



COMP4030

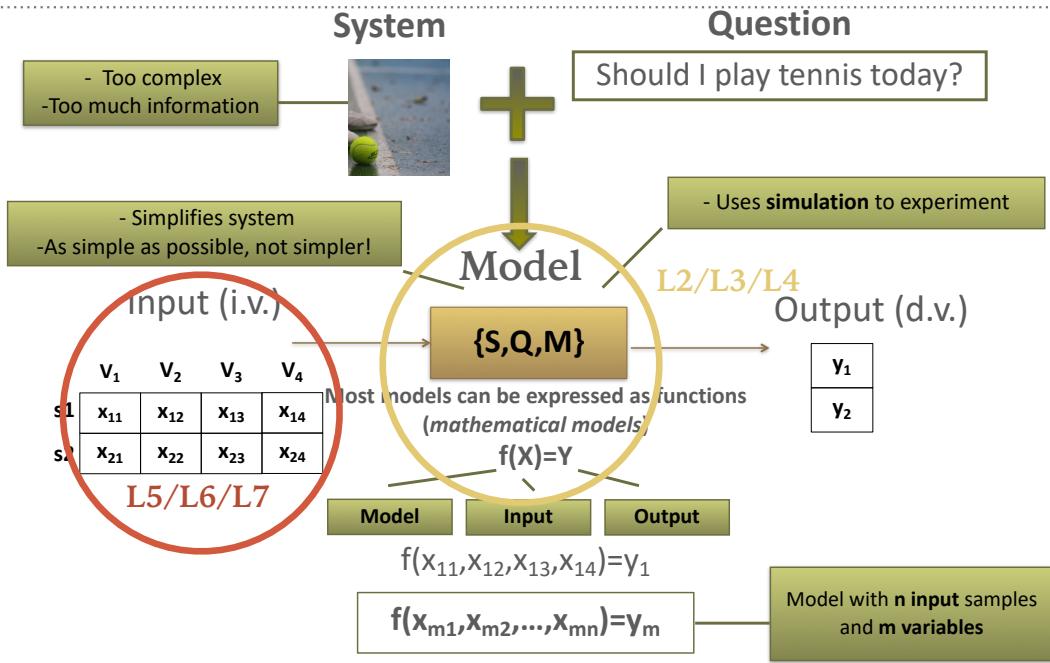
DATA MODELLING AND ANALYSIS

Lecture 5: Data Analysis

LECTURE OUTLINE

1. Summary
2. Lecture Outcomes
3. Modelling Data
4. Data vs knowledge
5. Describing data
 1. Types of data
 2. Statistics for: location, shape, dispersion, association
6. Modelling data
 1. Data overload
 2. Main challenges
 3. Data modelling pipeline

DMA: THE BIG PICTURE



LECTURE OUTCOMES



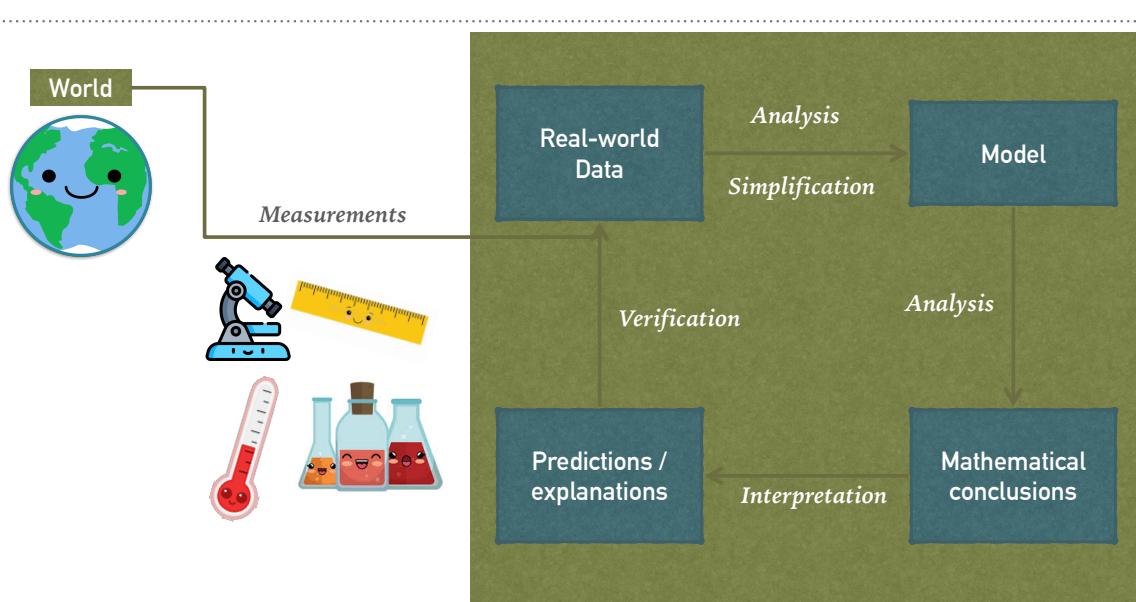
- At the end of this lecture, you should be able to answer these questions:
 - What is modelling data?
 - How is it done?
 - What is the difference between data and knowledge?
 - What are the main issues found when modelling data?
 - What are the types of data?
 - Which statistics are used to describe data?

