

Dr Bernard Butler

Department of Computing and Mathematics, WIT.

([bernard.butler@setu.ie](mailto:bernard.butler@setu.ie))

Preparation

Data Handling

Exploring Data 1

Exploring Data 2

Building Models

Autumn Semester, 2025

## Outline

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

Wrap up

# Data Mining (Week 1)

Introduction



Motivating Example

Preparation

Data Handling

Exploring Data 1

Exploring Data 2

Building Models

Prediction

Clustering

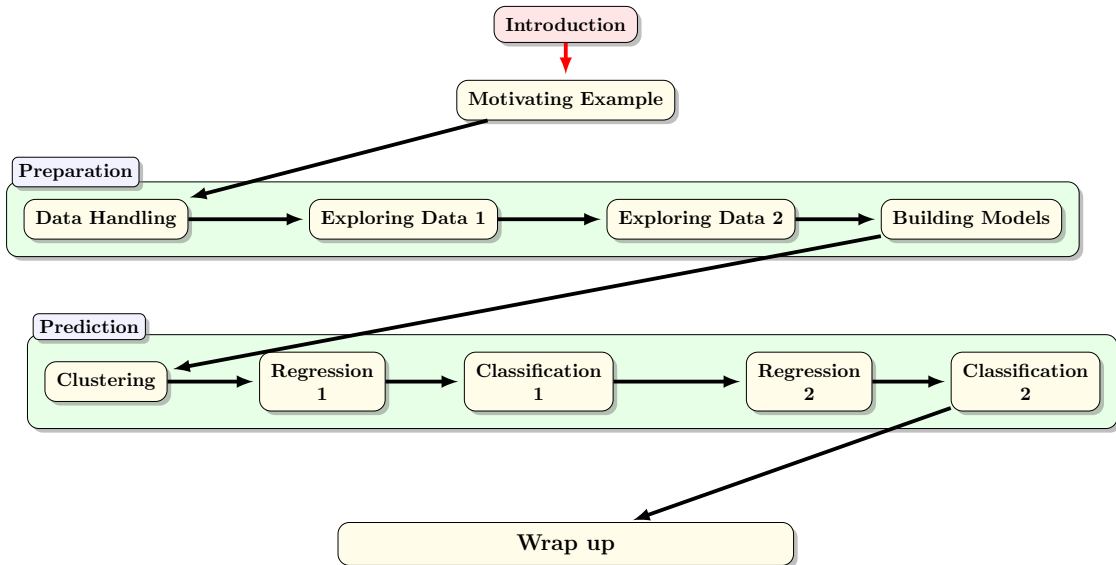
Regression  
1

Classification  
1

Regression  
2

Classification  
2

Wrap up



# Outline

---

1. Module Delivery	3
2. Module Introduction	6
3. Definitions of terms and their relationships	10
3.1. Exercise: Big Data Examples	13
4. The Context and History	14
5. Expertise and Roles	23
6. Is the IoT and Big Data always a force for good?	28
7. The DIKW pipeline	33
8. Data Mining Process Models	36
9. Overview of Machine Learning techniques	40
10. Resources for students	53

# How? — Delivery

## Contact hours - but subject to change!

- 2 hours of lectures per week (all students).

—Tuesday 12:15-13:15 (Room W13)

—Friday 14:15–15:15 (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.
- A 2-hour practical/lab session per week for each Group (Room D04).
  - Tuesday 13:15-15:15 - Group W1 (D04)
  - Friday 12:15-14:15 - Group W2 (D04)
  - Labs use python workbooks to define and implement data mining *workflows*.

## (Hardware) Requirements

- Use of own laptop (minimum of quad-core CPU and 8GB RAM is recommended)

# How? — Assessment Structure

## 100% Continuous Assessment

- 80% — two data investigations
  - Issued Week 3, submitted end Week 8 (Data Investigation, 40%)
  - Issued Week 8, submitted end Week 13 (Data Investigation, 40%)
- 20% — Labs and moodle quizzes

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

# Outline

---

1. Module Delivery	3
<b>2. Module Introduction</b>	<b>6</b>
3. Definitions of terms and their relationships	10
3.1. Exercise: Big Data Examples	13
4. The Context and History	14
5. Expertise and Roles	23
6. Is the IoT and Big Data always a force for good?	28
7. The DIKW pipeline	33
8. Data Mining Process Models	36
9. Overview of Machine Learning techniques	40
10. Resources for students	53

