

032-linear-regression-with-time-series-data

April 25, 2022

Linear Regression with Time Series Data

```
[1]: import matplotlib.pyplot as plt
import pandas as pd
import plotly.express as px
import pytz
from IPython.display import VimeoVideo
from pymongo import MongoClient
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_absolute_error
```

```
[2]: VimeoVideo("665412117", h="c39a50bd58", width=600)
```

```
[2]: <IPython.lib.display.VimeoVideo at 0x7f3a905397f0>
```

1 Prepare Data

1.1 Import

```
[3]: VimeoVideo("665412469", h="135f32c7da", width=600)
```

```
[3]: <IPython.lib.display.VimeoVideo at 0x7f3a90539cd0>
```

Task 3.2.1: Complete to the create a client to connect to the MongoDB server, assign the "air-quality" database to db, and assign the "nairobi" connection to nairobi.

- Create a client object for a MongoDB instance.
- Access a database using PyMongo.
- Access a collection in a database using PyMongo.

```
[4]: client = MongoClient(host="localhost", port=27017)
db = client["air-quality"]
nairobi = db["nairobi"]
```

```
[5]: VimeoVideo("665412480", h="c20ed3e570", width=600)
```

```
[5]: <IPython.lib.display.VimeoVideo at 0x7f39b4ad8910>
```

Task 3.2.2: Complete the wrangle function below so that the results from the database query are read into the DataFrame df. Be sure that the index of df is the "timestamp" from the results.

- Create a DataFrame from a dictionary using pandas.

```
[6]: def wrangle(collection):  
    #DB query  
    results = collection.find(  
        {"metadata.site": 29, "metadata.measurement": "P2"},  
        projection={"P2": 1, "timestamp": 1, "_id": 0},  
    )  
  
    df = pd.DataFrame(results).set_index("timestamp")  
  
    # Localize Timezone  
    df.index = df.index.tz_localize("UTC").tz_convert("Africa/Nairobi")  
  
    #Remove Outliers  
    df = df[df["P2"] < 500]  
  
    #Resample to 1H window , ffill missing values  
    df = df["P2"].resample("1H").mean().fillna(method="ffill").to_frame()  
  
    #Add Lag feature  
    df["P2.L1"] = df["P2"].shift(1)  
    #Drop NaN rows  
    df.dropna(inplace=True)  
  
    return df
```

```
[7]: VimeoVideo("665412496", h="d757475f7c", width=600)
```

```
[7]: <IPython.lib.display.VimeoVideo at 0x7f39b4ad8400>
```

Task 3.2.3: Use your wrangle function to read the data from the nairobi collection into the DataFrame df.

```
[8]: df = wrangle(nairobi)  
df.head(10)  
df.shape
```

```
[8]: (2927, 2)
```

```
[9]: # Check your work  
assert any([isinstance(df, pd.DataFrame), isinstance(df, pd.Series)])  
assert len(df) <= 32907  
assert isinstance(df.index, pd.DatetimeIndex)
```

```
[10]: VimeoVideo("665412520", h="e03eefff07", width=600)
```

```
[10]: <IPython.lib.display.VimeoVideo at 0x7f39b4ad89d0>
```

Task 3.2.4: Add to your `wrangle` function so that the `DatetimeIndex` for `df` is localized to the correct timezone, "Africa/Nairobi". Don't forget to re-run all the cells above after you change the function.

- Localize a timestamp to another timezone using pandas.

```
[11]: # Check your work
assert df.index.tzinfo == pytz.timezone("Africa/Nairobi")
```

1.2 Explore

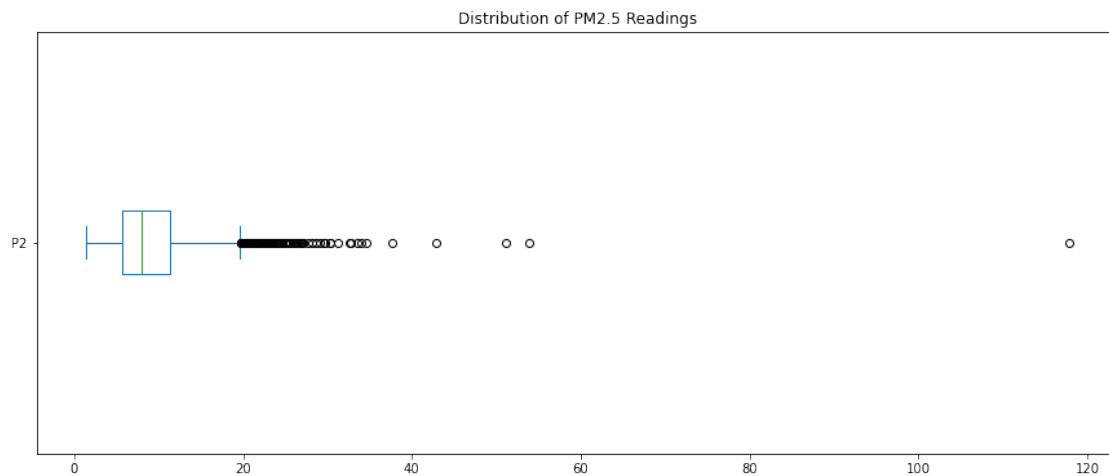
```
[12]: VimeoVideo("665412546", h="97792cb982", width=600)
```

```
[12]: <IPython.lib.display.VimeoVideo at 0x7f39ae35d2e0>
```

