

```
In [ ]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from statsmodels.tsa.stattools import adfuller
from sklearn.preprocessing import OrdinalEncoder
from sklearn.model_selection import train_test_split
from statsmodels.tsa.seasonal import seasonal_decompose

from sklearn.ensemble import RandomForestRegressor
from sklearn.model_selection import ParameterGrid
from sklearn.ensemble import GradientBoostingRegressor
```

```
In [ ]: df = pd.read_csv('train.csv', index_col = 'datetime', parse_dates = True)
df
```

| | ID | temperature | var1 | pressure | windspeed | var2 | electricity_consumption |
|---------------------|-------|-------------|-------|----------|-----------|------|-------------------------|
| datetime | | | | | | | |
| 2013-07-01 00:00:00 | 0 | -11.4 | -37.1 | 1003.0 | 571.910 | A | 216.0 |
| 2013-07-01 01:00:00 | 1 | -12.1 | -38.3 | 996.0 | 575.040 | A | 210.0 |
| 2013-07-01 02:00:00 | 2 | -12.9 | -20.0 | 1000.0 | 578.435 | A | 225.0 |
| 2013-07-01 03:00:00 | 3 | -11.4 | -37.1 | 996.0 | 582.580 | A | 216.0 |
| 2013-07-01 04:00:00 | 4 | -11.4 | -38.3 | 1005.0 | 586.600 | A | 222.0 |
| ... | ... | ... | ... | ... | ... | ... | ... |
| 2017-06-23 19:00:00 | 34895 | -0.7 | -15.0 | 1009.0 | 51.665 | A | 225.0 |
| 2017-06-23 20:00:00 | 34896 | -2.9 | -11.4 | 1005.0 | 56.105 | A | 213.0 |
| 2017-06-23 21:00:00 | 34897 | -1.4 | -12.9 | 996.0 | 56.375 | A | 210.0 |
| 2017-06-23 22:00:00 | 34898 | -2.9 | -11.4 | 996.0 | 67.210 | A | 210.0 |
| 2017-06-23 23:00:00 | 34899 | -2.1 | -11.4 | 1009.0 | 71.880 | A | 210.0 |

20490 rows x 7 columns

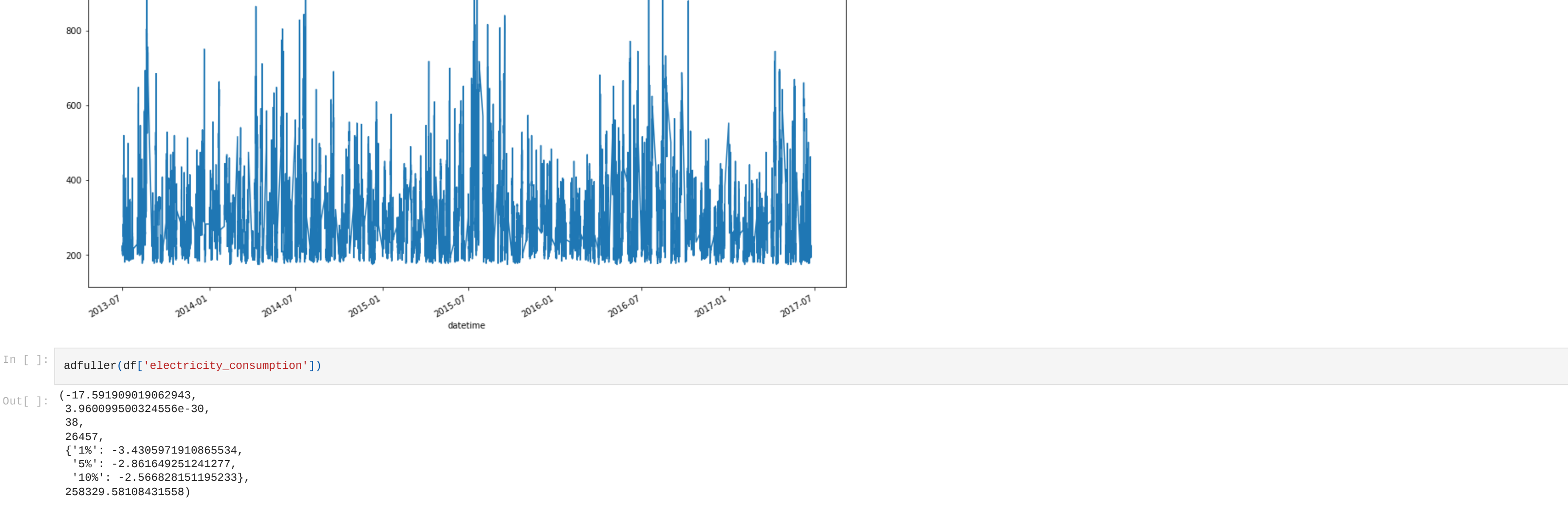
Changing var2 from letters to numbers

```
In [ ]: enc = OrdinalEncoder()
enc.fit(df[['var2']])
df['var2'] = enc.transform(df[['var2']])
df
```