# 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. . .

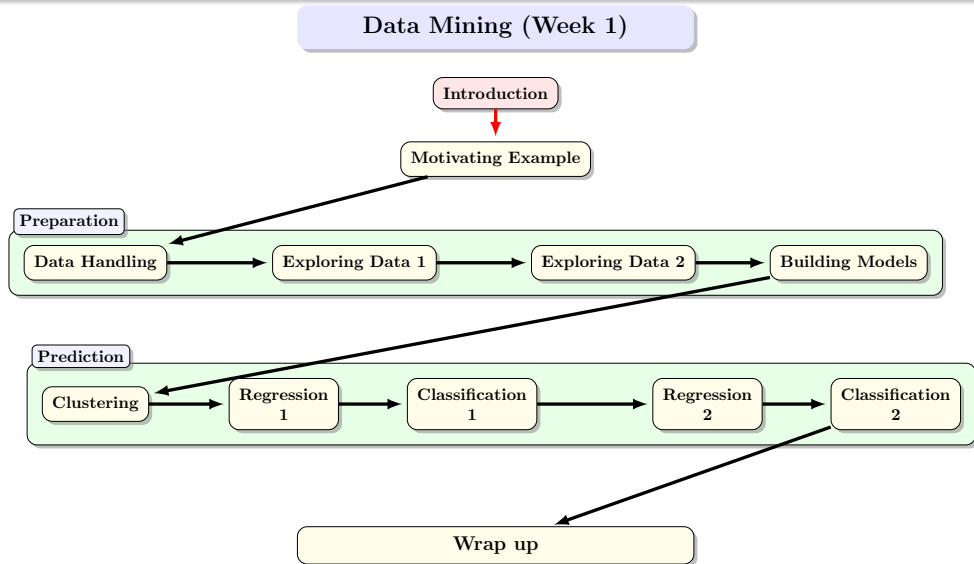
## Translation (Informal Aims)

- 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?





...

# Outline

---

1. Module Delivery	3
2. Module Introduction	6
<b>3. Definitions of terms and their relationships</b>	<b>10</b>
<b>3.1. Exercise: Big Data Examples</b>	<b>13</b>
4. The Context and History	14
5. Expertise and Roles	23
6. Is the IoT and Big Data always a force for good?	28
7. The DIKW pipeline	33
8. Data Mining Process Models	36
9. Overview of Machine Learning techniques	40
10. Resources for students	53

# Selected definitions

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

(McKinsey 2011)

**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)

**Generative AI:** describes algorithms that can be used to create (remix?) new content, including audio, code, images, text, simulations, and videos, (based on its training data) (McKinsey 2024)

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

**Artificial Narrow Intelligence (ANI):** also known as weak AI, (exists today) can be trained to perform a single or narrow task, often far faster and better than a human mind can (IBM 2025)

**Artificial General Intelligence (AGI):** also known as strong AI, is the (currently hypothetical) intelligence of a machine that can accomplish any intellectual task that a human can perform (Gartner 2024)

**Artificial Super Intelligence (ASI):** If ever realized, Super AI would think, reason, learn, make judgements and possess cognitive abilities that surpass those of human beings (IBM 2025)

# 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?

# Outline

---

1. Module Delivery	3
2. Module Introduction	6
3. Definitions of terms and their relationships	10
3.1. Exercise: Big Data Examples	13
<b>4. The Context and History</b>	<b>14</b>
5. Expertise and Roles	23
6. Is the IoT and Big Data always a force for good?	28
7. The DIKW pipeline	33
8. Data Mining Process Models	36
9. Overview of Machine Learning techniques	40
10. Resources for students	53

# 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: 2012–2016

## 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 third (machine) wave: 2017–2021

## 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

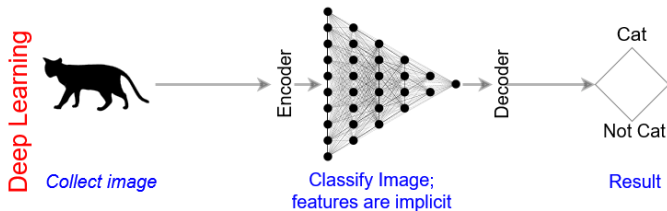
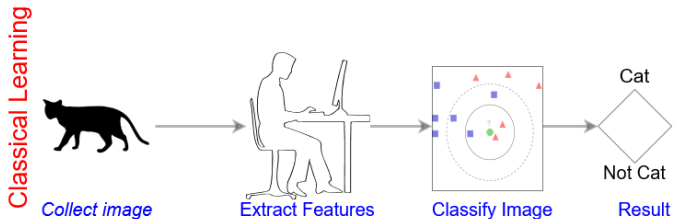
## Data Analysis