Task 3.2.5: Create a boxplot of the "P2" readings in `df`.

- Create a boxplot using pandas.

```
[13]: fig, ax = plt.subplots(figsize=(15, 6))
df["P2"].plot(kind="box", vert=False, title="Distribution of PM2.5 Readings",
→ax=ax);
```



```
[14]: VimeoVideo("665412573", h="b46049021b", width=600)
```

```
[14]: <IPython.lib.display.VimeoVideo at 0x7f39b415fcd0>
```

Task 3.2.6: Add to your `wrangle` function so that all "P2" readings above 500 are dropped from the dataset. Don't forget to re-run all the cells above after you change the function.

- Subset a DataFrame with a mask using pandas.

```
[15]: # Check your work
      assert len(df) <= 32906
```

```
[16]: VimeoVideo("665412594", h="e56c2f6839", width=600)
```

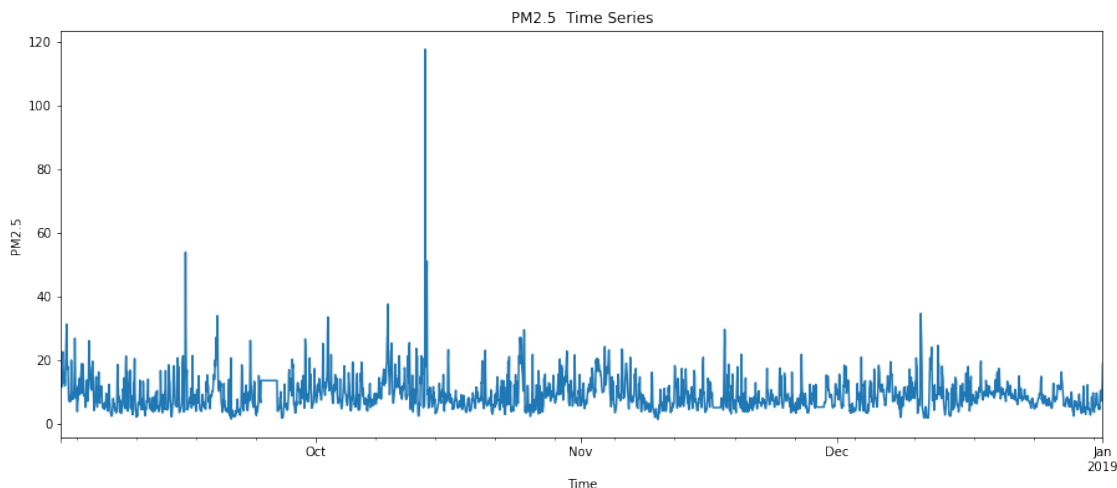
```
[16]: <IPython.lib.display.VimeoVideo at 0x7f39b41770d0>
```

Task 3.2.7: Create a time series plot of the "P2" readings in `df`.

- Create a line plot using pandas.

```
[17]: fig, ax = plt.subplots(figsize=(15, 6))
      df["P2"].plot(xlabel="Time", ylabel="PM2.5", title="PM2.5 Time Series", ax=ax)
```

```
[17]: <AxesSubplot:title={'center': 'PM2.5 Time Series'}, xlabel='Time',
      ylabel='PM2.5'>
```



```
[18]: VimeoVideo("665412601", h="a16c5a73fc", width=600)
```

```
[18]: <IPython.lib.display.VimeoVideo at 0x7f39b4057bb0>
```

Task 3.2.8: Add to your `wrangle` function to resample `df` to provide the mean "P2" reading for each hour. Use a forward fill to impute any missing values. Don't forget to re-run all the cells above after you change the function.

