



FIT3152 WEEK 1-8

Data Analytics (Monash University)

FIT3152

WEEK 1 – INTRODUCTION TO DATA SCIENCE

What is data science?

- Data Science is making a product out of data and using graphics to communicate the result
- Complex problems of societal interest/utility
- Large data sets, multiple data sets, messy data, incomplete data

Data scientist skills

- Understand problem from client's perspective
- Collect, cleanse and manage data
- Understand the data
- Analyse and model the data using statistical and machine learning techniques
- Communicate the results simply and effectively

Examples

- Road shooter found via mass data collection
- Food network
- Climate change time series
- Customer analytics

WEEK 2 – DATA VISUALISATION

Dimensions

- Graphics can show many dimensions
- Include colour, size, position, adjacency, connect, shape

Graphic types – Visualisation Zoo

- Time-series data
 - o Index charts
 - o Stacked graphs
 - o Small multiples
 - o Horizon Graphs
- Statistical distributions
 - o Stem-and-leaf plots
 - o Q-Q plots
 - o SPLOM
 - o Parallel coordinates
- Maps
 - o Flow maps
 - o Choropleth maps
 - o Graduated symbol maps
 - o Cartograms
- Hierarchies
 - o Node-link diagrams
 - o Adjacency diagrams
 - o Enclosure diagrams
- Networks
 - o Force-directed layouts
 - o Arc diagrams
 - o Matrix views

WEEK 3 – DATA MANIPULATION

Manipulation in R

- Making tables and summaries
- Applying functions and considering factors
- Transforming data
 - o Aggregate
 - o Cor
 - Correlation function only works on columns
 - o By
 - `By(data frame, column, function(df) function)`
 - Df is a temporary data frame created for each factor
 - Applied to a df and returns list
 - o As.table
 - List to a table
 - o As.data.frame
 - Table to a df
 - o Colnames
 - `colnames(data frame) = c("name1", "name2")`
 - Adds column name
 - o Merge
 - o Cbind
 - o Rbind
 - o Max
 - o Which.max
 - Determines the index of the maximum numeric vector
 - o Do.call
 - Calls a function on a list of arguments

WEEK 4 – REGRESSION MODELS

Fitting the regression

- Simple least squares regression assumes that the relationship is approximately linear and errors are approximately normally distributed
- Line of best fit

Regression diagnostics

- Look at median and coefficients
- P-value for variables and for overall model
 - o Small p-value = Is significant
- R-squared – Coefficient of determination
 - o High r-squared = How close the data is to the fitted regression line

Multiple linear regression

Regression with qualitative variables

- 0, 1 column

WEEK 5 – NETWORKS

Structure

- Nodes (vertices) and edges (lines/arcs)
 - o Can be directed with arrows or undirected
 - o Edges can be weighted
- Walk: Sequence of links

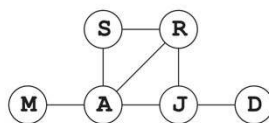
- Path: Walk with no repeated vertices
- Cycle: Walk that begins and ends at the same vertex
- Geodesic: Shortest path between two vertices
- Length: Number of links in walk/path
- Connected: Path between each pair of vertices
- Directed graphs: Travel on each edge is permitted in one direction only
- Link sequences
 - o Loop: An edge from a vertex to itself
 - o Complete: Graph where every vertex is joined to another vertex
 - o Subgraph: Subset of a graph
 - o Clique: Subgraph that is complete (super connected circle)
 - o Simple: Graph with no loops or multi-edges
- Degree: Number of edges connected to a vertex
 - o For directed, in-degree and out-degree
 - o `degree(graph object)`
- Diameter: Longest geodesic
- Average path length: Average geodesic
 - o Create distance matrix first
- Degree distribution: Probability distribution describing the magnitude of the vertices (weighted)

- *Diameter*: $= \max(\text{dist}(u,v)) = 3$
- *Degree distribution*:
- *Average path length*: 1.667