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

# The current (generative) wave: 2022–date

- Traditional ML (big data processing for predicting labels or numbers) becomes mainstream
  - ML skills become expected
  - ML is operationalised by incorporating it in other practices ( ML + Devops = MLops )
- *Generative AI* becomes practical
  - Large Language Models for generating textual responses to prompts, e.g., ChatGPT
  - Generate images and other rich media given textual descriptions, e.g., DALL-E
  - Beyond LLMs and next token prediction: RAG, Large Reasoning Models, AI Agents, ...
- AI is beginning to deliver on the promise identified by John McCarthy and Alan Turing and others in the 1950s.

# Classical versus Deep Learning - overview



Classical learning requires extensive setup before training.

# Classical versus Deep Learning - pros and cons

Classical Learning	Deep Learning
Can work well with less data (<10 <sup>6</sup> rows, say)	More training data gives more accuracy
Easier to interpret/explain models	Model is opaque and can be fragile
Training is relatively fast	Training can require many epochs
Training requires fewer resources	Training requires massive resources
Accuracy improvement falls off	Accuracy can improve with more training data
Requires feature engineering (by human)	Features are encoded implicitly in layers
Complex prediction requires complex model	With enough nodes can represent any function

➤ In conditions where one type of learning is weak, the other is often strong. ➤

# Emphasis of this module

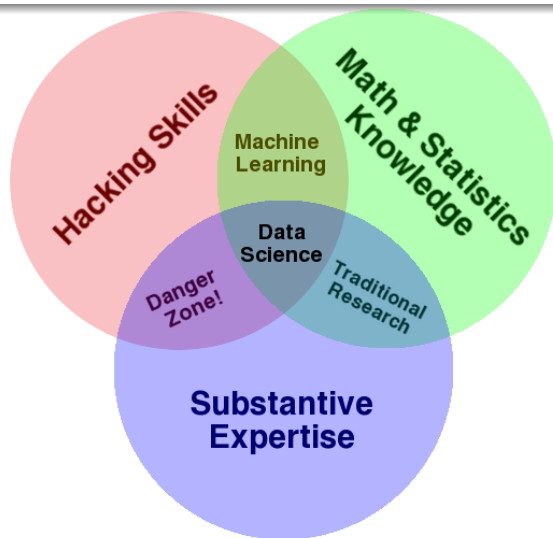
- This module covers foundations, classical models, some deep learning
- Foundations include EDA, notions of error and variance, training vs test data, ...
- Classical modelling includes feature selection and classical approaches to learning from data
- Neural network (see Data Mining 2) and Deep learning go straight to predicting based on (encoded) data
- Foundations are shared by both classical and deep learning
- **Deep learning is ideal for learning from labeled, web-scale big data.**
- Day-to-day, classical machine learning is often more suitable.

# Outline

---

1. Module Delivery	3
2. Module Introduction	6
3. Definitions of terms and their relationships	10
3.1. Exercise: Big Data Examples	13
4. The Context and History	14
<b>5. Expertise and Roles</b>	<b>23</b>
6. Is the IoT and Big Data always a force for good?	28
7. The DIKW pipeline	33
8. Data Mining Process Models	36
9. Overview of Machine Learning techniques	40
10. Resources for students	53

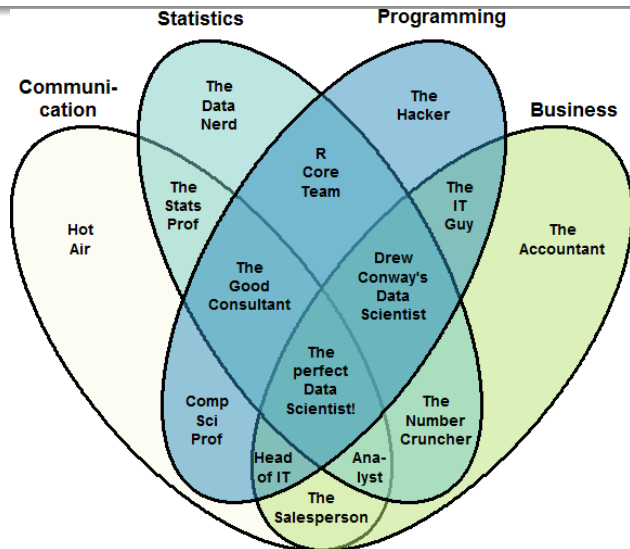
# Drew Conway's 3-set Venn Diagram of Data Science Expertise



