

```
1 from airflow import DAG
2 from airflow.operators.python import PythonOperator
3
4 from datetime import datetime, timedelta
5 import yfinance as yf
6 import pickle
7 import pandas as pd
8 from sklearn.linear_model import LinearRegression
9 from pathlib import Path
10 import logging
11
12
13 COMPANIES = ['AAPL', 'GOOGL', 'FB', 'MSFT', 'AMZN']
14
15 DF_DIRECTORY = "hw4/stock_dfs.pkl"
16 MODELS_DIRECTORY = "hw4/models.pkl"
17 HISTORY_DIRECT = 'hw4/histories.pkl'
18 RELATIVE_ERRORS_DIRECT = "hw4/relative_errors.csv"
19 PREDICTION_DIRECT = 'hw4/predictions.pkl'
20
21 FEATURE_COLS = ["Open", "Low", "Close", "Volume", '
    High']
22 PREDICTED_COL = ["High"]
23
24 default_args = {
25     'owner': 'yewen',
26     'depends_on_past': False,
27     'email': ['yz4175@columbia.edu'],
28     'email_on_failure': False,
29     'email_on_retry': False,
30     'retries': 1
31 }
32
33
34 def get_date_today():
35     """Get date today in str format such as 20201119
36     . """
37     return datetime.today().strftime('%Y%m%d')
38
39 def get_date_yesterday():
```

```

40     """Return yesterday in str such as 2021-11-19
41     ."""
42
43
44     return (datetime.today() - timedelta(days=1)).
45     strftime('%Y-%m-%d')
46
47
48 def get_data():
49     """Get new data in the past 10 days up to
50     including yesterday.
51     For example, if today is 11/22/2021, only collect
52     data up to
53     including 11/21/2021. Append as a list. Save to
54     pickle file.
55     """
56     stock_dfs = []
57     for i, company in enumerate(COMPANIES):
58         # because I run at 7 am everyday, using
59         period='10d'
60         # would get the data in the past 10 days
61         including
62         # yesterday
63         stock_df = yf.download(company, period='10d')
64         stock_dfs.append(stock_df)
65     with open(DF_DIRECTORY, 'wb') as f:
66         pickle.dump(stock_dfs, f)
67
68
69 def get_historical_data():
70     """Get all historical data up to including
71     yesterday. """
72     histories = []
73     for i, company in enumerate(COMPANIES):
74         # because I run at 7 am every day, this would
75         # actually just get me all the historic data
76         history = yf.download(company, period='max')
77         histories.append(history)
78     with open(HISTORY_DIREC, 'wb') as f:
79         pickle.dump(histories, f)
80
81
82 def compute_relative_error():

```

```

73         """Compute the relative error for today, update
CSV.
74         For example, today is 11/22, compute the
relative error between
75         prediction and high on 11/21.
76
77         On the 1st run, it will pass.
78         On the 2nd run, it will create a csv.
79         On the later runs, it will read the csv.
80         """
81         with open(DF_DIRECTORY, 'rb') as f:
82             # this data includes stock up to 11/21 in
the past
83             # 10 days
84             stock_dfs = pickle.load(f)
85
86         if Path(PREDICTION_DIREC).is_file():
87             with open(PREDICTION_DIREC, 'rb') as f:
88                 # a list of 5 predictions made on 11/21.
89                 predictions = pickle.load(f)
90
91         if Path(PREDICTION_DIREC).is_file():
92             relative_errors = []
93             for i, stock_df in enumerate(stock_dfs):
94                 # high price on 11/21.
95                 high = stock_df.iloc[-1, :]["High"]
96                 # prediction on 11/21.
97                 prediction = predictions[i]
98                 relative_error = (prediction - high) /
high
99                 relative_errors.append(relative_error)
100
101         if Path(RELATIVE_ERRORS_DIREC).is_file():
102             logging.info("Updating relative_errors
csv...")
103             relative_errors_df = pd.read_csv(
RELATIVE_ERRORS_DIREC)
104             relative_errors_df[get_date_yesterday
()] = relative_errors
105             relative_errors_df.to_csv(
RELATIVE_ERRORS_DIREC)

```

```

106         else:
107             # create a CSV
108             logging.info("Creating relative_errors
csv...")
109             relative_errors_df = pd.DataFrame(data=
relative_errors,
110             columns=[get_date_yesterday()],
111                                     index=
COMPANIES)
112             relative_errors_df.to_csv(
RELATIVE_ERRORS_DIREC)
113
114
115 def train_models():
116     """Make Linear Regression Models. """
117     with open(HISTORY_DIREC, 'rb') as f:
118         histories_df = pickle.load(f)
119
120     models = []
121     for history_df in histories_df:
122         n = history_df.shape[0]
123         logging.info(f"History shape: {history_df.
shape}")
124         reg = LinearRegression()
125         Xs = []
126         ys = []
127         for i in range(0, n-10):
128             slice = history_df.iloc[i:i + 10, :][
FEATURE_COLS]
129             X = slice.to_numpy().flatten().tolist()
130             Xs.append(X)
131             y = history_df.iloc[i + 10]['High']
132             ys.append(y)
133             reg.fit(Xs, ys)
134             models.append(reg)
135
136     with open(MODELS_DIRECTORY, 'wb') as f:
137         pickle.dump(models, f)
138
139

```

```

140 def make_prediction():
141     """Make a prediction for today. """
142     with open(MODELS_DIRECTORY, 'rb') as f:
143         models = pickle.load(f)
144
145     with open(DF_DIRECTORY, 'rb') as f:
146         stock_dfs = pickle.load(f)
147
148     predictions = []
149     for i, stock_df in enumerate(stock_dfs):
150         model = models[i]
151         X = stock_df[FEATURE_COLS].to_numpy().
152         flatten().reshape(1, -1)
153         logging.info(X)
154         prediction = model.predict(X)[0]
155         predictions.append(prediction)
156
157     with open(PREDICTION_DIREC, 'wb') as f:
158         pickle.dump(predictions, f)
159
160 with DAG('stock',
161         default_args=default_args,
162         description='Stock Price Analysis for HW4',
163         catchup=False,
164         start_date=datetime(2021, 1, 1),
165         schedule_interval='0 7 * * *',
166         ) as dag:
167     get_data_task = PythonOperator(task_id='get_data
168     ',
169                                     python_callable=
170                                     get_data)
171
172     get_historical_data_task = PythonOperator(
173     task_id='get_historic_data',
174     python_callable=get_historical_data)
175
176     make_prediction_task = PythonOperator(task_id="
177     make_prediction",

```

```
174 python_callable=make_prediction)
175
176     train_model_task = PythonOperator(task_id='
train_model',
177
178     python_callable=train_models)
178
179     compute_relative_error_task = PythonOperator(
180     task_id="compute_relative_error",
181
182     python_callable=compute_relative_error)
181
182     # task dependencies
183     # (get_data_task >> get_historical_data_task >>
184     train_model_task
185     # >> make_prediction_task >>
186     compute_relative_error_task)
185     (get_data_task >> get_historical_data_task >>
186     train_model_task >> make_prediction_task >>
187     compute_relative_error_task)
187
```