

Applied Data Science Lab | My Path Assessment Overview | work/ds-curriculum - JupyterLab | Data-Science-Lab/MongoDB.py | +

vm.wqu.edu/lab/tree/work/ds-curriculum/030-air-quality-in-nairobi/035-assignment.ipynb

File Edit View Run Kernel Tabs Settings Help

Filter files by name

Name	Last Modified
.ipynb, chec...	2 minutes ago
images	a month ago
031-data-wi...	7 minutes ago
032-linear...	a month ago
033-autore...	a month ago
034-arma...	a month ago
035-assign...	3 minutes ago

035-assignment.ipynb

Python 3 (pykernel)

- No re-sharing embedded videos with friends or colleagues.
- No adding this notebook to public or private repositories.
- No uploading this notebook (or screenshots of it) to other websites, including websites for study resources.

3.5. Air Quality in Dar es Salaam TZ

```
[49]: import warnings
import wqet_grader

warnings.simplefilter(action="ignore", category=FutureWarning)
wqet_grader.init("Project 3 Assessment")
```

```
[50]: # Import libraries here
import inspect
import time
from pprint import PrettyPrinter
import matplotlib.pyplot as plt
import pandas as pd
import plotly.express as px
import seaborn as sns
from pymongo import MongoClient
import pytz
from statsmodels.tsa.ar_model import AutoReg
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_absolute_error
from statsmodels.graphics.tsaplots import plot_acf, plot_pacf
from statsmodels.tsa.arima.model import ARIMA
```

1. Prepare Data

Simple | Python 3 (pykernel) | Idle | Mode: Command | Ln 1, Col 1 | English (United States) | 035-assignment.ipynb

35°C Smoke

Applied Data Science Lab | My Path Assessment Overview | work/ds-curriculum - JupyterLab | Data-Science-Lab/MongoDB.py | +

vm.wqu.edu/lab/tree/work/ds-curriculum/030-air-quality-in-nairobi/035-assignment.ipynb

File Edit View Run Kernel Tabs Settings Help

Filter files by name

Name	Last Modified
.ipynb, chec...	2 minutes ago
images	a month ago
031-data-wi...	7 minutes ago
032-linear...	a month ago
033-autore...	a month ago
034-arma...	a month ago
035-assign...	3 minutes ago

035-assignment.ipynb

Python 3 (pykernel)

```
from statsmodels.tsa.arima.model import ARIMA
```

1. Prepare Data

1.1. Connect

Task 3.5.1: Connect to MongoDB server running at host "localhost" on port "27017". Then connect to the "air-quality" database and assign the collection for Dar es Salaam to the variable name `dar`.

```
[52]: from pymongo import MongoClient

# Connect to server
client = MongoClient(host="localhost", port=27017)
# Connect to database
db = client["air-quality"]

# Get collection
dar = db["dar-es-salaam"]
```

```
[53]: wqet_grader.grade("Project 3 Assessment", "Task 3.5.1", [dar.name])
```

✓ Party time! 🎉 🎉
Score: 1

1.2. Explore

Task 3.5.2: Determine the numbers assigned to all the sensor sites in the Dar es Salaam collection. Your submission should be a list of integers.

Simple | Python 3 (pykernel) | Idle | Mode: Command | Ln 1, Col 1 | English (United States) | 035-assignment.ipynb

35°C Smoke

Applied Data Science Lab | My Path Assessment Overview | work/ds-curricu (2) - Jupyter | Data-Science-Lab/distinct | Introducing ChatGPT | Sensor sites in Dar Es Salaam | +

vm.wqu.edu/lab/tree/work/ds-curriculum/030-air-quality-in-nairobi/035-assignment.ipynb

File Edit View Run Kernel Tabs Settings Help

Filter files by name

030-air-quality-in-nairobi /

Name	Last Modified
030-air-quality-in-nairobi	6 minutes ago
031-data-wrangling-with-mc	a month ago
032-linear...	12 minutes ago
033-autore...	a month ago
034-arma...	a month ago
035-assign...	seconds ago

035-assignment.ipynb

Score: 1

Markdown

Python 3 (pykernel)

1.2. Explore

Task 3.5.2: Determine the numbers assigned to all the sensor sites in the Dar es Salaam collection. Your submission should be a list of integers.

[54]: sites = dar.distinct("metadata.site") # dar ---> variable holding collection sites

[54]: [11, 23]

[55]: wqet_grader.grade("Project 3 Assessment", "Task 3.5.2", sites)

Correct.

