

CS 643, Cloud Computing

Programming Assignment 2

Goal: The purpose of this individual assignment is to learn how to develop parallel machine learning (ML) applications in Amazon AWS cloud platform. Specifically, you will learn: (1) how to use [Apache Spark](#) to train an ML model in parallel on multiple EC2 instances; (2) how to use [Spark's MLlib](#) to develop and use an ML model in the cloud; (3) How to use [Docker](#) to create a container for your ML model to simplify model deployment.

Description: You have to build a wine quality prediction ML model in Spark over AWS. The model must be trained in parallel using 4 EC2 instances. Then, you need to save and load the model in a Spark application that will perform wine quality prediction; this application will run on one EC2 instance. The assignment must be implemented in Java on Ubuntu Linux. The details of the assignment are presented below:

- Input for model training: we share 2 datasets with you for your ML model. Both datasets are available in Canvas, under Programming Assignment 2.
 - TrainingDataset.csv: you will use this dataset to train the model in parallel on multiple EC2 instances.
 - ValidationDataset.csv: you will use this dataset to validate the model and optimize its performance (i.e., select the best values for the model parameters).
- Input for prediction testing: TestDataset.csv. We will use this file, which has a similar structure with the two datasets above, to test the functionality and performance of your prediction application. Your prediction application should take such a file as input. This file is not shared with you, but you can use the validation dataset to make sure your application works.
- Output: The output of your application will be a measure of the prediction performance, specifically the F1 score, which is available in MLlib.
- Model Implementation: You have to develop a Spark application that uses MLlib to train for wine quality prediction using the training dataset. You will use the validation dataset check the performance of your trained model and to potentially tune your ML model parameters for best performance. You should start with a simple linear regression or logistic regression model from MLlib, but you can try multiple ML models to see which one leads to better performance. For classification models, you can use 10 classes (the wine scores are from 1 to 10). Note: there will be extra-credit for the top 5 applications/students in terms of prediction performance (see below under grading).
- Docker container: You have to build a Docker container for your prediction application. In this way, the prediction model can be quickly deployed across many different environments.
- The model training is done in parallel on 4 EC2 instances.
- The prediction with or without Docker is done on a single EC2 instance.

Submission: You will submit in Canvas, under Programming Assignment 2, text/Word/pdf file that contains:

- A link to your code in [GitHub](#). The code includes the code for parallel model training and the code for prediction application.
- A link to your container in [Docker Hub](#).

This file must also describe step-by-step how to set-up the cloud environment and run the model training and the application prediction. For the application prediction, you should provide instructions on how to run it with and without Docker.

Grading:

- | | |
|---|-------------|
| - Parallel training implementation | – 50 points |
| - Single machine prediction application | – 25 points |
| - Docker container for prediction application | – 25 points |
| - Extra-credit for top 5 prediction performance | – 20 points |