



Cryptocurrency Prediction using Social Media Sentiment & Volume

SPRING 2019 | DATA ENGINEERING & PLATFORMS FOR ANALYTICS

GROUP 5

Outline

- ▶ Executive Summary
- ▶ Business Use Case
- ▶ Tech Stack – Data and Tools
- ▶ Design Considerations
- ▶ Enhanced Entity Relationship Model
- ▶ Dimensional Model

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.

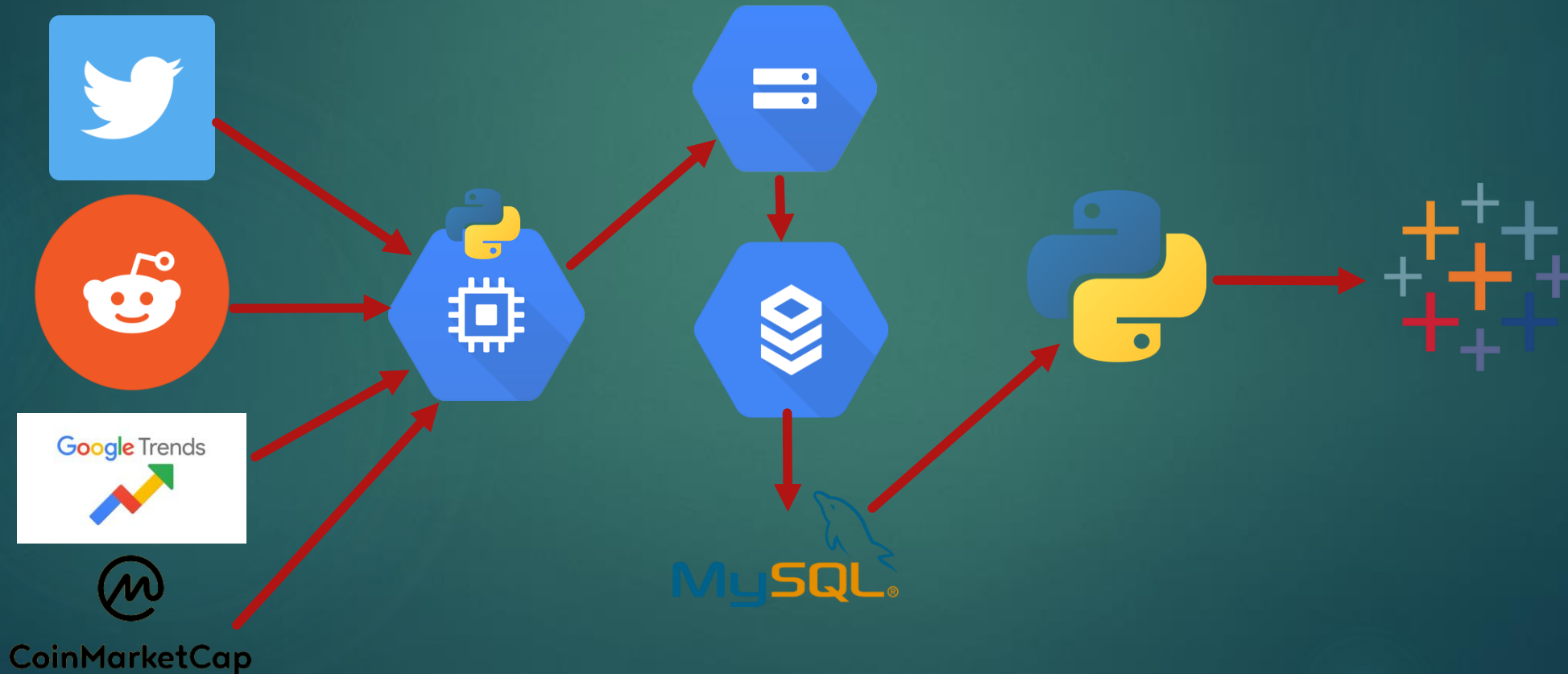
Business Use Case

A price prediction model based on metrics and sentiments from multiple online resources and incorporate it into larger system that automatically and intelligently manages a cryptocurrency portfolio.