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



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



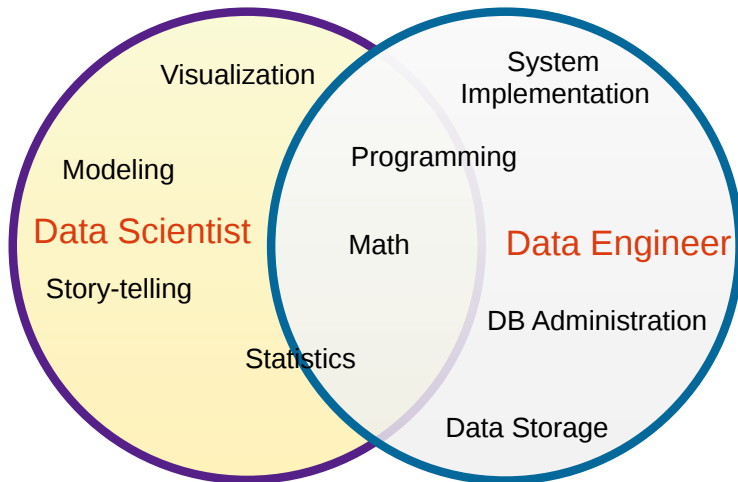
*Source:* <https://datascience.stackexchange.com/a/2406>

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



Source: <http://www.kdnuggets.com/2016/03/cartoon-citizen-data-scientist.html>

# 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*...

# Outline

---

1. Module Delivery	3
2. Module Introduction	6
3. Definitions of terms and their relationships	10
3.1. Exercise: Big Data Examples	13
4. The Context and History	14
5. Expertise and Roles	23
<b>6. Is the IoT and Big Data always a force for good?</b>	<b>28</b>
7. The DIKW pipeline	33
8. Data Mining Process Models	36
9. Overview of Machine Learning techniques	40
10. Resources for students	53

# 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)

# Ethical frameworks and legislation

## UNESCO human rights approach

- ① Proportionality and Do No Harm
- ② Safety and Security
- ③ Right to Privacy and Data Protection
- ④ Multi-stakeholder and Adaptive Governance & Collaboration
- ⑤ Responsibility and Accountability
- ⑥ Transparency and Explainability
- ⑦ Human Oversight and Determination
- ⑧ Sustainability (c.f., SDGs)
- ⑨ Awareness & Literacy
- ⑩ Fairness and Non-Discrimination

## EU AI Act 2025

- Rules for each risk class
- Need for transparency
- Encouraging innovation

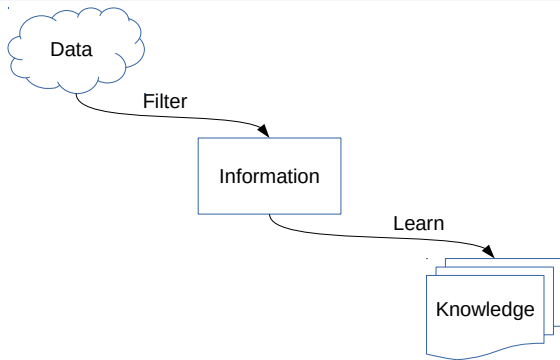


# Outline

---

1. Module Delivery	3
2. Module Introduction	6
3. Definitions of terms and their relationships	10
3.1. Exercise: Big Data Examples	13
4. The Context and History	14
5. Expertise and Roles	23
6. Is the IoT and Big Data always a force for good?	28
<b>7. The DIKW pipeline</b>	<b>33</b>
8. Data Mining Process Models	36
9. Overview of Machine Learning techniques	40
10. Resources for students	53

