Cryptocurrency Prediction using Social Media Indicators

SPRING 2019 | DATA ENGINEERING PLATFORMS FOR ANALYTICS

GROUP 5 – AKARSH SAHU, NAZIH KALO, ANGELA TENG, SOWMYA NALLAPANENI

Outline

- ▶ Executive Summary
- Business Use Case
- ► Tech Stack Data and Tools
- Design Considerations
- Enhanced Entity Relationship Model
- Dimensional Model
- Visualizations and Demo
- ▶ Future Recommendations



Executive Summary

- The cryptocurrency market is very nascent and volatile. This makes the asset class very risky but also potentially lucrative.
- •Behavioral sciences and related scientific literature provide evidence that there is a close relationship between social media and price fluctuations of cryptocurrencies.
- **Hypothesis:** Social media indicators are strongly correlated with prices and can be a good prediction indicator for short-term investments in cryptocurrencies

Project
Background/Hypothesis

Project Objective

- Automated Data Integration Connect Data Sources, Ingestion, and Consolidation
- Efficient and Scalable Data Storage: Normalized EER Model, OLAP Dimensional Model, and Cloud SQL Integration
- •Generate real-time insights on cryptocurrency price: Leveraging Data Visualization, Statistical Analyses, and Forecasting.

- Strong Correlation b/w
 Cryptocurrency price and Social media indicators. But Correlation doesn't imply causation!
- Build a regression model based on metrics available to **test the hypothesis** further.
- •Leverage NLP and Sentiment Analysis to generate sentiments from text sources.
- Build a **real-time predictive model** utilizing all the metrics , sentiment, and demographics.

Conclusion/Future Recommendations

Business Use Case



Day-traders can use our database to build complex trading models based on social media sentiment and volume



Long-term investors can find established/safer assets to diversify their portfolio



▶ Industry experts can use the dataset to analyze the public's knowledge and interest in their specific cryptocurrency/blockchain technology



Data & Tools

Storage& Data Ingestion Analysis/Visualization Prepare database sources Google Trends **SQL**

Data Source



<u>Data source</u> www.bitinfocharts.com

Description

Number of daily tweets for each cryptocurrency

<u>Time Period</u> 2014-Today

Number of Records 48.598



<u>Data source</u> www.reddit.com

Description

- All posts to /r/cryptocurrency subreddit
- 2. Total daily subscribers to all cryptocurrency subreddits

<u>Time Period</u>

2014-Today

Number of Records

- 1. 442,038
- 2. 121,138



<u>Data source</u> https://trends.google.com

<u>Description</u>
Google trend data

Time Period 2014-Today

Number of Records 120.975



<u>Data source</u> <u>www.cryptocompare.com</u>

<u>Description</u>
OHLCV for all cryptocurrencies

Time Period 2009-Today

Number of Records 82,637

Ingestion



Beautifulsoup & Requests API
Scraped # of tweets per currency





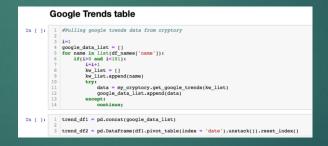
Praw & Psaw APIs
Scraped all posts
from/r/cryptocurrency







Cryptory API
Scraped google
trend data per
currency





Requests API
Scraped pricing data for every currency

```
Pricing Table

| import requests | import datetime | import negress | import datetime | import pandan as pd |
| idef daily price historical(symbol, comparison.symbol, all_dateTrue, linit=1, aggregate=1, exchange="i"):
| datef daily price historical(symbol) | idef exchange | ide
```

Prepare/Clean





Cleaning Tweets & Merging with coin ids



	Date	No. of Tweets	coin_id	date_id
0	2018-01-08	7	0	3294
1	2018-01-09	20	0	3295
2	2018-01-10	21	0	3296
3	2018-01-11	4	0	3297
4	2018-01-12	6	0	3298
5	2018-01-13	14	0	3299
6	2018-01-14	3	0	3300
7	2018-01-15	18	0	3301
8	2018-01-16	1	0	3302



