



The Open University of Sri Lanka  
Faculty of Engineering Technology  
Department of Electrical and Computer Engineering

186

Study Programme	: Bachelor of Software Engineering Honours
Name of the Examination	: Final Examination
<b>Course Code and Title</b>	<b>: EEX4373 Data Science</b>
Academic Year	: 2022/23
Date	: 04 <sup>th</sup> November 2023
Time	: 0930-1230hrs
Duration	: <b>3 hours</b>

### General Instructions

1. Read all instructions carefully before answering the questions.
2. This question paper consists of **Five (5) questions** in **five (5) pages**.
3. Answer **ALL** questions. All questions carry equal marks.
4. This is a Closed Book Test (**CBT**).
5. Answers should be in clear hand writing.
6. Do not use Red colour pen.

---

### Question 1

Machine learning automates the process of data analysis. It helps predictions based on data in real time without human intervention. There are three (03) basic types of Machine learning methods namely, supervised, unsupervised and reinforcement learning methods.

- (a) Describe the key features in reinforcement learning with an example. [04 marks]
- (b) Describe the clustering and association learning methods in unsupervised machine learning methods with an example for each. [06 marks]
- (c) Figure 1 contains the process from preparing the data to implementation of a new model. Copy the labels into the answer book and name them and briefly explain the task carried out in each step. [07 marks]

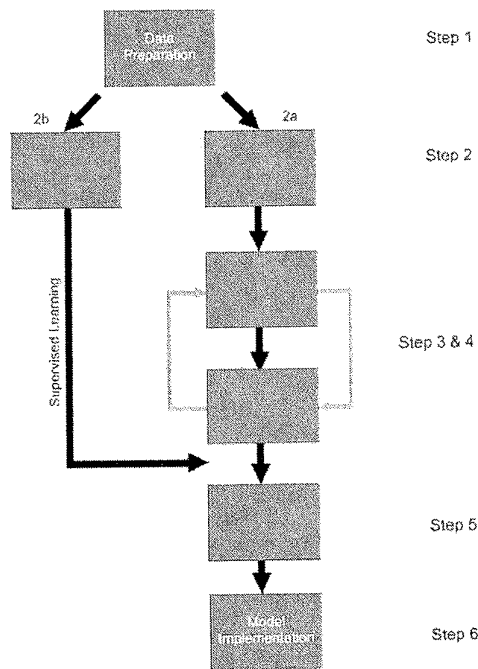


Figure 1- Machine Learning process in Data Science

- (d) When developing a machine learning model, the dataset is split into training and testing sets. Explain why this step is important in machine learning. [03 marks]

### Question 2

Assume that you are a data scientist working for a retail company, and you've been provided with a dataset containing information about the company's sales and customer behavior. The dataset includes various attributes such as customer demographics, purchase history, and product details etc. Your goal is to perform exploratory data analysis to gain insights from the data.

- (a) Describe the process of Exploratory Data Analysis (EDA) that you can carry out for the above-described company data set. [04 marks]
- (b) Describe four (04) key statistics, including their key characteristics, that can be used in exploratory data analysis. [04 marks]
- (c) Clearly describe three (03) techniques that can be used to handle missing data with advantages and disadvantages of each of the methods. [06 marks]
- (b) Assume that you need to create a few charts for each attribute for better understanding. Assuming that you have identified histograms, scatter plots and bar charts as the most suitable types of charts. Compare and contrast each of the selected charts taking into consideration the attributes from the given scenario. [06 marks]

### Question 3

Imagine you are working as a data scientist for a social media company, and your task is to build a machine learning model to classify user posts into two categories: "Spam" and "Not Spam." The company is facing a growing issue of getting more and more spam posts, which are being manually identified at the present. The company wants an automated solution to filter out spam content. You decide to use the k-Nearest Neighbors (KNN) algorithm for this classification task.

- (a) Explain the steps of the KNN algorithm in the context of this classification problem. [06 marks]
- (b) Describe how you would choose the optimal value of 'k' for the KNN model. [04 marks]
- (c) What are the advantages and disadvantages of using KNN for this classification task? [04 marks]
- (d) During model development, it is noted that the model is overfitted.
  - (i) Explain the concept of regularization in the context of model fitting. [02 marks]
  - (ii) Describe two specific techniques that can be used to mitigate overfitting problem in this scenario. [04 marks]

### Question 4

Assume that you are a data scientist working for a large super market chain, and you are tasked with analyzing customer purchase behavior. The marketing team is interested in segmenting customers based on their purchasing habits to create targeted marketing campaigns. You decide to use the k-Nearest Neighbors (KNN) clustering algorithm for this purpose.

- (a) Explain the steps on how the KNN clustering algorithm works to segment the customers in the above given context. [04 marks]
- (b) Describe the key factors to consider when choosing the optimal value of 'k' for KNN clustering. [02 marks]
- (c) Write two (02) advantages and two (02) disadvantages of using KNN clustering for customer segmentation? [04 marks]
- (d) Assume that you are given a set of data points to form two clusters using K-means clustering method. The data points are,

(0,0), (0,1), (1,1), (1,0), (0.5,0.5), (5,5), (5,6), (6,6), (6,5), (5.5,5.5)

Choose the first two data points as seed points and initialize the clusters as  $C_1 = \{(0,0)\}$  and  $C_2 = \{(0,1)\}$ .