	J	A	S	D	R	M
J		1	2	1	1	2
A			1	2	1	1
S				3	1	2
D					2	3
R						2
M						

Distance Matrix

Degree	N
0	0
1	2
2	1
3	2
4	1



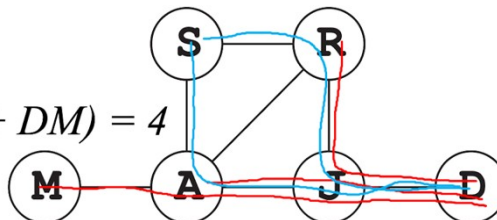
Vertex characteristics

- Importance of a vertex is based on degree and centrality
- Centrality of a vertex
 - o Betweenness
 - Indicates the degree to which the vertex is between other vertices
 - Calculate by hand
 - Measures the hub potential of a node
 - `betweenness(graph object)`

$$c_B(v) = \sum_{s \neq t \neq v \in V} \frac{\sigma(s, t|v)}{\sigma(s, t)}$$

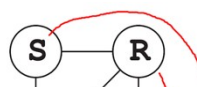
$$c_B(J) = (DR + DA + DS + DM) = 4$$

$$c_B(M) = 0$$



- o Closeness
 - Vertex is close if there is a small total distance between it and all the other vertices
 - Calculate by hand
 - Measures how well a node is connected locally
 - `closeness(graph object)`

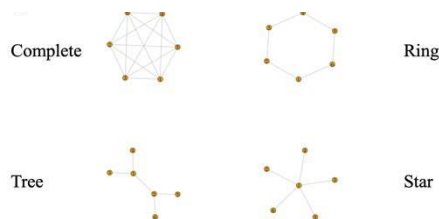
$$c_{cl}(v) = \frac{1}{\sum_{u \in V} \text{dist}(u, v)}$$



- Eigenvector
 - Higher weight to vertices with neighbours that are more central
 - Weights a node according to the quality of its connection where nodes connected to important nodes are ranked higher
 - `evcent(graph object)`

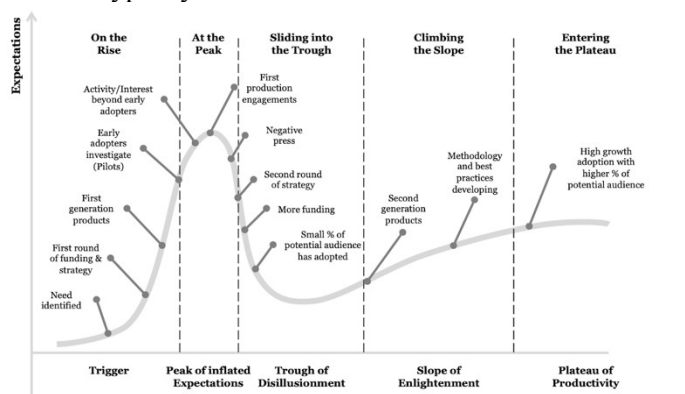
Four network topologies

- Complete (most robust)
- Ring
- Tree
- Star (most fragile)