| | ID | temperature | var1 | pressure | windspeed | var2 | electricity_consumption |
|---------------------|-------|-------------|-------|----------|-----------|------|-------------------------|
| datetime | | | | | | | |
| 2013-07-01 00:00:00 | 0 | -11.4 | -37.1 | 1003.0 | 571.910 | 0.0 | 216.0 |
| 2013-07-01 01:00:00 | 1 | -12.1 | -38.3 | 996.0 | 575.040 | 0.0 | 210.0 |
| 2013-07-01 02:00:00 | 2 | -12.9 | -20.0 | 1000.0 | 578.435 | 0.0 | 225.0 |
| 2013-07-01 03:00:00 | 3 | -11.4 | -37.1 | 996.0 | 582.580 | 0.0 | 216.0 |
| 2013-07-01 04:00:00 | 4 | -11.4 | -38.3 | 1005.0 | 586.600 | 0.0 | 222.0 |
| ... | ... | ... | ... | ... | ... | ... | ... |
| 2017-06-23 19:00:00 | 34895 | -0.7 | -15.0 | 1009.0 | 51.665 | 0.0 | 225.0 |
| 2017-06-23 20:00:00 | 34896 | -2.9 | -11.4 | 1005.0 | 56.105 | 0.0 | 213.0 |
| 2017-06-23 21:00:00 | 34897 | -1.4 | -12.9 | 996.0 | 61.375 | 0.0 | 210.0 |
| 2017-06-23 22:00:00 | 34898 | -2.9 | -11.4 | 996.0 | 67.210 | 0.0 | 210.0 |
| 2017-06-23 23:00:00 | 34899 | -2.1 | -11.4 | 1009.0 | 71.880 | 0.0 | 210.0 |

20490 rows x 7 columns

```
In [ ]: df['electricity_consumption'].plot(figsize = (16, 12))
Out [ ]: <AxesSubplot: xlabel='datetime'>
```



```
In [ ]: adfuller(df['electricity_consumption'])
Out [ ]: (1.37, 0.0000000000000000, 3.9609999999324556e+36, 18, 26457, ['n': -3.48899710108000524, 'dn': -2.86168251241277, 'dm': -2.56628131192293], 256329, [6108431558])
```

Time series is stationary

Using datetime as variables

```
In [ ]: df['hour'] = df.index.hour
df['month'] = df.index.month
df['day_of_month'] = df.index.day
df['day_of_week'] = df.index.day_of_week
df
```

