

# BSc - Data Mining 1

## Topic 01 : Module Overview

---

### Part 03 : Module Introduction

Dr Bernard Butler

Department of Computing and Mathematics, WIT.  
([bernard.butler@waltoninstitute.ie](mailto:bernard.butler@waltoninstitute.ie))

Autumn Semester, 2021

#### Outline

- Introduction, definitions and context
- Roles, expertise and ethics
- Workflow and process models
- Overview of Machine Learning Algorithms
- Delivery and Assessment
- Resources

# What is the AIM of the module?

## Aim, as per Module Descriptor\*...

The purpose of this module is to introduce the student to the fundamental concepts and techniques of Data Mining. The student will become familiar with Data Mining approaches (such as prediction, classification, clustering) and their typical solution techniques (methods and algorithms) to datasets...

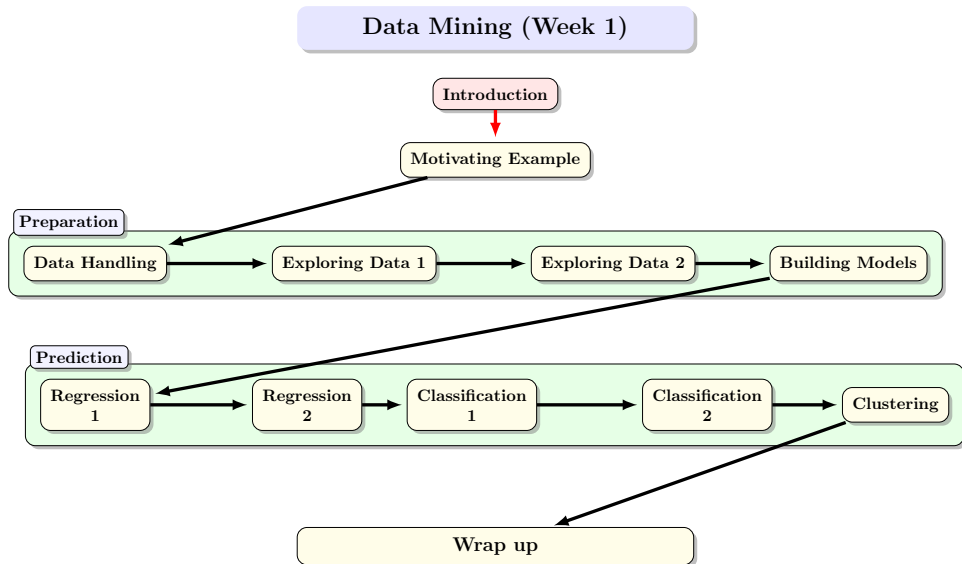
## Translation (Informal Aims)

- 1 Collect observations from a variety of processes, yielding large amounts of data.
- 2 Preprocess this data, selecting relevant features only.
- 3 Use data-intensive analysis techniques to obtain insights.
- 4 Postprocess analysis results, validate, visualise and refine the process.

---

\*Also, see the A11350 module descriptor for the learning outcomes for a more formal description of this module.

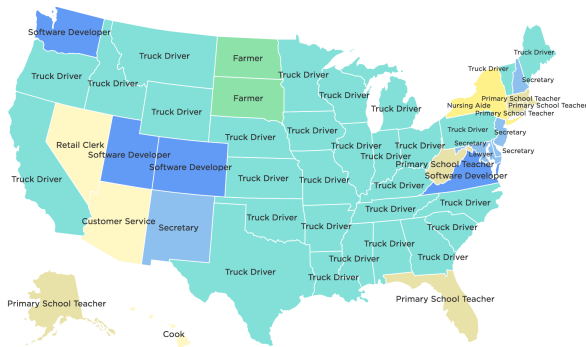
# What topics does it contain?



...

# Why are Data Mining (and automation) so important?

- Most common job by state in USA (2014)<sup>†</sup> ...



- By 2035 autonomous end-to-end delivery can be achieved.

