks
- Height
- Weight

**Categorical**

- Gender
- Play Cricket

**Continuous**

- Prev_Exam_Marks
- Height
- Weight

# Data exploration
## Variable analysis

- Univariate analysis

  - Continuous variables

  - Categorical variables


- Bivariate analysis

  - Continuous & Continuous

  - Categorical & Categorical

  - Categorical & Continuous



Strong positive correlation
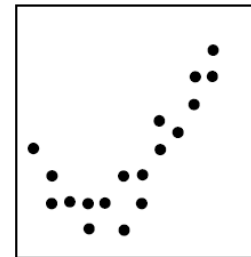
Moderate positive correlation

No correlation

Moderate negative correlation

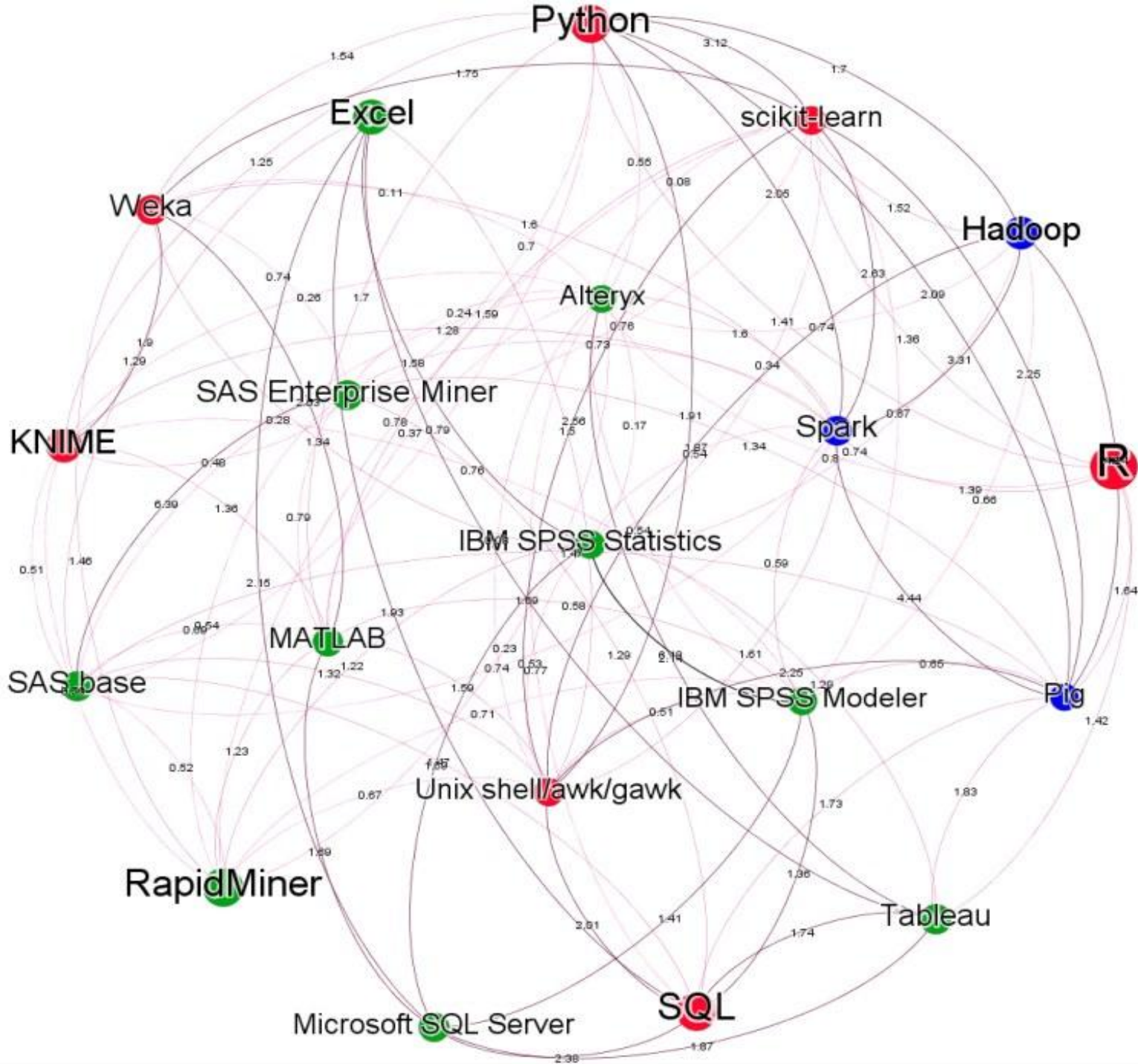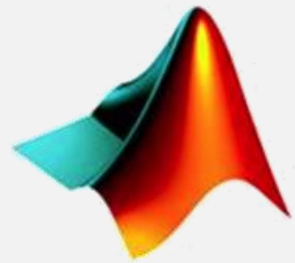Strong negative correlation