MODELLING DATA

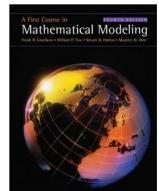


- This module: **Data Modelling and Analysis**
 - Why are we talking about data?
 - Data: motivation behind modelling and analysis
- **Analysis:** organise and understand
- **Modelling:** describe and predict
- **Optimisation:** build & improve models
- Mathematical Models allow for description and prediction of data:
 - Description
 - e.g.: Summarisation, Clustering, Visualisation
 - Prediction (generalise to new cases)
 - e.g.: Classification, Regression, Forecasting

MODELLING DATA: DATA IN RELATION TO THE WORLD



Adapted from:





MODELLING DATA: DATA OVERLOAD

- Information/data overload
- Data is being generated at an ever increasing rate
- Rough estimate:
 - Data volume in the world doubles every 20 months
- What does the data mean?
- What can we do with the data?
- Our ability to analyse data (e.g. through computational tools) is developing slower than our ability to generate and store data



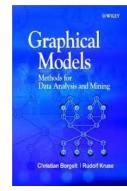
MODELLING DATA: DATA OVERLOAD

- Data overload comes in many forms
- Many sources of data with many applications, e.g:
 - Images → classification of satellite imagery
 - Sounds → music recognition, speech recognition
 - Text → automated translation
 - DNA → gene identification
 - Demographics → marketing decision support
 - Commercial transactions → inventory optimization
 - Weather → weather forecasting
 - etc.



MODELLING DATA: DATA VS KNOWLEDGE

- For data to be useful, we need to extract **meaning**
 - Data → knowledge
- Differences between data and knowledge

Data	Knowledge	
Refers to instances (e.g. objects and events)	Refers to classes (e.g. sets of objects and events)	Adapted from
Describes individual properties	Describes general patterns, principles, etc.	
Usually in huge amounts	Aim: to extract a minimal set of statements	
Usually easy to obtain	Usually hard to obtain	
Can't be used directly for generating predictions	Can be used directly for generating predictions	



BUT...WHAT IS DATA?

- When we talk about data, what are we talking about?
- Factual information used as a basis for reasoning, discussion, or calculation
 - Includes measurements or statistics
- Information in digital form that can be transmitted or processed



NOT THIS



5.94,66755.39,0,
59.12,42826.99,0,0,0,0,30.09
35.64,50656.8,0,0,0,0,30.10
115.94,67905.07,0,0,0,0,30.12
115.94,66938.9,0,0,0,0,30.13
119.24,49,86421.04,0,0,0,0,30.14
129.98,5.0,0,0,0,0,30.15

BUT THIS!

