

# COURSE OUTLINE

Course Title	<b>Introduction to Data Science and Artificial Intelligence</b>		
Course Code	<b>PS0002</b>		
Offered	Study Year 1, Semester 2		
Course Coordinator	Xiang Liming (Assoc Prof)	lmxiang@ntu.edu.sg	6513 7451
Pre-requisites	PS0001 or CZ1003 or CZ1103		
AU	3		
Contact hours	Lectures: 26, Tutorials: 13, Laboratories: 10		
Approved for delivery from			
Last revised	08 Sept 2023		

## Course Aims

In this era of information, vast amounts of new data are produced every day from various fields including scientific research, healthcare, industry and service processes. Using data effectively and extracting meaningful insights from data can significantly improve efficiencies, cut costs and add more value to organizations. This course aims to provide you with an understanding of basic techniques for data analysis, machine learning and dimension reduction for big data, and expose you to hands-on computational tools that are fundamental for data science. Besides supervised and unsupervised learning, another fundamental technology of Artificial Intelligence – reinforcement learning will also be introduced, including Markov decision process and Q-learning.

Suited for anyone from different backgrounds, this course will show you how you could apply various methods to data examples and case studies from both research and industrial sources in the Singapore context.

## Intended Learning Outcomes

Upon successfully completing this course, you should be able to:

1. Identify interesting data-driven questions in your respective field of study.
2. Formulate meaningful study problems that you want to explore.
3. Collect/extract relevant data, visualize and perform exploratory analysis on data.
4. Perform machine learning models to extract meaningful insights from data.
5. Implement above techniques with R.
6. Present your analysis results and problem solution via an engaging written communication.

## Course Content

What is data science? Why it is essential? Analytic thinking

Study problem formulation, Data acquisition, Data wrangling

Coding essentials with R

Data visualization, Exploratory data analysis

Statistical insight

Machine learning models for prediction, Linear regression

Machine learning models for classification: Linear classifier, Logistic regression

Support vector machine

Machine learning for clustering: K-means method

Other clustering technics

Dimension reduction: Curse of dimensionality, Principal component analysis

The state of art, Neural networks, Deep learning

Written communication for data analysis results

## Assessment

Component	Course ILOs tested	Weighting	Team / Individual	Assessment Rubrics
<b>Continuous Assessment</b>				
<b>Laboratories</b>				
Assignment	3, 4, 5, 6	20	team	See Appendix for rubric
<b>Tutorials</b>				
Quiz 1	1, 2, 3, 4, 5, 6	15	individual	See Appendix for rubric
Quiz 2	2, 3, 4, 5	15	individual	See Appendix for rubric
<b>Examination (2 hours)</b>				
Short Answer Questions	2, 3, 4, 5, 6	50	individual	See Appendix for rubric
<b>Total</b>		<b>100%</b>		

## Formative Feedback

You will receive written and verbal feedback from the lecturer for lab assignments and written examinations.

You will receive written and verbal feedback from the lecturer on the quizzes.

## Learning and Teaching Approach

<b>Lectures</b> (26 hours)	Lectures will deliver the theoretical knowledge required to understand various components involved in data analysis and machine learning process.
<b>Tutorials</b> (13 hours)	Labs sessions will: <ul style="list-style-type: none"><li>• Demonstrate practical applications of data science techniques in various application fields.</li><li>• Enable you to code using R and address any coding issues.</li></ul>
<b>Laboratories</b> (10 hours)	Labs sessions will: <ul style="list-style-type: none"><li>• Demonstrate practical applications of data science techniques in various application fields.</li><li>• Enable you to code using R and address any coding issues.</li></ul>

## Reading and References

The main references/textbooks relevant to the course materials are:

1. Gareth James, Daniela Witten, Trevor Hastie and Robert Tibshirani (2014). An Introduction to Statistical Learning: with Applications in R, Springer New York. Available at <http://www-bcf.usc.edu/~gareth/ISL/> .
2. Hastie, T., Tibshirani, R. and Friedman, J. (2017). The Elements of Statistical Learning: Data Mining, Inference and Prediction, Springer. <https://web.stanford.edu/~hastie/ElemStatLearn>
3. Kabacoff, Robert (2011). R in Action: Data analysis and graphics with R, Manning.
4. Stuart Russell and Peter Norvig (2016) Artificial Intelligence: A Modern Approach, 3rd edition. Pearson

## Course Policies and Student Responsibilities

### (1) General

You are expected to complete all assigned pre-class readings and activities, attend all seminar classes punctually and take all scheduled assignments and tests by due dates. You are expected to take responsibility to follow up with course notes, assignments and course related announcements for seminar sessions they have missed. You are expected to participate in all seminar discussions and activities.

### (2) Absenteeism

Absence from class without a valid reason will affect your overall course grade. Valid reasons include falling sick supported by a medical certificate and participation in NTU's approved activities supported by an excuse letter from the relevant bodies.

If you miss a lecture, you must inform the course instructor via email prior to the start of the class.

### (3) Absence Due to Medical or Other Reasons

If you are sick and not able to attend a quiz or midterm, you have to submit the Medical Certificate (or another relevant document) to the administration to obtain official leave. In this case, the missed assessment component will not be counted towards the final grade. There are no make-up quizzes or make-up midterm.