- Current situation:

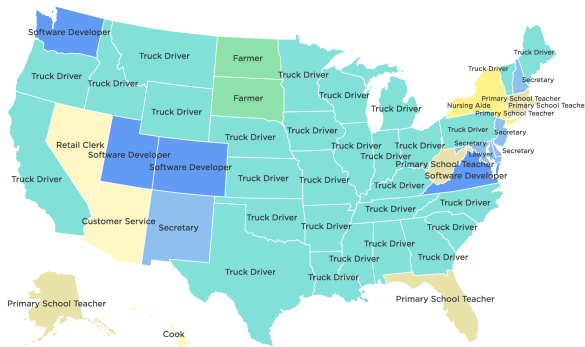
Any cognitive task that requires less than 2 seconds to perform can be automated in the short term.

See: Andrew Ng: Why 'Deep Learning' Is a Mandate for Humans, Not Just Machines

<sup>†</sup><https://www.npr.org/sections/money/2015/02/05/382664837/map-the-most-common-job-in-every-state>

# Why are Data Mining (and automation) so important?

- Most common job by state in USA (2014)<sup>†</sup> ...
- By 2035 autonomous end-to-end delivery can be achieved.
- Current situation:



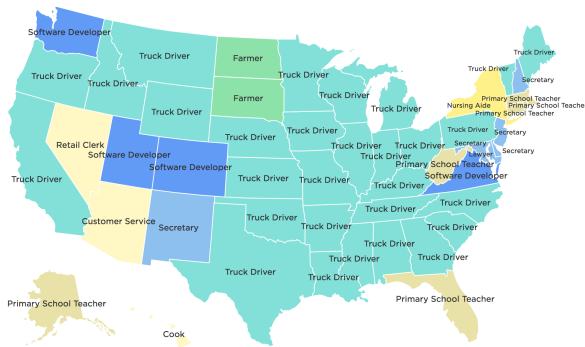
Any cognitive task that requires less than 2 seconds to perform can be automated in the short term.

See: Andrew Ng: Why 'Deep Learning' Is a Mandate for Humans, Not Just Machines

<sup>†</sup><https://www.npr.org/sections/money/2015/02/05/382664837/map-the-most-common-job-in-every-state>

# Why are Data Mining (and automation) so important?

- Most common job by state in USA (2014)<sup>†</sup> ...
- By 2035 autonomous end-to-end delivery can be achieved.
- Current situation:



Any cognitive task that requires less than 2 seconds to perform can be automated in the short term.

*See: Andrew Ng: Why 'Deep Learning' Is a Mandate for Humans, Not Just Machines*

<sup>†</sup><https://www.npr.org/sections/money/2015/02/05/382664837/map-the-most-common-job-in-every-state>

## Selected definitions

---

**IoT:** A network of pervasive connected objects able to collect and exchange data from embedded sensors, with the infrastructure and services to support them. (Various)

**Big Data:** High volume, high velocity, and/or high variety information assets that require new forms of processing to enable enhanced decision making, insight discovery and process optimization. (Gartner 2012)

**Data Scientist:** can ask the right questions, {generate} and consume the results of analysis of Big Data effectively. (McKinsey 2011)

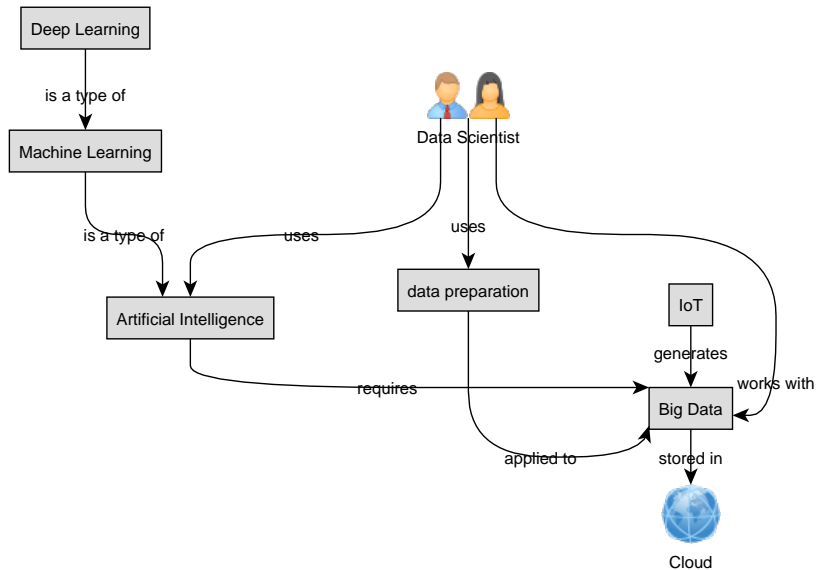
**Artificial Intelligence:** the capability of a machine to imitate intelligent human behavior (Webster 2017)