Creating Reddit-Coin Join Table

```
# # Create reddit_coins Join Table
# #### Matches the reddit posts in reddit_data table to the coins in coins table
#Zip the names into tuple of (name, slug, symbol)
zipped_names = list(zip(df_names['name'],df_names['slug'],df_names['symbol']))
#Create a search list - seperating the three terms with OR operator (|)
search_list = []
for (name, slug, symbol) in zipped_names:
     listt = [name, slug, symbol]
pat = '|'.join(listt)
search_list.append(pat)
#Add boolean/dummy columns
dummy_df = pd.DataFrame(dict((name, reddit_data.title.str.contains(name, re.IGNORECASE))
                                    for name in search list))
#Convert dummy columns into rows with post id as the index
i, j = np.where(dummy_df)
coins_mentioned_series = pd.Series(dict(zip(zip(i, j), dummy_df.columns[j])))
#Create final join table between reddit and coins table. Rename columns of dataframe.
coin_reddit_join = pd.DataFrame(coins_mentioned_series).reset_index()
coin_reddit_join.columns = ['post_id', 'coin_id', 'coin_name']
coin_reddit_join = coin_reddit_join.drop('coin_name', axis = 1)
coin_reddit_join2 = coin_reddit_join.set_index(['post_id', 'coin_id'],)
coin reddit join2
```

title	subreddit	score	num_comments	date	author	
						post_id
Swap (XWP) added to CNPool.cc today	CryptoCurrency	1	0	2019-05-31 23:56:44	Mahaprajapati	0
Millionaire Alexander Amado Johnson and his wi	CryptoCurrency	1	0	2019-05-31 23:55:15	webnowcompany	1
How is OmiseGOing? By Kasima Tharnpipitchai, D	CryptoCurrency	29	4	2019-05-31 23:48:29	sebikun	2
Gemini is down??	CryptoCurrency	1	0	2019-05-31 23:47:14	albuquerquetulio	3
Analyst: Regardless of Newest Bitcoin Pullback	CryptoCurrency	1	0	2019-05-31 23:46:36	cryptotradinglife	4

#Concatenate DF's of name-based subreddits reddit_subscribers_name = pd.concat(subreddit_list2) #Add the date_id column reddit_subscribers_name1 = reddit_subscribers_name.merge(dates_df, how = 'left', left_on = 'date', right_on = 'date', ri

total_subscribers subreddit date coin_id 2015-01-01 882 AION AION 2015-01-02 884 2015-01-03 AION 2015-01-04 887 AION AION 2015-01-05 896



Merging coin ids with Trends

		trend
date	coin_id	
2015-01-01	0	4.145594
2015-01-02	0	20.267349
2015-01-03	0	27.176672
2015-01-04	0	24.873564
2015-01-05	0	26.255429



Cleaned/Parsed HTTP request & merged with coin ids

```
// Accord list of dataframes
pricing_df = pd.concat(pricing_list)

// Arge with coin id
pricing_df2 = pricing_df.merge(df_names[['symbol', 'coin_id']], how = 'left', left_on = 'noting_df3 = pricing_df2.drop(['name', 'symbol', 'time'], axis = 1)

// Arke it a datetime object
pricing_df3['timestamp'] = pd.to_datetime(pricing_df3['timestamp'])

// Areaame timestamp to date
pricing_df3.rename(columns = {'timestamp':'date'}, inplace = True)

// Areaame timestamp to date
pricing_df3.rename(columns = {'timestamp':'date'}, inplace = True)

// Areaame timestamp to date
```

		close	high	low	open	volumefrom	volumeto
date	coin_id						
2017-08-11	0	0.06687	0.4444	0.06687	0.06950	39437.24	2994.43
2017-08-12	0	0.16250	0.1830	0.06687	0.06687	6232.81	765.38
2017-08-13	0	0.19000	0.3880	0.12500	0.16250	10382.08	1885.11
2017-08-14	0	0.32000	0.3300	0.13100	0.19000	7936.82	1788.88
2017-08-15	0	1.00000	1.8000	0.32000	0.32000	3467.80	2905.58