Users:

- ▶ 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



Database Design Considerations

- Application Type

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

- Data Format and Size

- Cryptocurrency Financial Data

- Data Maintenance and Support

- Open Source Project

- Hardware and Software

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

- Elements relevant for pricing model analysis will be stored

- ACID properties and use cases were considered when developing the constraints, datatypes and ranges of the database

Data Extraction



praw/psaw API

Reddit Table

Reddit posts from /r/Cryptocurrency from 2017-2019

```
In [5]: 1 import praw
2 from praw import PushshiftAPI

In [6]: 1 #Connect to reddit API
2 reddit = praw.Reddit(client_id='GjFUMQ8A7nXlg', client_secret='PfbhtsXGAUNiKy8PRGPuP3Joro', user_agent='DEPA_Pr
3 api = PushshiftAPI()

In [ ]: 1 #Extract all posts to /r/Cryptocurrency from 2017-2019
2
3 start_epoch=int(dt.datetime(2017, 1, 1).timestamp())
4 end_epoch = int(dt.datetime(2019, 1, 1).timestamp())
5
6 reddit_data = pd.DataFrame(api.search_submissions(after=start_epoch,
7 before=end_epoch,
8 subreddit='CryptoCurrency',
9 filter=['author', 'title', 'subreddit', 'num_comments', 'created', 'score']))
```

| | 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 buddies. ... | 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 |

Scrapped all reddit posts to /r/cryptocurrency for our timeframe



beautifulsoup

Twitter Table

```
In [129]: 1 import re
2 import csv
3 import requests
4 from bs4 import BeautifulSoup
5 from IPython.display import HTML

In [182]: 1 import pandas as pd
2 tweets_df = pd.DataFrame({'Date': [], 'No. of Tweets': [], 'Coin': []})
3
4 for i in range(0, len(df_names['symbol'])-1):
5     #len(coin_name['symbol'])-1
6
7     coin = df_names['symbol'][i]
8     url = 'https://bitfinexdata.com/compasion/tweets?coin=lower(' + coin + ')&html'
9     headers = {'User-Agent': 'Chrome/54.0.2840.90'}
10    response = requests.get(url, headers=headers)
11    html = response.text
12
13    from bs4 import BeautifulSoup
14    soup = BeautifulSoup(html, 'html.parser')
15
16    x = soup.find_all('script')
17
18    data_1 = re.findall(r'\\(lowerData\\.\\[\\])', str(x))
19    data_1 = str(data_1)
20
21
22    if data_1 == '[]':
23        continue
24    data_2 = data_1.split(',')
25    data_2[0] = data_2[0][3:]
26    data_2[1:len(data_2)-1] = data_2[1:len(data_2)-1][1:-4]
27    data_pd = pd.DataFrame(data_2)
28    data_clean = data_pd[0].str.split(',')
29    data_clean = data_clean[0:2]
30    data_clean['Date'] = data_clean['Date'].str.slice(0,20)
31    data_clean['Coin'] = data_clean['Coin'].str.upper()
32    tweets_df = tweets_df.append(data_clean, ignore_index = True)
33
```

| | 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 |

Scrapped # of tweets per currency. Granularity = day

Google Trends



cryptory API

Google Trends table

```
In [ ]: 1 #Pulling google trends data from cryptory
2
3 i=1
4 google_data_list = []
5 for name in list(df_names['name']):
6     if(i>0 and i<10):
7         i=i+1
8         kw_list = []
9         kw_list.append(name)
10        try:
11            data = my_cryptory.get_google_trends(kw_list)
12            google_data_list.append(data)
13        except:
14            continue

In [ ]: 1 trend_df1 = pd.concat(google_data_list)
2
3 trend_df2 = pd.DataFrame(df1.pivot_table(index = 'date').unstack().reset_index())
```