**Machine Learning:** Branch of computer science {and related fields} that gives computers the ability to learn without being explicitly programmed. (Samuel 1959)

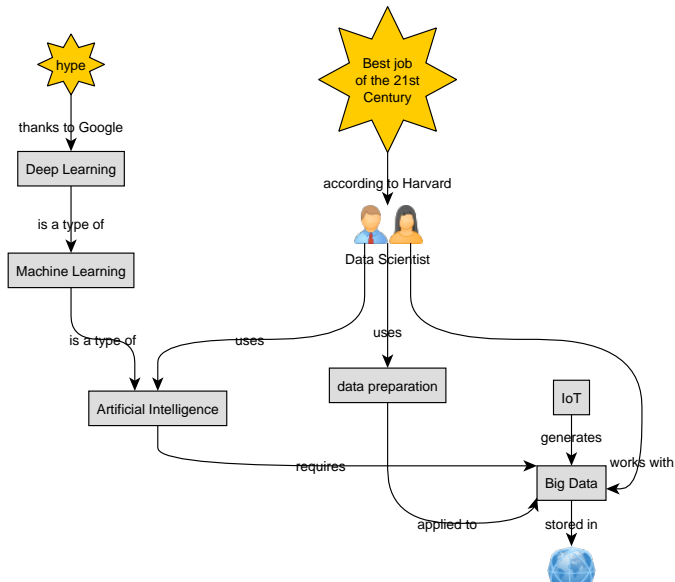
**Deep Learning:** Use of very large neural networks with many layers of “neurons” that can be trained to generate robust models of their input, whose classification performance scales with the amount of data supplied. (Various)



# Relationships between terms



# Annotated Relationships between terms!



# What is data mining and how does it relate to similar terms?

## Operational Definitions

- deriving knowledge from large and/or complex datasets, with *guidance* from the data scientist
- “Data mining is the study of efficiently finding structures and patterns in large data sets. It draws from and influences the disciplines of programming, mathematics/statistics, database management and machine learning.”

## Primary goals

- From messy and noisy raw data, deriving structure and context
- Applying scalable learning algorithms to these higher value data sets

## Secondary goals

- Modelling and understanding the error and other consequences of the modelling process.
- Building data-driven processes, architectures & frameworks: *Big Data*

# Interlude: Examples of Big Data

## Exercise

please consider (real world) processes generating *Big Data*.  
Can you come up with 3 examples in 3 minutes?

# Prehistory, or before 2007...

## Data Generation

- Transactions (bank, retail)
- Activity, e.g., texts
- Basic e-commerce

## Data Processing

- Databases, SQL, stored procedures
- Consultants, system integrators
- Proprietary statistical software

## Data Analysis

- Reporting: looking back
- Descriptive statistics
- Simple plots

# The first (batch) wave: 2007–2011

## Data Generation

- As before...
- Web activity: comments, etc.
- 360degree view

## Data Processing

- As before...
- NoSQL
- hadoop ecosystem (batch analytics)

## Data Analysis

- As before...
- Personalisation and recommendation
- Predictive Analytics

# The second (streaming) wave: 2012–2015

## Data Generation

- As before...
- Social Media!
- IoT (early adopters)

## Data Processing

- As before...
- Apache Spark
- R vs. python

## Data Analysis

- As before...
- Data understanding
- Weak AI: assistants, etc.