Storage



Compute Engine VM

 Run all python web scraping and transformation scripts in a VM instance





Storage Buckets

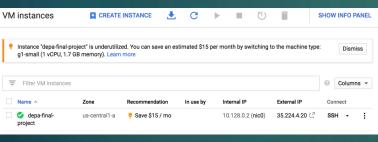
- Store data as .csv files
- Store raw data files for potential future use-cases



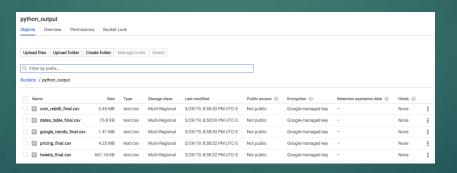


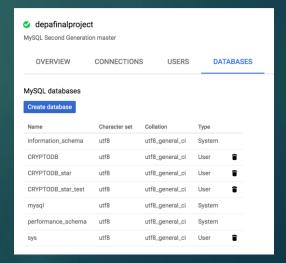
CloudSQL Instance

- Create DB through MySQL connection using DDL scripts
- Import Data from storage bucket csv's
- Second DDL/DMLfor starschema









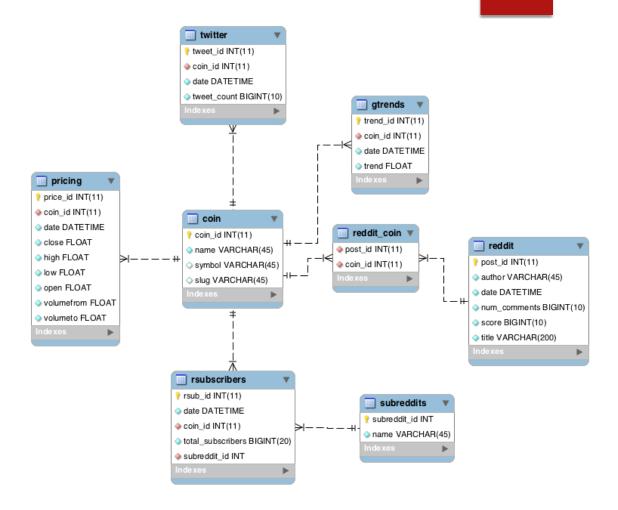
Enhanced Entity Relationship Diagram

▶ Entities:

- ▶ reddit: /r/cryptocurrency subreddit post data
- ▶ rsubscribers: subscriber count for over 150 coin subreddits
- ▶ pricing: cryptocurrency market OHLCV data
- ▶ twitter: tweet count by coin mentioned
- ▶ gtrends: Google trend data by coin
- ▶ reddit_coin: join-table linking reddit post data to coin data
- ▶ **subreddits**: name of subreddit and subreddit id
- ► coin: coin information containing coin names

► Normalized Data:

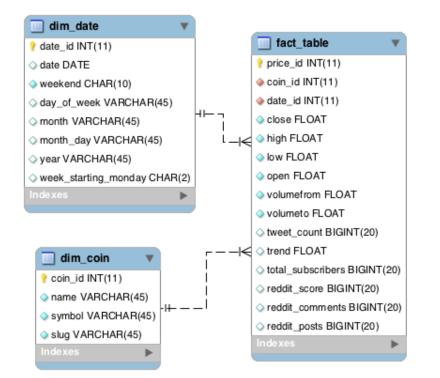
- ▶No functional dependencies
- ▶No transitive dependencies



Dimensional Model (Star Schema)

- ▶ <u>4 Step Dimensional Design Process</u>
- Business Process
 - Trader/User-oriented
- Grain
 - Currency-date level detail
- Dimensions
 - Determined by grain statement
- Facts
 - Social media aggregates

▶ Balancing user requirements with data realities





NoSQL and Graph Databases



Inappropriate Model

- MongoDB cannot model or store relationships if no complexity exists
- Our datasets are sourced through webscraping methods via APIs

Degraded Performance

•Use speed plummets as the amount of data increases and the number of join keys grows

Inappropriate Language

•"almost SQL" languages are inappropriate for SQL "joins" and similar functionality

