# COMP6237 Data Mining:
# Introduction to Data Mining

## Sarah Hewitt and
## Markus Brede

Sarah.Hewitt@soton.ac.uk

Brede.Markus@gmail.com

# Teaching Staff

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


- Sarah Hewitt – WAIS
  - Sarah.Hewitt@soton.ac.uk
  - 32/3017


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

# Module Overview

- Not quite so new module, run for the 5th 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 28 lectures + additional tutorials
  - Wide range of data mining topics


- Assessment
  - 50% 2 hour examination
  - 20% Individual coursework
  - 30% Group coursework

The lecture slots are as follows:

| Day | Time | Room |
| --- | --- | --- |
| Mondays | 4PM | 58/1067 (Murray) |
| Tuesdays | 11AM | 29/1101 (Graham Hills) |
| Thursdays | 4PM | 58/1067 (Murray) |
| Fridays | 11AM | 29/1101 (Graham Hills) |

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

When we are not using sessions for formal teaching, the rooms are available for you to use for the group project. At those times both Sarah & Markus will endeavour to be in their respective offices should you wish to get assistance with any aspects of the course (it is advisable to email us before to give us a heads-up that you're coming though).

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 |
| --- | --- | --- | --- |
| 27-Jan | 1 | Sarah & Markus | Intro to data mining |
| 28-Jan | | | |
| 30-Jan | Markus | | Linear Regression |
| 31-Jan | Markus | Markus | Maximum Likelihood Estimation |
| 03-Feb | 2 | Markus | Tutorial/seminar: linear regression and MLE; CW1 set |
| 04-Feb | | | |
| 06-Feb | | Markus | Logistic regression |
| 07-Feb | | Sarah & Markus | Group coursework Q & A |
| 10-Feb | 3 | Sarah & Markus | Group coursework Q & A |
| 11-Feb | | Sarah & Markus | Group coursework Q & A (if needed) |
| 13-Feb | | Markus | Dealing with non-linear data |
| 14-Feb | | Markus | Tutorial/seminar: logistic regression |
| 17-Feb | 4 | Sarah | Making recommendations; CW2 set |
| 18-Feb | | Sarah | Discovering Groups |

# Module Timetable

- We have 4 slots timetabled for every week
  - Will **not use all slots** every week (some weeks we'll use all of them, in other weeks only 2 of them)
  - 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 – Sarah – Markus – Revisions

# Coursework Timetable

- Group coursework
  - Set next week; report submission at the end of the term (May 15)
  - Will have presentation sessions before Easter
  - More in CW Q & A session Feb 7 in which we want you to have formed groups
- Individual coursework
  - Set week 4 (Feb 17)
  - Due before the Easter break (March 20)

# 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

# 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?
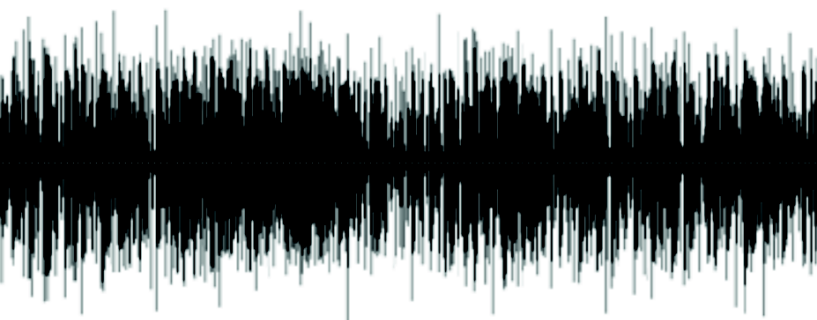
# Categorizing data: Structured/ Unstructured

# Categorizing data: Dynamic/static

# Categorizing data: Unimodal/multimodal

# Typical Data Mining Pipeline

# 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
  - Ranking search results
- 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 8 (+/- 1)
  - As a group, you need to choose a **predictive** 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 Thursday/Friday slots next week 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/COMP6237-1920-classlist

# Key Dates

- Each team needs to submit a 1-page project brief by the end of the day of the 14th of Feb.

- In week 8 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 15.