

Individual reflection

Junfan Huang

June 21, 2019

1 Summary

This time, we had an exploration of deal with pyspark and parallelisation. To begin with, we assumed that we have a system where contains "lots of" messages. In order to make user have a clear world, we have to make sure the system is able to detect spam, ham and some fraud messages. After that, the system should be able to encrypt the message parallelly. During this exploration, we learned how to use pyspark to do machine learning by using series classification algorithms (like Logistic Regression, Random Forests) and had a experience of doing vectorized computation. In the end, we measure the time of en/decryting some artificial unbalanced datasets to compare the linear speedup curves.

The dataset is not easy to build when we decided to do this, so we took some efforts to build the dataset by using the previous project code and other online resources (can be find in **documentation/preprocess/**). The dataset consists of books, spam emails, ham emails. Our first step was to distinguish this out.

2 Usage of pyspark

Learning pyspark is a non-trivial task. There are a number of troubles, for example, setup: how to make it executable on our own lapto; dealing with data format of pyspark; the csv reader or dataframe creator for pyspark sometimes didn't do the right job; a bad datapoint will influence the whole processing. After learning some basic operation of the pyspark, we were doing word counts for the dataset by using pyspark and got some cool graphs though the frequencies of the words.

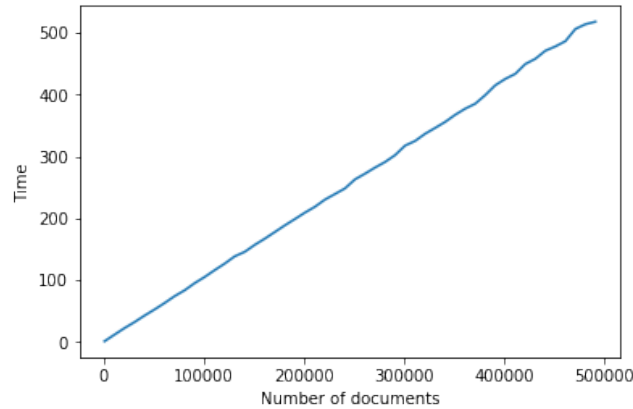


Figure 2: Linear speed up

4 Parallelisability and speed

During the exploration of this assignment, we understood that some processes are not parallelisable which limited the max speedup. For example the work dependent on the previous process is not parallelisable. I also implement the previous **Jaccard distance matrix** example by using vectorized computation which is 17 times faster than before.

In addition, we found that the communication speed is a serious issue for dealing with big data when we doing our summer project. Reading data from the hard disk is much slower (possibly hundreds of times) when compared to reading the same data from RAM directly. It is necessary to keep in mind when we deal with big data in our future work.

5 Unfinished things

Since pyspark contains a lot of machine learning tools, it could be interesting to play around with kdd data and previous project and doing a comparison of the speed, especially for the algorithms can be parallelised.

Another thing is because of the out of memory issue, I cannot train the whole dataset(got bad accuracy) and limit to the time I am unable to figure out the reason of this issue and got a whole better prediction.

Although we got some graphs, we don't have enough time to write a beautiful report to explain it and setup a good environment to run. Moreover, we didn't play around with the number of cores as we planned.

Thanks for data science toolbox and your detailed feedbacks! And thanks to everyone who helped me improving this year! Hope this course will go better in the future!