# The current (machine) wave: 2016–?

## Data Generation

- As before...
- Machine-generated (e.g., fake news)
- IoT (mainstream)

## Data Processing

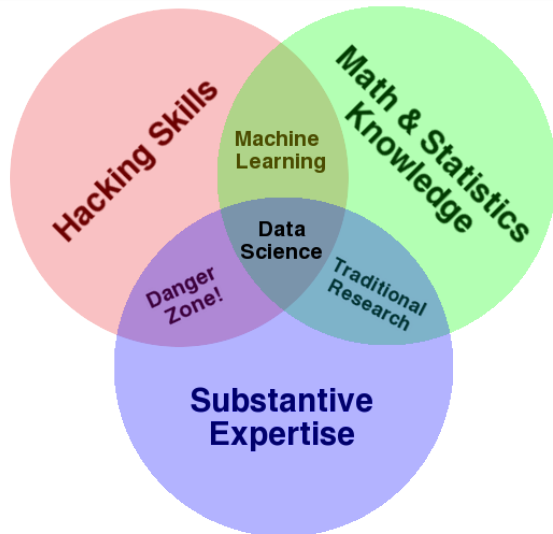
- As before...
- Microservices: move function to data
- Decoupled databases with schema-on-read

## Data Analysis

- As before...
- Deep learning inflection point
- Visualisation

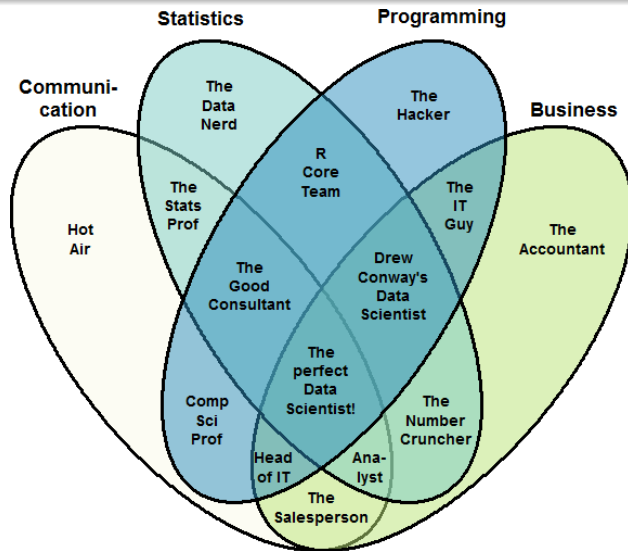


# Drew Conway's 3-set Venn Diagram of Data Science Expertise



*Source:* <http://drewconway.com/zia/2013/3/26/the-data-science-venn-diagram>

# Stephan Kolassa's 4-set Venn Diagram of Data Science Expertise

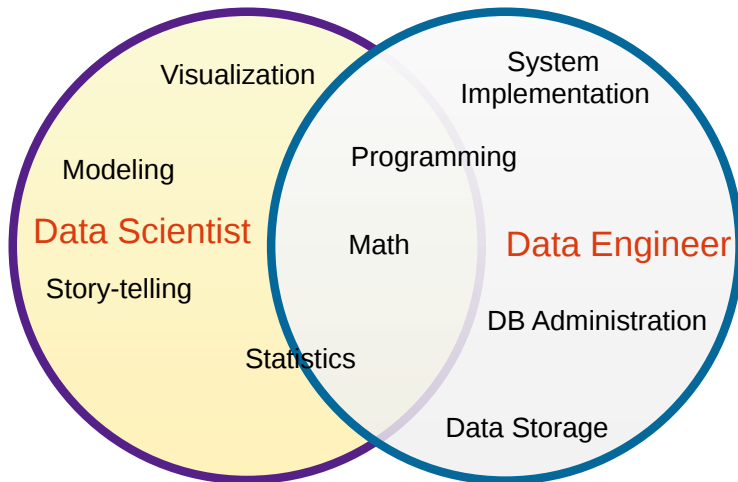


*Source:* <https://datascience.stackexchange.com/a/2406>

# Gartner suggests the need for a *Citizen Data Scientist*



# Data Scientist vs Data Engineer



*Source:* <https://ryanswanstrom.com/2014/07/08/data-scientist-vs-data-engineer/>  
Also the traditional roles of *Data Analyst* and *Software Engineer*...

