# COMP6237 Data Mining:
# Introduction to Data Mining

## Jo Grundy and
## Markus Brede

J.Grundy@soton.ac.uk

Markus.Brede@soton.ac.uk

# Funded PhD Position

## PhD Studentship: Data Driven Modelling of Supply Chain Constraints in Maritime Transport

**University of Southampton**

| | | | |
|---|---|---|---|
| **Qualification Type:** | PhD | **Placed On:** | 17th January 2022 |
| **Location:** | Southampton | **Closes:** | 25th February 2022 |
| **Funding for:** | UK Students | | |
| **Funding amount:** | For UK students, Tuition Fees and a stipend of £15,609 tax-free per annum for up to 3.5 years. | | |
| **Hours:** | Full Time | | |

★ View Employer Profile

**Apply**

**Supervisory Team:** Adam Sobey (CMEE), Markus Brede (ECS) and Amy Parkes (Industry)

**Project description**

Recent incidents, such as the blocking of the Suez Canal or lack of workers at Chinese ports due to COVID, demonstrate the importance of maritime trade. Ships carry >90% of world trade, moving 3.5 billion tonnes of cargo (expected to triple by 2050) and 350 million passengers through 1,200 ports every year. These goods and people are carried by ~56,000 merchant ships ranging from small domestic ferries to 400m long containerships which transport 23,000 shipping containers at a time. The result of this is a complex flow network defined by port capacities, trade requirements, and routing choices of a large number of diverse ships. It is vital that we become more efficient in how we move goods around the world and that the resilience in these shipping networks is maintained and increased.

# Teaching Staff

- Credit goes to Jon Hare who developed a large part of the module


- Jo Grundy – VLC
    - Jo.Grundy@soton.ac.uk
    - 32/4053


- Markus Brede – AIC
    - Markus.Brede@soton.ac.uk
    - 32/4033

# Module Overview

- ## Not quite so new module, run for the 7$^{th}$ time

  - ### See feedback from last year

- ## Created to fill a gap

  - ### Data mining is almost synonymous with advanced machine learning

    - #### Inevitably some overlaps with COMP3206/COMP6208

      - Should be complementary and offer different views

    - #### Slightly more applied pragmatic focus

      - How do you work with real world data?
      - How do you solve real problems?

# Module Structure

- Around 26 lectures + additional tutorials
  - Wide range of data mining topics

- Assessment
  - 70% 2 hour examination
  - 30% Group coursework

| Day | Time | Room |
|-----|------|------|
| Mondays | 9AM | B46 2003 |
| Tuesdasy | 9AM | B46 2003 |
| Thursdays | 9AM | B44 1041 |
| Fridays | 5PM | B02 1039 |

There will generally be three lectures each week, making use of the Monday, first Thursday slot and the Friday slot. For some of the weeks we will also use the second Thursday slot.

The current timetable is shown below - be aware that this might change (especially if you ask us to add additional tutorial sessions):

| Date | Semester Week | Lecturer(s) | Topic/Title |
|------|---------------|-------------|-------------|
| 01-Feb | 1 | Jo & Markus | Intro to data mining |
| 03-Feb | | Markus | Linear Regression |
| 04-Feb | | Markus | Maximum Likelihood Estimation |
| 08-Feb | 2 | Markus | Tutorial/seminar: linear regression and MLE; CW set |
| 10-Feb | | Markus | Logistic regression |
| 11-Feb | | Markus | Dealing with non-linear data |
| 14-Feb | 3 | Jo & Markus | Group coursework Q & A |
| 15-Feb | | Jo & Markus | Group coursework Q & A |
| 17-Feb | | Jo & Markus | Group coursework Q & A |
| 18-Feb | | Jo & Markus | Group coursework Q & A (if needed) |
| 22-Feb | 4 | Markus | Tutorial/seminar: logistic regression |
| 23-Feb | | Markus | Intro to information theory |
| 28-Feb | 5 | Jo | Making Recommendations |
| 01-Mar | | Jo | Finding Groups, CW brief due |
| 02-Mar | | Jo | Covariance recap |
| 07-Mar | 6 | Jo | Embedding Data |
| 08-Mar | | Jo | Search |

# Module Timetable

- We have 4 slots timetabled for every week
    - Mon 9am
    - Tue 9am
    - Thu 9am.
    - Fri 5pm

- All sessions will be recorded.