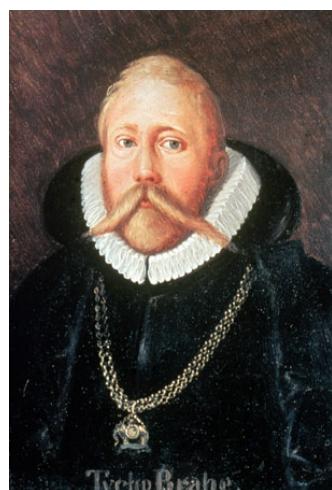
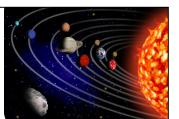
DESCRIBING DATA: DATA REPRESENTATION



- Different formats
- Common representation: tables
 - Columns: **attributes, features**
 - Measurable characteristics
 - Last column: **Output, label, class**
 - rows (tuples): data objects, **instances**
 - case, point, record, sample, entity
 - Input(s) / **output(s)**

	Outlook	Temp.	Humid.	Windy	Play
1	sunny	hot	high	false	no
2	sunny	hot	high	true	no
3	overcast	hot	high	false	yes
4	rainy	mild	high	false	yes
5	rainy	cool	normal	false	yes
6	rainy	cool	normal	true	no
7	overcast	cool	normal	true	yes
8	sunny	mild	high	false	no
9	sunny	cool	normal	false	yes
10	rainy	mild	normal	false	yes
11	sunny	mild	normal	true	yes
12	overcast	mild	high	true	yes
13	overcast	hot	normal	false	yes
14	rainy	mild	high	true	no

MODELLING DATA: DATA VS KNOWLEDGE – AN EXAMPLE



Tycho Brahe
(1546 – 1601)



Johannes Kepler
(1571-1630)

MODELLING DATA: DATA VS KNOWLEDGE

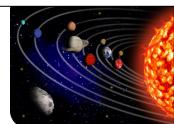


Tycho Brahe

(1546 – 1601)

Tabularum Rudolphini									
Tabula Aequationum MARTIS.									
Anomalia	Interven-	Anomalia	Interven-	Anomalia	Interven-	Anomalia	Interven-	Anomalia	Interven-
Eccentria,	Contra-	Excentrica,	Contra-	Excentrica,	Contra-	Excentrica,	Contra-	Excentrica,	Contra-
Contra-	Contra-	Contra-	Contra-	Contra-	Contra-	Contra-	Contra-	Contra-	Contra-
Mercurij	Mercurij	Veneris	Veneris	Terrestriis	Terrestriis	Mercurij	Veneris	Terrestriis	Solis
120	1. 5. 16	115.17.11	37717	14.5.02	170	1.19.14	1.16.35	147.13.45	149.15.2
121	1. 5. 20	116.19.15	145030	151	1.14.23	1.04.45	146.18.43	140.00.1	
122	1. 5. 29	117.22.39	144371	152	2.10.49	1.04.50	149.23.44	139.88.7	
123	1. 6. 7	118.25.13	174642	153	1.09.23	1.07.52	150.28.45	137.44.4	
124	1. 6. 16	119.28.59	144358	154	2.19.16	1.09.10	151.33.46	139.66.6	
125	1. 6. 25	120.31.53	144075	155	1.14.13	1.11.12	151.37.47	139.50.1	
126	1. 6. 34	121.34.42	143903	156	1.09.10	1.07.57	152.43.48	138.57.9	
127	1. 6. 43	122.37.50	143855	157	2. 9.30	1.03.43	153.44.49	138.53.1	
128	1. 7. 1	123.41.16	143661	158	1.11.13	1.04.80	154.49.49	138.51.3	
129	1. 7. 10	124.44.37	143645	159	1.19.11	1.05.55	155.55.5	138.26.5	
130	1. 7. 19	125.48.0	143278	160	1.14.8	1.07.13	157.0.13	139.17.7	
131	1. 7. 28	126.51.40	143093	161	1.12.39	1.08.00	158.5.49	138.99.9	
132	1. 7. 37	127.55.19	142995	162	1.13.16	1.11.43	159.11.47	139.93.1	
133	1. 8. 5	128.59.3	142726	163	1.09.00	1.10.00	160.16.47	138.52.7	
134	1. 8. 14	129.63.0	142614	164	1.13.22	1.11.12	161.22.51	138.50.9	
135	1. 8. 23	130.67.4	142410	165	1.12.48	1.11.57	162.27.53	137.77.3	
136	1. 8. 32	131.71.1	142370	166	1.21.39	1.13.60	163.31.53	137.16.1	
137	1. 8. 41	132.75.7	142357	167	1.13.0	1.12.50	164.37.53	136.64.1	
138	1. 8. 50	133.79.3	142353	168	1.12.4	1.12.39	165.43.59	136.64.1	
139	1. 9. 8	134.83.0	142357	169	1.09.47	1.12.19	166.50.59	136.62.1	
140	1. 9. 17	135.87.6	142353	170	1.11.47	1.12.19	167.56.5	136.49.5	
141	1. 9. 26	136.91.3	142157	171	1.12.48	1.13.21	168.63.59	136.47.6	
142	1. 9. 35	137.95.0	142154	172	1.21.39	1.13.21	169.7.57	136.45.4	
143	1. 9. 44	138.98.7	142150	173	1.13.2	1.13.21	170.18.53	136.43.4	
144	1. 9. 53	139.10.4	142031	174	1.12.13	1.12.20	171.15.0	136.41.3	
145	1. 9. 62	140.10.7	142031	175	1.12.13	1.12.20	172.12.0	136.39.2	
146	1. 9. 71	141.10.4	141942	176	1.12.13	1.12.20	173.19.49	136.37.1	
147	1. 9. 80	142.10.1	141942	177	1.12.13	1.12.20	174.17.53	136.35.0	
148	1. 9. 89	143.10.8	140659	178	1.12.43	1.12.20	175.15.59	136.32.9	
149	1. 9. 98	144.10.5	140657	179	1.12.43	1.12.20	176.7.57	136.30.8	
150	1. 9. 107	145.10.2	140657	180	1.0.0	1.12.32	177.63.59	136.28.7	
151	1. 9. 116	146.10.9	140654	181	1.12.32	1.12.32	178.6.59	136.26.6	

