# Course project - Big Data Concepts

Yash Pratap Solanky

Github link - https://github.com/ysolanky/big-data-project

## Introduction

Cryptocurrencies have recently become a widely recognised way for many individuals to invest their money. Cryptocurrency values have increased dramatically during the last several years. As a result, I tried numerous ways to find an efficient and accurate model to predict the price using machine learning algorithms. For the course project, I analysed four cryptocurrencies and attempted to predict the price of bitcoin. The CSV time series data is uploaded to GitHub. The time-series data in CSV format is posted to GitHub. The GitHub repository is then cloned into a Jupyter notebook under the Google cloud platform's AI platform. The time-series graph and the prediction graph are then saved as png files and uploaded to Github and to a bucket on the Google cloud platform that has been created separately.

## Background

I attempted to predict cryptocurrencies due to the fluctuating nature and fluctuation in the values of cryptocurrencies since the pandemic has piqued the interest of people all around the world. We can completely compute, store, analyse, and predict this project on the Google cloud platform thanks to the usage of virtual machines. Although this project could have been completed on a local system, the usage of a virtual machine makes it extremely simple for anybody to replicate it for any time series data. It might be for stocks, cryptocurrency, or something else.

## Methodology

**Step 1:** Downloading the time series data of various cryptocurrencies from https://www.cryptodatadownload.com/.

I downloaded the daily prices of Bitcoin, Litecoin, Ethereum and Dogecoin from this website. The range of time series of all the four coins was different.

**Step 2:** Creating a repository and uploading the downloaded data to the repository



**Step 3:** Creating a Jupyter notebook on the AI platform on the Google cloud platform and then cloning the repository to this Jupyter environment.

The AI platform supports Jupyter notebooks and has inbuilt support for git repositories. Making it easy to link the notebook with repository,

**Step 4:** Performing exploratory data analysis on the data.

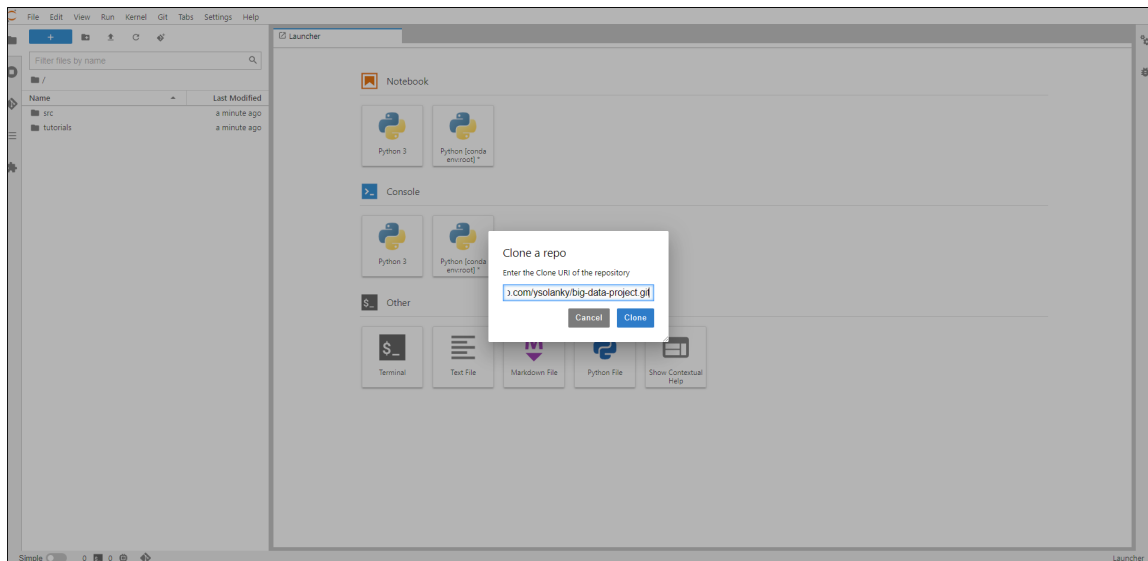We will performing data engineering (cleaning data, transforming data), and then using the final dataset to make predictions using the autoregressive integrated moving average (ARIMA) algorithm.We then will also be using facebook's prophet library to make predictions.

For ARIMA we will be using 2364 rows for training data and will use the last 30 days for testing. For the hyperparameters we will be using (10,1,5). For prophet we will be using 2300 days for training and testing on the final 94 days.