- Resample time series data in pandas.
- Impute missing time series values using pandas.

```
[19]: # Check your work
      assert len(df) <= 2928
```

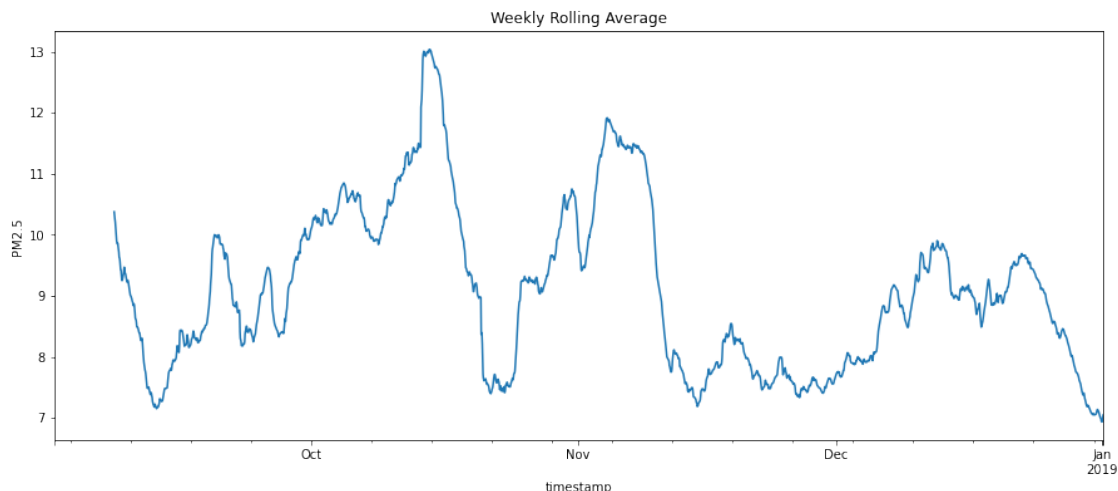
```
[20]: VimeoVideo("665412649", h="d2e99d2e75", width=600)
```

```
[20]: <IPython.lib.display.VimeoVideo at 0x7f39b4057e80>
```

Task 3.2.9: Plot the rolling average of the "P2" readings in `df`. Use a window size of 168 (the number of hours in a week).

- What's a rolling average?
- Calculate a rolling average in pandas.
- Create a line plot using pandas.

```
[21]: fig, ax = plt.subplots(figsize=(15, 6))
      df["P2"].rolling(168).mean().plot(ax=ax, ylabel="PM2.5", title="Weekly Rolling_
      ↪Average");
```



```
[22]: VimeoVideo("665412693", h="c3bca16aff", width=600)
```

```
[22]: <IPython.lib.display.VimeoVideo at 0x7f39aeded400>
```

Task 3.2.10: Add to your `wrangle` function to create a column called "P2.L1" that contains the mean "P2" reading from the previous hour. Since this new feature will create NaN values in your DataFrame, be sure to also drop null rows from `df`.

- Shift the index of a Series in pandas.
- Drop rows with missing values from a DataFrame using pandas.

```
[23]: # Check your work
      assert len(df) <= 11686
      assert df.shape[1] == 2
```

```
[24]: VimeoVideo("665412732", h="059e4088c5", width=600)
```

```
[24]: <IPython.lib.display.VimeoVideo at 0x7f39b405f160>
```

Task 3.2.11: Create a correlation matrix for df.

- [Create a correlation matrix in pandas.](#)

```
[25]: df.corr()
```

```
[25]:
```

	P2	P2.L1
P2	1.000000	0.650679
P2.L1	0.650679	1.000000

```
[26]: VimeoVideo("665412741", h="7439cb107c", width=600)
```

