


Project 2: ETL with Cryptocurrencies

Jacob Cortez, Julian Loudon, Dan Tinarwo, Paola Mateo

SMU-Data Science Bootcamp

Overview

A brief look at what we will discuss on this presentation

- Purpose of Project
 - Methods
 - Potential Analysis
 - Data Limitations
- 

Purpose and Goals

Extract, Transform, and Load data

Data files were derived to observe trends of 3 cryptocurrency types (bitcoin, dogecoin, and ethereum) over the span of 3 years and imported into MongoDB

Could make future predictions

An analysis can easily be observed when data is imported into a database where all nine files may be accessed

View ongoing trends for a specific type of cryptocurrency

Keeping track of vital information can ensure the success of the company, and even the potential customer.



Guide on how Goal Was Reached

01

Collect Data Sources

Data sources were found on Kaggle.com with ~475,000 columns each dataset.

02

Data Cleaning

Files were imported into Jupyter by using the Pandas library. Null variables and unnecessary columns were deleted. Columns were renamed and reorganized.

03

Specific Data Transformation

Because the "timestamp" column was in epoch time, we changed it so it can have a readable date and time.

04

Load Data onto MongoDB

Once the CSV files were uploaded with the updated and relevant data, they were transferred to MongoDB in order to collectively view data at a time.

05

Filter Desired Information and Transfer Back to Jupyter

Specific data can be sent to Jupyter to view relationships between variables and make further analysis.

```
#Format the columns accordingly
dogecoin19_df['Number of Trades'] = dogecoin19_df['Number of Trades'].map('{:,.0f}'.format)
dogecoin19_df['Open'] = dogecoin19_df['Open'].map('{:,.4f}'.format)
dogecoin19_df['High'] = dogecoin19_df['High'].map('{:,.4f}'.format)
dogecoin19_df['Low'] = dogecoin19_df['Low'].map('{:,.4f}'.format)
dogecoin19_df['Close'] = dogecoin19_df['Close'].map('{:,.4f}'.format)
dogecoin19_df['Volume'] = dogecoin19_df['Volume'].map('{:,.4f}'.format)
dogecoin19_df['VWAP'] = dogecoin19_df['VWAP'].map('{:,.4f}'.format)
dogecoin19_df.head()
```

Formatting the values to the hundredth or ten-thousandth place

	Timestamp	Crypto Currency Type	Number of Trades	Open	High	Low	Close	Volume	VWAP
0	1555079640	4	10	\$0.0028	\$0.0028	\$0.0028	\$0.0028	\$192,109.0000	\$0.0028
1	1555079700	4	14	\$0.0028	\$0.0028	\$0.0028	\$0.0028	\$70,399.0000	\$0.0028
2	1555079760	4	7	\$0.0028	\$0.0028	\$0.0028	\$0.0028	\$217,483.0000	\$0.0028
3	1555079820	4	7	\$0.0028	\$0.0028	\$0.0028	\$0.0028	\$4,247,299.0000	\$0.0028
4	1555079880	4	21	\$0.0028	\$0.0028	\$0.0028	\$0.0028	\$173,376.0000	\$0.0028

```
#Rewrite content in Crypto Currency Type
dogecoin_rename= {4:"Dogecoin"}

dogecoin19_df["Crypto Currency Type"]=dogecoin19_df["Crypto Currency Type"].map(dogecoin_rename)
dogecoin19_df
```