- Will **not use all slots** every week (some weeks we'll use all of them, in other weeks only 2 of them)
    - Will typically use Tue, Thu, Fri
    - Have a look at the course webpage!
    - This may sometimes also change – we'll update you by email (check ECS module page)

- Roughly the plan is:

    Markus – Jo – Markus – Revisions

# Coursework Timetable

- Group coursework
  - Set next week; report submission at the end of the term (May 13)
  - Will have presentation sessions after Easter
  - More in CW Q & A sessions in week 3; by that time we want you to have formed groups
  - Once you have formed a group, please enter into this wiki:

    https://secure.ecs.soton.ac.uk/student/wiki/w/COMP6237 -2022-classlist

# Resources

- Course website [handouts, slides, interactive demos]
  - http://comp6237.ecs.soton.ac.uk
- ECS module pages [syllabus, announcements]
  - https://secure.ecs.soton.ac.uk:/module/comp6237
- Reading material
  - Toby Segaran. Programming Collective Intelligence: Building Smart Web 2.0 Applications. O'Reilly, 2007
  - Aurélien Géron. Hands-On Machine Learning with Scikit-Learn and TensorFlow: Concepts, Tools, and Techniques to Build Intelligent Systems. O'Reilly Media. March 2017
  - J. Leskovec et al. Mining of Massive Datasets. Third Edition. Cambridge University Press. 2020
  - M. J. Zaki and W. Meira, Data Mining and Machine Learning: Fundamenal Concepts and Algorithms. Cambridge University Press. 2020.

# What is Data Mining?

"Data mining is an *interdisciplinary* subfield of computer science. It is the computational process of discovering patterns in large data sets involving methods at the intersection of *artificial intelligence, machine learning, statistics, and database systems*. The overall goal of the data mining process is to <u>extract information from a data set and transform it into an understandable</u> structure for further use."

– wikipedia

# What is Data Mining?

"Generally, data mining (sometimes called data or knowledge discovery) is the process of analyzing data from different perspectives and summarizing it into useful information - information that can be used to increase revenue, cuts costs, or both."

– Bill Palace, Anderson Graduate School of Management at UCLA, 1996

DATA
(input)

Data Mining

INFORMATION
(output)

# What is Data?

- Data is any sequence of one or more symbols given meaning by specific act(s) of interpretation.

- Data (or datum - a single unit of data) is not information.
  - Data requires interpretation to become information.
  - To translate data to information, there must be several known factors considered. The factors involved are  determined by the creator of the data and the desired information.

# What is Information?

- There is a formal definition → Information theory … will have a bit of a look at this later.

- "Actionable knowledge"
  - **Prediction**
    - Christoph Adami (Michigan State) defines information as: 'the ability to make predictions with a likelihood better than chance'.
  - **Understanding**
    - Making sense of the data

# What is Data Mining?

- Given lots of data …

- **Discover patterns and models** that are:
  - **Valid**: hold on new data with some certainty
  - **Useful**: should be possible to act on the item
  - **Unexpected**: non-obvious to the system
  - **Understandable**: humans should be able to interpret the pattern

# Two Complementary Goals of Data Mining

Use some variables to predict unknown
or future values of other variables



Prediction

Understanding

Find human-interpretable patterns that
describe the data

# What kinds of data are we interested in mining?

# Categorizing data: Structured/ Unstructured

# Categorizing data: Dynamic/static/stream

# Categorizing data: Unimodal/multimodal

# Typical Data Mining Pipeline

DATA
(input)

Data Mining

INFORMATION
(output)

DATA (input) → Preprocessing → Feature Extraction → Intelligent Algorithm → INFORMATION (output)

Data Mining

# Descriptive Techniques

*PCA*
*ICA*
*MDS*
*Clustering*
*Anomaly Detection*
*…*

Intelligent Algorithm

# Predictive Techniques

*Classification*
*Ranking*
*Regression*
*Matrix Completion*
*…*

# The Plan for the Next 12 Weeks

- You will learn to solve real-world problems – e.g.:
  - Recommender systems
  - Market Basket Analysis
  - Document filtering and spam detection
  - Duplicate document detection
  - Link prediction
  - Community detection
  - Ranking search results
  - Social network analysis
- You will also learn various tools & techniques - e.g.:
  - Linear algebra (SVD, Eigendecomposition & PCA, NNMF, etc.)
  - Optimisation (e.g. stochastic gradient descent)
  - Dynamic programming (frequent itemsets)
  - Hashing (LSH, Sketching, Bloom Filters)
  - Statistics of regression analysis
  - Information theory
  - Network theory

# The Group Coursework

- You need to form groups
  - Target size is 6 (**strictly**)
  - As a group, you need to choose a data mining problem to work on
    - (You'll need to train and evaluate models and compare their performance [possibly against approaches from others])

- Come along to the slots in week 3 to discuss your ideas for problems to work on with us

- Enter your team name and team members on the student wiki:

  https://secure.ecs.soton.ac.uk/student/wiki/w/COMP6237-2022-classlist

# Key Dates

- Each team needs to submit a 1-page project brief by the end of the day of week 3 (18th of Feb).

- After Easter groups must present their idea and approaches to the class.

  – Teams should be prepared to present in the first slot; to ensure fairness we will pick teams at random

- Teams must submit a conference paper by 4pm on May 13.