Score: 1

Task 3.5.3: Determine which site in the Dar es Salaam collection has the most sensor readings (of any type, not just PM2.5 readings). Your submission readings_per_site should be a list of dictionaries that follows this format:

[['_id': 6, 'count': 70360], ['_id': 29, 'count': 131852]]

Note that the values here are from the Nairobi collection, so your values will look different.

[59]: result = dar.aggregate([{'\$group': {'_id': '\$metadata.site', 'count': {'\$count': {}}}]

Simple | Python 3 (pykernel) | Idle | Mode: Command | Ln 1, Col 1 | English (United States) | 035-assignment.ipynb

Data-Science-Lab-...zip | Show all

35°C Smoke | Search | ENG | 13:28 10-03-2023

Applied Data Science Lab | My Path Assessment Overview | work/ds-curricu - JupyterLab | Data-Science-Lab/MongoDB.py | +

vm.wqu.edu/lab/tree/work/ds-curriculum/030-air-quality-in-nairobi/035-assignment.ipynb

File Edit View Run Kernel Tabs Settings Help

Filter files by name

030-air-quality-in-nairobi /

Name	Last Modified
030-air-quality-in-nairobi	2 minutes ago
031-data-wrangling-with-mc	a month ago
032-linear...	7 minutes ago
033-autore...	a month ago
034-arma...	a month ago
035-assign...	3 minutes ago

035-assignment.ipynb

Markdown

Python 3 (pykernel)

1. Prepare Data

1.1. Connect

Task 3.5.1: Connect to MongoDB server running at host "localhost" on port 27017. Then connect to the "air-quality" database and assign the collection for Dar es Salaam to the variable name dar.

[52]: from pymongo import MongoClient

Connect to server

client = MongoClient(host="localhost", port=27017)

Connect to database

db = client["air-quality"]

Get collection

dar = db["dar-es-salaam"]

[53]: wqet_grader.grade("Project 3 Assessment", "Task 3.5.1", [dar.name])

Party time!

Score: 1

1.2. Explore

Task 3.5.2: Determine the numbers assigned to all the sensor sites in the Dar es Salaam collection. Your submission should be a list of integers.

Simple | Python 3 (pykernel) | Idle | Mode: Command | Ln 1, Col 1 | English (United States) | 035-assignment.ipynb

Data-Science-Lab-...zip | Show all

35°C Smoke | Search | ENG | 13:23 10-03-2023

Applied Data Science Lab x My Path Assessment Overview x work/ds-curricu (2) - Jupyter x Data-Science-Lab/distinct x Introducing ChatGPT x Sensor sites in Dar Es Salaam x +

vm.wqu.edu/lab/tree/work/ds-curriculum/030-air-quality-in-nairobi/035-assignment.ipynb

File Edit View Run Kernel Tabs Settings Help

Launcher x 035-assignment.ipynb x 031-data-wrangling-with-mc x

Filter files by name

/ / ds-curriculum / 030-air-quality-in-nairobi /

Name	Last Modified
031-data-w...	12 minutes ago
032-linear...	a month ago
033-aure...	a month ago
034-arma...	a month ago
035-assign...	seconds ago

1.2. Explore

Task 3.5.2: Determine the numbers assigned to all the sensor sites in the Dar es Salaam collection. Your submission should be a list of integers.

```
[54]: sites = dar.distinct("metadata.site") # dar --> variable holding collection sites
```

```
[54]: [11, 23]
```

```
[55]: wqet_grader.grade("Project 3 Assessment", "Task 3.5.2", sites)
```

Correct.
Score: 1

Task 3.5.3: Determine which site in the Dar es Salaam collection has the most sensor readings (of any type, not just PM2.5 readings). Your submission `readings_per_site` should be a list of dictionaries that follows this format:

```
[{"_id": 6, "count": 70360}, {"_id": 29, "count": 131852}]
```

Note that the values here are from the Nairobi collection, so your values will look different.

```
[59]: result = dar.aggregate([{"$group": {"_id": "$metadata.site", "count": {"$count": {}}}])
```

Simple 0 9 Python 3 (ipykernel) Idle Mode Command Ln 1, Col 1 English (United States) 035-assignment.ipynb