# 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.

# Outline

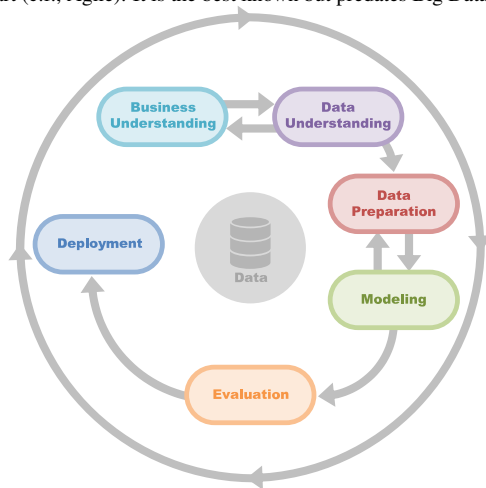
---

1. Module Delivery	3
2. Module Introduction	6
3. Definitions of terms and their relationships	10
3.1. Exercise: Big Data Examples	13
4. The Context and History	14
5. Expertise and Roles	23
6. Is the IoT and Big Data always a force for good?	28
7. The DIKW pipeline	33
<b>8. Data Mining Process Models</b>	<b>36</b>
9. Overview of Machine Learning techniques	40
10. Resources for students	53

# Cross Industry Standard Process (for) Data Mining

## CRISP-DM

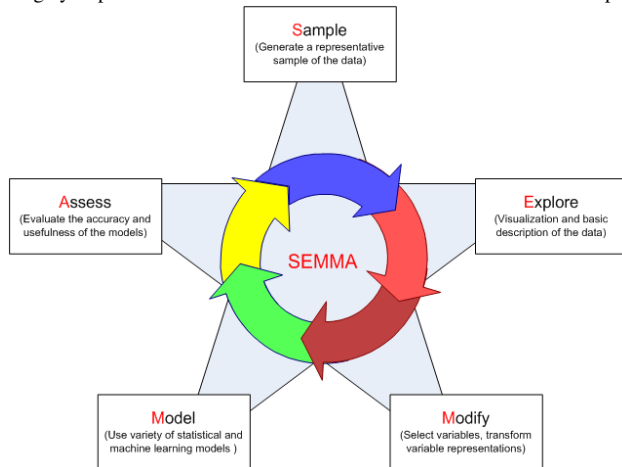
CRISP-DM is a high-level iterative process model. It gives much weight to data understanding and preprocessing and involves the data and problem owners from the start (c.f., Agile). It is the best known but predates Big Data, etc.



# 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.

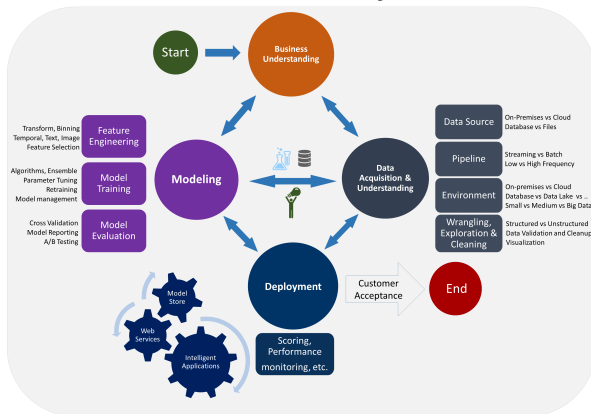


# Microsoft Team Data Science Process

## 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



# Outline

---

1. Module Delivery	3
2. Module Introduction	6
3. Definitions of terms and their relationships	10
3.1. Exercise: Big Data Examples	13
4. The Context and History	14
5. Expertise and Roles	23
6. Is the IoT and Big Data always a force for good?	28
7. The DIKW pipeline	33
8. Data Mining Process Models	36
<b>9. Overview of Machine Learning techniques</b>	<b>40</b>
10. Resources for students	53