Cleaning data: renaming the specified number to the actual crypto-currency name

	Timestamp	Crypto Currency Type	Number of Trades	Open	High	Low	Close	Volume	VWAP
0	1555079640	Dogecoin	10	\$0.0028	\$0.0028	\$0.0028	\$0.0028	\$192,109.0000	\$0.0028
1	1555079700	Dogecoin	14	\$0.0028	\$0.0028	\$0.0028	\$0.0028	\$70,399.0000	\$0.0028
2	1555079760	Dogecoin	7	\$0.0028	\$0.0028	\$0.0028	\$0.0028	\$217,483.0000	\$0.0028
3	1555079820	Dogecoin	7	\$0.0028	\$0.0028	\$0.0028	\$0.0028	\$4,247,299.0000	\$0.0028
4	1555079880	Dogecoin	21	\$0.0028	\$0.0028	\$0.0028	\$0.0028	\$173,376.0000	\$0.0028
...
306307	1565823300	Dogecoin	0	\$0.0028	\$0.0028	\$0.0028	\$0.0028	\$0.0000	\$nan
306308	1565823960	Dogecoin	0	\$0.0028	\$0.0028	\$0.0028	\$0.0028	\$0.0000	\$nan
306309	1565824680	Dogecoin	0	\$0.0027	\$0.0027	\$0.0027	\$0.0027	\$0.0000	\$nan
306310	1565824740	Dogecoin	0	\$0.0027	\$0.0027	\$0.0027	\$0.0027	\$0.0000	\$nan
306311	1565826000	Dogecoin	0	\$0.0027	\$0.0027	\$0.0027	\$0.0027	\$0.0000	\$nan

Replacing the renamed column information onto the column

Function Used to Change Epoch Time to a Readable Timestamp

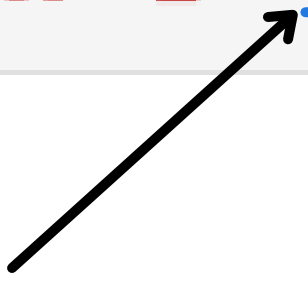
```
def epoch(inputdate):  
    return datetime.datetime.fromtimestamp(inputdate).strftime('%Y-%m-%d %H:%M:%S')
```

```
timestamp = bitcoin19_df['Timestamp']  
timestamp1 = bitcoin20_df['timestamp']  
timestamp2 = bitcoin21_df['timestamp']
```

```
bitcoin19_df['Date Time'] = bitcoin19_df['Timestamp'].apply(lambda x:epoch(x))  
bitcoin20_df['Date Time'] = bitcoin20_df['timestamp'].apply(lambda x:epoch(x))  
bitcoin21_df['Date Time'] = bitcoin21_df['timestamp'].apply(lambda x:epoch(x))
```

```
bitcoin21_df['Date Time']
```

```
#Export file as a CSV with a header  
dogecoin19_df.to_csv("../Project2_Crypto_db/DB_dogecoin19.csv", index=False, header=True)  
dogecoin20_df.to_csv("../Project2_Crypto_db/DB_dogecoin20.csv", index=False, header=True)  
dogecoin21_df.to_csv("../Project2_Crypto_db/DB_dogecoin21.csv", index=False, header=True)
```