35°C Smoke

Applied Data Science Lab x My Path Assessment Overview x work/ds-curricu (2) - Jupyter x Data-Science-Lab/wrangle x Introducing ChatGPT x Sensor sites in Dar Es Salaam x +

vm.wqu.edu/lab/tree/work/ds-curriculum/030-air-quality-in-nairobi/035-assignment.ipynb

File Edit View Run Kernel Tabs Settings Help

Launcher x 035-assignment.ipynb x 031-data-wrangling-with-mc x

Filter files by name

/ / ds-curriculum / 030-air-quality-in-nairobi /

Name	Last Modified
031-data-w...	14 minutes ago
032-linear...	a month ago
033-aure...	a month ago
034-arma...	a month ago
035-assign...	in a few seconds

```
# Localize time
df.index = df.index.tz_localize("UTC").tz_convert("Africa/Dar_es_Salaam")

# Remove outliers
df = df[df["P2"] < 100]

# Resample to 1hour period, fill in missing values
y = df["P2"].resample("1H").mean().fillna(method="ffill")

return y
```

Use your `wrangle` function to query the `dar` collection and return your cleaned results.

```
[62]: y = wrangle(dar)
y.head()
```

```
[62]: timestamp
2018-01-01 03:00:00+03:00    9.456327
2018-01-01 04:00:00+03:00    9.400033
2018-01-01 05:00:00+03:00    9.331658
2018-01-01 06:00:00+03:00    9.528776
2018-01-01 07:00:00+03:00    8.861250
Freq: H, Name: P2, dtype: float64
```

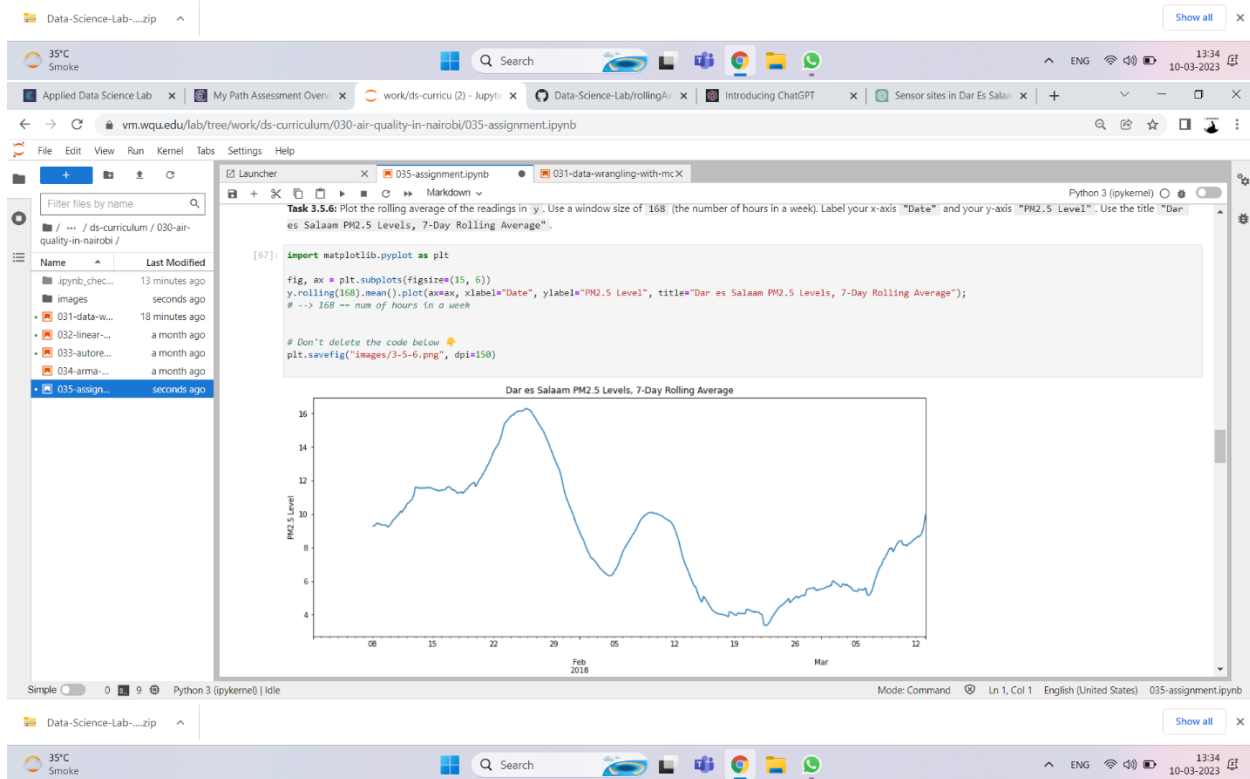
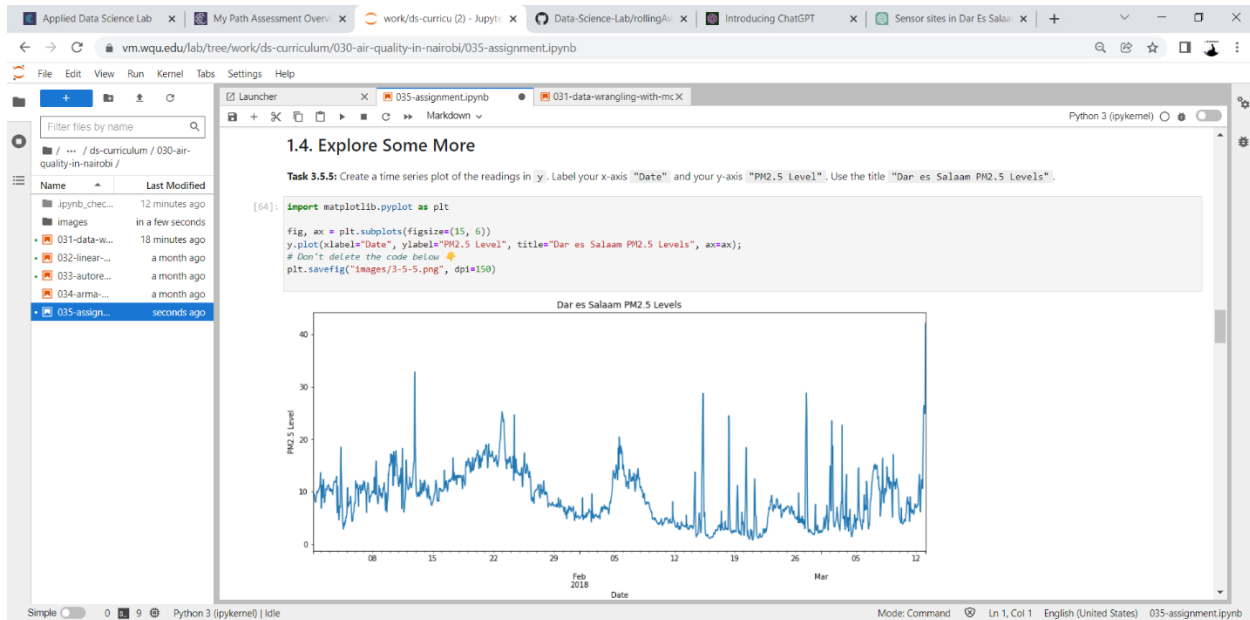
```
[63]: wqet_grader.grade("Project 3 Assessment", "Task 3.5.4", wrangle(dar))
```

