# ISYS30221 Artificial Intelligence:
# TOPIC A: Machine Learning Methods for Data Mining

## I. Overview

The digital revolution has made data easy to digitally capture and inexpensive to store. Subsequently, many organisations and businesses are increasingly collecting data about various aspects of their operation, services, customers, suppliers and competitors with the hope of finding useful information with which to gain leverage or competitive advantage. As a result, **many organisations are storing between 20-100 terabytes of data _finding themselves data rich yet information poor_** with a recent study by IBM showing that 1:3 business leaders do not trust their data. This is likely to continue with IBM estimating that 2.3 trillion gigabytes of data is being created per day, expecting this to increase to 43 Trillion GB (40 Zetta bytes) by 2020. This has _led to an exponential demand (912% increase!) for professionals with expertise in managing and analysing big data sets_. This is particularly true for the following domains:

- **Business and Finance:** identifying and interpreting buyer/seller behaviour patterns, risk assessment data, product distribution patterns, foreign exchange rates, stock indexes and prices, interest rates, credit card data, fraud;
- **Government:** identifying behaviours and data patterns relating to: policy changes, tax evasion, benefit fraud, terrorism, social problems, law and order;
- **Internet:** extracting useful information to support personal or business decisions;
- **Health Care:** interpreting diverse diagnostic and treatment/patient care information stored by hospital management systems and biological/biometric data to inform patient care and drug development;
- **Telecommunication network:** Calling patterns and fault management systems.

The aim of this part of the module is therefore to **introduce you to a number of 'machine learning' techniques to process and discover patterns in data**. Machine learning allows computational systems to adaptively improve their performance with experience accumulated from the observed data rather than to use pre-defined rules.

The content is themed around the notion of data as a resource; how it is acquired, prepared for analysis and finally how we can learn from it. A practical approach to teaching machine learning will be adopted and you will therefore use the freely available WEKA data mining package ([www.cs.waikato.ac.nz/ml/weka](www.cs.waikato.ac.nz/ml/weka)) to put theory into practice. WEKA is full industrial-strength implementation of the techniques taught on this module and will enable you to better understand the techniques taught, appreciate their strengths and applicability to various problem domains.

The **essential reading** for this section of the module is as follows:

- Witten, I., Frank, E., Hall, M and Pal, C (2016), _Data Mining: Practical Machine Learning Tools and Techniques,_ Morgan Kaufmann; 4th edition (22nd Dec 2016), ISBN-13: 978-0128042915.

  _NB we will make use of Witten's youtube videos (links given in this document and the lab sheets). You are strongly encouraged to watch ALL of the youtube videos recommended_

You will also find the following textbooks <u>very</u> useful:

- Samarasinghe, S. (2007) <u>Neural Networks for Applied Sciences and Engineering</u>, First edition, Auerback Publishing, ISBN: 9780849333750

Online copies of <u>both</u> textbooks are available **for free** in the Resource List on NOW. Also, a number of supporting case studies and research papers will be made available on NOW for each of the taught topics.  Note that professionals in the data mining field typically follow: http://www.kdnuggets.com/ - you too may find it useful too!

The content will be delivered over 9 weeks and organised into the following four parts:

       Part 1. Introduction to Knowledge Discovery in Databases
       Part 2. The Basic Components of Machine Learning Models
       Part 3. Supervised Machine Learning Methods for Classification and Prediction
       Part 4. Unsupervised Machine Learning Methods for Clustering

It is important to state that you will be exposed to a broad array of approaches to data analysis in a very short space of time, which will initially appear rather daunting! However, your assessment will enable you to focus your learning and thus develop a deeper understanding of a specific set of methods, models and techniques as appropriate to your chosen data set(s) and analytics task. You are therefore NOT expected to become an expert in every machine learning or data analytics technique!

The remainder of this document provides an indicative teaching and learning schedule for the 9-week teaching period for this section of the module. You should read it carefully and ensure you follow the directed reading. If you have any queries about this section of the module, then please contact me.