Not ACID

 Non-relational databases like MongoDB do not follow ACID properties



Inappropriate database

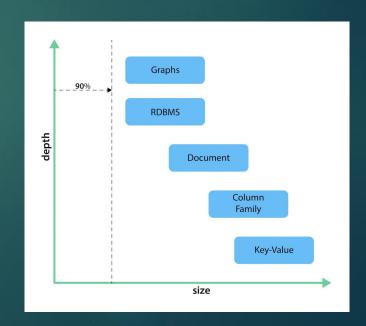
- Graph databases are inappropriate for traditional relational data.
- Difficult to perform mass analytics queries across all the relationships and records.

Not efficient for high volumes of transactions

•Graph databases are not as useful for operational use cases since they are not good at handling large queries that span the entire database.

Not an independent master data management solution

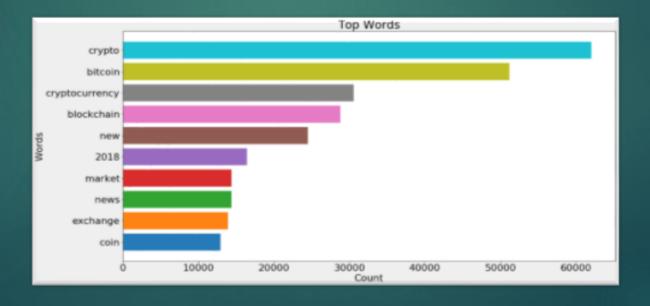
- A graph database is just a data store and doesn't give you a business-facing user interface to query or manage relationships.
- •It will not provide advanced match and survivorship functionality or data quality capabilities.



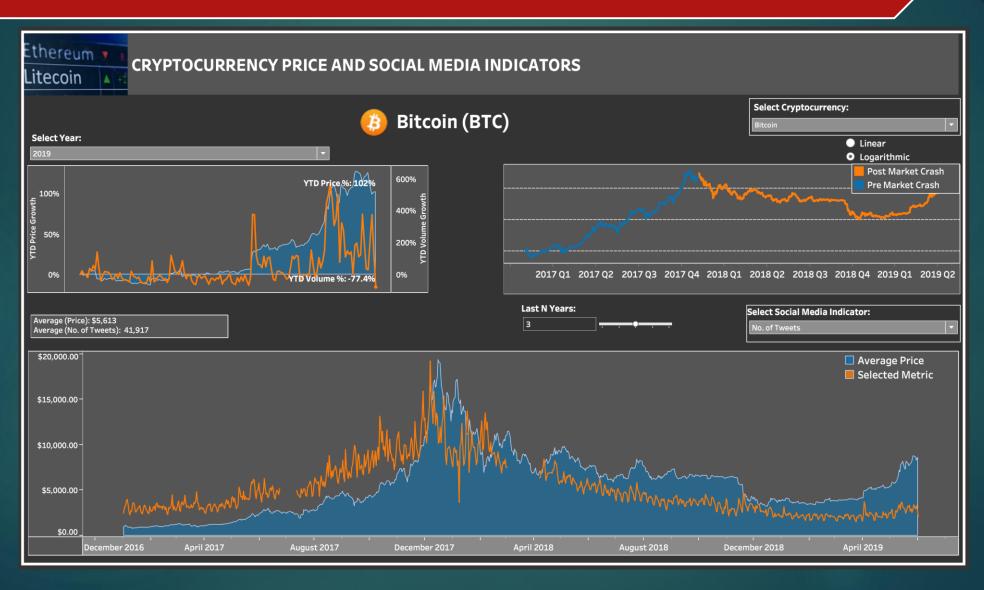
Analysis/Visualization (*)