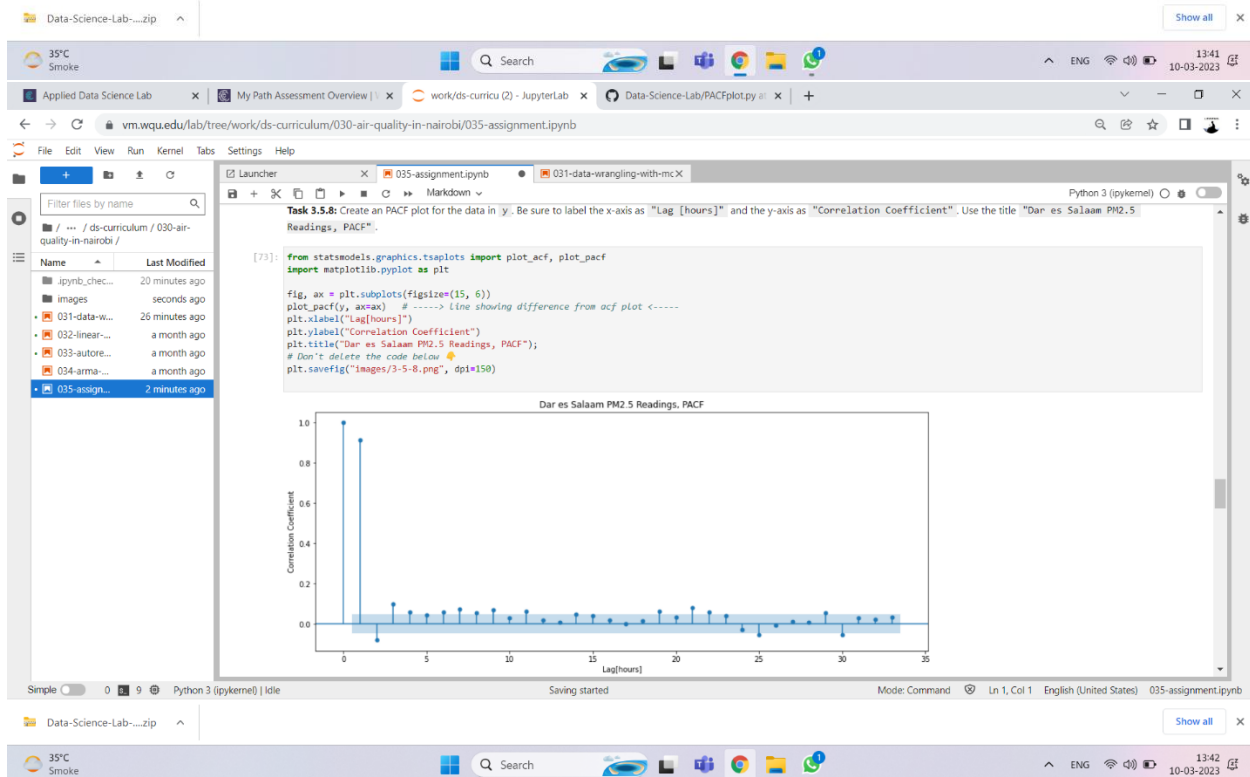
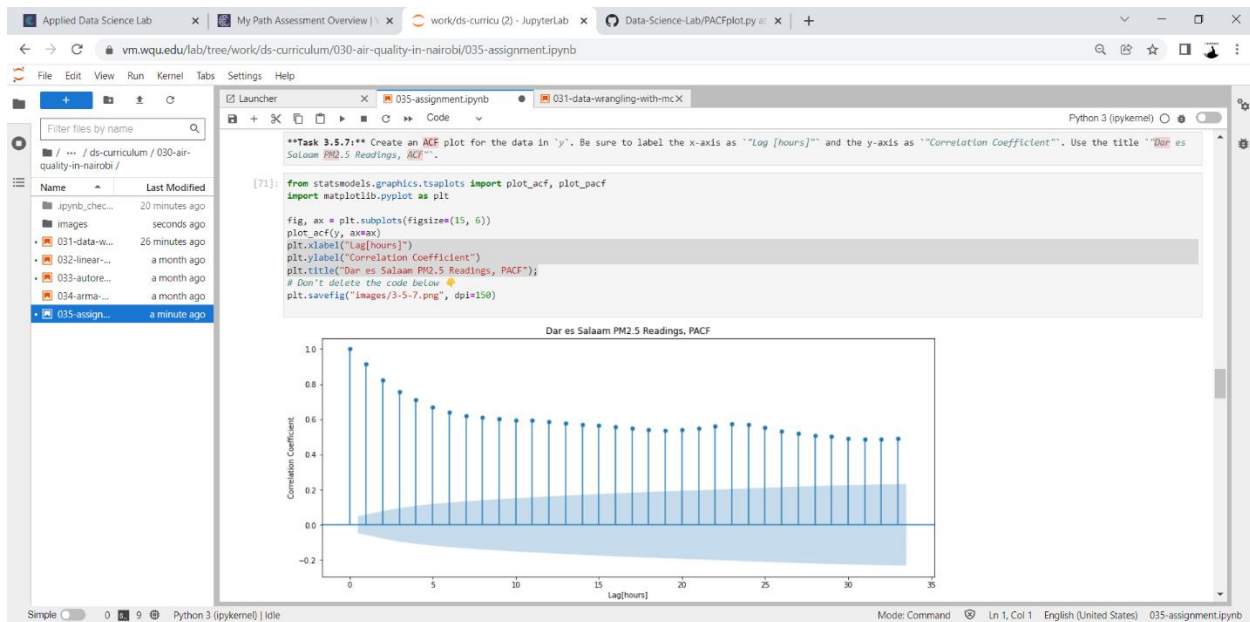
Party time!
Score: 1