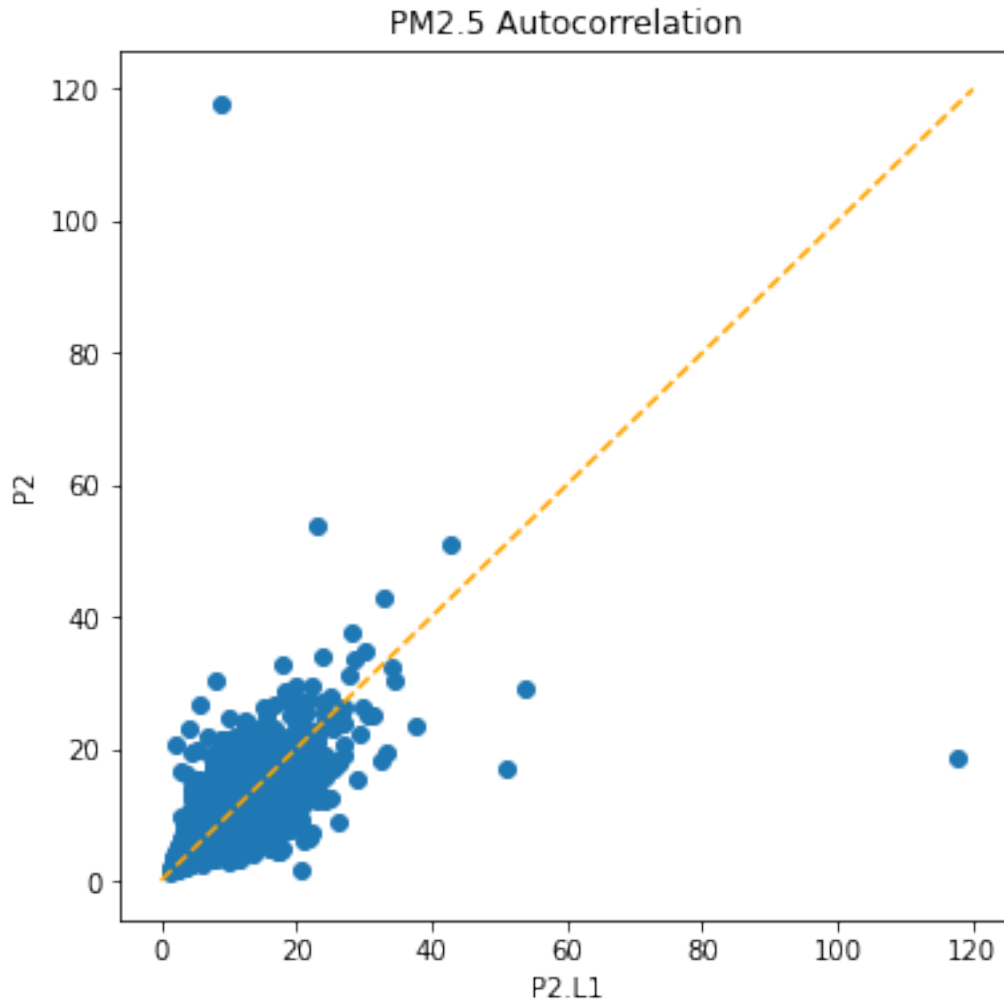
```
[26]: <IPython.lib.display.VimeoVideo at 0x7f39aee0c4f0>
```

Task 3.2.12: Create a scatter plot that shows PM 2.5 mean reading for each our as a function of the mean reading from the previous hour. In other words, "P2.L1" should be on the x-axis, and "P2" should be on the y-axis. Don't forget to label your axes!

- [Create a scatter plot using Matplotlib.](#)

```
[27]: fig, ax = plt.subplots(figsize=(6, 6))
ax.scatter(x=df["P2.L1"], y=df["P2"])
ax.plot([0,120],[0,120], linestyle="--", color="orange")
plt.xlabel("P2.L1")
plt.ylabel("P2")
plt.title("PM2.5 Autocorrelation")
```

```
[27]: Text(0.5, 1.0, 'PM2.5 Autocorrelation')
```



1.3 Split

```
[28]: VimeoVideo("665412762", h="a5eba496f7", width=600)
```

```
[28]: <IPython.lib.display.VimeoVideo at 0x7f39aed80370>
```

Task 3.2.13: Split the DataFrame `df` into the feature matrix `X` and the target vector `y`. Your target is "P2".

- Subset a DataFrame by selecting one or more columns in pandas.
- Select a Series from a DataFrame in pandas.

```
[29]: target = "P2"  
y = df[target]  
  
# features=["P2.L1"]
```

```
# X_train= df[features]
# X_train.head()

X = df.drop(columns=target)
```

```
[30]: VimeoVideo("665412785", h="03118eda71", width=600)
```

```
[30]: <IPython.lib.display.VimeoVideo at 0x7f39aed737f0>
```

Task 3.2.14: Split X and y into training and test sets. The first 80% of the data should be in your training set. The remaining 20% should be in the test set.

- [Divide data into training and test sets in pandas.](#)

```
[31]: cutoff = int(len(X) * 0.8)

X_train, y_train = X.iloc[:cutoff], y.iloc[:cutoff]
X_test, y_test = X.iloc[cutoff:], y.iloc[cutoff:]

#len(X_train) + len(X_test) == len(X)
```

2 Build Model

2.1 Baseline

Task 3.2.15: Calculate the baseline mean absolute error for your model.