Basic Reddit Word Analysis



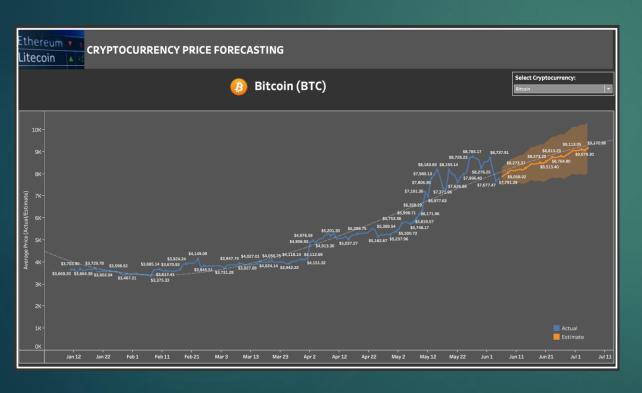


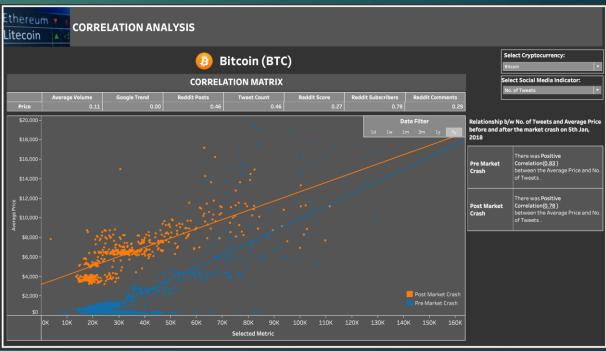
Analysis/Visualization (**)



Analysis/Visualization

Interactive Demo: Tableau Public Website





Future Work & Improvements

- Acquire user demographic data
- Include additional market data
 - ▶ Market cap, circulating supply, spread, etc.
- Develop scheduler for real-time data ingestion
- Build predictive model
 - Using NLP on twitter and reddit data
- Revisit dimensional design process and level of granularity

Thank you!

ANY QUESTIONS?



Database Design Considerations

Requirements



Conceptual



Physical model

Application Type

Online Analytical
Processing for Business
Intelligence and end
users (OLAP)

Data Format and Size

Cryptocurrency Financial Data

Business Requirements

Optimal user experience

Modeling Steps

- Elements relevant for pricing model analysis will be stored
- ACID properties used to set constraints, datatypes and ranges of the database
- technical,
 performance, and
 process requirements

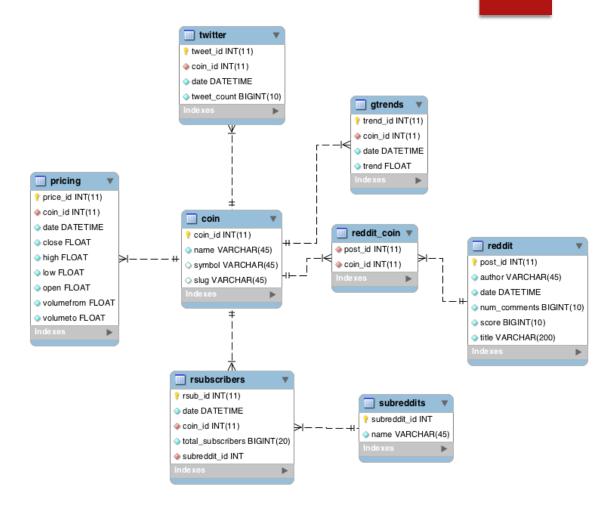
Hardware and Software

Google Cloud compute and virtual machines, and Google Cloud buckets for storage

- Financial constraints, usability constraints

Enhanced Entity Relationship Diagram

- ▶ Entities:
- reddit: /r/cryptocurrency sub post data
- reddit_subs: subscriber count for >150 coin subreddits
- Pricing: cryptocurrency market OHLCV data
- ► Twitter: tweet count by coin mentioned
- gtrends: Google trend data by coin
- reddit_Coin: join-table linking reddit post datato coin ID
- date:dateinformation
- ▶ Relationship & Cardinality
- ▶ Normalization Achieved:
 - ▶No functional dependencies
 - ▶No transitive dependencies



Graph Databases/ Neo4j

For our particular database, we also decided not to use Neo4j for the following reasons:

Inappropriate database

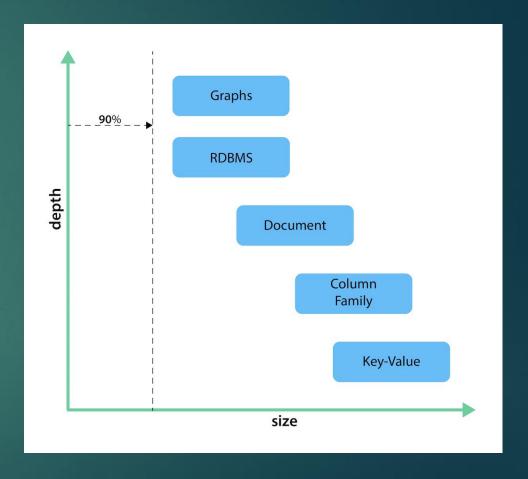
- Graph databases are inappropriate for traditional relational data.
- Difficult to perform mass analytics queries across all the relationships and records.

Not efficient for high volumes of transactions

•Graph databases are not as useful for operational use cases since they are not good at handling large queries that span the entire database.

Not an independent master data management solution

- A graph database is just a data store and doesn't give you a business-facing user interface to query or manage relationships.
- •It will not provide advanced match and survivorship functionality or data quality capabilities.



APPENDIX

Summary of Data Analysis Pipeline

Extract



Conclusion

- ► Through this project we:
 - Created a database that

Ingestion



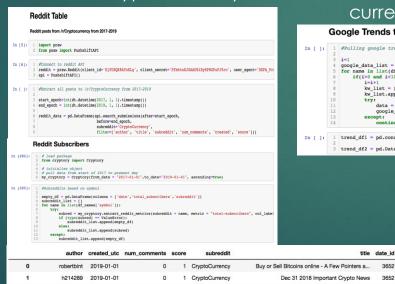
Beautifulsoup API Scraped # of tweets per currency



	Date	No. of Tweets	coin_id	date_id
0	2018-01-08	7	0	3294
1	2018-01-09	20	0	3295
2	2018-01-10	21	0	3296
3	2018-01-11	4	0	3297
4	2018-01-12	6	0	3298
5	2018-01-13	14	0	3299
6	2018-01-14	3	0	3300
7	2018-01-15	18	0	3301
8	2018-01-16	1	0	3302



Praw & Psaw APIs Scraped all posts from/r/cryptocurrency



1 CryptoCurrency

CryptoCurrency

CryptoCurrency

1 CryptoCurrency

1 CryptoCurrency

coinmarshal 2019-01-01

h1121900 2019-01-01

instasmarter 2019-01-01

h2022395 2019-01-01



Cryptory API

Scraped google trend data per currency

title date id

I am sharing Crypto with my WhatsApp buddies. ...

Dec 31 2018 Important Crypto News

Dec 31 2018 Important Crypto News

Best places to spend Bitcoin

Dec 31 2018 Important Crypto News 3652



uace	coin_iu	
2015-01-01	0	4.145594
2015-01-02	0	20.267349
2015-01-03	0	27.176672
2015-01-04	0	24.873564
2015-01-05	0	26.255429



Cryptocompare api Scraped pricing data for every currency

```
Pricing Table
         import pandas as pd
       def daily_price_historical(symbol, comparison_symbol, all_data-True, limit-1, aggregate-1, exchange-''):
    url = 'https://min-apl.oryproceapare.com/data/histoday?fayme()sizint-()saggregate-()'\
    if exchange-()' format(expholupper(), comparison_symbol.upper(), limit, aggregate)
    url += 'sec()'.format(exchange)
    if all_data;
               if all_datar
url + 'kallota=true'
page requests.get(url)
data - page.jen(||'Data'|)
df = pd.DataFrame(data)
df 'timestamp' = (datetime.fromtimestamp(d) for d in df.time]
         for symbol in list(df_names['symbol']):
                       df = daily_price_historical(symbol, 'USD')
                      df['name'] = symbol
pricing_list.append(df)
                except:

df2 = pd.DataFrame(columns = df.columns)
                       pricing list.append(df2)
```