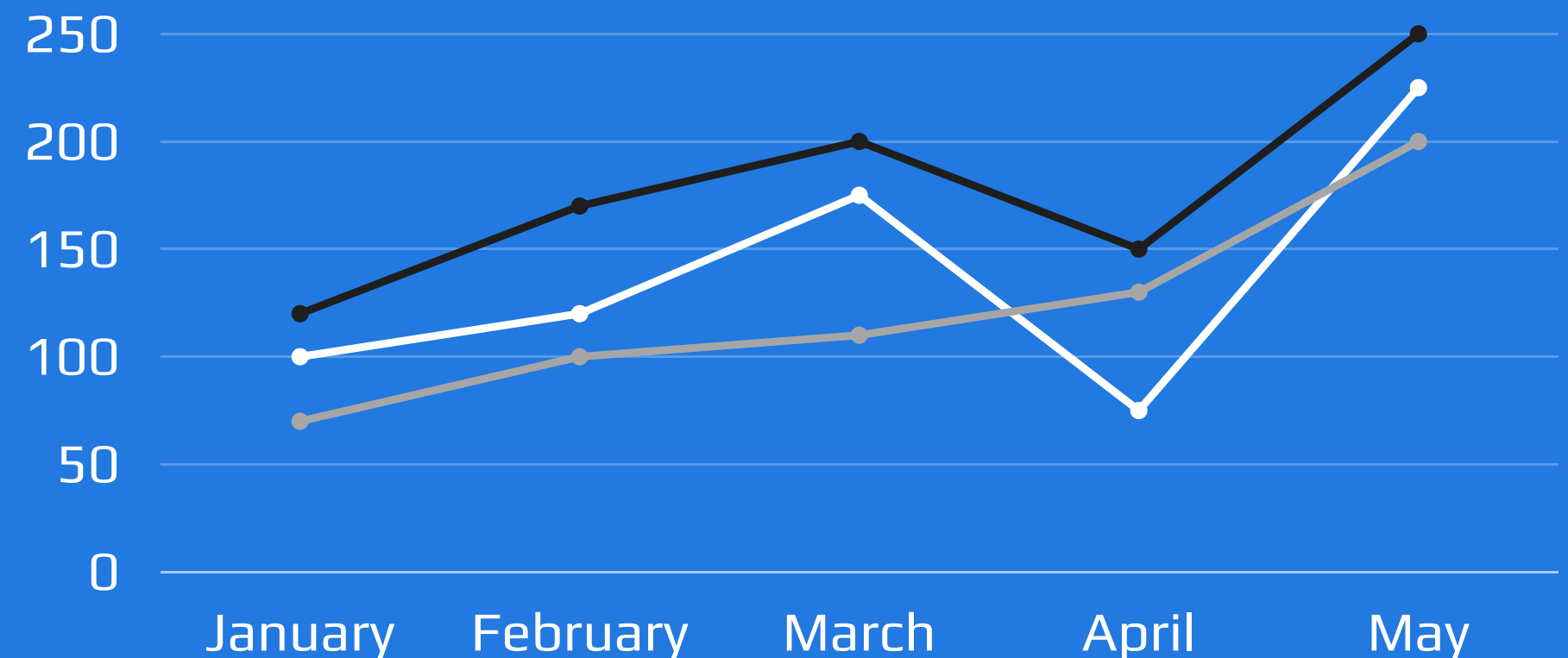


Exporting files as a new CSV file. File had to be created, and path was copied so that we can have a new dataset with the clean data

Potential Analysis Based on Data

WE WOULD BE ABLE TO DETERMINE/ANSWER THE FOLLOWING:

- The most overall successful type of cryptocurrency
- The month in which the count (amount of trades) is the highest.
- Which calendar quarter had the most volume
- The average amount of trades in Bitcoin.
- Observing the trend for "Target" (residualized returns) in Ethereum



GRAPH WAS INCLUDED FOR VISUAL PURPOSES. **NOT** ACCURATE DATA

Data Limitations

Data Files Were Too Large

Once new data files were uploaded to MongoDB, attempting to push data back to Jupyter took roughly about 25-30 minutes, and sometimes even crashing the computer. Filtering is needed to further analyze data.

Inconsistencies in Timestamps

There are certain times where a specific crypto does not have trades. For example, Dogecoin may not have any trades for 10 minutes total one day, but trades for every minute the next day. This creates the data to be slightly inconsistent and difficult to create an accurate conclusion.

Conclusion

- 9 total CSV files were used for this project
- MongoDB was the database where CSV files were imported. We can load the data from this database to filter it out by the relationships that will be observed
- Potential analysis could include results regarding the best month for a trade in any of the 3 cryptocurrencies.
- Allowing room for error may be a solution to the listed data limitations.
- Overall, the purpose is to begin the process for a data analysis with different CSV files to create future predictions or data reports for ongoing trends.

Utilized Data Files:

Bitcoin Files:

- https://www.kaggle.com/yamqwe/cryptocurrency-extra-data-bitcoin?select=full_data__1__2019.csv
- https://www.kaggle.com/yamqwe/cryptocurrency-extra-data-bitcoin?select=full_data__1__2020.csv
- https://www.kaggle.com/yamqwe/cryptocurrency-extra-data-bitcoin?select=full_data__1__2021.csv

Dogecoin Files:

- https://www.kaggle.com/yamqwe/cryptocurrency-extra-data-dogecoin?select=full_data__4__2019.csv
- https://www.kaggle.com/yamqwe/cryptocurrency-extra-data-dogecoin?select=full_data__4__2020.csv
- https://www.kaggle.com/yamqwe/cryptocurrency-extra-data-dogecoin?select=full_data__4__2021.csv

Ethereum Files:

- https://www.kaggle.com/yamqwe/cryptocurrency-extra-data-ethereum?select=full_data__6__2019.csv
- https://www.kaggle.com/yamqwe/cryptocurrency-extra-data-ethereum?select=full_data__6__2020.csv
- https://www.kaggle.com/yamqwe/cryptocurrency-extra-data-ethereum?select=full_data__6__2021.csv