# Complete the following disadvantages of IoT and Big Data

m\_\_\_\_ s\_\_\_\_v\_\_\_\_l\_\_\_\_\_

i\_\_\_\_t\_\_\_\_y \_h\_f\_

d\_\_\_\_c\_ b\_\_\_\_n\_\_\_\_

d\_\_\_\_\_l \_f \_r\_\_\_\_c\_

b\_\_\_\_s

l\_\_\_\_ o\_ t\_\_\_\_s\_\_\_\_r\_\_\_\_y

## And those disadvantages are...

---

mass surveillance

identity theft

device botnets

denial of service

bias

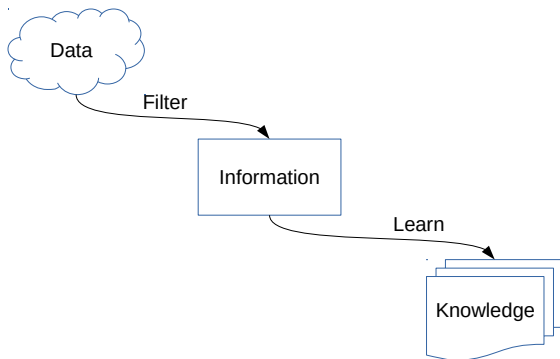
lack of transparency

## Ethical Concerns

---

- protecting privacy (informed consent; undoing pseudonymisation)
- ensuring transparency of decisions (how was the decision made?)
- breaking cycles of bias (biased data leads to biased results)
- enabling validation (ensuring correct usage of techniques)
- enabling decisions to be challenged (openness and due process)

# The Data to Knowledge Pipeline



## Data Filtering

- Clean (drop unwanted observations)
- Summarise (remove observation detail)
- Reduce (remove/transform variables)

## Learning

- Derive models
- Validate models
- Analyse discordance



# Data - Information - Knowledge - Wisdom

## Example of the DIKW chain

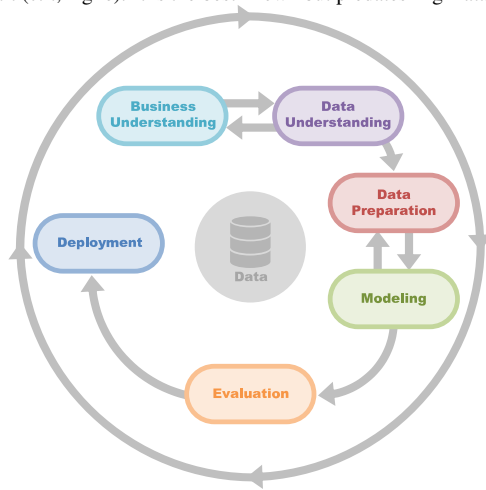
- Servers and applications log events in files and/or databases [DATA]
- *Collector* agents select specific events, in context [INFORMATION]
- Machine learning *classifiers* learn system behaviour and identify anomalies [KNOWLEDGE]
- Humans and software use this knowledge to prevent future problems [WISDOM]

Note that the DIKW chain is often represented as a pyramid.

# Cross Industry Standard Process (for) Data Mining

## CRISP-DM

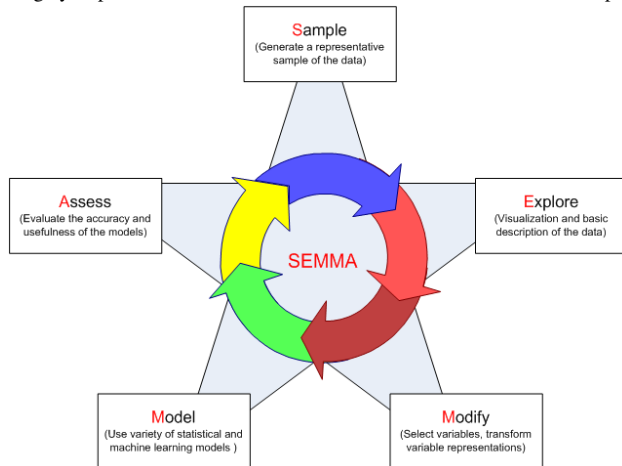
CRISP-DM is a high-level iterative process model. It gives much weight to data understanding and preprocessing and involves the data and problem owners from the start (c.f., Agile). It is the best known but predates Big Data, etc.



# Sample, Explore, Model, Modify, Assess

## SEMMA

SEMMA is promoted by SAS and takes a more operational view of data mining, using a (statistical) *model-building* metaphor. Business input is essential but largely implicit. It is more concrete than CRISP-DM so it tends to map well to DM tool workflows.