WEEK 6 – TIDY AND DIRTY DATA

Gartner Hype Cycle

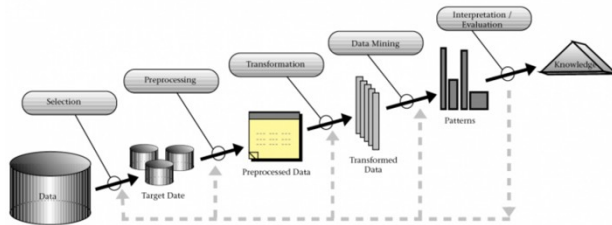


Five stages of understanding

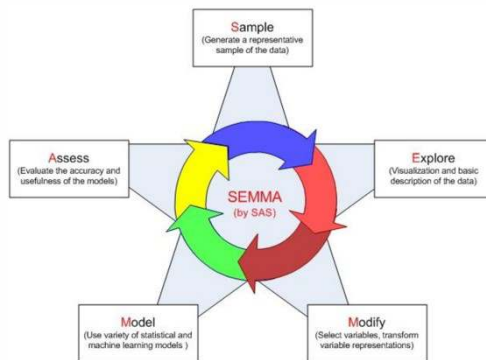
- Descriptive
 - What happened?
 - Condense data into smaller, more useful pieces of information
- Diagnostic
 - Why did it happen?
 - Data analysis by employing predefined criteria
- Explorative
 - What might be interesting?
 - Uncover underlying structure, patterns, and anomalies
 - Manual discovery and automated discovery (machine learning & clustering)
- Predictive
 - What is likely to happen?
 - Predictive modelling on historical data to produce future likelihood of events
 - E.g. Random forest, bagging, & boosting
- Prescriptive
 - What can we do about it?
 - Suggests the best option for handling a future scenario

Data science methodologies (NE)

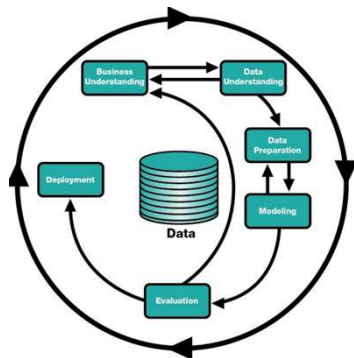
- KDD (Knowledge Discovery in Databases)
 - o Broad process of finding knowledge in data
 - o Emphasises the high-level use of data-mining methods
 - o Selection -> preprocessing -> transformation -> data mining -> interpretation/evaluation



- SEMMA (Sample, Explore, Modify, Model, and Assess)



- o Methodology for data mining processes
- CRISP-DM (Cross-Industry Standard Process for DM)
 - o Preferred due to inclusion of business aspects
 - o Business understanding
 - Understanding the project objectives and requirements from a business perspective
 - o Data understanding
 - Initial data collection and activities to get familiar with the data
 - o Data preparation
 - Activities required to construct the final dataset from the initial raw data
 - Performed repeatedly and not in any prescribed order
 - o Modelling
 - Modelling techniques are selected and applied
 - o Evaluation
 - Evaluate model and review steps taken to ensure it properly achieves the business objectives
 - o Deployment
 - Knowledge will need to be organised and presented in a way that the customer can use



Business Understanding	Data Understanding	Data Preparation	Modeling	Evaluation	Deployment
Determine Business Objective Background Business Objectives Business Success Criteria Assess Situation Inventory of Resources Requirements, Assumptions, and Constraints Risks and Contingencies Terminology Costs and Benefits Determine Data Mining Goals Data Mining Goals Data Mining Success Criteria Produce Project Plan Project Plan Initial Assessment of Tools and Techniques	Collect Initial Data Initial Data Collection Report Describe Data Data Description Report Explore Data Data Exploration Report Verify Data Quality Data Quality Report	Select Data Rationale for Inclusion/Exclusion Clean Data Data Cleaning Report Construct Data Derived Attributes Generated Records Integrate Data Merged Data Format Data Reformatted Data Dataset Dataset Description	Select Modeling Techniques Modeling Technique Modeling Assumptions Generate Test Design Test Design Build Model Parameter Settings Models Model Descriptions Assess Model Model Assessment Revised Parameter Settings	Evaluate Results Assessment of Data Mining Results w.r.t. Business Success Criteria Approved Models Review Process Review of Process Determine Next Steps List of Possible Actions Decision	Plan Deployment Deployment Plan Plan Monitoring and Maintenance Monitoring and Maintenance Plan Produce Final Report Final Report Final Presentation Review Project Experience Documentation