# The “5 Tribes of Machine Learning”

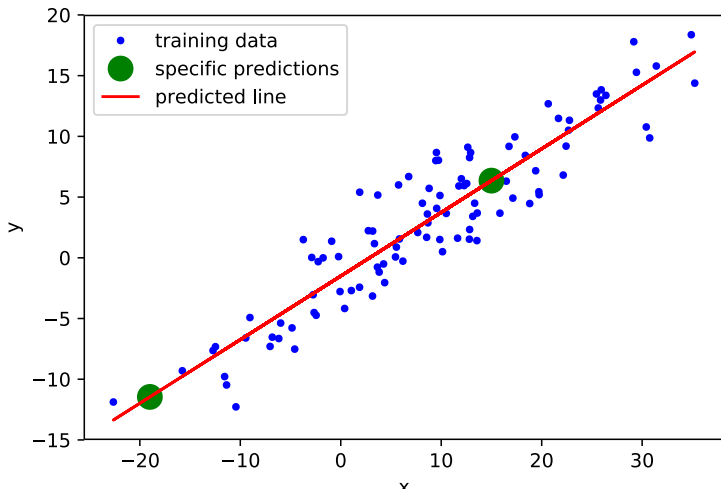
Tribe	Origins	Learning Algorithm
Symbolists	Logic, Philosophy	Inverse Deduction
Connectionists	Neuroscience	Backpropagation
Evolutionaries	Mathematical Biology	Genetic Programming
Bayesians	Statistics	Probabilistic Inference
Analogizers	Psychology	Kernel Machines

*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  $\{\mathbf{x}_i, \mathbf{y}_i\}$ .

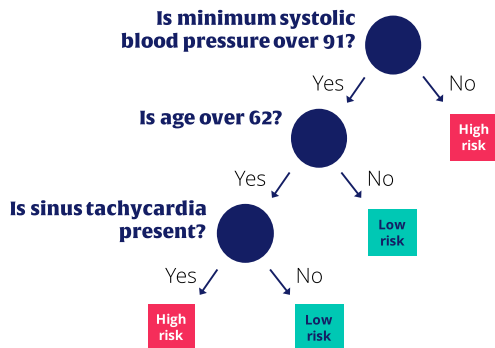


# Classification

## 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.

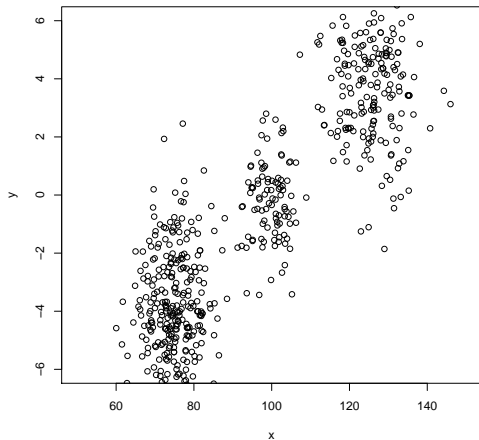


# Clustering

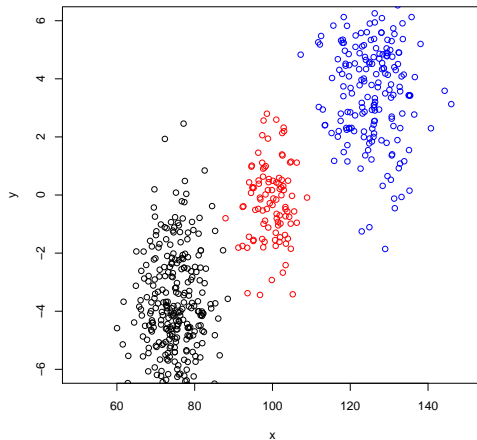
## Definition

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



## Clustering Example 2: Image Analysis

### Le et al (2012)

- Google Brain simulator crawled the web, looking for patterns in photographs on the web.
- *It was **not** looking for anything in particular!*
- One pattern came up strongly...

Supervised learning question: what does this represent?

## Clustering Example 2: Image Analysis

### Le et al (2012)

- Google Brain simulator crawled the web, looking for patterns in photographs on the web.
- It was *not* looking for anything in particular!
- One pattern came up strongly...

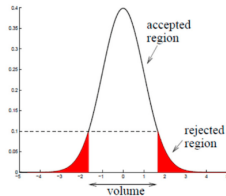
Supervised learning question: what does this represent?



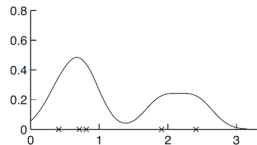
# Anomaly Detection

## Definition

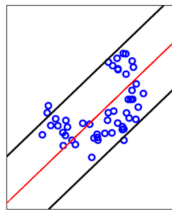
Anomaly detection identifies data points, events, and/or observations that depart from a dataset's normal behavior. Anomalous data can indicate problems, such as fraud, or opportunities, like a surge in demand for a product.