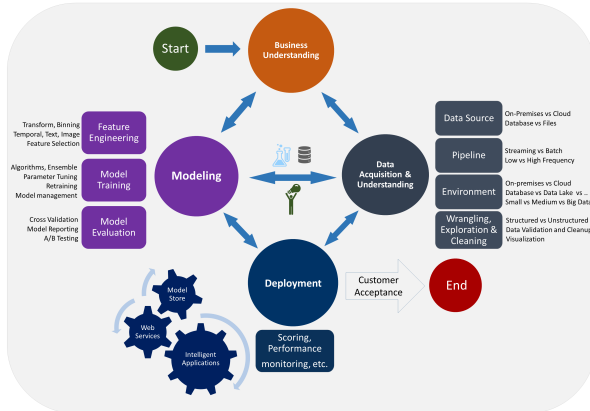


# Microsoft Team Data Science Process

## TDSP

TDSP is the most detailed process model of the 3. It is much more recent. It is cloud-aware and directly references Azure and other Microsoft technologies. Typically there are two main cycles, one involving the Business, the other involving Deployment. Interestingly, there is a Start and End, so it is more project-focused.

### Data Science Lifecycle



# The “5 Tribes of Machine Learning”

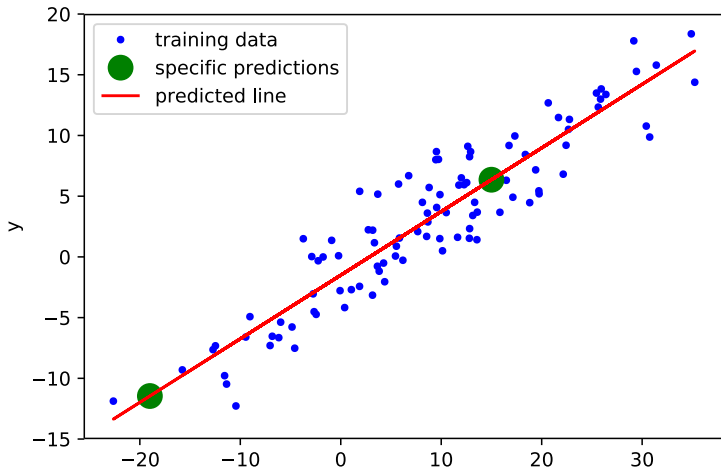
Tribe	Origins	Learning Algorithm
Symbolists	Logic, Philosophy	Inverse Deduction
Connectionists	Neuroscience	Backpropagation
Evolutionaries	Mathematical Biology	Genetic Programming
Bayesians	Statistics	Probabilistic Inference
Analogizers	Psychology	Kernel Machines

*Summarised from Domingos (2015) “The Master Algorithm”*

# Regression

## Definition

Given data comprising a set of independent variables (of any type  $\mathbf{x}$ ) with a set of dependent variables (numeric only  $\mathbf{y}$ ), find the relationship  $\mathbf{y} = f(\mathbf{x})$  having the maximum likelihood given the available observations  $\{\mathbf{x}_i, \mathbf{y}_i\}$ .

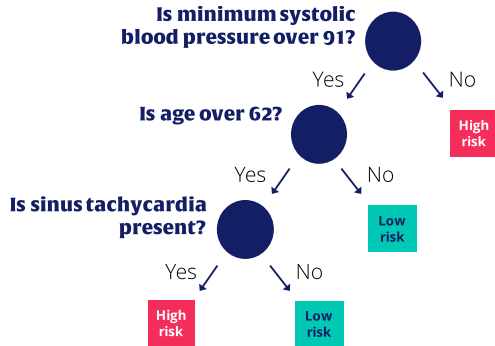


# Classification

## Definition

Given data comprising a set of independent variables (of any type  $\mathbf{x}$ ) with a set of dependent variables (categorical  $\mathbf{y}$  (labels)), find the relationship  $\mathbf{y} = f(\mathbf{x})$  having the maximum likelihood given the available observations  $\{\mathbf{x}_i, \mathbf{y}_i\}$ .

There are many ways of representing  $f$ : a classification tree is shown here.

