

```

1 # modules.project_helper.py
2 # https://github.com/QuantCS109/TrumpTweets/blob/master/modules/project_helper.py
3 # This is a .py file used in the notebooks
4
5
6 import pandas as pd
7 import numpy as np
8 from collections import Counter
9 from pytz import timezone
10 from datetime import timedelta
11 import pickle
12 import re
13 from sklearn.ensemble import RandomForestClassifier
14
15
16
17 class TweetData:
18
19     def __init__(self, file='../data/input_data/trump_archive_db.csv'):
20         self.file = file
21         self.raw_data = []
22         self.error_tweets = {}
23
24         self.read_data()
25         self.raw_tweets = self.parse()
26         self.clean_tweets = self.clean()
27         self.text = self.create_text()
28         self.words = self.tokenize_text()
29         self.vocab_to_int, self.int_to_vocab = self.create_lookup_tables()
30         self.int_words = self.create_int_words()
31
32         self.daily_tweets = None
33         self.get_daily_tweets()
34
35     def read_data(self):
36         with open(self.file, mode='r', errors='ignore') as f:
37             for row in f:
38                 self.raw_data.append(row)
39
40     def parse(self):
41         timestamps = []
42         tweets = []
43         raw_tweets = pd.DataFrame(columns=['tweets'])
44         for i, tweet in enumerate(self.raw_data):
45             try:
46                 timestamps.append(timezone('US/Central').localize(pd.to_datetime((tweet
47                     [-21:-2]))))
48                 tweets.append(tweet[:-22])
49             except:
50                 self.error_tweets[i] = tweet
51         raw_tweets['tweets'] = tweets
52         raw_tweets.index = timestamps
53         raw_tweets.index.name = 'timestamp'
54
55         return raw_tweets
56
57     def clean_step_1(self, tweet):
58
59         # Remove whitespace before and after tweet
60         tweet = tweet.strip(' ')
61         tweet = tweet.lstrip('\n')
62         tweet = tweet + '\n\n'
63         return tweet
64
65     def clean_step_2(self, tweet):
66         # Remove http links
67         tweet = re.sub(r"http\S+", '', tweet)

```

```

67     # Remove hash tags
68     tweet = re.sub(r"#\S+", '', tweet)
69     # Remove twitter handles
70     tweet = re.sub(r"@\S+", '', tweet)
71     # Turn everything to lower case
72     tweet = tweet.lower()
73     # Remove symbols other than letters in the alphabet and numbers
74     tweet = re.sub(r"\'", '', tweet)
75     tweet = re.sub(r'[^a-zA-Z0-9]', ' ', tweet)
76     # Remove whitespace before and after tweet, add one white space
77     tweet = re.sub(r'[^a-zA-Z]', ' ', tweet)
78     tweet = ' '.join(tweet.split())
79     tweet = tweet + ' '
80     return tweet
81
82 def clean(self):
83     clean_tweets = pd.DataFrame(columns=['tweets'])
84     clean_tweets['tweets'] = self.raw_tweets['tweets'].apply(self.clean_step_1)
85     clean_tweets['tweets'] = clean_tweets['tweets'].apply(self.clean_step_2)
86     clean_tweets.index = self.raw_tweets.index
87     return clean_tweets
88
89 def create_text(self):
90     text = ''.join(self.clean_tweets.tweets)
91     return text
92
93 def tokenize_text(self):
94     words = self.text.split()
95     # Remove all words with 5 or fewer occurrences
96     word_counts = Counter(words)
97     return [word for word in words if word_counts[word] > 5]
98
99 def create_lookup_tables(self):
100
101     word_counts = Counter(self.words)
102     # words sorted in descending frequency
103     sorted_vocab = sorted(word_counts, key=word_counts.get, reverse=True)
104     int_to_vocab = {ii: word for ii, word in enumerate(sorted_vocab)}
105     vocab_to_int = {word: ii for ii, word in int_to_vocab.items()}
106     return vocab_to_int, int_to_vocab
107
108 def create_int_words(self):
109     return [self.vocab_to_int[word] for word in self.words]
110
111 def get_daily_tweets(self):
112     self.clean_tweets['timestamp'] = self.clean_tweets.index
113     after_4_tweets = self.clean_tweets.timestamp.dt.hour >= 15
114     self.clean_tweets['after4_date'] = self.clean_tweets.timestamp.dt.date
115     self.clean_tweets.loc[after_4_tweets, 'after4_date'] = self.clean_tweets.
116                                     timestamp[after_4_tweets].dt.date\
117                                     + timedelta(days=1)
118     self.daily_tweets = self.clean_tweets.groupby('after4_date')['tweets'].apply(
119         lambda x: ' '.join(x))
120     self.daily_tweets = self.daily_tweets.to_frame('tweets')
121     self.daily_tweets.index.name = 'date'
122
123 class APIData(TweetData):
124
125     def __init__(self, file='../data/input_data/trumptwits.csv'):
126         super().__init__(file=file)
127
128     def read_data(self):
129         self.raw_data = pd.read_csv(self.file)
130
131     def parse(self):
132         raw_tweets = pd.DataFrame(columns=['tweets'])