**Dr Jon Tepper**
**Room ERD284/Office hr: Tuesday, 1100-1230/E: Jonathan.Tepper@ntu.ac.uk/ T: 0115 8488363**

## II. Indicative Teaching and Learning Schedule

The indicative teaching and learning schedule for Topic A of the ISYS30221 Artificial Intelligence module is given in table 1 below. (Note that this schedule may change under exceptional circumstances. Every effort, however, will be made to notify you beforehand).

| Uni Week | Week Comm. | Indicative Content and Reading | | |
|---|---|---|---|---|
| | | **Lecture** | **Laboratory** | **Core Reading/viewing** |
| 11 | 09/10/17 | **Part 1. Introduction to Data Analytics and CRISP-DM**<br>a) The CRISP-DM process<br>b) Getting to know your data | N/A | Witten et al (2016): pg 3-65<br>See NOW for additional research papers and case studies including: BBC Horizon programme: The Age of Big Data (watch it here…you will have to register and may have to search for the programme):<br>http://bobnational.net/record/155036/media_id/158510 )<br>Also listen to the following Radio4 programme:<br>http://www.bbc.co.uk/programmes/b03kqfzx<br>..and watch the Click programme on AI:<br>http://www.bbc.co.uk/programmes/b06dst0x<br><br>**Useful articles/sites:**<br>https://www.youtube.com/watch?v=hZxnzfnt5v8 (types of data)<br>http://changingminds.org/explanations/research/sampling/sampling_terminology.htm#frac (quick ref to relevant stats sampling terms)<br>https://www.analyticsvidhya.com/blog/2017/09/6-probability-distributions-data-science/ (common distributions we should know) |
| 12 | 16/10/17 | **Part 2. The Basic Components of Machine Learning Models**<br>a) Defining the learning problem<br>b) Components of any machine learning system | N/A | Witten et al (2016): pg 43-88 (See NOW for additional research papers and case studies)<br>http://www.youtube.com/watch?v=b4zr9Zx5WiE (Peter Norvig from Google, Machine learning hottest tech trends in 3-5 years)<br>http://www.tutorialspoint.com/data_mining/index.htm (general introduction to data mining and machine learning)<br>https://www.youtube.com/watch?v=yDLKJtOVx5c (Machine learning and its applications) |
| 13 | 23/10/17 | c) Model Fitting, Selection and Evaluation<br>• Training and testing<br>• Predicting performance<br>• Cross-validation and bootstrapping<br>• Confusion matrix / ROC curves<br>• Calculating Recall and precision | 1. The Explorer<br>• Preparing and loading given data set(s)<br>• Reviewing data filtering algorithms<br>• Reviewing the Learning Algorithms in Weka – noting different types of algorithms (Bayes, Trees, Rules, Functions, Lazy, MI and MIsc) | Witten et al (2016): pg 161-202, 479-501<br>http://www.youtube.com/watch?v=hihuMBCuSIU (Model Selection with Cross Validation)<br>https://www.youtube.com/watch?v=jiQamxz2ZcQ&nohtml5=False (Bias-variance trade-off)<br>https://www.youtube.com/watch?v=Z5TtopYX1Gc (Sensitivity vs Specificity)<br>http://www.youtube.com/watch?v=lHa1UYAxGxs (ROC curve #1)<br>http://www.youtube.com/watch?v=sWAJsiVh1Gg (ROC curve #2)<br>https://www.youtube.com/watch?v=vbpm9-WYuzU (Precision and Recall)<br>https://www.youtube.com/watch?v=0zZYBALbZgg (Disproving through Hypothesis |