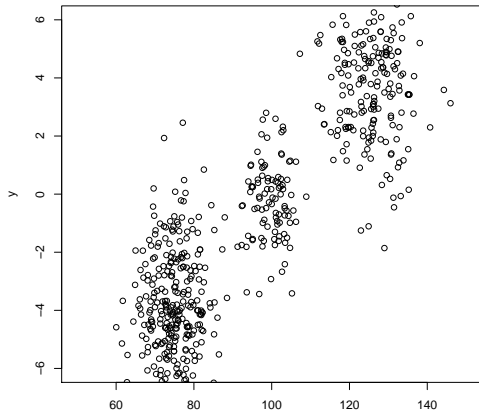


# Clustering

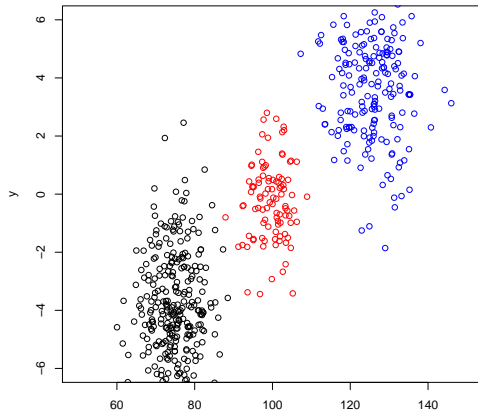
## Definition

Clustering is the process of grouping data into classes or clusters, so that objects within a cluster have high similarity with each other but are dissimilar to objects in other clusters. Different similarity measures and/or algorithms result in different cluster arrangements.

Single cluster



Three clusters

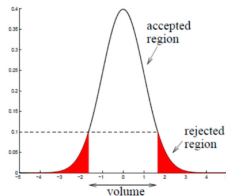




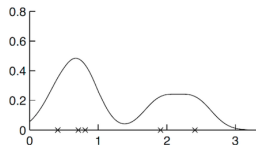
# Anomaly Detection

## Definition

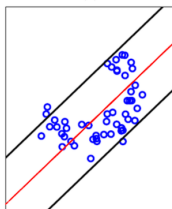
Anomaly detection identifies data points, events, and/or observations that depart from a dataset's normal behavior. Anomalous data can indicate problems, such as fraud, or opportunities, like a surge in demand for a product.



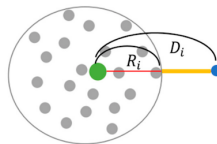
(a)



(b)



(c)

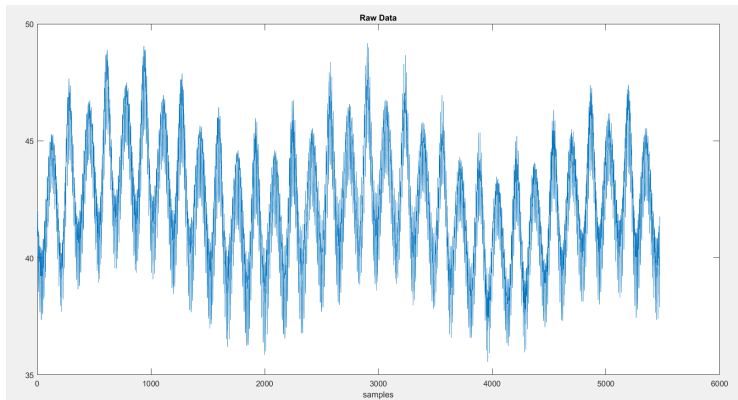


(d)

# Time Series Analysis

## Definition

Time series data is a sequence of observations on the values that a variable taken at regularly-spaced time intervals. This data is sequentially correlated and techniques are needed to determine seasonality, trends, anomalies etc.



Source: <https://stats.stackexchange.com/q/458491>

# Association Rules Mining

## Definition

*Frequent itemset mining* looks for associations and correlations among items in large data sets. Associations are expressed as rules and quantified in terms of their *support* and *confidence*. The classical example is market basket analysis and the famous rule about buying diapers and beer together. See example transaction data below