| | ID | temperature | var1 | pressure | windspeed | var2 | electricity_consumption | hour | month | day_of_month | day_of_week |
|---------------------|-------|-------------|-------|----------|-----------|------|-------------------------|------|-------|--------------|-------------|
| datetime | | | | | | | | | | | |
| 2013-07-01 00:00:00 | 0 | -11.4 | -37.1 | 1003.0 | 571.910 | 0.0 | 216.0 | 0 | 7 | 1 | 0 |
| 2013-07-01 01:00:00 | 1 | -12.1 | -38.3 | 996.0 | 575.040 | 0.0 | 210.0 | 1 | 7 | 1 | 0 |
| 2013-07-01 02:00:00 | 2 | -12.9 | -20.0 | 1000.0 | 578.435 | 0.0 | 225.0 | 2 | 7 | 1 | 0 |
| 2013-07-01 03:00:00 | 3 | -11.4 | -37.1 | 996.0 | 582.580 | 0.0 | 216.0 | 3 | 7 | 1 | 0 |
| 2013-07-01 04:00:00 | 4 | -11.4 | -38.3 | 1005.0 | 586.600 | 0.0 | 222.0 | 4 | 7 | 1 | 0 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 2017-06-23 19:00:00 | 34895 | -0.7 | -15.0 | 1009.0 | 51.665 | 0.0 | 225.0 | 19 | 6 | 23 | 4 |
| 2017-06-23 20:00:00 | 34896 | -2.9 | -11.4 | 1005.0 | 56.105 | 0.0 | 213.0 | 20 | 6 | 23 | 4 |
| 2017-06-23 21:00:00 | 34897 | -1.4 | -12.9 | 996.0 | 61.375 | 0.0 | 210.0 | 21 | 6 | 23 | 4 |
| 2017-06-23 22:00:00 | 34898 | -2.9 | -11.4 | 996.0 | 67.210 | 0.0 | 210.0 | 22 | 6 | 23 | 4 |
| 2017-06-23 23:00:00 | 34899 | -2.1 | -11.4 | 1009.0 | 71.880 | 0.0 | 210.0 | 23 | 6 | 23 | 4 |

20490 rows x 12 columns

Splitting dataset into predictor values and what we're looking for

```
In [ ]: x = df.iloc[:, [1, 2, 3, 4, 5, 7, 8, 9, 10]]
X
Out [ ]: array([[1.140e+01, -1.730e+01, 1.003e+03, ..., 7.000e+00, 1.000e+00, 6.000e+00],
       [1.230e+01, -1.930e+01, 9.960e+02, ..., 7.000e+00, 1.000e+00, 6.000e+00],
       [1.230e+01, -1.930e+01, 1.000e+03, ..., 7.000e+00, 1.000e+00, 6.000e+00],
       [1.140e+00, -1.200e+01, 9.956e+02, ..., 6.000e+00, 2.300e+01, 4.000e+00],
       [2.000e+00, -1.140e+01, 9.960e+02, ..., 6.000e+00, 2.300e+01, 4.000e+00],
       [2.100e+00, -1.140e+01, 1.009e+03, ..., 6.000e+00, 2.300e+01, 4.000e+00]])
```

```
In [ ]: y = df.iloc[:, 6].values
Y
Out [ ]: array([216., 210., 225., ..., 213., 210., 210.])
```

```
In [ ]: X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size = 0.2, random_state = 42)
```

P1: Random Forest Regressor

Grid search

```
In [ ]: grid = {'n_estimators': [200], 'max_depth': [2, 3, 4, 5, 6, 7, 8, 9, 10], 'max_features': [2, 3, 4, 5, 6, 7, 8, 9], 'random_state': [17, 18, 19]}
rf = RandomForestRegressor()
test_scores = []
for g in ParameterGrid(grid):
    rf.set_params(**g) #unpacking parameters from dictionary
    rf.fit(X_train, Y_train) #fitting data
    test_scores.append(rf.score(X_test, Y_test)) #getting score
```

```
In [ ]: best_idx = np.argmax(test_scores)
print(test_scores[best_idx], ParameterGrid(grid)[best_idx])
0.5765848931450903 {'random_state': 19, 'n_estimators': 200, 'max_features': 8, 'max_depth': 18}
```

```
In [ ]: rfr = RandomForestRegressor(random_state = 19, n_estimators = 200, max_features = 8, max_depth = 18)
rfr.fit(X_train, Y_train)
print(rfr.score(X_train, Y_train))
print(rfr.score(X_test, Y_test))
0.680580415972389
0.5765848931450903
```