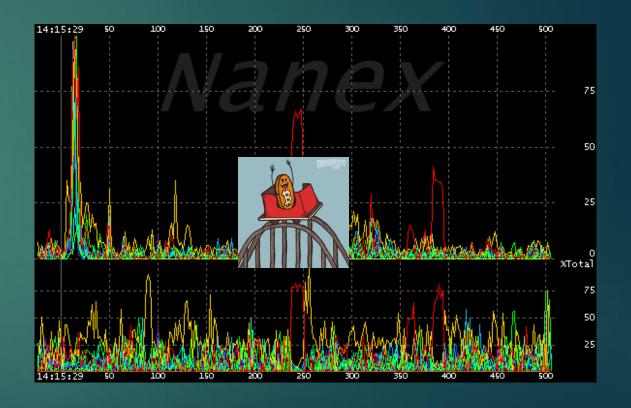
		close	high	low	open	volumefrom	volumeto
date	coin_id						
2017-08-11	0	0.06687	0.4444	0.06687	0.06950	39437.24	2994.43
2017-08-12	0	0.16250	0.1830	0.06687	0.06687	6232.81	765.38
2017-08-13	0	0.19000	0.3880	0.12500	0.16250	10382.08	1885.11
2017-08-14	0	0.32000	0.3300	0.13100	0.19000	7936.82	1788.88
2017-08-15	0	1.00000	1.8000	0.32000	0.32000	3467.80	2905.58

Executive Summary

Background:

- ► The cryptocurrency market is very nascent and volatile. This makes the asset class very risky but also potentially lucrative.
- Understanding the impact of social media on price direction can provide a cryptocurrency trader a better trading strategy.

Objective: To collect market information and social media engagement data for the major cryptocurrencies and tokens to create a platform for short-term trading strategy.



Executive Summary

- •The cryptocurrency market is very volatile.
- Understanding the impact of social media on price direction can provide a cryptocurrency trader a better trading strategy.

Project Background

Objective

- Automated Data Integration Connect Data Sources, Ingestion, and Consolidation
- Efficient and Scalable Data Storage:
 Normalized EER Model, OLAP Dimensional Model, and Cloud SQL Integration
- •Generate real-time insights on cryptocurrency price: Leveraging Data Visualization, Statistical Analyses, and Forecasting.

- Leverage NLP and Sentiment Analysis
- Build a real-time predictive model utilizing all the metrcs, sentiment, and demography.
- Build a regression model based on metrics available to test the hypothesis further.

Future Recommendations

Prepare/Clean

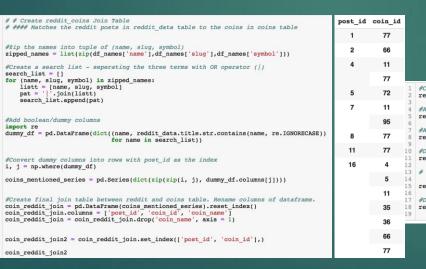


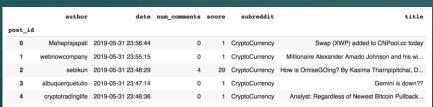
Cleaning Tweets & Merging with coin id





Creating Reddit-Coin Join Table









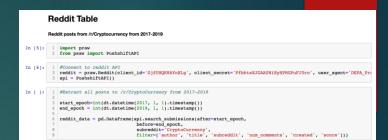
Cleaning Reddit Subscribers

```
#Concat list of dataframes
                                                                      pricing_df = pd.concat(pricing_list)
#Concatenat
                                                                                                                  ol', 'coin_id']], how = 'left',
                #Concatenate list of dataframes
reddit subs
                                                                                                                 l', 'time'], axis = 1)
               trend_df1 = pd.concat(google_data_list)
#Add the da
reddit subs
               trend_df2 = pd.DataFrame(trend_df1.pivot_table(index = 'date').unstack()).reset index()
                                                                                                                  ing_df3['timestamp'])
#Add the co
reddit subs
               trend_df3 = trend_df2.merge(df_names[['name', 'coin_id']], how = 'left', left_on = 'level_0',
#Drop the n
                                                                                                                  e'}, inplace = True)
reddit_subs
# CONCATENA 11 trend_df3 = trend_df3.drop(['level_0', 'name'], axis = 1)
                                                                                                                  lace = True)
               trend_df3.columns = ['date', 'trend', 'coin_id']
           16 trend_df3.set_index(['date', 'coin_id'], inplace = True)
```