## Academic Integrity

Good academic work depends on honesty and ethical behaviour. The quality of your work as a student relies on adhering to the principles of academic integrity and to the NTU Honour Code, a set of values shared by the whole university community. Truth, Trust and Justice are at the core of NTU's shared values.

As a student, it is important that you recognize your responsibilities in understanding and applying the principles of academic integrity in all the work you do at NTU. Not knowing what is involved in maintaining academic integrity does not excuse academic dishonesty. You need to actively equip yourself with strategies to avoid all forms of academic dishonesty, including plagiarism, academic fraud, collusion and cheating. If you are uncertain of the definitions of any of these terms, you should go to the [Academic Integrity website](#) for more information.

On the use of technological tools (such as Generative AI tools), different courses / assignments have different intended learning outcomes. Students should refer to the specific assignment instructions on their use and requirements and/or consult your instructors on how you can use these tools to help your learning. Consult your instructor(s) if you need any clarification about the requirements of academic integrity in the course.

## Course Instructors

Instructor	Office Location	Phone	Email
Xiang Liming (Assoc Prof)	SPMS MAS04-11	6513 7451	lmxiang@ntu.edu.sg

## Planned Weekly Schedule

Week	Topic	Course ILO	Readings/ Activities
1	Coding essentials with R	1	Readings: Lecture notes, reference book of Kabacoff (2011).  Demo of R programming essentials.
2	What is data science? Why it is essential? Analytic thinking Study problem formulation, Data acquisition, Data wrangling	1, 2, 3	Readings: Lecture notes, reference book of Kabacoff (2011), including multiple case studies with data science applications in various fields of study.  Activity: after reading, students will be asked to describe the application scenarios aligned with their major.  Demo of R package dplyr for data wangling and preparation.
3	Data visualization	1, 2, 3	Readings: Lecture notes, reference book of Kabacoff (2011).  Show visualization tools using packages dplyr and ggplot2 for visualizing and exploring data.

4	Exploratory data analysis	1, 2, 3,4	<p>Readings: Lecture notes, reference book of Kabacoff (2011).</p> <p>Demo of more methods for visualizing data and exploratory data analysis.</p> <p>Activity: Lab Assignment 1</p>
5	Machine learning models for prediction, K-nearest neighbors regression	1, 2, 3, 4, 5	<p>Readings: Lecture notes, reference books of James et al (2014) and Hasti et al (2017).</p> <p>Demo of kNN regression in R</p>
6	Machine learning models for prediction, Linear regression	1, 2, 3, 4, 5	<p>Readings: Lecture notes, reference books of James et al (2014) and Hasti et al (2017).</p> <p>Demo of linear regression analysis using lm() and model diagnostics</p>
7	Machine learning models for classification: Logistic regression	1, 2, 3, 4, 5	<p>Readings: Lecture notes, reference books of James et al (2014) and Hasti et al (2017).</p> <p>Show the use of function glm() for classification based on logistic regression.</p> <p>Quiz 1 will be conducted.</p>
8	Machine learning models for classification: k-nearest neighbors	1, 2, 3, 4, 5	<p>Readings: Lecture notes, reference books of James et al (2014) and Hasti et al (2017).</p> <p>Show the use of R package MASS, GLM kNN for classification</p>
9	Machine learning models for classification: Support vector machine	1, 2, 3, 4, 5	<p>Readings: Lecture notes, reference books of James et al (2014) and Hasti et al (2017).</p> <p>Demo of the SVR model with R package e1071.</p>
10	Machine learning models for classification: Support vector machine	1, 2, 3, 4, 5	<p>Demo of the SVR with other kernel functions for improving classification performance.</p> <p>Activity: Lab Assignment 2</p>
11	Unsupervised machine learning for clustering: K-means method	1, 2, 3, 4, 5,6	<p>Readings: Lecture notes, reference books of James et al (2014) and Hasti et al (2017).</p> <p>Demo of K-means clustering kmeans() and compare with classification</p>
12	Unsupervised machine learning for clustering: Hierarchical clustering	1, 2, 3, 4, 5, 6	<p>Readings: Lecture notes, reference books of James et al (2014) and Hasti et al (2017).</p> <p>Demo of more examples for clustering using R function hclust().</p>

			Quiz 2 will be conducted.
13	Review	1, 2, 3, 4, 5, 6	Readings: Lecture notes; various case studies and corresponding sample papers

## Appendix 1: Assessment Rubrics

### *Rubric for Laboratories: Assignment (20%)*

Levels of Performance	Criteria Description
A+ (Exceptional) A (Excellent)	Provides clear, efficient, working and well-documented code; evidence of programming understanding and concern for code efficiency beyond getting correct solution. Demonstrated ability to develop multiple approaches to programming task, and understanding of their respective advantages.
A- (Very good) B+ (Good)	Provides clear, efficient, working and well-documented code; evidence of programming understanding.
B (Average) B- (Satisfactory) C+ (Marginally satisfactory)	Working but limited documentation of code.
C (Bordering unsatisfactory) C- (Unsatisfactory)	Write the code with lots of help from TA and instructor. Limited code documentation or demonstration of conceptual understand.
D (Deeply unsatisfactory) F (0-44)	Lack of demonstrated conceptual understanding. Non-functional code.

### *Rubric for Tutorials: Two Quizzes (30%)*

Point-based marking (not rubrics based)

### *Rubric for Examination: Short Answer Questions (50%)*

Point-based marking (not rubrics based)