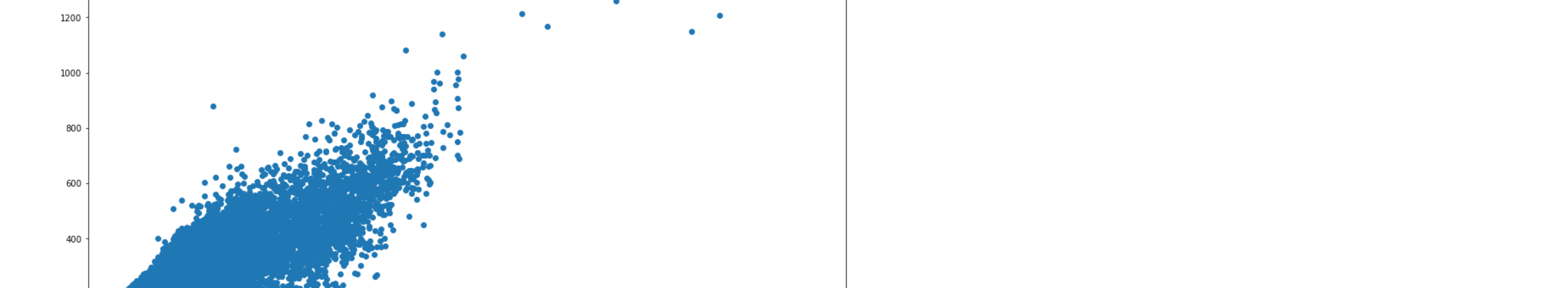
Making predictions

```
In [ ]: train_prediction = rfr.predict(X_train)
test_prediction = rfr.predict(X_test)
```

```
In [ ]: plt.rc('figure', figsize = (16, 8))
plt.scatter(train_prediction, Y_train)
```



```
In [ ]: plt.scatter(test_prediction, Y_test)
```

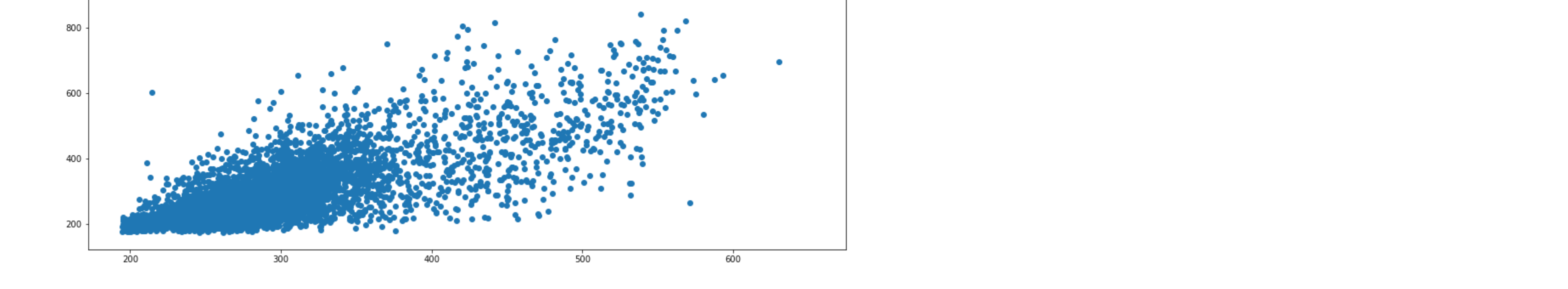


Checking which features are the most important

```
In [ ]: feature_importanceRFR = rfr.feature_importances_
sorted_index = np.argsort(feature_importanceRFR[::-1]) #sorting the index by descending order
x1 = range(len(feature_importanceRFR))
feature_names = df.columns[df.columns.tolist()[1:]] #names are the same as names of columns
labels = np.array(feature_names[sorted_index])
```

```
In [ ]: plt.bar(x1, feature_importanceRFR[sorted_index], tick_label = labels)
```

```
Out [ ]: <BarContainer object of 9 artists>
```