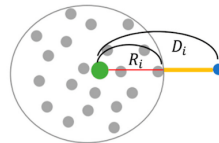
(a)



(b)



(c)



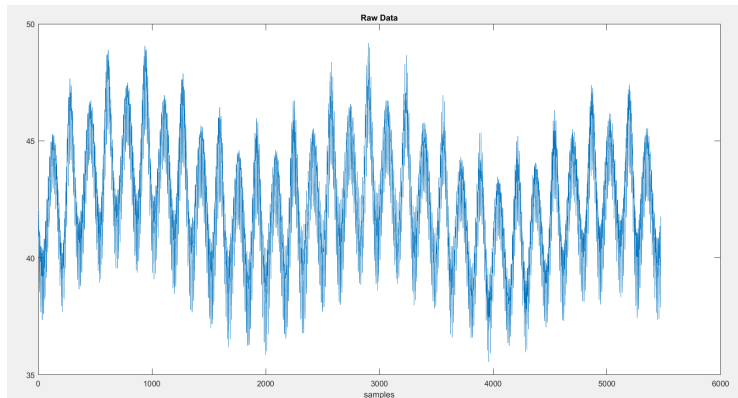
(d)

Source: Appl. Sci. 2019, 9, 4018; doi:10.3390/app9194018

# 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

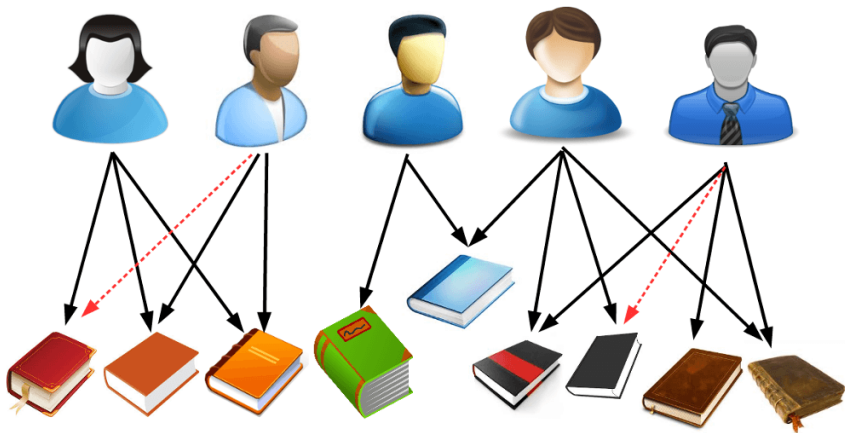


	Beer	Bread	Milk	Diaper	Eggs	Coke
$T_1$	0	1	1	0	0	0
$T_2$	1	1	0	1	1	0
$T_3$	1	0	1	1	0	1
$T_4$	1	1	1	1	0	0
$T_5$	0	1	1	1	0	1

# Recommender Systems

## Definition

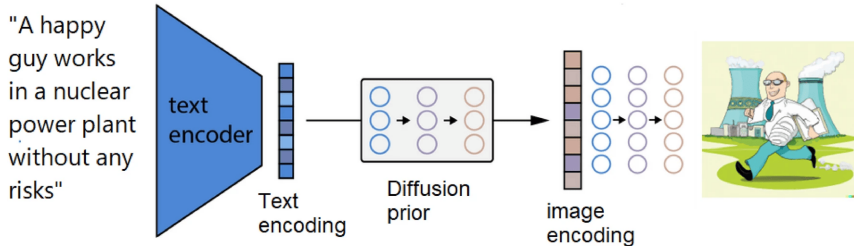
Collaborative filtering is an extension of item-based association rules mining to consider relationships between users and items. Generally, if two users have similar behaviour and/or preferences, items favoured by one user will also be favoured by the other. This can be used to recommend “new” items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



# Generating text and images

## Definition

Text, images and other rich media can be encoded as numbers and used to train deep learning models that discover hidden relationships and can be used to predict next tokens in a sequence of words, and with suitable mapping techniques such as stable diffusion, can use stable diffusion to generate new images based on the learnt *latent space* instead of just pixels.