		total_subscribers	subreddit
date	coin_id		
2015-01-01	2	882	AION
2015-01-02	2	884	AION
2015-01-03	2	885	AION
2015-01-04	2	887	AION
2015-01-05	2	896	AION



- Praw/psaw API
- Scraped all reddit posts from/r/cryptocurrency for our timeframe of 2 years





	author	created_utc	num_comments	score	subreddit	title	date_id
0	robertbint	2019-01-01	0	1	CryptoCurrency	Buy or Sell Bitcoins online - A Few Pointers a	3652
1	h214289	2019-01-01	0	1	CryptoCurrency	Dec 31 2018 Important Crypto News	3652
2	coinmarshal	2019-01-01	5	1	CryptoCurrency	I am sharing Crypto with my WhatsApp buddles	3652
3	h1121900	2019-01-01	0	1	CryptoCurrency	Dec 31 2018 Important Crypto News	3652
4	instasmarter	2019-01-01	13	1	CryptoCurrency	Best places to spend Bitcoin	3652
5	h1121900	2019-01-01	0	1	CryptoCurrency	Dec 31 2018 Important Crypto News	3652
6	h2022395	2019-01-01	0	1	CryptoCurrency	Dec 31 2018 Important Crypto News	3652

Data
Extraction and
Web Scraping

Twitter

- Beautifulsoup API
- Scraped # of tweets per currency.
 Granularity at the day level



	Twitter Table	
n [125]:	impact re impact cov impact cover impact cover impact coperate from bet impact impact impact from Drybnand impact impact from Drybnand impact from Drybnand	
n [182]=	impers points as pl impers points as pl threat, pl points	

	Date	No. of Tweets	coin_id	date_id
0	2018-01-08	7	0	3294
1	2018-01-09	20	0	3295
2	2018-01-10	21	0	3296
3	2018-01-11	4	0	3297
4	2018-01-12	6	0	3298
5	2018-01-13	14	0	3299
6	2018-01-14	3	0	3300
7	2018-01-15	18	0	3301
8	2018-01-16	1	0	3302

Google Trends

- Cryptory API
- Scraped google trend data for every currency in our timeframe of interest



	Go	pogle Trends table
In []:	1	#Pulling google trends data from cryptory
	3	i=1
	4	google_data_list = []
	5	for name in list(df_names['name']):
	6	if(i>0 and i<101): i=i+1
	8	kw list = []
	9	kw list.append(name)
	10	try:
	11	<pre>data = my_cryptory.get_google_trends(kw_list) google_data_list.append(data)</pre>
	13	<pre>google_data_list.append(data) except:</pre>
	14	continue;
		Annual Affi or and annual Afficiance Annual Afficiance
n []:	2	trend_df1 = pd.concat(google_data_list)
	3	<pre>trend_df2 = pd.DataFrame(df1.pivot_table(index = 'date').unstack()).reset_index()</pre>

×	date	trend	coin_id
0	2017-01-01	14.939345	С
1	2017-01-02	14.939345	C
2	2017-01-03	15.686312	C
3	2017-01-04	24.276435	0
4	2017-01-05	22.035534	C
5	2017-01-06	17.553730	C
6	2017-01-07	13.071927	C
7	2017-01-08	10.457541	C
8	2017-01-09	8.216640	C
9	2017-01-10	18.300697	C

Data Storage

Compute Engine VM

- Run all webscraping scripts in python 3
- Transform and normalize tables
- Push data to buckets

Storage Buckets

- Store data as .csv files
- Store raw data files for potential future use-cases

CloudSQL Instance

- Cryptodb Database created through MySQL connection
- Run DDL scripts
- Import DML/Data from storage bucket csv's