Final forecast

```
In [ ]: df1 = pd.read_csv('test.csv', index_col = 'datetime', parse_dates = True)

#turning var2 into numbers
enc2 = OrdinalEncoder()
enc2.fit(df1[['var2']])
df1['var2'] = enc2.transform(df1[['var2']])

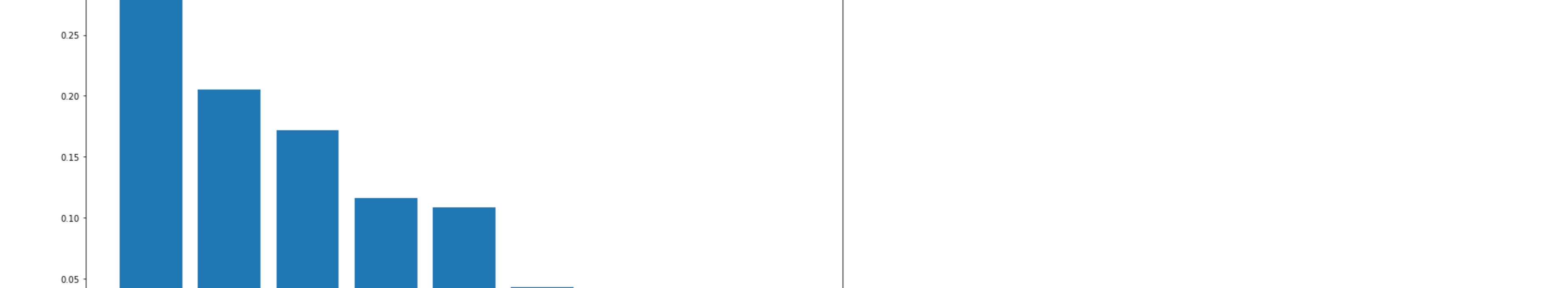
#datetime as variables
df1['hour'] = df1.index.hour
df1['month'] = df1.index.month
df1['day_of_month'] = df1.index.day
df1['day_of_week'] = df1.index.day_of_week

#X values for prediction
X_pred = df1.iloc[:, 1:11].values

#model
rfr = RandomForestRegressor(random_state = 19, n_estimators = 200, max_features = 8, max_depth = 18)
rfr.fit(X, Y)
prediction = rfr.predict(X_pred)
```

```
In [ ]: df1['electricity_consumption'] = prediction
df1['electricity_consumption'].plot()
```

```
Out [ ]: <AxesSubplot: xlabel='datetime'>
```



Saving predictions to file

```
In [ ]: df2 = pd.read_csv('sample_submission.csv')
df2['electricity_consumption'] = prediction
df2.set_index('ID', inplace = True)
df2.to_csv('sample_prediction.csv')
```

Score: lower is better

```
In [ ]: import matplotlib.image as mpimg
img = mpimg.imread('successcontestresult.png')
imgplot = plt.imshow(img)
plt.show()
```



P2: Gradient Boosted Regressor

Grid search to tune parameters

```
In [ ]: grid = {'n_estimators': [50, 75, 100], 'learning_rate': [0.1, 0.15, 0.2], 'min_samples_split': [200], 'min_samples_leaf': [30], \
       'max_depth': [5, 6, 7, 8], 'max_features': [4], 'subsample': [0.8], 'random_state': [18]}
gbr = GradientBoostingRegressor()
test_scores = []
for g in ParameterGrid(grid):
    gbr.set_params(**g) #unpacking parameters from dictionary
    gbr.fit(X_train, Y_train) #fitting data
    test_scores.append(gbr.score(X_test, Y_test)) #getting score
```

```
best_idx = np.argmax(test_scores)
print(test_scores[best_idx], ParameterGrid(grid)[best_idx])
0.5642548931450903 {'subsample': 0.8, 'random_state': 18, 'n_estimators': 100, 'min_samples_split': 250, 'min_samples_leaf': 50, 'max_features': 4, 'max_depth': 5, 'learning_rate': 0.2}
```

```
In [ ]: grid = {'n_estimators': [100], 'learning_rate': [0.2], 'min_samples_split': [200, 250, 300], 'min_samples_leaf': [30, 40, 50, 60], \
       'max_depth': [5], 'max_features': [4], 'subsample': [0.8], 'random_state': [19]}
gbr = GradientBoostingRegressor()
test_scores = []
for g in ParameterGrid(grid):
    gbr.set_params(**g) #unpacking parameters from dictionary
    gbr.fit(X_train, Y_train) #fitting data
    test_scores.append(gbr.score(X_test, Y_test)) #getting score
```

```
best_idx = np.argmax(test_scores)
print(test_scores[best_idx], ParameterGrid(grid)[best_idx])
0.5737097599748987 {'subsample': 0.8, 'random_state': 18, 'n_estimators': 100, 'min_samples_split': 200, 'min_samples_leaf': 30, 'max_features': 4, 'max_depth': 5, 'learning_rate': 0.2}
```

```
In [ ]: grid = {'n_estimators': [200], 'learning_rate': [0.2], 'min_samples_split': [200], 'min_samples_leaf': [30], \
       'max_depth': [5, 6, 7, 8], 'max_features': [7], 'subsample': [0.8], 'random_state': [19]}
gbr = GradientBoostingRegressor()
test_scores = []
for g in ParameterGrid(grid):
    gbr.set_params(**g) #unpacking parameters from dictionary
    gbr.fit(X_train, Y_train) #fitting data
    test_scores.append(gbr.score(X_test, Y_test)) #getting score
```

```
best_idx = np.argmax(test_scores)
print(test_scores[best_idx], ParameterGrid(grid)[best_idx])
0.645190487020191 {'subsample': 0.8, 'random_state': 18, 'n_estimators': 100, 'min_samples_split': 200, 'min_samples_leaf': 30, 'max_features': 7, 'max_depth': 8, 'learning_rate': 0.2}
```

```
In [ ]: grid = {'n_estimators': [100], 'learning_rate': [0.2], 'min_samples_split': [200], 'min_samples_leaf': [30], \
       'max_depth': [5], 'max_features': [7], 'subsample': [0.8, 0.7, 0.6, 0.5], 'random_state': [20]}
gbr = GradientBoostingRegressor()
test_scores = []
for g in ParameterGrid(grid):
    gbr.set_params(**g) #unpacking parameters from dictionary
    gbr.fit(X_train, Y_train) #fitting data
    test_scores.append(gbr.score(X_test, Y_test)) #getting score
```

```
best_idx = np.argmax(test_scores)
print(test_scores[best_idx], ParameterGrid(grid)[best_idx])
0.645190487020191 {'subsample': 0.8, 'random_state': 18, 'n_estimators': 100, 'min_samples_split': 200, 'min_samples_leaf': 30, 'max_features': 7, 'max_depth': 8, 'learning_rate': 0.2}
```

