**WHAT WENT WELL?**

Setting up Hadoop environment in local machine seemed intimidating at first, but I found couple resources online such as YouTube videos and Stack Overflow that offered adequate content. Throughout this project, I had done a bunch of trial and error and lot of trouble shooting. My curiosity also peaked due to a couple success HDFS installation over YouTube videos. This had given me a strong motivation that I could deployed the system successfully on my local machine.

Besides, I documented the project properly will help me to pick up latter where I left off. For example, there were couple errors when installing HDFS, I screenshotted and put a side note of the solution. This would remind instantly to pick up where I left off. Another example, I noticed working with virtual machine can be goofy specifically the IP address may change overtime. So, before running Hive, I may want to make sure that the each VMs is pinging

Lastly, I noticed that SQL and HiveQL have language similarities. I found a cheat sheet that I really liked (http://hortonworks.com/wp-content/uploads/2016/05/Hortonworks.CheatSheet.SQLtoHive.pdf)    that explicitly compares the difference between the two. I brushed up my SQL skills from Code Academy and helped me tremendously when creating the JOIN function.


**WHAT DID NOT GO WELL?**

The first data set that I downloaded from Kaggle apparently has poor data quality specifically no ID. This means if I would like to perform JOIN function, it is nearly impossible. The only way to solve it is to create a random table. I tried my best to solve it, but apparently no luck. I decided to change my dataset because I was afraid that I was running out of time.

Even though the second dataset worked well, I failed to process the timestamp. My dataset timestamp format has a mix between DD-MM-YYYY hh:mm and MM/DD/YYYY hh:mm. According to Hive Language Manual, there are two ways to process timestamp:

1.   YYYY-MM-DD hh:mm:ss
2.   Unix timestamp function

I tried to custom format on Microsoft Excel to YYYY-MM-DD hh:mm, but for some of the reasons some rows that have MM/DD/YYYY hh:mm format was not converted. I decided to create my own unix

timestamp function with a help of different forums, but still failed. Apparently, when I tried to process the time in HiveQL, it still works; however, I highly doubt it shows all the data correctly.



**WHAT WOULD I TRY DIFFERENTLY NEXT TIME?**

There are three main takeaways from this project:

First, I would make a deeper research between Hive and Spark. After having a conversation with a senior data scientist, he mentioned that nowadays most companies are using Spark rather Hive and have better

compatibility with machine learning programs such as Python. In fact, Spark has better documentation that Hive.

Second, I would pick three different datasets that I would be interested to work on. The main purpose is to have some back up in case if I choose poor dataset quality (e.g. No ID). It sounds like that I am irresilient the fact that I hate to jump over if I could not solve a problem.

Third, when creating the unix timestamp, I should have documented every code that I tried. This would give me a comparison between codes that I have written and not repeating the same codes that I might have written.