# Introduction to Data Mining

Jonathon Hare
jsh2@ecs.soton.ac.uk

Markus Brede
mb8@ecs.soton.ac.uk

# Teaching Staff

- **Jonathon Hare**

  - jsh2@ecs.soton.ac.uk

  - 1/2003


- **Markus Brede**

  - mb8@ecs.soton.ac.uk

  - 32/4033

# Module Overview

- Fairly new module - this is the second time it has run!

  - See feedback from last year

- Created to fill a gap:

  - Data mining is almost synonymous with applied machine learning

    - Inevitably some overlap in topics with COMP3206/COMP6208

      - Should be complementary & offer different views

    - Much more applied/pragmatic focus

      - How do you work with real world data?

      - How do you solve real problems?

# Module Structure

- Around 24 lectures + additional tutorials

  - Wide range of data mining topics


- Assessment:

  - 50% 2 Hour Final Exam

  - 20% Individual Coursework

  - 30% Group Coursework

# Coursework Timetable

- Group Coursework

  - Set today; report submission on the 10th May; presentations following that.

  - More info at the end of the lecture!

- Individual coursework

  - Set 24th Feb (week 4); due 22th March (just before Easter break)

# Resources

- Course web site (handouts, slides [inc 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.
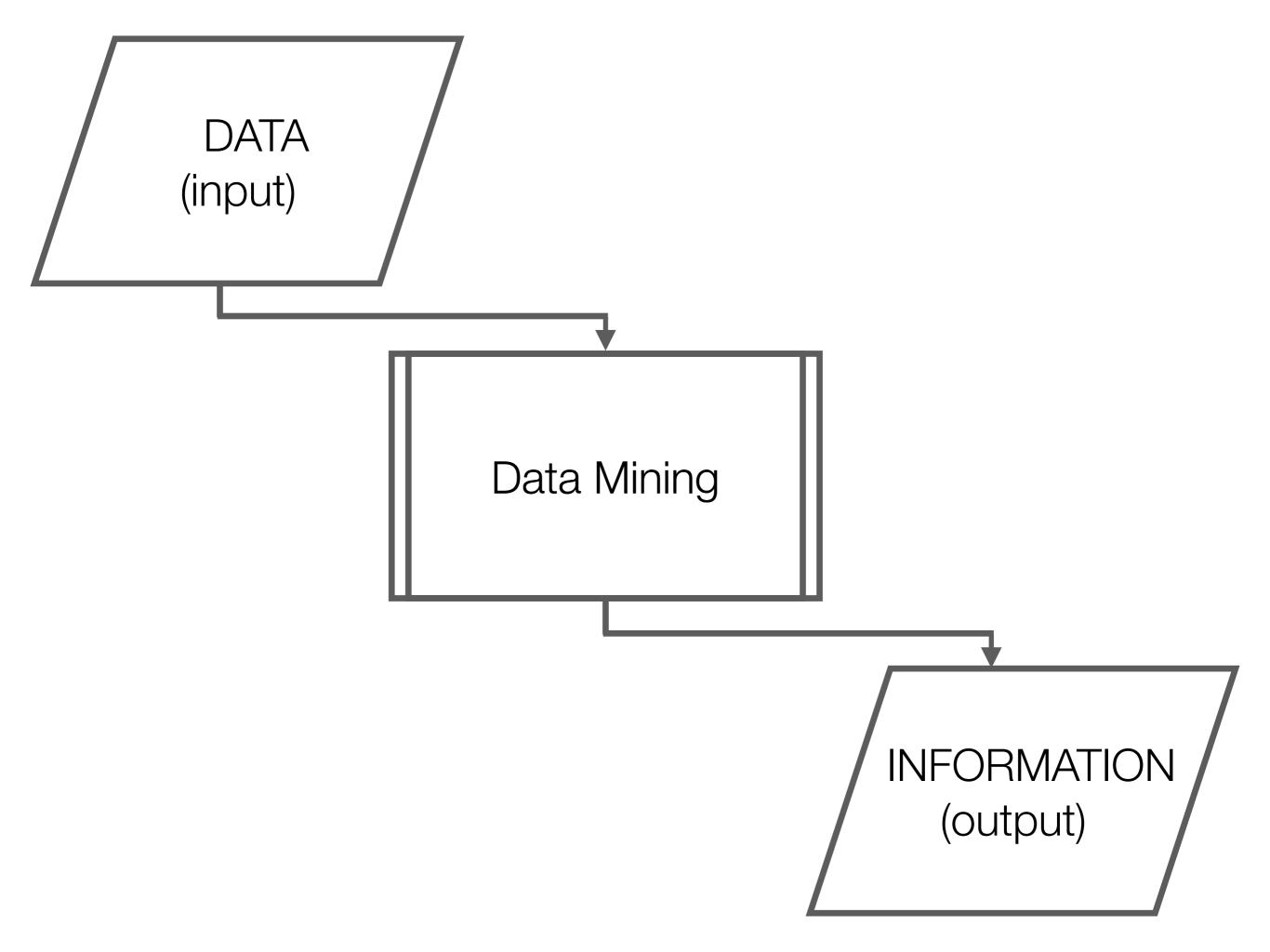
# What is "Data Mining"?