Decreasing learning rate

```
In [ ]: gbr = GradientBoostingRegressor(n_estimators = 400, learning_rate = 0.1, min_samples_split = 200, min_samples_leaf = 30, max_depth = 8, max_features = 7, subsample = 0.8, random_state = 18)
gbr.fit(X_train, Y_train) #fitting data
print(gbr.score(X_test, Y_test))
0.650891050157375
```

```
In [ ]: gbr = GradientBoostingRegressor(n_estimators = 800, learning_rate = 0.05, min_samples_split = 200, min_samples_leaf = 30, max_depth = 8, max_features = 7, subsample = 0.8, random_state = 18)
gbr.fit(X_train, Y_train) #fitting data
print(gbr.score(X_test, Y_test))
0.667684339984897
```

Final parameters for best model:

n_estimators = 800,

learning_rate = 0.05,

min_samples_split = 200,

min_samples_leaf = 30,

max_depth = 8,

max_features = 7,

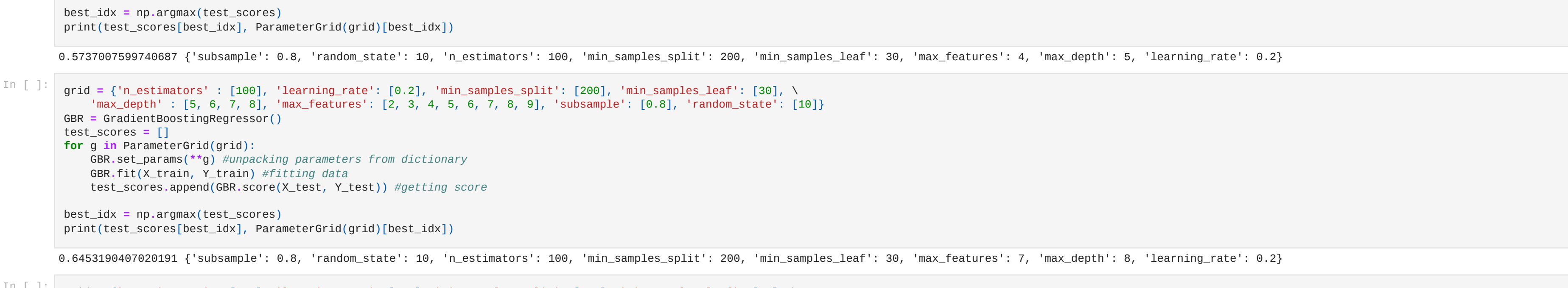
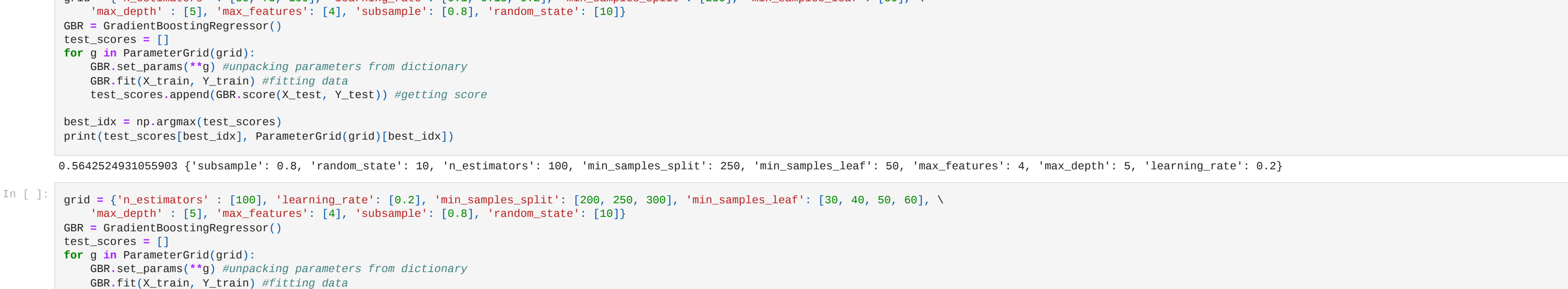
subsample = 0.8,