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

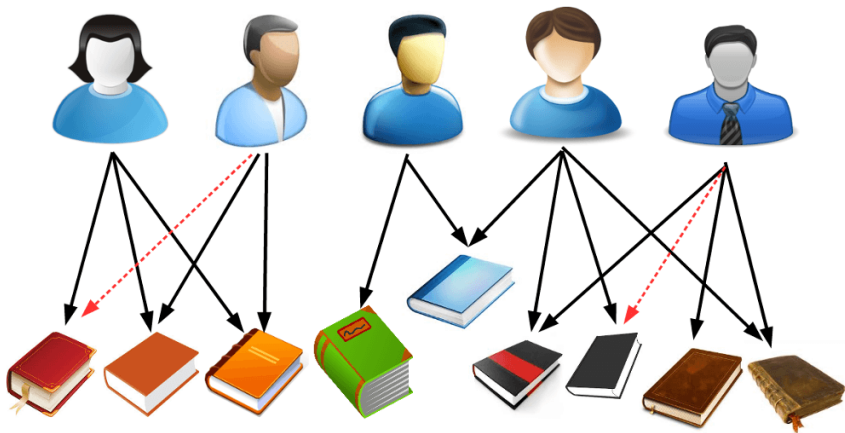


	Beer	Bread	Milk	Diaper	Eggs	Coke
$T_1$	0	1	1	0	0	0
$T_2$	1	1	0	1	1	0
$T_3$	1	0	1	1	0	1
$T_4$	1	1	1	1	0	0
$T_5$	0	1	1	1	0	1

# Recommender Systems

## Definition

Collaborative filtering is an extension of item-based association rules mining to consider relationships between users and items. Generally, if two users have similar behaviour and/or preferences, items favoured by one user will also be favoured by the other. This can be used to recommend “new” items to users. Used very commonly on ecommerce websites for cross-selling.



# How? — Delivery

## Contact hours

- One 2-hour lecture per week.

—Friday 13:15 (GMT/IST; Room E15)

- We cover concepts, definitions, examples, etc.
- Lecture objective is to improve understanding of the topic.
- Practical/Lab is used to develop practical experience of an integrated set of topics
- Feel free to ask questions.
- One 2-hour practical/lab session per week on Tuesday afternoons (Room D05).
  - 13:15-15:00 - W2 Group
  - 15:15-17:00 - W1 Group
  - Labs use python workbooks to define and implement data mining *workflows*.

## (Hardware) Requirements

- Use of own (moderately powerful: multi-core CPU, min 8GB RAM) laptop is recommended!

# How? — Assessment Structure

## 50% Continuous Assessment; 50% Written Exam

- 45% — two data investigations
  - Issued Week 3, submitted end Week 7 (Data Investigation, 20%)
  - Issued Week 7, submitted end Week 12 (Data Investigation, 25%)
- 5% — Engagement metric, based on class activity and moodle quizzes.

These are indicative and might change (in terms of when assignments are issued and their relative weightings)...

# Resources



- URL: [Moodle: Data Mining I-89898-\[2021-2022\]](#)
- Used for all notices, assignment briefs and practical work submissions.



- URL: [Data Mining 1 \(2021 Semester 1\) pages on github.io](#)
- Used for content delivery (lecture notes and labs).

## Software

All software used during this module is open source or freely available for non-commercial use (full details given in notes). Primarily

- Anaconda (**Python 3.8 and later, 64 bit**)
- scikit-learn
- pandas

[www.anaconda.com](http://www.anaconda.com)

[scikit-learn.org](http://scikit-learn.org)

[pandas.pydata.org](http://pandas.pydata.org)

# Further Reading

Please note that the notes and labs we provide should be sufficient to pass to pass this module, so the books below are intended as *further reading*, not *recommended reading*.

## **Data Mining, Concepts and Techniques**

by *Jiawei Han, Michelline Kamber and Jian Pei* (Available in the library)

Broad selection of topics, looks at the entire data mining process including how to collect and preprocess data, discusses selected algorithms in depth.

## **Mining of Massive Data Sets**

by *Jure Leskovec, Anand Rajaraman and Jeffrey David Ullman* (Available as PDF)

Good mix of mathematical rigour and a treatment of the *Big Data* aspects for data mining at scale.

## **Python for Everybody: Exploring Data using Python 3**

by *Charles R Severance* (Available as PDF)

Useful summary of basic python concepts and a good introduction to the type of data manipulation can can be performed using core python data structures and idioms.

## **Python Data Science Handbook**

by *Jake vanderPlas* (Available from website), is both a textbook and a set of freely-available Jupyter notebooks that go into more detail on implementing some of the material in this module using pandas etc.