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

Scrapped google trend data for every currency for our timeframe of interest

Data Storage

STEP 1 Compute Engine VM

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

STEP 2 Storage Buckets

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

STEP 3 CloudSQL Instance

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

VM instances

CREATE INSTANCE

Instance "depa-final-project" is underutilized. You can save an estimated \$15 per month by switching to the machine type: g1-small (1 vCPU, 1.7 GB memory). [Learn more](#)

Filter VM instances

| Name | Zone | Recommendation | In use by | Internal IP | External IP | Connect |
|--------------------|---------------|----------------|-----------|-------------------|-------------|---------|
| depa-final-project | us-central1-a | Save \$15 / mo | | 10.128.0.2 (nic0) | 35.224.4.20 | SSH |

python_output

Objects Overview Permissions Bucket Lock

Upload files Upload folder Create folder Manage holds Delete

Filter by prefix...

Buckets / python_output

| Name | Size | Type | Storage class | Last modified | Public access | Encryption | Retention expiration date | Holds |
|-------------------------|-----------|----------|----------------|---------------------------|---------------|--------------------|---------------------------|-------|
| coin_reddit_final.csv | 2.45 MB | text/csv | Multi-Regional | 5/29/19, 8:38:30 PM UTC-5 | Not public | Google-managed key | - | None |
| dates_table_final.csv | 75.8 KB | text/csv | Multi-Regional | 5/29/19, 8:50:09 PM UTC-5 | Not public | Google-managed key | - | None |
| google_trends_final.csv | 1.41 MB | text/csv | Multi-Regional | 5/29/19, 8:38:30 PM UTC-5 | Not public | Google-managed key | - | None |
| pricing_final.csv | 4.25 MB | text/csv | Multi-Regional | 5/29/19, 8:38:32 PM UTC-5 | Not public | Google-managed key | - | None |
| tweets_final.csv | 661.74 KB | text/csv | Multi-Regional | 5/29/19, 8:38:32 PM UTC-5 | Not public | Google-managed key | - | None |

depa-final-project

MySQL Second Generation master

OVERVIEW CONNECTIONS USERS DATABASES

MySQL databases

Create database

| Name | Character set | Collation | Type |
|--------------------|---------------|-----------------|--------|
| information_schema | utf8 | utf8_general_ci | System |
| CRYPTODB | utf8 | utf8_general_ci | User |
| mysql | utf8 | utf8_general_ci | System |
| performance_schema | utf8 | utf8_general_ci | System |
| sys | utf8 | utf8_general_ci | User |

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

VM instances

[CREATE INSTANCE](#) [Download](#) [Refresh](#) [Play](#) [Stop](#) [Delete](#) [SHOW INFO PANEL](#)

Instance "depa-final-project" is underutilized. You can save an estimated \$15 per month by switching to the machine type: g1-small (1 vCPU, 1.7 GB memory). [Learn more](#) [Dismiss](#)

Filter VM instances

| <input type="checkbox"/> Name ^ | Zone | Recommendation | In use by | Internal IP | External IP | Connect |
|--|---------------|----------------|-----------|-------------------|-------------|---------|
| <input checked="" type="checkbox"/> depa-final-project | us-central1-a | Save \$15 / mo | | 10.128.0.2 (nic0) | 35.224.4.20 | SSH |

python_output

[Objects](#) [Overview](#) [Permissions](#) [Bucket Lock](#)

[Upload files](#) [Upload folder](#) [Create folder](#) [Manage holds](#) [Delete](#)

Filter by prefix...

Buckets / python_output

| <input type="checkbox"/> Name | Size | Type | Storage class | Last modified | Public access | Encryption | Retention expiration date | Holds |
|--|-----------|----------|----------------|---------------------------|---------------|--------------------|---------------------------|-------|
| <input type="checkbox"/> coin_reddit_final.csv | 2.45 MB | text/csv | Multi-Regional | 5/29/19, 8:38:30 PM UTC-5 | Not public | Google-managed key | - | None |
| <input type="checkbox"/> dates_table_final.csv | 75.8 KB | text/csv | Multi-Regional | 5/29/19, 8:50:09 PM UTC-5 | Not public | Google-managed key | - | None |
| <input type="checkbox"/> google_trends_final.csv | 1.41 MB | text/csv | Multi-Regional | 5/29/19, 8:38:30 PM UTC-5 | Not public | Google-managed key | - | None |
| <input type="checkbox"/> pricing_final.csv | 4.25 MB | text/csv | Multi-Regional | 5/29/19, 8:38:32 PM UTC-5 | Not public | Google-managed key | - | None |
| <input type="checkbox"/> tweets_final.csv | 661.74 KB | text/csv | Multi-Regional | 5/29/19, 8:38:32 PM UTC-5 | Not public | Google-managed key | - | None |