# Traditional computing vs. Machine learning

Models transform input(s) to outputs. Example: *Compound Interest*

## Traditional Computing: Implement a known model

- Given input values  $P = \text{Principal (amount invested)}$ ,  $r = \text{interest rate (as percentage)}$  and  $n = \text{term (number of years invested)}$
- Output value is  $I = P \left(1 + \frac{r}{100}\right)^n - P$  is the *compound interest*.
- Note that the model is provided *explicitly* and is programmed as such.
- To validate the implementation of the model: unit testing on specific examples.

# Traditional computing vs. Machine learning

Models transform input(s) to outputs. Example: *Compound Interest*

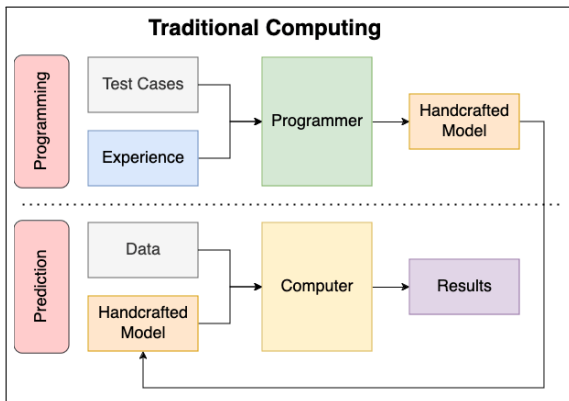
## Traditional Computing: Implement a known model

- Given input values  $P = \text{Principal (amount invested)}$ ,  $r = \text{interest rate (as percentage)}$  and  $n = \text{term (number of years invested)}$
- Output value is  $I = P \left(1 + \frac{r}{100}\right)^n - P$  is the *compound interest*.
- Note that the model is provided *explicitly* and is programmed as such.
- To validate the implementation of the model: unit testing on specific examples.

## Machine learning: Derive a model from data

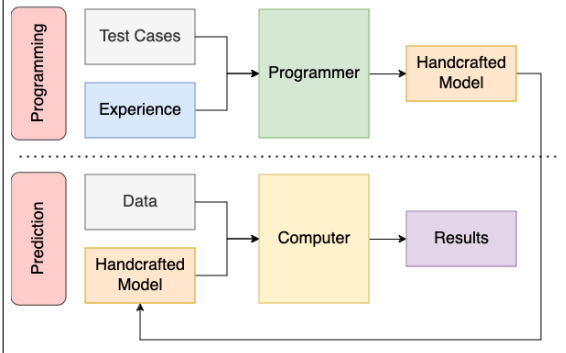
- Given inputs  $P$ ,  $r$  and  $n$  and associated compound interest output  $I$ , interpreted as before, but the formula for  $I$  is not known
- ... with several examples of inputs and associated output (training data)
- We *search* for a model that relates the given inputs to the outputs.
- To validate: use the model to predict the output from previously unseen (test) data.

# Comparing the two paradigms

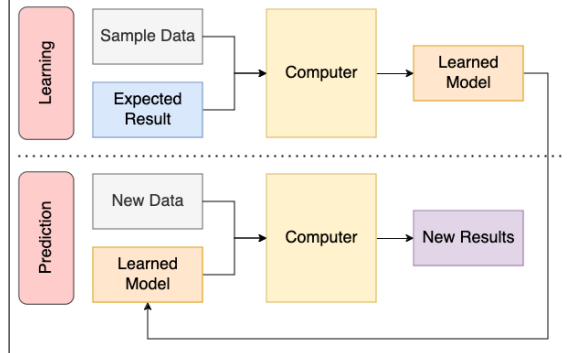


# Comparing the two paradigms

## Traditional Computing



## Machine Learning



# Outline

---

1. Module Delivery	3
2. Module Introduction	6
3. Definitions of terms and their relationships	10
3.1. Exercise: Big Data Examples	13
4. The Context and History	14
5. Expertise and Roles	23
6. Is the IoT and Big Data always a force for good?	28
7. The DIKW pipeline	33
8. Data Mining Process Models	36
9. Overview of Machine Learning techniques	40
10. Resources for students	53



# Resources



- URL: [Moodle: link to be added when module is created](#)
- Used for all notices, assignment briefs and practical work submissions.



- URL: [Data Mining 1 \(2025 Semester 1\) pages on github.io](#)
- Used for content delivery (lecture notes and labs).



- URL: [Slack workspace for comments and Q&A relating to the module.](#)

## 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 3.12 and later, 64 bit**)
- scikit-learn
- pandas

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

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

[pandas.pydata.org](https://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.