Curvilinear relationship

Ulster University

# Data exploration
## Treatment of missing values

- Reasons for missing data

    - Data extraction

    - Data collection


- Methods for treating missing values

    - Deletion

    - Mean/median/mode imputation

    - Predictive model

    - kNN imputation

Ulster University

# Data exploration

- Feature selection

- Classification

Ulster
University

# Software tools

# Software Tools

- Overview

- Learning curve

- Trouble shooting / Debugging

- User community

- Cost

- Available libraries

- Data analytics process

Ulster
University

# Software Tools
## R - Overview

- Open source
- Focus towards statistics
- Strong visualization tools





Based on estimates from:
Abel and Sander (2014) *Science* Vol. 343 no. 6178 pp. 1520 – 1522

# Software Tools
## R - Overview

- Rstudio

- Packages

    - zoo – work with time-series data

    - lattice – to visualize data

    - caret – machine learning package

Ulster
University

# Software Tools
## Weka – Overview
**http://www.cs.waikato.ac.nz/ml/weka/downloading.html**

- Open source
- GUI / Java Library
- Strong visualization tools for data exploration
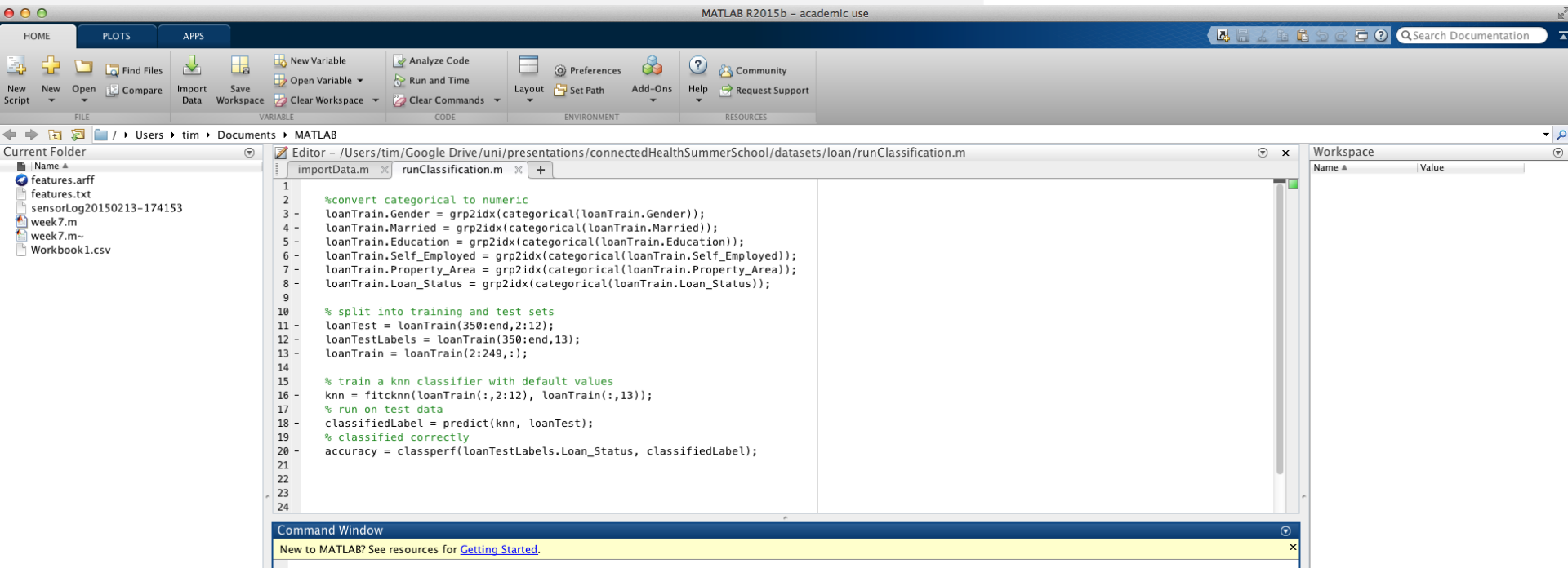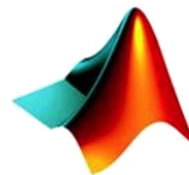
# Software Tools
## Matlab - Overview