depafinalproject

MySQL Second Generation master

[OVERVIEW](#) [CONNECTIONS](#) [USERS](#) [DATABASES](#)

MySQL databases

[Create database](#)

| Name | Character set | Collation | Type |
|--------------------|---------------|-----------------|--------|
| information_schema | utf8 | utf8_general_ci | System |
| CRYPTODB | utf8 | utf8_general_ci | User |
| mysql | utf8 | utf8_general_ci | System |
| performance_schema | utf8 | utf8_general_ci | System |
| sys | utf8 | utf8_general_ci | User |

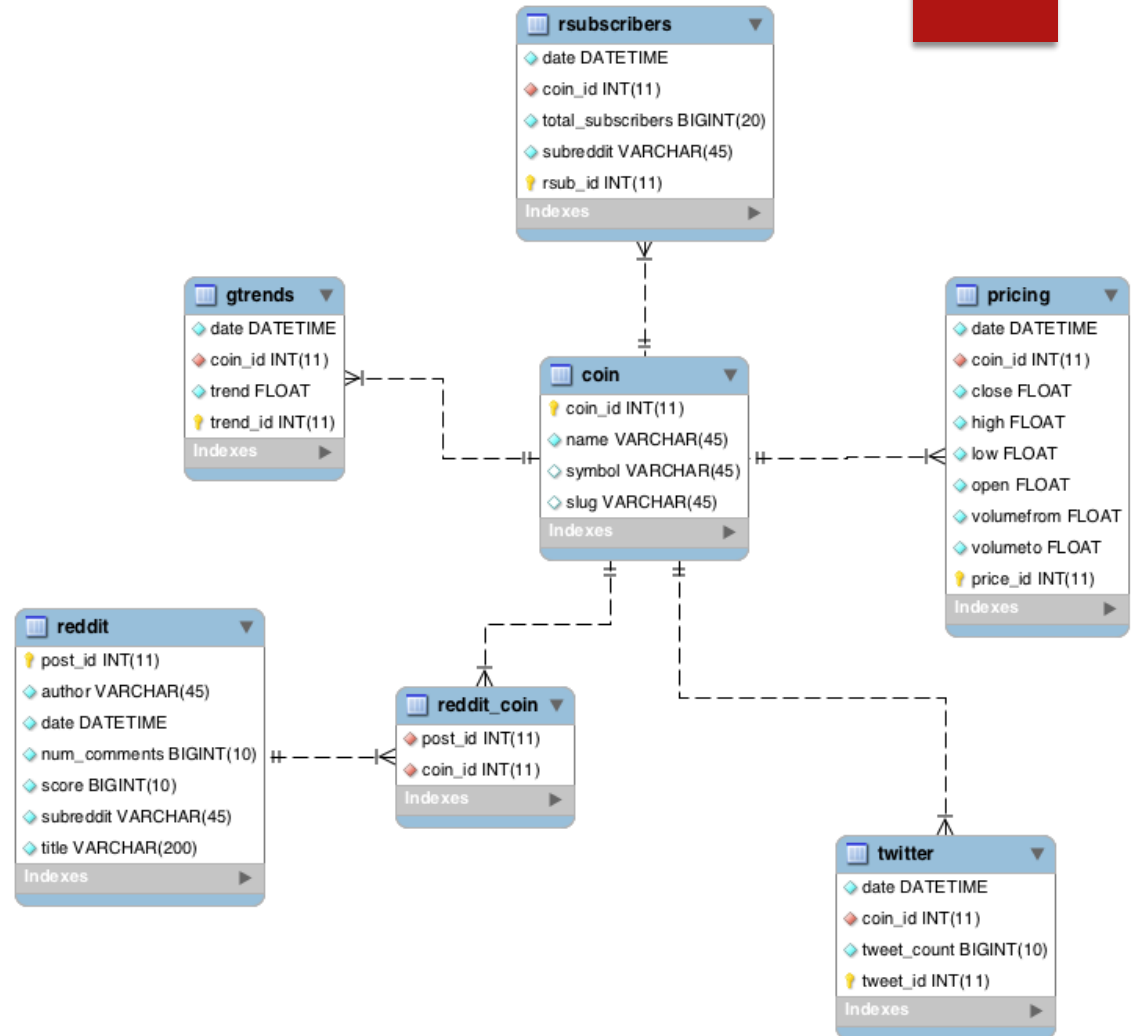