Tab. Lat.



MODELLING DATA: DATA VS KNOWLEDGE



KEPLER'S LAWS OF PLANETARY MOTION

1. The orbit of a planet is an ellipse with the Sun at one of the two foci.
2. A line segment joining a planet and the Sun sweeps out equal areas during equal intervals of time.
3. The square of the orbital period of a planet is proportional to the cube of the semi-major axis of its orbit.



Johannes Kepler

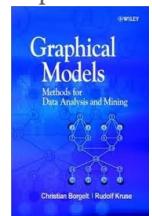
(1571-1630)



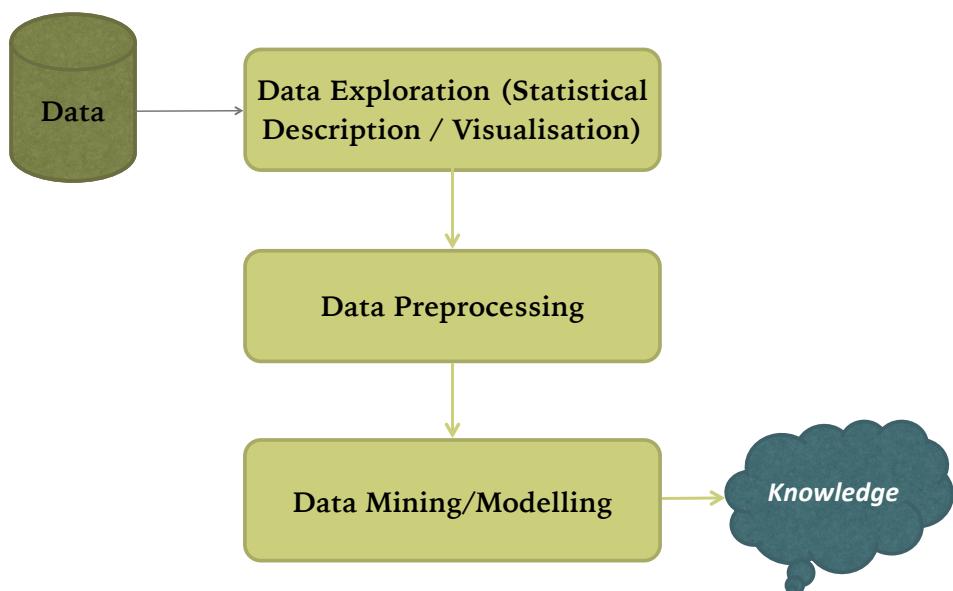
MODELLING DATA: DATA VS KNOWLEDGE - ASSESSMENT

