

YZV211E - INTRODUCTION TO DATA SCIENCE AND ENGINEERING

Emircan EROL, 150200324, Istanbul Technical University

Bitcoin Data

I could not manage to use an API because historical data is not free and after research I found a website and downloaded historical data(I did not want to use ready data). Also I wrote code for Crypto Compare API.

Tweet Data

My TwitterAPI application is under review, twint is used in this project. In order to initialize twint configurations are declared and because of twitter a problem in twint which is global and recent. Code automatically detects last date and goes on. At the end file is saved as space separated values.

Code Snippet 1: get tweets

```
1 # allows nested use of asyncio.run
2 nest_asyncio.apply()
3
4 c = twint.Config()
5 c.Username = 'elonmusk'
6 c.Since = '2019-09-08 17:57:00'
7
8 with open('elonMuskTweets.csv') as f:
9     for line in f:
10         pass
11         last_line = line
12 date = last_line[20:39]
13
14 c.Until = date
15 c.Output_pandas = True
16 c.Output = 'elonMuskTweets.csv'
17 c.Show_cashtags = True
18 c.Show_hashtags = True
19
20 # twint is searching
21 twint.run.Search(c)
```

Beautify Tweets

Because tweets have space characters and file is also space separated, division is a must(5th and 6th rows). Timezone is changed from UTC+03:00 to UTC±00:00. Unix time is created and data frame is returned.

Code Snippet 2: edit tweet

```
1 with open('elonMuskTweets.csv', 'r') as f:
2     titles = f.readline().strip().split(' ')
3     df = pd.DataFrame(columns=["unix", "time", "tweet"])
4     for line in f:
5         infos = line[:56].split(' ')
6         tweet = line[56:].strip()
7
8         date = infos[1]
9         hour = infos[2]
10
11         three_hours = datetime.timedelta(hours=3)
12         time_obj = datetime.datetime.strptime(date + " " + hour, '%Y-%m-%d %H:%M:%S') - three_hours
13
14         # unix timestamp is created
15         unix = time.mktime(time_obj.timetuple())
16
17         # new rows are added to df
18         df = df.append({'unix': unix, 'time': time_obj, 'tweet': tweet}, ignore_index=True)
```

Control Missed Value

If values are not ordered (missing values) this function warns and prints first distorted unix timestamp. It does not effect anything just control if there is some missing value.

Code Snippet 3: edges

```
1 def controlMissedValue(btc:pd.DataFrame):
2     i = int(btc['unix'][0] / 60)
3     for x in btc['unix']:
4         x = int(x/60)
5         if i != x:
6             print(x*60)
7             break
8     i -= 1
```

Fill Empty Data

Code Snippet 4: fill empty data

```
1 with open('BTCUSDnew.csv', 'r') as f:
2     # file created to write
3     not_empty = open('BTCUSD_min.csv', 'w')
4     titles = f.readline().strip().split(',')
5     prev_line = f.readline().strip().split(',')
6     not_empty.write(",".join(titles[:-1]) + '\n' + ",".join(←
7         prev_line) + '\n')
8     prev_unix = prev_line[0]
9     for line in f:
10         line = line.strip().split(',')
11         unix = line[0]
12         if int(unix) == (int(prev_unix) - 60):
13             prev_line = line
14             prev_unix = unix
15             not_empty.write(",".join(line) + '\n')
16             continue
17
18         space_len = int((int(prev_unix) - int(unix)) / 60) - 1
19         for i in range(space_len, 0, -1):
20             new_unix = int(unix) + (60 * i)
21             new_line = prev_line
22             new_line[0] = str(new_unix)
23             not_empty.write(",".join(new_line) + '\n')
24             prev_unix = unix
25             prev_unix = unix
26             not_empty.write(",".join(line) + '\n')
27
28 not_empty.close()
```

Join Tweet and Prices

My choice is to reach price values by their indices. Location on prices data frame created based on unix data and index is calculated. Hourly, day start, day end unix timestamps are also derived indices and unix timestamp.

Code Snippet 5: join tweet prices

```
1 cols = ['unix', 'date', 'tweet', 'day start', 'tweet time price', 'ten←
2     minute', 'hour', 'day end']
3
4 df = pd.DataFrame(columns=cols) # create a dataframe to save data
```

```

3     for line in tw.iloc: # every line in tweets DataFrame .iloc is
4
5         temp = line['unix'] # to save time value because t will change
6         t=line['unix'] # better for calculations
7         t = t-t%60 # round to minute
8         date = line[1] # better to represent as date
9         tweet = line['tweet']
10
11        int_location = int((btc['unix'][0] - t) / 60) # calculate ←
            where the value is
12        ds = t - t%86400 # day start
13
14        #prices
15        nw = btc['unix'][int_location] # get open price for tweet time
16        nw_open = btc['open'][int_location]
17        ten_minute = btc['open'][int_location - 10] # int_loc + 10 ←
            minutes
18        hour = btc['open'][int_location - 60]
19        day_loc = int((btc['unix'][0]-ds)/60)
20        day_start = btc['open'][day_loc]
21        day_end = btc['open'][day_loc - 1440] # day_loc - 1440: go to ←
            1 day later
22
23        # create a temporary dataframe and append to df
24        final_line = pd.DataFrame({'unix':[temp],
25                                   'date':[date],
26                                   'tweet':[tweet],
27                                   'day start':[day_start],
28                                   'tweet time price':[nw_open],
29                                   'ten minute':[ten_minute],
30                                   'hour':[hour],
31                                   'day end':[day_end]})
32
33        df = df.append(final_line, ignore_index=True)

```