Enhanced Entity Relationship Diagram

► 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 data to coin ID
- date : date information

► Relationship & Cardinality

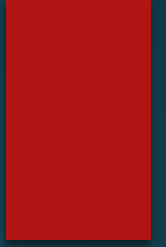
► Datatypes



Dimensional Model (Snowflake)

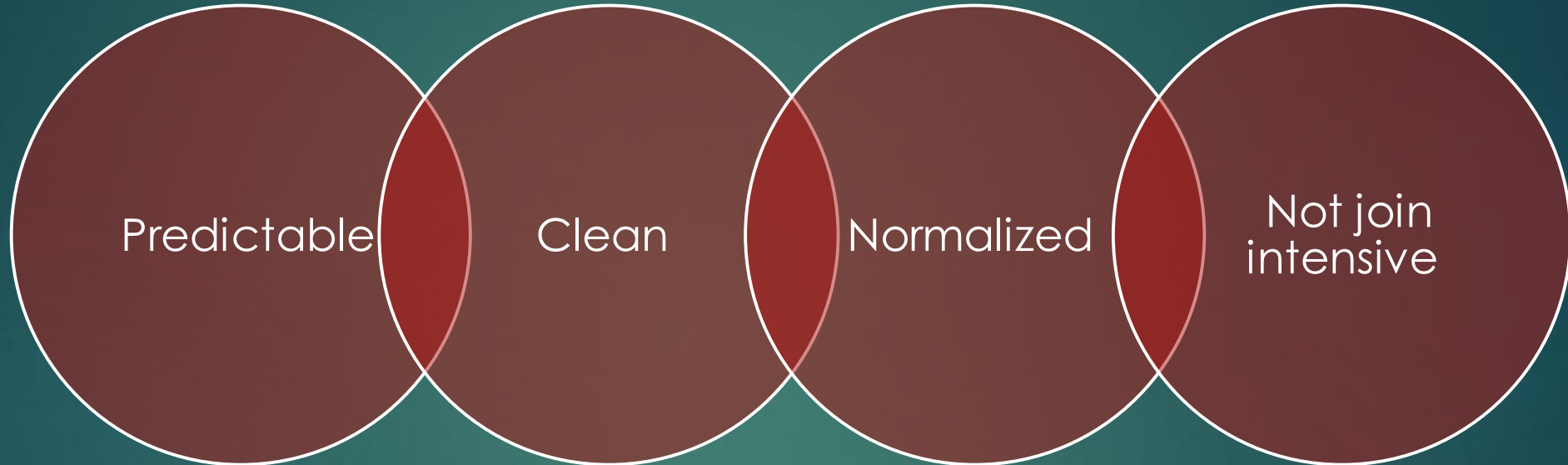
TBD

Tableau



MongoDB

Our data is:



Thus, for our particular database, we decided not use MongoDB.

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

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.

