

IEMS 5780
Building and Deploying Scalable Machine Learning Services
2019-2020 Term 1

Final Examination

Release date : 2019-12-09 (Monday)

Deadline : 2019-12-14 (Saturday) Noon (12:00)

Instructions

1. The full mark of this examination is 100.
2. There are altogether FOUR questions in this paper, each with sub-questions.
3. You should attempt all questions.
4. Give concise answers and avoid lengthy paragraphs.
5. When working on the exam paper, you can refer to any materials including lecture slides, books, and online resources.
6. You should finish the examination on your own, and should NOT discuss with your classmates during the examination period.
7. You should submit a TEXT document (.txt file) containing your answers to the questions to Blackboard. Refer to the template released together with this paper.

Question 1 - Machine Learning Application (24 Marks)

You are working in a bank in Hong Kong. The bank is interested in predicting which customers are more likely to borrow money, so that the salespersons can be more focused when calling customers to promote the bank's products. Your supervisor has assigned you to study how to achieve this task using machine learning.

- a) **(3 Marks)** List at least THREE types of data that are available inside the bank that will be useful in this task.
- b) **(16 Marks)** You have decided that supervised learning can be used to train a model to predict whether a customer will borrow money. Describe how you would implement such a model. You can focus on the following issues when giving your answer:
 - i) How would you prepare the data for training the model?
 - ii) Name some important features that will be used to represent the customers
 - iii) What will be the input and output of your model?
 - iv) How will you evaluate the model? (e.g. What experiments will you do? What metric(s) will you use?)
- c) **(5 Marks)** Can you suggest any additional data publicly available that can contribute to this task? Describe how they can be used.

Question 2 - Text Classification (20 Marks)

You are working for a e-commerce Website. Customers of the Website can purchase various kinds of products online, and they can leave comments on the products, which will be publicly available to any user of the Website. Your supervisor would like you to develop a text classification model to automatically category the comments of the users. Examples of categories include “product quality”, “usability of the Website”, “quality of delivery service” and “product price”.

- a) **(3 Marks)** You are given a small sample of user comments collected on the Website. Suggest THREE pieces of information you would like to extract from this dataset (by performing some analysis), which will be useful for implementing the machine learning model later.
- b) **(3 Marks)** Suggest some preprocessing steps that you will apply to the data before using them in your model training.
- c) **(4 Marks)** Describe how you can convert a comment (represented as a string of characters) into a feature vector to be fed into the machine learning model.
- d) **(4 Marks)** Name TWO different machine learning models that can be used to perform supervised learning in this task. Briefly describe how inference is done for these two models.
- e) **(6 Marks)** You trained a model and it achieved 90% accuracy on the test dataset.
 - i) What other metric(s) would you check to ensure that your model's performance is good enough?
 - ii) Give an example to explain why the model is bad even when it achieved an overall accuracy of 90%.

Question 3 - Recommendation System (20 Marks)

You are working in a company that is operating an app in which users can share book reviews (comments on books that they have read). As the numbers of users of the app and books with reviews have increased significantly recently, the company has decided to implement a recommendation system to recommend books to the users. You are assigned this task.

- a) **(6 Marks)** When you review the data collected in the app, you realize that the app has not allowed users to rate the books they have read. Users can only write reviews and submit to the app. But in order to understand the tastes of the users, you need to know whether the users think positively or negatively about the books. How would you solve this problem using machine learning?
- b) **(12 Marks)** Now assume that you now have ratings on books given by the users (e.g. a user may assign a rating of 10 to books he likes very much, and a rating of 1 if he dislikes the book).

- i) Name two different algorithms that can be used to generate recommendations for the users. Describe how they work.
 - ii) Assume that your system will generate a predicted rating for each (user, book) pair. How would you evaluate the performance of your system? What metric will you use? Give some examples of true and predicted ratings and show how the metric is calculated.
- c) **(2 Marks)** For users who have never written any reviews and rated any book in your app, how would you generate a list of recommended books to him/her?

Question 4 - Network Programming & Concurrent Programming (36 Marks)

- a) **(6 Marks)** Explain the difference between TCP (Transmission Control Protocol) and UDP (User Datagram Protocol).
- b) **(4 Marks)** Explain what is the usage of a port. Why is a port number necessary?
- c) **(8 Marks)** When sending a message using TCP from one program to another program running on another computer, we need to agree on the method of determining when a message is fully received on the other side. Describe two ways of delimiting a message, and discuss their disadvantages.
- d) **(6 Marks)** Imagine a system in which there are multiple server programs running on different server machines. Over time, some machines will be disabled while new machines will be added to run additional server programs. All server programs offer the same function or service to the client programs. Suggest a method for implementing this system (the clients and the servers) such that requests from the clients will be evenly distributed to the active server programs.
- e) **(6 Marks)** Consider a computer with four CPUs. You have a program for processing a large number of text documents stored in the local hard disk. Explain why using multiprocessing in your program will be a better choice than using multi-threading, assuming that you are using Python to develop your program.
- f) **(6 Marks)** Explain why asynchronous programming will make a program run faster in some cases than using sequential programming. Give an example to illustrate your points.