```

```

132         raw_tweets['tweets'] = self.raw_data['text']
133         raw_tweets.index = pd.to_datetime(self.raw_data['time'])
134         raw_tweets.index.name = 'timestamp'
135         raw_tweets = raw_tweets.sort_index()
136         raw_tweets = raw_tweets.tz_convert('US/Central')
137         return raw_tweets
138
139     def clean_step_1(self, tweet):
140         tweet = tweet.rstrip('b\\')
141         tweet.rstrip('\\n\\n\\n\\n')
142         tweet.rstrip('\\')
143         return tweet
144
145
146 class IntradayData:
147
148     def __init__(self, file='../data/input_data/ES_intraday.csv'):
149         self.file = file
150         self.raw_data = self.read_data()
151
152     def read_data(self):
153         fin_data = pd.read_csv(self.file)
154         fin_data.index = pd.to_datetime(fin_data['Date'] + ' ' + fin_data['Time']).dt.
            tz_localize('US/Central')
155         fin_data.index.name = 'timestamp'
156         fin_data = fin_data.drop(columns=['Date', 'Time'])
157         return fin_data
158
159     def get_data(self):
160         return self.raw_data[['Open', 'Close']]
161
162
163 class FuturesCloseData:
164     def __init__(self, path='../data/input_data/futures_close.csv'):
165         self.instrument_list = ['ES', 'NQ', 'CD', 'EC', 'JY', 'MP', 'TY', 'US', 'C', 'S',
            , 'W', 'CL', 'GC']
166         self.df = self.load(path)
167
168     def load(self, path):
169         df = pd.read_csv(path)
170         df.set_index('date', inplace=True)
171         df.index = pd.to_datetime(df.index)
172         return df
173
174     def features(self, inst):
175         return self.momentum(inst)
176
177     def price(self, inst):
178         return self.df[inst]
179
180     def returns(self, inst, start=0, end=1):
181         returns = (self.df[inst].shift(-end) - self.df[inst].shift(-start)) / self.df[
            inst].shift(-start)
182         return returns
183
184     def momentum(self, inst, lag=60):
185         momo = pd.DataFrame((self.df[inst] - self.df[inst].shift(lag)) / self.df[inst])
186         momo = momo.dropna()
187         momo.columns += '_{}D'.format(lag)
188         return momo
189
190     def log_returns(self, start=0, end=1):
191         return np.log(self.df.shift(-end)) - np.log(self.df.shift(-start))
192
193     def single_log_returns(self, inst, start=0, end=1):
194         return np.log(self.df[inst].shift(-end)) - np.log(self.df[inst].shift(-start))
195

```

```

196
197 class VolFeatures:
198     def __init__(self, path='../data/features/vol_features.pkl'):
199         self.instrument_list = ['ES', 'NQ', 'CD', 'EC', 'JY', 'MP', 'TY', 'US', 'C', 'S'
200             , 'W', 'CL', 'GC']
201         self.df = self.load(path)
202         self.col_dict = {inst: [key for key in self.df.columns if re.match(r"{}_+".
203             format(inst), key)]
204             for inst in self.instrument_list}
205
206     def features(self, inst):
207         return self.df[self.col_dict[inst]]
208
209     def load(self, path):
210         pickle_in = open(path, "rb")
211         vol_pd = pickle.load(pickle_in)
212         pickle_in.close()
213         return vol_pd.fillna(vol_pd.mean())
214
215 class TweetReturnsFeatures(VolFeatures):
216     def __init__(self, path='../data/features/tweet_returns_features.csv'):
217         super().__init__(path)
218
219     def load(self, path):
220         tweet_returns = pd.read_csv(path)
221         tweet_returns.set_index('date', inplace=True)
222         tweet_returns.index = pd.to_datetime(tweet_returns.index)
223         return tweet_returns
224
225 class MarketFeatures(TweetReturnsFeatures):
226     def __init__(self, path='../data/features/market_features.csv'):
227         super().__init__(path)
228
229 class TradeModel:
230
231     def __init__(self, model=RandomForestClassifier, *args, **kwargs):
232         self.model = model(*args, **kwargs)
233
234     def fit(self, X, y):
235         self.model.fit(X, y)
236
237     def position(self, X, cutoff=0.55):
238         # converting predictions from {0,1} to {-1,1}, short/long
239         position = 2 * self.model.predict(X) - 1
240         position[self.model.predict_proba(X).max(axis=1) <= cutoff] = 0
241         return position
242
243     def _strategy_returns(self, x, y):
244         strat_ret = x[:-2] * y[:-2]
245         strat_ret_cum = strat_ret.cumsum()
246         return strat_ret, strat_ret_cum
247
248     def strategy_returns(self, X, returns, cutoff=0.55):
249         return self._strategy_returns(returns, self.position(X, cutoff))
250
251     def sharpe(self, X, returns, cutoff=0.55):
252         rets = self.strategy_returns(X, returns, cutoff)[0]
253         return np.mean(rets) / np.std(rets)
254
255
256
257
258
259
260

```