| Uni Week | Week Comm. | Indicative Content and Reading | | |
|---|---|---|---|---|
| | | **Lecture** | **Laboratory** | **Core Reading/viewing** |
| | | • Evaluating numeric prediction (e.g. mean squared error (MSE), root mean squared error (RMSE), correlation coefficients) | • Pre-processing data and creating your own data sets in Arrf format<br><br>2. The Knowledge Flow Interface<br>• Getting started – add data, select attribute and tree method<br>• Reviewing Knowledge Flow components: Visualisation and Evaluation<br>• Configuring and connecting components<br>• Running the experiment and interpreting the output | testing)<br>https://www.youtube.com/watch?v=eyknGvncKLw (p values)<br>https://www.youtube.com/watch?v=lwpobQmUTd8&nohtml5=False (confidence intervals)<br>https://www.youtube.com/watch?v=fXOS4Q3nJQY (z scores)<br>https://www.youtube.com/watch?v=0Pd3dc1GcHc&nohtml5=False (t-tests)<br>https://www.youtube.com/watch?v=UmAJJtEo6cQ (z or t scores?)<br>https://www.youtube.com/watch?v=1Ldl5Zfcm1Y&nohtml5=False (Chi squared test)<br>https://www.youtube.com/watch?v=t2ryZyytW5w (two means t- test in Excel)<br>https://www.youtube.com/watch?v=c68JLu1Nfkw&nohtml5=False (Evaluating numerical prediction models)<br>https://www.youtube.com/watch?v=JxweGr03W9g&list=PLea0WJq13cnCZZ3sXVEZ2OE5CLeZUlCmm&nohtml5=False (Model Evaluation) |
| 14 | 30/10/17 | **Part 3. Supervised Machine Learning Methods for Classification and Prediction**<br>a) Keep it Simple Stupid…probably!<br>• Occams Razor<br>• 1R - inferring 1 rule from data<br>• Naive Bayes<br>– Quick tour of probability theory<br>– Bayes Theorem<br>– Naive Independence assumption<br>– Example classification<br>– Limitations<br><br>**Surgery Session #1 (Getting familiar with the coursework:**Mark sample coursework submissions**)** | | Witten et al (2016): pg 96-105<br>http://www.youtube.com/watch?v=-8eSOmTPUbk (Intro to probability)<br>https://www.youtube.com/watch?v=yi97YB8EyDg&nohtml5=False (Conditional prob)<br>https://www.youtube.com/watch?v=3XL6w4_EKrg (Intro to Bayes Theorem)<br>http://www.youtube.com/watch?v=wxmNvH7XBOA (Mario Kart and Bayes!) |
| 15 | 06/11/17 | b) Divide and Conquer with Decision Trees<br>• What is a decision tree?<br>• Learning from Claude Shannon: Father of the | 3. The Experimenter<br>• Getting started – add data and models to evaluate<br>• Running an experiment<br>• Analysing the results | Witten et al (2016): pg 105-112, 209-221<br>http://www.youtube.com/watch?v=wL9aogTuZw8 (overview of ID3 algorithm)<br>http://www.youtube.com/watch?v=-dCtJjIEEgM (Lecture from Uni of Columbia)<br>https://www.youtube.com/watch?v=z2Whj_nL-x8&noredirect=1 (The genius of Claude Shannon) |