Dirty data

- Can be:
 - o Incorrect
 - Valid data (e.g. dates)
 - o Inaccurate
 - Correct but not accurate (e.g. Sydney, VIC)
 - o Incomplete
 - Missing/empty fields
 - o Duplicate
 - o Violates business rules
 - E.g. date should precede expiration date
 - o Inconsistent
 - Due to uncontrolled data redundancy (e.g. Customer listed three times)
 - o Non-integrated
 - Data is stored redundantly and inconsistently across many systems
 - Primary keys don't match or are not unique

Tidy data

- Consistent format that has:
 - o Each variable in its own column
 - o Each observation in its own row
 - o Each value in its own cell
- Actions:
 - o Replace missing values
 - o Standardisation
 - o Normalisation

WEEK 7 – DECISION TREES

Machine Learning

- Automated learning of a concept given some examples of data
- How can a model learn a concept?
 - o Descriptive – captures training data
 - o Predictive – generalises to unseen data
 - o Explanatory – describes the concept to be learned

Classification

- Find a model to predict class as a function of the other attributes
- Goal is for unseen records to be assigned a class as accurately as possible
- Test set is used to determine the accuracy
- Training set is used to build the model

Decision trees

- Can be used to profile existing data and classify new instances
- Robust to noise and missing values
- Leaf nodes (class) and non-leaf nodes
- Branches corresponding to the values of the decision attributes
- Applications
 - o Profiling
 - o Classifying unseen data
 - o Example areas:
 - Consumer credit evaluation
 - Spam filtering
 - Predicting customer churn
- Pros
 - o Easy to interpret
 - o Can handle discrete and continuous input
 - o Robust to outliers and missing values
- Cons
 - o Unstable
 - o Can become overly complex (overfitting)

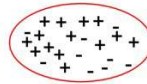
Top-down induction: ID3

- Greedy approach
 - o At each step, determine the best decision for the decision attribute
 - o May not be the best decision overall but stick with it
 - o May not result in best overall decision tree
 - o Goal is to increase the homogeneity (same) of the resulting datasets with respect to the class variable
- Information gain
 - o ID3 uses information gain as splitting criteria
 - o Measures how well a given attribute separates the training set into homogenous groups
 - o Expected reduction in entropy
- Entropy
 - o Measures the uncertainty in a random variable
 - o Ordered data = low entropy
 - o Entropy is 0 if all members belong to the same class
 - o For a two class problem
$$\begin{aligned}\text{Entropy}(S) &= -P_{c1} \log_2(P_{c1}) - P_{c2} \log_2(P_{c2}) \\ &= -\frac{N_{c1}}{N} \log_2\left(\frac{N_{c1}}{N}\right) - \frac{N_{c2}}{N} \log_2\left(\frac{N_{c2}}{N}\right)\end{aligned}$$
 - o For a multi-class problem

$$\begin{aligned}\text{Entropy}(S) &= - \sum_{i=1}^C P_i \log_2(P_i) \\ &= - \sum_{i=1}^C \frac{N_i}{N} \log_2\left(\frac{N_i}{N}\right)\end{aligned}$$

○ Example:

E.g. Suppose S is a collection of 14 examples, 9 positive and 5 negative → [9+,5-]



E.g. suppose S has all positive or all negative

$$\begin{aligned}\text{Entropy}(S) &= -P_{c1} \log_2(P_{c1}) - P_{c2} \log_2(P_{c2}) \\ &= -\frac{9}{14} \log_2\left(\frac{9}{14}\right) - \frac{5}{14} \log_2\left(\frac{5}{14}\right) \\ &= 0.940\end{aligned}$$

Entropy

Metrics for performance evaluation

- Calculate accuracy using confusion matrix
- Accuracy is trues/everything

	PREDICTED CLASS	
	Class=Yes	Class=No
	Class=Yes	Class=No
ACTUAL CLASS	a	b
	c	d

a: TP (true positive)

b: FN (false negative)

c: FP (false positive)

d: TN (true negative)

	PREDICTED CLASS	
	Class=Yes	Class=No
	Class=Yes	Class=No
ACTUAL CLASS	a (TP)	b (FN)
	c (FP)	d (TN)