**Step 5:** Creating a google cloud platform bucket

In addition to saving the outputs from our exploratory data analysis and the ARIMA prediction to Github, we wish to publish the results on a google cloud platform bucket.



**Step 6:** Saving the graphs from the project to the bucket.

Using the google.cloud library in Python, we can link our pre-existing GCP bucket with the Jupyter notebook. And save the outputs to the bucket.

**Step 7:** Making the graphs public from the bucket

We can now make the contents of the bucket publicly available through a URL. This can be done by editing the access of the files individually.



## Results

Through our exploratory analysis we observed that towards the end of 2020, all the cryptocurrencies experienced a rapid increase in prices. There is still some debate as to what caused it, as the rise in crypto prices is in contrast to the world economy which due to the COVID-19 pandemic has been in a decline.

Using Arima, we were not able to accurately anticipate the price of Bitcoin. The algorithm predicted the price of bitcoin to remain stable, yet it plummeted significantly. This is rather expected as Bitcoin and all other cryptocurrencies follow an erratic pattern.

Above, you can see the prediction (in red line) significantly different from the actual data (blue line). Hence, for any real world application, this prediction is not acceptable.



Surprisingly, even worse results were obtained using the prophet library. Above you can see the first graph with prediction (in red line) significantly different from the actual data (blue line). The second graph is a forecast of the time series.

**Discussion**

An interpretation of the results

Even though ARIMA is a very sophisticated algorithm and is considered one of the best ways to make predictions on time series data, it proved to be relatively poor for predicting cryptocurrencies. The same result was obtained when prediction was done using facebook's prophet library. This leads to the conclusion that cryptocurrencies are hard to predict using the measures that have been proven to work in great effect for stocks.

A discussion of how you employed the technologies/skills from this course.

For this project, I did the following:

- Implemented a pipeline (e.g., download - transform - summarize - visualize) - The pipeline was implemented on a jupyter notebook in google cloud platform. The data was downloaded in the form of csv files, it was transformed by changing variables to datetime objects and removing the unnecessary columns. It was summarised and then visualised.
- Investigated a big data cloud platform environment - Google cloud platforms AI platform and storage bucket were used in the scope of this project.
- Developed a data management plan that describes how this dataset is being collected, stored, preserved, and shared - The dataset is being collected through a cryptocurrency historical archive and  stored,preserved and shared on github
- (1) Pipelines, (2) Established techniques for data cleaning/quality assurance (3) Publishing data in a repository (automatically or manually). All three of the above stated points were implemented in the scope of this project. The pipeline was created

in a Jupyter notebook on the AI platform on GCP, the technique for data cleaning and quality assurance has been demonstrated in the notebook. The results in the form of graphs have been published to GitHub.

Any barriers or failures you encountered?

The traditional and the modern techniques used to predict stock prices have resulted in fruitless results for prediction of cryptocurrencies. Very high value for the Root mean square error (RMSE) was obtained in these experiments, hence these results are far from acceptable in real world application. I think that the sudden rise and drop in price make it very difficult to predict using parametric linear model constructions, yet if with more information of weekly, seasonal and yearly trends, the prediction accuracy may be improved.



Above is a 365 day forecast of BTC. It predicts a somewhat steady growth.

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

## **Conclusion**

Cryptocurrency prices are hard to predict using traditional methods for stock prediction because of the lack of trends that stocks follow, as cryptocurrencies are more volatile than stocks. Recently it has been observed that a single event could trigger an unprecedented rise or fall in the price of crypto or stock. For example, Elon Musk tweeted that everyone should be buying GameStop stocks and Dogecoin cryptocurrency which led to a huge increase in their respective price. Events like this are hard to be integrated as a part of input information for forecasting.

But using the GCP and github, this project can be modeleted to implement prediction on any time series data. I expect it to perform much better for traditional time series like stocks.

# References

1. https://www.qwiklabs.com/focuses/1846?catalog_rank=%7B%22rank%22%3A1%2C%22num_filters%22%3A0%2C%22has_search%22%3Atrue%7D&parent=catalog&search_id=7008005

2. https://cloud.google.com/ai-platform/training/docs/working-with-cloud-storage

3. https://www.cryptodatadownload.com

4. https://github.com/googleapis/python-storage

5. https://stackoverflow.com/questions/37003862/how-to-upload-a-file-to-google-cloud-storage-on-python-3