| | | Indicative Content and Reading | | |
|---|---|---|---|---|
| Uni Week | Week Comm. | Lecture | Laboratory | Core Reading/viewing |
| | | information age <br> • Learning decision trees from examples <br> – Quinlan's ID3 algorithm <br> – Limitations and solutions <br> – Applications | • Revisiting The Explorer <br> • Training, testing and visualising different machine learning schemes | |
| 16 | 13/11/17 | c) Simple Linear Regression for Numerical Prediction <br> • What is linear regression? <br> • Simple linear regression <br> – Ordinary least squares technique <br> – Approach 1: using original data <br> – Approach 2: using redefined data <br> – Making predictions <br> • Calculating correlation coefficients <br> – Coefficient of determination (R2) <br> – Pearson's coefficient of correlation (R) <br> – Spearman's coefficient of rank correlation (Rs) <br> • Evaluating the fit <br> Other forms of linear regression <br><br> **Surgery Session #2 (Share and discuss coursework ideas)** | 4. Numerical Prediction <br> • Reviewing numerical data sets <br> • Using Weka for numerical prediction using simple linear regression <br> • Evaluating what's been learnt | Witten et al (2016): pg 128-133 <br> https://www.youtube.com/watch?v=ZkjP5RJLQF4 (Brandon Foltz Lec #1 - Basics) <br> https://www.youtube.com/watch?v=iAgYLRy7e20 (Brandon Foltz Lec #2 - Algebra) <br> https://www.youtube.com/watch?v=Qa2APhWjQPc&nohtml5=False (Brandon Foltz Lec #3 – Least Squares) <br> https://www.youtube.com/watch?v=kHZBy1uVNnM&nohtml5=False (Brandon Foltz Lec #4 – Evaluating fit) <br> https://www.youtube.com/watch?v=dQNpSa-bq4M&list=PLIeGtxpvyG-IqjoU8IiF0Yu1WtxNq_4z-&nohtml5=False (Brandon Foltz Lec on Multiple Regression) <br> https://www.youtube.com/watch?v=Y2khrpVo6qI&nohtml5=False (Time series processing using ARIMA models) |
| 17 | 20/11/17 | d) Artificial Neural Networks (supervised learning) <br> • Basic principles <br> • The Perceptron <br> • Problems of Linear Separability <br> • Multi-layered Perceptrons (MLP) trained with Back- | 5 & 5S Artificial Neural Networks <br> • Basic matrix algebra for understanding artificial neural networks <br> • Working with the perceptron for simple logic gates <br> • Understanding the basics Back-propagation | Witten et al (2016): pg  260-272, 417-466 <br> Samarasinghe (2007): 1-151; 195-220; 245-261. <br> http://www.youtube.com/watch?v=DG5-UyRBQD4 (Understanding neural nets without the maths) <br> https://www.youtube.com/watch?v=S3iQgcoQVbc&nohtml5=False (Perceptrons for simple classification) <br> https://www.youtube.com/watch?v=0qVOUD76JOg&nohtml5=False (TED talk on neural nets) |
| 18 | 27/11/17 | | | |

| Uni Week | Week Comm. | Indicative Content and Reading | | |
|---|---|---|---|---|
| | | **Lecture** | **Laboratory** | **Core Reading/viewing** |
| | | propagation<br>• Applications<br>**Surgery Session #3 (Getting familiar with coursework:** Remark sample coursework scripts**)** | • Getting back to Weka<br>  Using MLP with Back-propagation to classify and perform prediction (to be completed as homework if not during the lab). | https://www.youtube.com/watch?v=u5GAVdLQyIg    and https://www.youtube.com/watch?v=IlmNhFxre0w (Engaging lectures on multi-layered perceptrons)<br>https://www.youtube.com/watch?v=aVId8KMsdUU&list=PL29C61214F2146796&nohtml5=False (Advanced - derivation of the Backpropagation Learning algorithm for the interested student!) |
| 19 | 04/12/17 | **Part 4. Unsupervised Machine Learning Methods for Clustering**<br>a) Simple K-means<br><br>**Independent study for the interested student:**<br>b) Kohonen's Self-organising Maps (SOMs) | 6. Clustering in WEKA<br>• Clustering a simple data set using K-means and other algorithms<br>• Evaluating the clusters | K-means:<br>Witten et al (2016): pg 141-156.<br>https://www.youtube.com/watch?v=zHbxbb2ye3E&nohtml5=False (Fantastic illustration/visualisation of the K-means algorithm)<br><br>For interested student (SOMs):<br>Samarasinghe (2007): 337-436.<br>https://www.youtube.com/watch?v=H9H6s-x-0YE&nohtml5=False (Understanding SOMs within 6 mins!)<br>https://www.youtube.com/watch?v=HKj2ASG0DKQ&ebc=ANyPxKqz60EIWrF9Dg4V0Fzy-XE58j2vEEdXiIaGdMpyRpJ8fNFjbsQCkIYQbWzU_GqXyd3pX8tWDQGOT2Z7u9D7gfbPPX5rYQ&nohtml5=False (More detailed treatment of SOMs)<br>https://www.youtube.com/watch?v=iWPhGKniTew (Lecture by the master himself, Teuvo Kohonen) |

**Table 1.** Indicative content and reading schedule for Topic A. Machine Learning Methods for Data Mining. You should make every effort to view ALL of the recommended youtube videos.