- Not all knowledge is equal: how do we assess it?
- Assessment Criteria
 - Correctness (probability, success in tests)
 - Generality (range of validity, conditions of validity)
 - Usefulness (relevance, predictive power)
 - Comprehensibility (simplicity, clarity, parsimony)
 - Novelty (previously unknown, unexpected)
- Priority
 - Science: correctness, generality, comprehensibility
 - Economy: usefulness, comprehensibility, novelty

Adapted from



MODELLING DATA: DATA MODELLING PIPELINE



DESCRIBING DATA TYPES OF DATA: CATEGORICAL (DISCRETE) DATA



A. Nominal

- Data characterised by names, labels or categories
- There is no order
 - Blood type (A, B, AB, O)
 - Gender
 - Country of origin
 - Attributes with yes/no values
- Data is **distinct**
- valid statistics: **mode**

DESCRIBING DATA TYPES OF DATA: CATEGORICAL (DISCRETE) DATA



B. Ordinal

- Data can be arranged in some **order** but differences between values are meaningless.
- It may be linguistic, but also numeric!
 - Product rating: poor, good, excellent.
 - Gymnastic scores.
 - Satisfaction level (very unsatisfied, somewhat unsatisfied, neutral, somewhat satisfied, very satisfied).
- **Distinctness** and **order**
- valid statistics: **median, percentile**



DESCRIBING DATA TYPES OF DATA: NUMERIC (METRIC) DATA

C. Interval

- Data are ordered and differences are meaningful.
- There is no true 0.
- Multiplication and division are meaningless, so ratios have no meaning
 - Temperature: $50^\circ - 40^\circ = 60^\circ - 50^\circ$, but 100° is not $2 * 50^\circ$.
- **Distinctness, order** and **addition (subtraction)**
- valid statistics: mean, st. dev., correlation, regression



DESCRIBING DATA TYPES OF DATA: NUMERIC (METRIC) DATA

D. Ratio

- data is ordered and all arithmetic operators are applicable; ratios are meaningful
 - e.g. litres of petrol consumed per day per individual
- **distinctness, order, addition** and **multiplication (division)**
- valid statistics: **all**

DESCRIBING DATA: EXERCISE: DATA CATEGORISATION



1. IQ scores? (100 = ‘average’; 140 = ‘genius’)
2. Sport team squad numbers? (1 to 22)
3. Golf score?
4. Age?
5. Date of birth?
6. Political party (Labour, Lib-Dem, Conservative)?
7. Exam grades (A, B, C, D, E, F)?
8. pH (acidity = $-\log[H^+]$)?
9. Number of children

DESCRIBING DATA: EXERCISE: DATA CATEGORISATION



10. Income?
11. Eye colour?
12. Degrees in Fahrenheit?
13. Degrees in Kelvin?
14. Number of sales?
15. Level of spiciness?
16. Time?
17. Constellations
18. Zodiac sign



DATA ANALYSIS: DESCRIBING DATA – LOCATION



A. Mode (nominal)

- Attribute value that is most frequent in the sample
- Not unique
 - more than one value may have the largest frequency
- The most general localisation measure
 - applies to all scales of measurement

DATA ANALYSIS: DESCRIBING DATA – LOCATION



B. Median (Ordinal):

- If the data elements x_1, x_2, \dots, x_n are sorted, the median is:

$$\tilde{x} = \begin{cases} x_{\left(\frac{n+1}{2}\right)} & \text{if } n \text{ is odd} \\ \frac{1}{2}(x_{\left(\frac{n}{2}\right)} + x_{\left(\frac{n}{2}+1\right)}) & \text{if } n \text{ is even} \end{cases}$$