random_state = 18.

Making predictions

```
In [ ]: train_prediction = gbr.predict(X_train)
test_prediction = gbr.predict(X_test)
```

```
In [ ]: plt.rc('figure', figsize = (16, 8))
plt.scatter(train_prediction, Y_train)
```

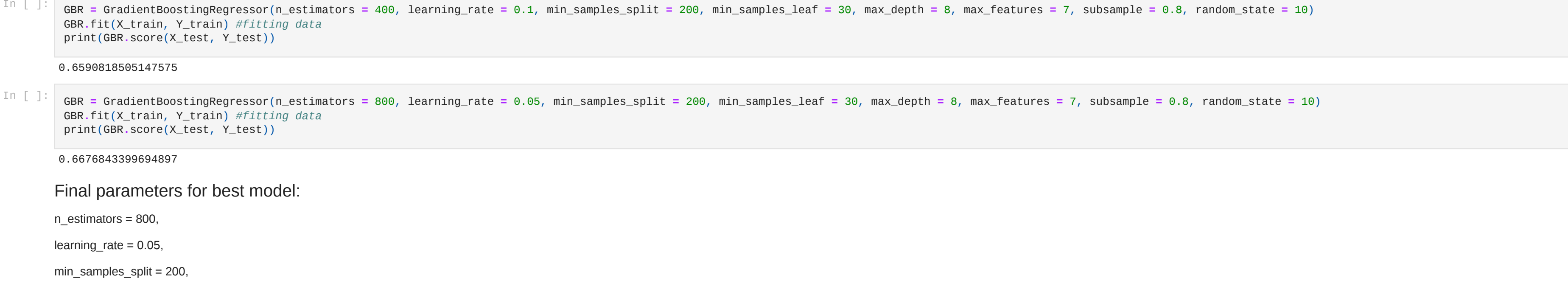


Finding most important features

```
In [ ]: feature_importanceGBR = gbr.feature_importances_
sorted_index = np.argsort(feature_importanceGBR[::-1]) #sorting the index by descending order
x1 = range(len(feature_importanceGBR))
feature_names = df.columns[df.columns.tolist()[1:]] #names are the same as names of columns
labels = np.array(feature_names[sorted_index])
```

```
In [ ]: plt.bar(x1, feature_importanceGBR[sorted_index], tick_label = labels)
```

```
Out [ ]: <BarContainer object of 9 artists>
```