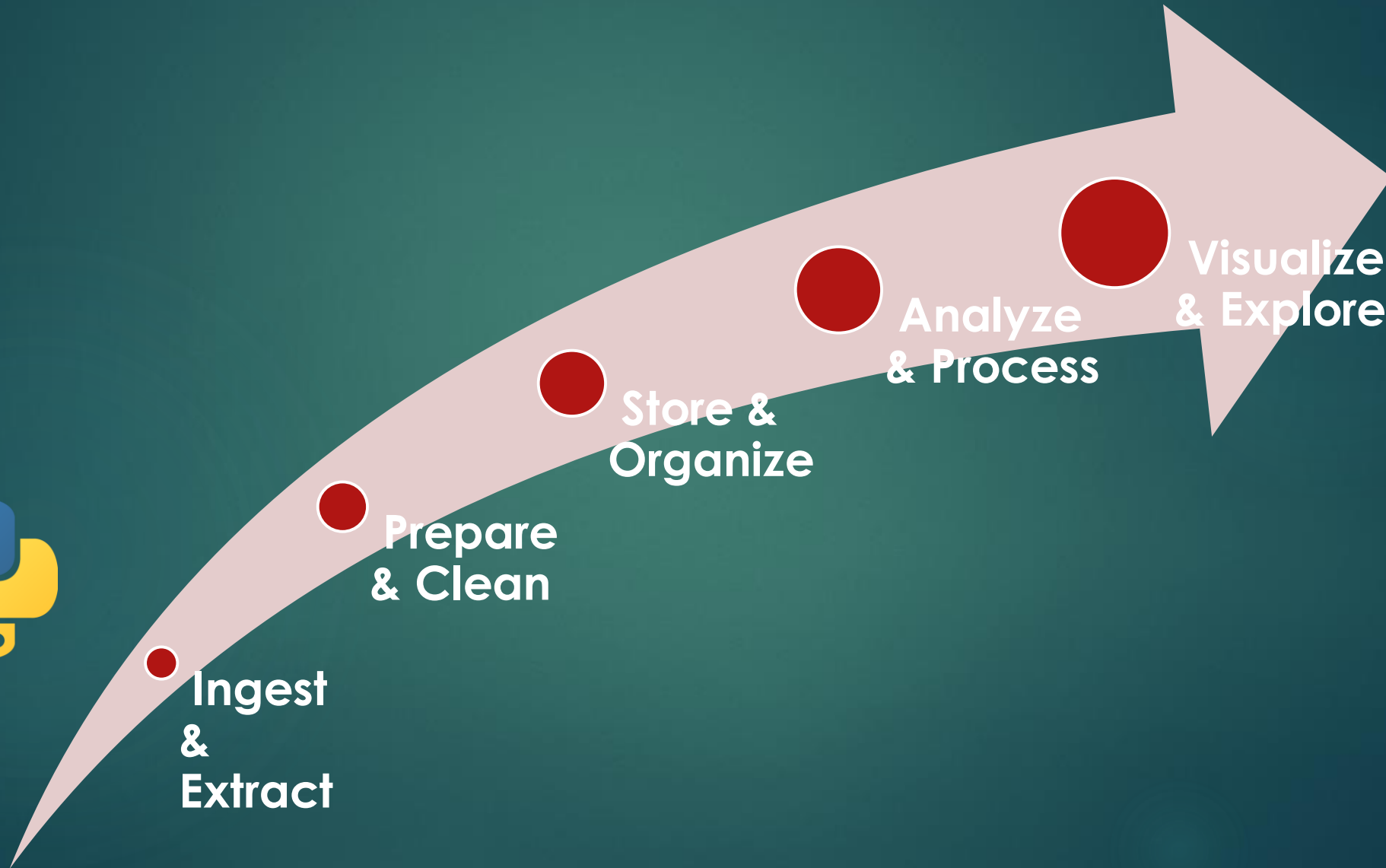
Not efficient for high volumes of transactions

- 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.

Summary of Data Analysis Pipeline



Ingest
&
Extract

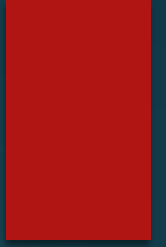
Prepare
&
Clean

Store &
Organize

Analyze
&
Process

Visualize
&
Explore

Future Recommendations



Conclusion