C. Arithmetic mean (interval, ratio):

- Only applicable to metric attributes

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

DATA ANALYSIS: DESCRIBING DATA – DISPERSION



D. Range:

- A. difference between the maximum and minimum value

E. Inter-quantile range:

- p -quantile: the value relative to which the fraction p of the data is smaller
- The $(1/2)$ -quantile value is the median
- the p -inter-quantile range (IQR) is the difference between the $(1-p)$ -quantile value and the p -quantile value
- note: $0 < p < 1/2$
- It is common to use $p = 1/4$, *inter-quartile* range
- In which case: $IQR = Q3 - Q1$

DATA ANALYSIS: DESCRIBING DATA – LOCATION



F. Average absolute deviation

$$d_{\bar{x}} = \frac{1}{n} \sum_{i=1}^n |x_i - \bar{x}|$$

G. Variance

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

H. Standard deviation

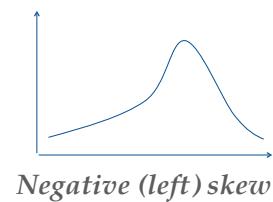
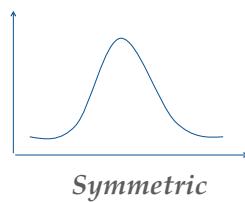
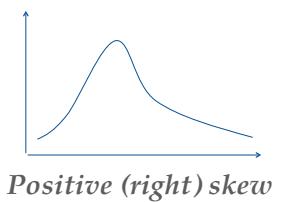
$$s = \sqrt{s^2} = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$$

DATA ANALYSIS: DESCRIBING DATA – SHAPE



► Skewness:

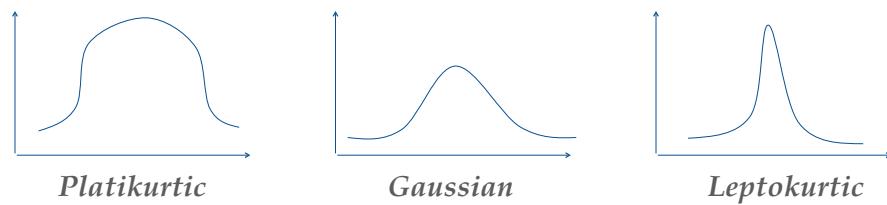
- Measures the asymmetry of probability distribution of a real-valued random variable about its mean.
- It can be positive, negative and even undefined.



DATA ANALYSIS: DESCRIBING DATA – SHAPE



- Kurtosis:
 - Measures how peaked a distribution is
 - Measure is usually computed relative to the Gaussian



DATA ANALYSIS: DESCRIBING DATA – RELATIONSHIPS



- Measures a statistical relationship between two variables.
- **Pearson coefficient:** measures strength and direction of the linear relationship between two variables.
- Ranges between $[-1, +1]$

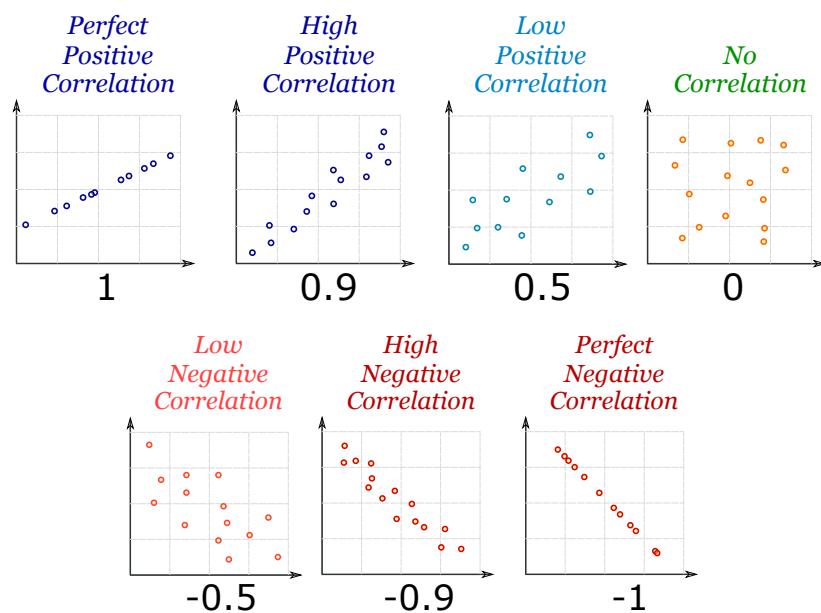
$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