- [Calculate summary statistics for a DataFrame or Series in pandas.](#)

```
[32]: y_pred_baseline = [(y_train.mean())] * len(y_train)
mae_baseline = mean_absolute_error(y_train, y_pred_baseline)

print("Mean P2 Reading:", round(y_train.mean(), 2))
print("Baseline MAE:", round(mae_baseline, 2))
```

Mean P2 Reading: 9.27

Baseline MAE: 3.89

2.2 Iterate

Task 3.2.16: Instantiate a [LinearRegression](#) model named `model`, and fit it to your training data.

- [Instantiate a predictor in scikit-learn.](#)
- [Fit a model to training data in scikit-learn.](#)

```
[33]: #Instantiate
model = LinearRegression()
#Fit
```



```
model.fit(X_train, y_train)
```

[33]: LinearRegression()

2.3 Evaluate

[34]: VimeoVideo("665412844", h="129865775d", width=600)

[34]: <IPython.lib.display.VimeoVideo at 0x7f39aed249d0>

Task 3.2.17: Calculate the training and test mean absolute error for your model.

- [Generate predictions using a trained model in scikit-learn.](#)
- [Calculate the mean absolute error for a list of predictions in scikit-learn.](#)

```
[35]: training_mae = mean_absolute_error(y_train, model.predict(X_train))
test_mae = mean_absolute_error(y_test, model.predict(X_test))
print("Training MAE:", round(training_mae, 2))
print("Test MAE:", round(test_mae, 2))
```

Training MAE: 2.46

Test MAE: 1.8

3 Communicate Results

Task 3.2.18: Extract the intercept and coefficient from your model.

- [Access an object in a pipeline in scikit-learn](#)

```
[36]: intercept = round(model.intercept_, 2)
coefficient = round(model.coef_[0], 2)

print(f"P2 = {intercept} + ({coefficient} * P2.L1)")
```

P2 = 3.36 + (0.64 * P2.L1)

[37]: VimeoVideo("665412870", h="318d69683e", width=600)

[37]: <IPython.lib.display.VimeoVideo at 0x7f39aed33190>

Task 3.2.19: Create a DataFrame `df_pred_test` that has two columns: "y_test" and "y_pred". The first should contain the true values for your test set, and the second should contain your model's predictions. Be sure the index of `df_pred_test` matches the index of `y_test`.

- [Create a DataFrame from a dictionary using pandas.](#)

```
[38]: df_pred_test = pd.DataFrame(
    {
        "y_test": y_test,
```

```

        "y_pred": model.predict(X_test)
    }
)
df_pred_test.head()

```

```

[38]:
      timestamp      y_test      y_pred
2018-12-07 17:00:00+03:00    7.070000    8.478927
2018-12-07 18:00:00+03:00    8.968333    7.865485
2018-12-07 19:00:00+03:00   11.630833    9.076421
2018-12-07 20:00:00+03:00   11.525833   10.774814
2018-12-07 21:00:00+03:00    9.533333   10.707836

```

```

[39]: VimeoVideo("665412891", h="39d7356a26", width=600)

```

```

[39]: <IPython.lib.display.VimeoVideo at 0x7f39aed333d0>

```

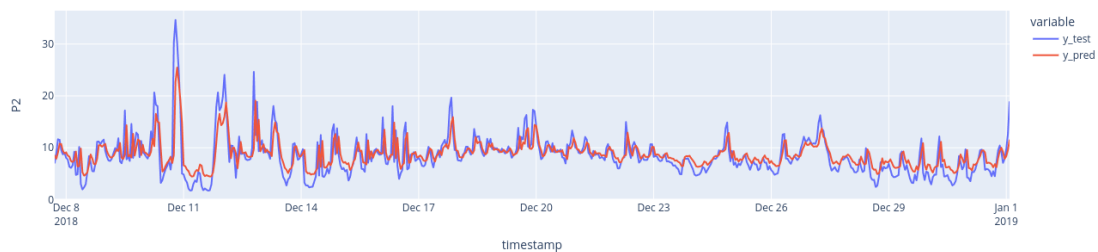
Task 3.2.20: Create a time series line plot for the values in `test_predictions` using `plotly` express. Be sure that the y-axis is properly labeled as "P2".

- Create a line plot using `plotly` express.

```

[40]: fig = px.line(df_pred_test, labels={"value": "P2"})
      fig.show()

```



Copyright © 2022 WorldQuant University. This content is licensed solely for personal use. Redistribution or publication of this material is strictly prohibited.