Final forecast

```
In [ ]: df1 = pd.read_csv('test.csv', index_col = 'datetime', parse_dates = True)

#turning var2 into numbers
enc2 = OrdinalEncoder()
enc2.fit(df1[['var2']])
df1['var2'] = enc2.transform(df1[['var2']])

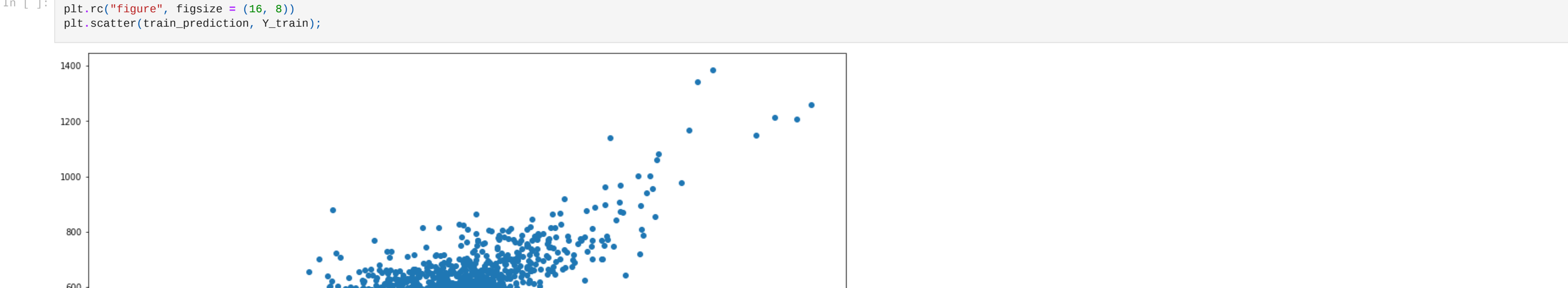
#datetime as variables
df1['hour'] = df1.index.hour
df1['month'] = df1.index.month
df1['day_of_month'] = df1.index.day
df1['day_of_week'] = df1.index.day_of_week

#X values for prediction
X_pred = df1.iloc[:, 1:11].values

#model
gbr = GradientBoostingRegressor(n_estimators = 800, learning_rate = 0.05, min_samples_split = 200, min_samples_leaf = 30, max_depth = 8, max_features = 7, subsample = 0.8, random_state = 18)
gbr.fit(X, Y)
prediction = gbr.predict(X_pred)
```

```
In [ ]: df1['electricity_consumption'] = prediction
df1['electricity_consumption'].plot()
```

```
Out [ ]: <AxesSubplot: xlabel='datetime'>
```

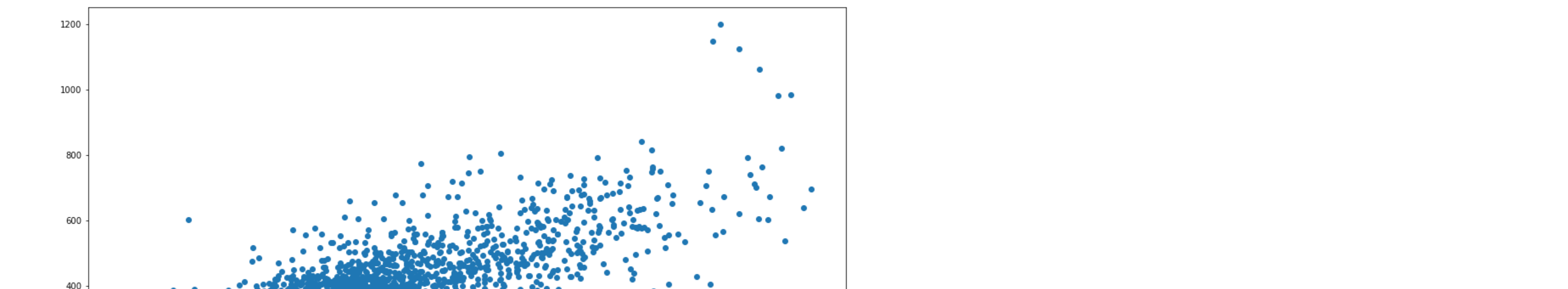


Saving prediction to file

```
In [ ]: df4 = pd.read_csv('sample_submission.csv')
df4['electricity_consumption'] = prediction
df4.set_index('ID', inplace = True)
df4.to_csv('sample_prediction_GBR.csv')
```

Score: lower is better

```
In [ ]: img = mpimg.imread('GBR.png')
imgplot = plt.imshow(img)
plt.show()
```



```
In [ ]:
```