Most widely-used metric:

$$\text{Accuracy} = \frac{a + d}{a + b + c + d} = \frac{TP + TN}{TP + TN + FP + FN}$$

Also:

$$\text{Precision} = \frac{tp}{(tp + fp)}$$

$$\text{Sensitivity} = \frac{tp}{(tp + fn)}$$

WEEK 8 – OTHER CLASSIFIERS

Metrics

- Precision and recall

Precision is:

$$\frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}}$$

Recall is:

$$\frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$

Overfitting/underfitting

- Overfitting
 - Model does not generalise well and is excessively complex
 - Performs well on training data but not on testing data
 - Low recall
 - To avoid
 - Pre-pruning: Stop growing tree earlier
 - Post-pruning: Allow tree to overfit and then post-prune
 - More effective in practice
 - Generally, remove rules if they have little effect on the error rate
- Underfitting
 - Model is too simple to give accurate labels
 - Performs poorly on training and testing data
 - Low precision

Cross validation

- K-fold cross-validation steps
 - o Partition data into k disjoint subsets
 - o Train on k-1 partitions and test remaining one
- | | | |
|----|----|----|
| D1 | D2 | D3 |
|----|----|----|

 → D1 + D2 for Training, D3 for Test

D1	D2	D3
----	----	----

 → D1 + D3 for Training, D2 for Test

D1	D2	D3
----	----	----

 → D2 + D3 for Training, D1 for Test
- Make sure that all folds have the same distribution of classes
 - Stratification ensures that classes are properly represented across partitions
 - Leave-one-out cross-validation
 - o Each case is left out and the model is trained on all of the remaining instances

Bayesian classifiers

- Bayes' Theorem is conditional probability
- Naïve Bayes classifier
 - o Robust to isolated noise points
 - o Can adapt quickly to new data
 - o Handles missing values by ignoring them
 - o Robust to irrelevant attributes

Classifier model evaluation

- Confusion matrix and performance measures
 - ROC (Receiver operating characteristics) for binary classifiers
 - o Tradeoff between positive hits and false alarms
 - o TPR on y axis and FPR on x axis
- $$\text{True Positive Rate: } TPR = \frac{TP}{TP + FN}$$
- $$\text{False Positive Rate: } FPR = \frac{FP}{FP + TN}$$
- o Where FPR is 0 on the x-axis is the perfect classification
 - o On y = x is 50/50, no better than random
- AUC
 - o Calculate area under ROC curve
 - Lift charts
 - o Measure of effectiveness of a predictive model
 - o Lift factor = success rate with model/success rate without model
- The probability of any chocolate bar chosen at random being "good" is 3/7.
 - If we just sample the top 2 (28%) we're most confident of then the probability of "good" is 2/2.

$$\text{Lift} = \frac{\left(\frac{2}{2}\right)}{\left(\frac{3}{7}\right)} = \frac{7}{3} = 2.33 \text{ or } 233\%$$

- For top 3 (42%), P(good) = 2/3

$$\text{Lift} = \frac{\left(\frac{2}{3}\right)}{\left(\frac{3}{7}\right)} = \frac{14}{9} = 1.56 \text{ or } 156\%$$

Brand	Actual Class	Conf
Aero	0	0.3
Cherry Ripe	0	0.4
Kitkat	0	0.4
Bounty	1	0.6
Snickers	0	0.7
Flake	1	0.8
Violet Crumble	1	1.0

WEEK 9 – ENSEMBLE METHODS AND ANN

Ensemble methods

- Original training data – Create multiple data sets – Build multiple classifiers – Combine classifiers

- Works best when:
 - o Individual classifiers are moderately accurate
 - o Individual classifiers are not correlated
 - o Decision trees work well as individual classifiers
- Bagging
 - o Multiple replicates of original data by sampling, with replacing, from the training set
 - o Combine the classifiers by taking a majority vote to produce the final version