- Owned/maintained/developed by Mathworks
- Costs (base product, add-on-products)
    - Individual £1600, add-on-products e.g. Statistics and Machine Learning £800
    - Home £85, add-on-products £25
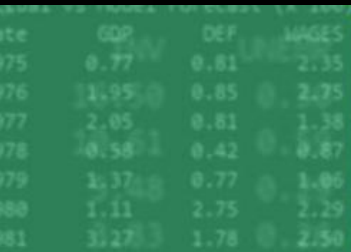    - Student £29, add-on-products £16
- Good for rapid prototyping

Ulster University

# Software Tools
## Matlab - Overview

- Intuitive GUI

- Debugging capabilities

- User Community

- Learning curve

# Software Tools
## Data analytics processes

- Importing data
    - Drag and drop, UI for importing and dealing with missing values
- Data exploration
    - *>>histogram(categorical(loanTrain.Property_Area));*
    - *>>scatter(loanTrain.ApplicantIncome,loanTrain.LoanAmount)*
    - Max, min, standard deviation, mean etc.
    - *>>a=categorical(loanTrain.Education)*
    - *>>summary(a)*

Ulster
University

# Software Tools
## Data analytics processes

- Feature selection

    - Cannot be performed in UI [5]

- Classification

    - Can use UI

    - Or code

        - *>>knn=ClassificationKNN.fit(Xtrain, Ytrain)*

**Ulster University**

[5] http://uk.mathworks.com/help/stats/select-subset-of-features-with-comparative-predictive-power.html

# Software Tools
## Python

- General purpose programming language

- Open source

- Large user community

- 2 versions, 2.x and 3.x available from
  https://www.python.org/downloads/

  - Currently more library support for 2.x

[5] http://uk.mathworks.com/help/stats/select-subset-of-features-with-comparative-predictive-power.html

Ulster
University

# Software Tools
## Python

- Popular Libraries

    - NumPy – Numerical Python

    - Matplotlib – Plotting graphs

    - Pandas – Structured data operations

    - Scikit learn – For machine learning

    - OS – Operating system and file operations

    - BeautifulSoup – Scrape webpages

Ulster
University

# Software Tools
## Python

https://try.jupyter.org/

Ulster University

# Software Tools
## Python

Pandas

- *df['columnName'].hist(bins=n)*

- *df.boxplot(column='columnName')*

- *df['columnName'].plot(kind='bar')*
    *df['Property_Area'].value_counts().plot(kind='bar')*


Matplotlib

- *plt.pyplot.scatter(df['ApplicantIncome'], df['LoanAmount'])*

- *plt.pyplot.boxplot(x=df['ApplicantIncome'])*

*where df is a DataFrame*

Ulster
University