n is the sample size

x_i, y_i are the single samples indexed with i

$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ is the sample mean of X

DATA ANALYSIS: DESCRIBING DATA – RELATIONSHIPS



DATA ANALYSIS: DESCRIBING DATA – RELATIONSHIPS

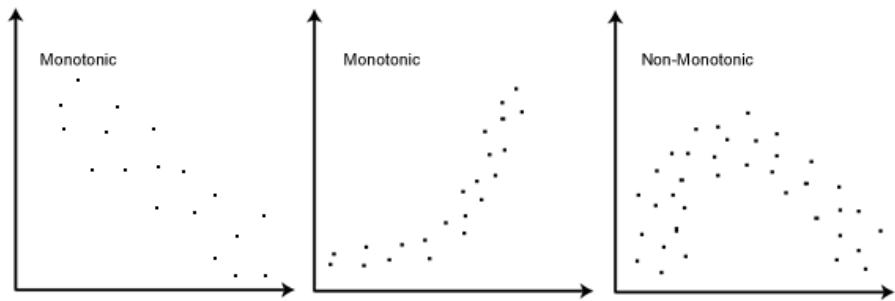


- The Pearson coefficient is not applicable to ordinal data.
- **Spearman coefficient:** non-parametric rank-based correlation coefficient
- It measures monotonic relationship
- Values range between [-1,1]

$$\rho = \frac{\sum (x'_i - m_x)(y'_i - m_y)}{\sqrt{\sum (x'_i - m_x)^2 \sum (y'_i - m_y)^2}}$$

- where m is the mean, $x' = \text{rank}(x)$, and $y' = \text{rank}(y)$

DATA ANALYSIS: DESCRIBING DATA – RELATIONSHIPS





MODELLING DATA: MAIN CHALLENGES

- In this lecture, we will cover three challenges [1]:
 - Too much data
 - Too little data
 - Fractured data
- Analysis will help us understand which of these challenges, if any, are present in our problems.
- Lecture 6 and 7 will delve in how to “fix” them, if possible.

A. Famili, W.M. Shen, R. Weber, E. Simoudis, *Data preprocessing and intelligent data analysis*, Intell. Data Anal. 1 (1) (1997) 1–28.



MODELLING DATA: MAIN CHALLENGES

- Too much data
- Corrupt and noisy data
- Irrelevant data
- Very large data-sizes
- Very large number of attributes
- Different types of data

A. Famili, W.M. Shen, R. Weber, E. Simoudis, *Data preprocessing and intelligent data analysis*, Intell. Data Anal. 1 (1) (1997) 1–28.



MODELLING DATA: TOO MUCH DATA

- Some domains (e.g. satellite imagery, medical imagery, telecommunications) have very large:
 - Volumes of data
 - Rates of data production
- Real-time data analysis: very difficult if not impossible (depending on hardware available)
- Example solutions:
 - Sampling
 - Scaling down the resolution of the data (effectively attribute and instance selection of some kind)
 - Dimensionality reduction (for attributes)
 - High-performance computing, etc.