"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 extract information from a data set and transform it into an understandable structure for further use.*"

**–Wikipedia**

"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?

# 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?

# What is information?

- "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

# Two complementary goals of data mining



Prediction

Understanding

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

# Two complementary goals of data mining



Prediction

Understanding

Find human-interpretable patterns that describe the data

# What kinds of data are we interested in mining?

# Categorising data: Structured/unstructured

# Categorising data: Dynamic/static

# Categorising *data mining*: Unimodal/multimodal

What is the *typical* data mining pipeline?

```
┌─────────────────┐
│  DATA           │
│  (input)        │
└─────────────────┘
        │
        ▼
┌─────────────────┐
║  Data Mining    ║
└─────────────────┘
        │
        ▼
        ┌─────────────────┐
        │  INFORMATION    │
        │  (output)       │
        └─────────────────┘
```

DATA
(input)

Preprocessing → Feature Extraction

Intelligent
Algorithm

Data Mining

INFORMATION
(output)

**Descriptive Techniques**
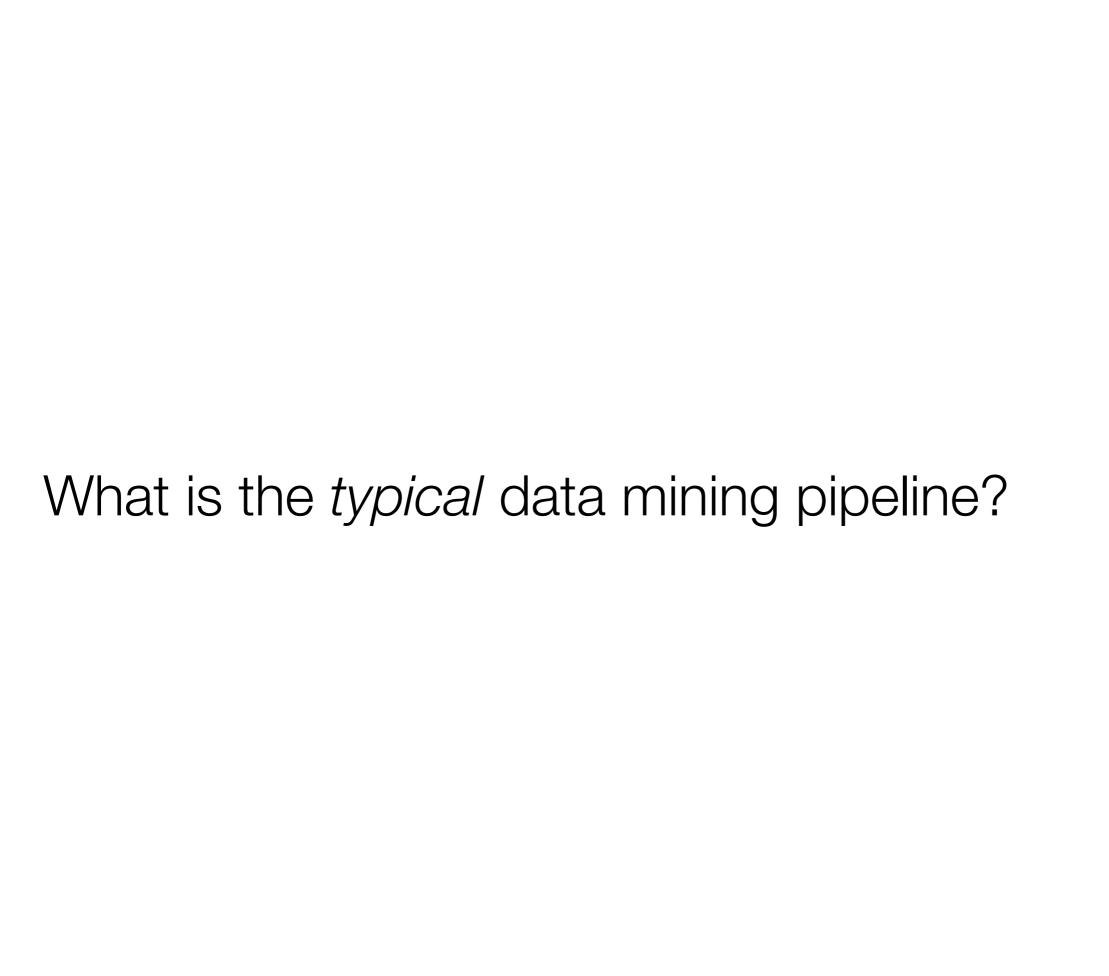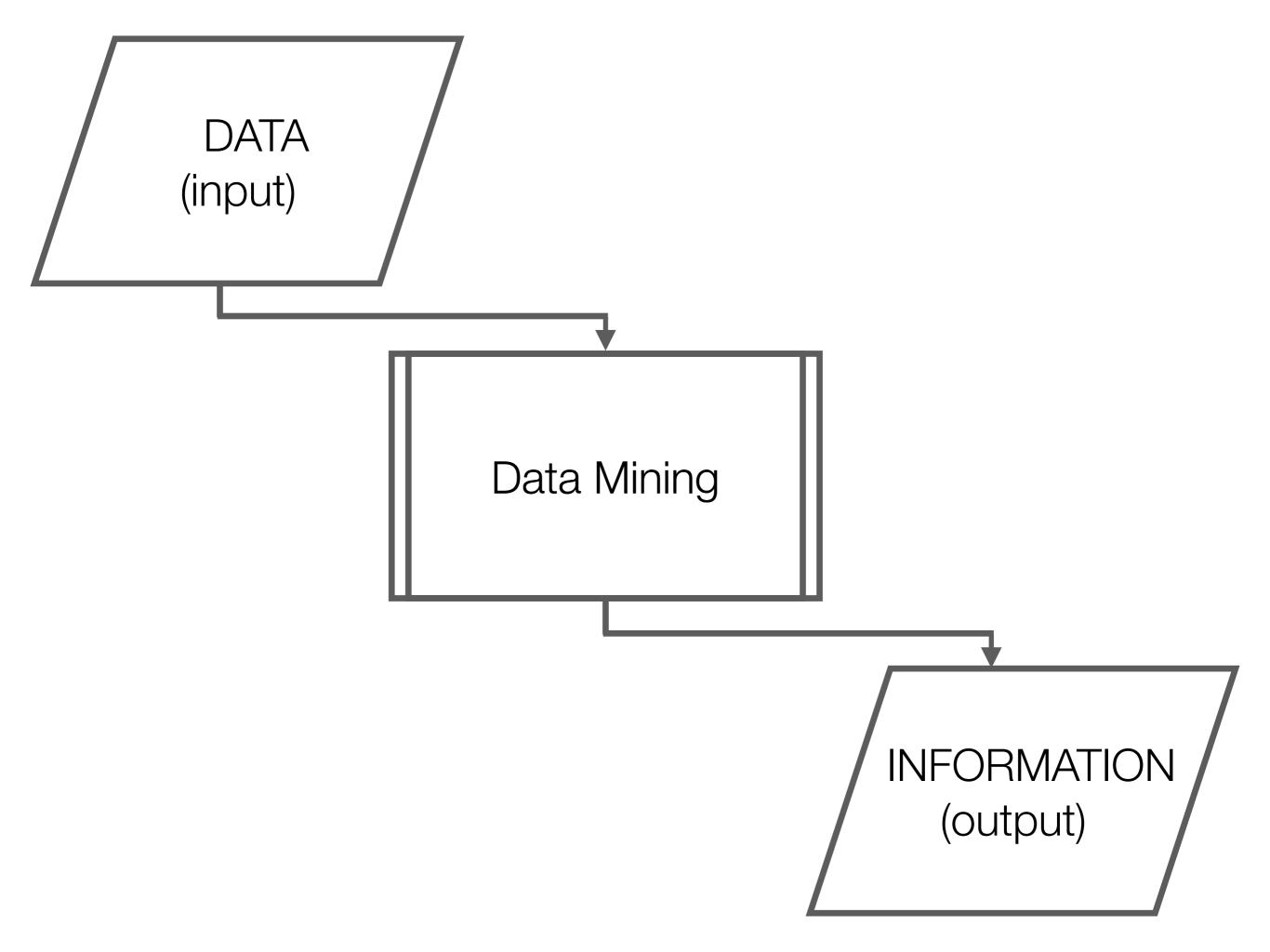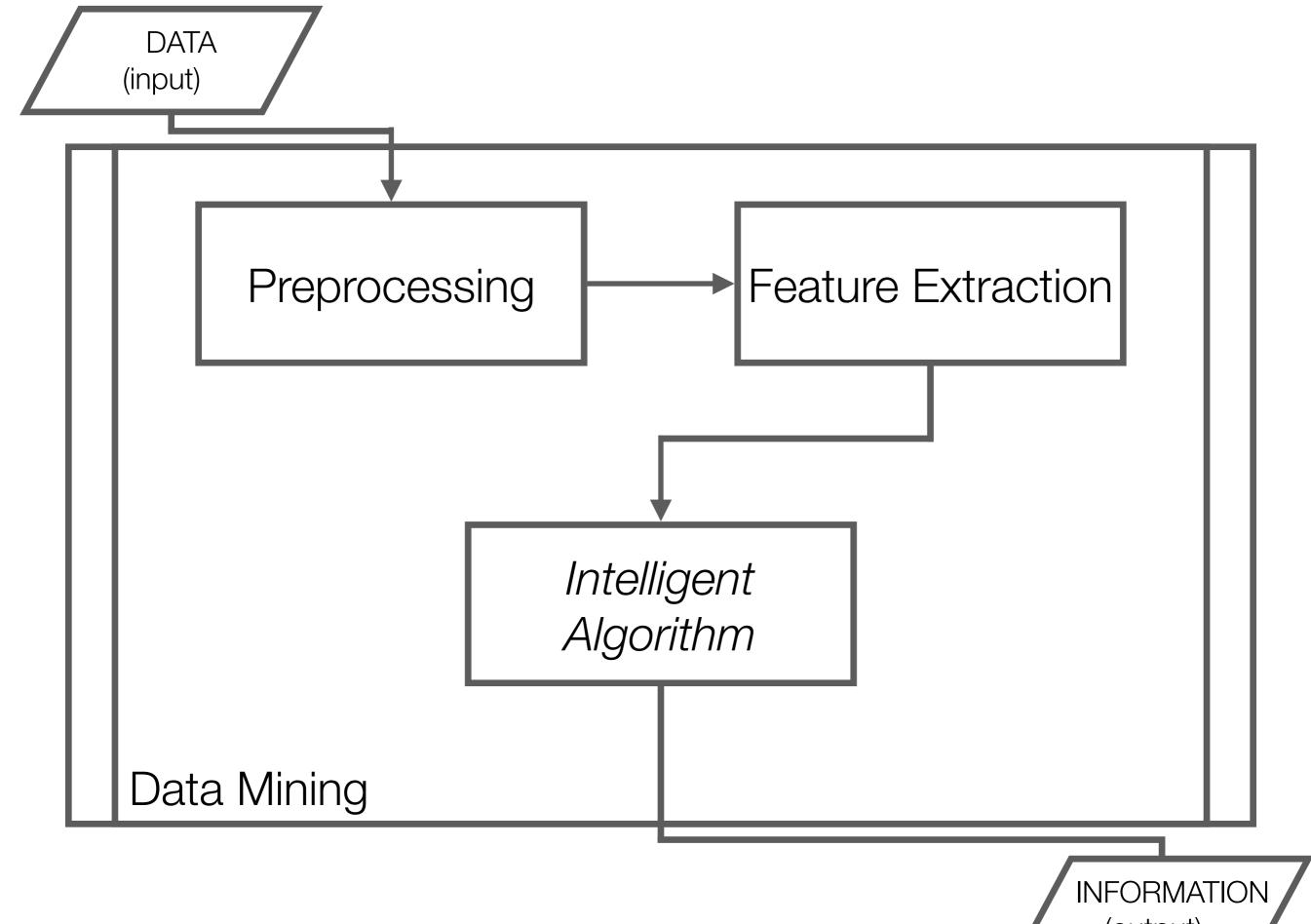
**Predictive Techniques**

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

*…*

*Intelligent Algorithm*

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

*…*

The plan for the next 12 weeks

# …we're going to look at a range of topics…

- **You will learn to solve real-world problems** - e.g.:

  - Recommender systems

  - Market Basket Analysis

  - Document filtering and spam detection

  - Duplicate document detection

- **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)

# The Group Coursework

- You need to form groups

  - Target size is 5-6 people

  - As a group you need to chose 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 the Friday afternoon slot this week to discuss your ideas for problems to work on with us

- Key dates:

  - Each team must submit a 1-page project brief by the end of the day on the 17th Feb.

  - On the 24th Feb (afternoon) each team must pitch their project to the class (2 minutes to pitch; 3 for Q&A)

  - Teams must submit a conference paper by 4PM on the 10th May

  - Projects will be presented in the lecture/tutorial slots in week 11 & 12

    - Teams should be prepared to present in the first slot; to ensure fairness we'll pick teams at random.