MODELLING DATA: MAIN CHALLENGES

- Too little data
- Depending on the problem, a much more serious issue than having too much data
- Issues include:
 - Missing attributes
 - Missing attribute values
 - Small amount of data



MODELLING DATA: TOO LITTLE DATA

- Even if no missing values, data set may be too small
- Some techniques require large data sets
 - Machine learning models: require a lot of examples to be able to distinguish between classes
- Reliability of predictions might be too low
 - Examples are not enough to generalise
- Statistical tests and sample size
 - Confidence of results



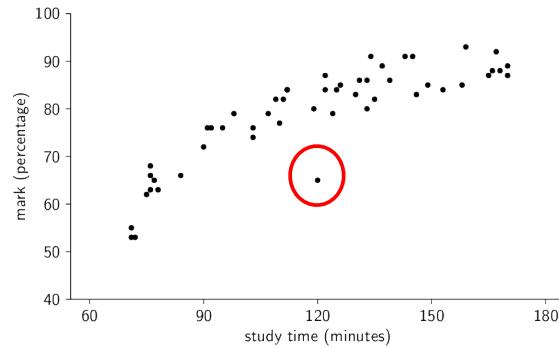
MODELLING DATA: MAIN CHALLENGES

- Fractured data
- Fractured data is a very common issue in problems where fusion of data is necessary
- Issues include:
 - Outliers
 - Incompatible data
 - Multiple sources of data
 - Data from multiple levels of granularity

MODELLING DATA: FRACTURED DATA



- Values lies far away from other population values
- Causes:
 - Noise
 - Mixed distributions
 - Natural exceptions
- Problem
 - Distort predictions
- Measures
 - Test abnormality
 - Eliminate (carefully!)



MODELLING DATA: FRACTURED DATA – OUTLIER

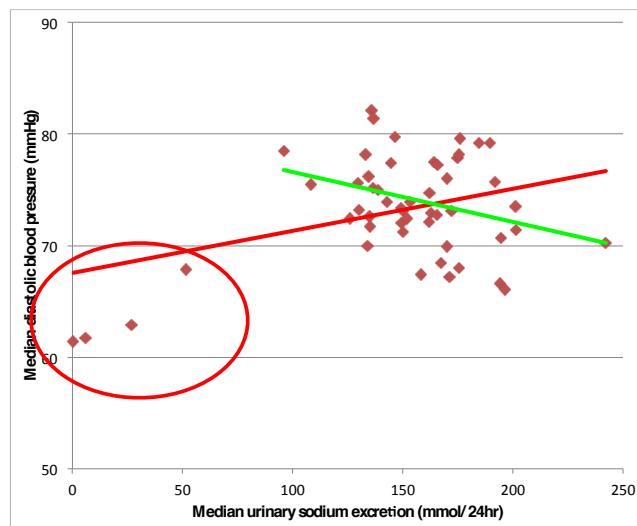


- The **identification** and **elimination** of outliers has to be done very carefully.

Intersalt: an international study of electrolyte excretion and blood pressure. Results for 24 hour urinary sodium and potassium excretion.

British Medical Journal;
297: 319-328, 1988.

If we exclude these four ‘outliers’, which are non-industrialised countries with non-salt diets, we get a quite different result!





MODELLING DATA: FRACTURED DATA – MULTIPLE SOURCES

- Data collected by different groups can often lead to data compatibility issues
- Data collected by different teams/projects at different times
- Examples of specific causes of incompatibility:
 - different experimental designs
 - different choices of attributes
 - different hardware



MODELLING DATA: FRACTURED DATA – MULTIPLE SOURCES

- In large enterprises it may be that:
 - data is scattered in different departments and platforms
 - data is acquired and maintained using different software
 - the goal, depth and standard of data collection may vary
- These subtle differences can cause problems when the data from the different sources is combined



MODELLING DATA: FRACTURED DATA – MULTIPLE SOURCES

- In real world applications data can come from different levels of granularity
- Levels of granularity: aerospace example
 - fleet level
 - e.g. all aircraft of a particular type
 - aircraft level
 - e.g. a particular fin-number
 - system level
 - e.g. a particular engine of an aircraft
 - system operation level
 - e.g. operation of the engine for a particular duration or cycle
- The finer the granularity of the data, the more difficult it will be to collect it and to model it.

LECTURE OUTCOMES



- Now, you should be able to answer these questions:
 - What is modelling data?
 - How is it done?
 - What is the difference between data and knowledge?
 - What are the main issues found when modelling data?
 - What are the types of data?
 - Which statistics are used to describe data?

THE END

Questions 