1.4. Explore Some More

Task 3.5.5: Create a time series plot of the readings in `y`. Label your x-axis "Date" and your y-axis "PM2.5 Level". Use the title "Dar es Salaam PM2.5 Levels".

Simple 0 9 Python 3 (ipykernel) Idle Mode Command Ln 1, Col 1 English (United States) 035-assignment.ipynb





Applied Data Science Lab | My Path Assessment Overview | work/ds-curricu (2) - JupyterLab | Data-Science-Lab/ARmodel.py | +

vm.wqu.edu/lab/tree/work/ds-curriculum/030-air-quality-in-nairobi/035-assignment.ipynb

File Edit View Run Kernel Tabs Settings Help

Launcher 035-assignment.ipynb 031-data-wrangling-with-mcX Python 3 (pykernel)

Filter files by name

/ / ds-curriculum / 030-air-quality-in-nairobi /

Name	Last Modified
031-data-w...	30 minutes ago
032-linear...	a month ago
033-autore...	a month ago
034-arma...	a month ago
035-assign...	seconds ago

Task 3.5.11: You're going to use an AR model to predict PM2.5 readings, but which hyperparameter settings will give you the best performance? Use a `for` loop to train your AR model on using settings for `p` from 1 to 30. Each time you train a new model, calculate its mean absolute error and append the result to the list `maes`. Then store your results in the Series `mae_series`.

```
[85]: from statsmodels.tsa.ar_model import AutoReg
      from sklearn.metrics import mean_absolute_error

      # Use AR model to predict PM2.5 readings
      # Hyperparameter --> p
      p_params = range(1, 31)
      maes = []
      for p in p_params:
          #Train model
          model = AutoReg(y_train, lags=p).fit()

          #Generate in-sample pred
          y_pred = model.predict().dropna()

          #Calculate mae
          mae = mean_absolute_error(y_train.iloc[p:], y_pred)
          maes.append(mae)

      mae_series = pd.Series(maes, name="mae", index=p_params)
      mae_series.head()

[85]: 1    0.947888
      2    0.933894
      3    0.920850
      4    0.920153
      5    0.919519
      Name: mae, dtype: float64

[86]: wqet_grader.grade("Project 3 Assessment", "Task 3.5.11", mae_series)
```

Simple 0 9 Python 3 (pykernel) Idle Mode Command Ln 1, Col 1 English (United States) 035-assignment.ipynb

35°C Smoke

Applied Data Science Lab | My Path Assessment Overview | work/ds-curricu (2) - JupyterLab | Data-Science-Lab/finalModel.py | +

vm.wqu.edu/lab/tree/work/ds-curriculum/030-air-quality-in-nairobi/035-assignment.ipynb

File Edit View Run Kernel Tabs Settings Help

Launcher 035-assignment.ipynb 031-data-wrangling-with-mcX Python 3 (pykernel)

Filter files by name

/ / ds-curriculum / 030-air-quality-in-nairobi /

Name	Last Modified
031-data-w...	31 minutes ago
032-linear...	a month ago
033-autore...	a month ago
034-arma...	a month ago
035-assign...	a minute ago

Task 3.5.12: Look through the results in `mae_series` and determine what value for `p` provides the best performance. Then build and train `final_model` using the best hyperparameter value.

Note: Make sure that you build and train your model in one line of code, and that the data type of `best_model` is `statsmodels.tsa.ar_model.AutoRegResultsWrapper`.

```
[87]: from statsmodels.tsa.ar_model import AutoReg
      from statsmodels.tsa.arima.model import ARIMA

      mae_series # locate best_p
      best_p = 28

      # build and train model
      best_model = AutoReg(y_train, lags=best_p).fit()

      # calculate training residuals for best_model
      y_train_resid = best_model.resid
      y_train_resid.name = "residuals"
      y_train_resid.head()

[87]: timestamp
2018-01-02 07:00:00+03:00    1.732488
2018-01-02 08:00:00+03:00   -0.381568
2018-01-02 09:00:00+03:00   -0.560971
2018-01-02 10:00:00+03:00   -2.215760
2018-01-02 11:00:00+03:00    0.006468
Freq: H, Name: residuals, dtype: float64

[88]: wqet_grader.grade(
      "Project 3 Assessment", "Task 3.5.12", [isinstance(best_model.model, AutoReg)]
      )
```

✓ Score: 1

Simple 0 9 Python 3 (pykernel) Idle Mode Command Ln 1, Col 1 English (United States) 035-assignment.ipynb

Applied Data Science Lab | My Path Assessment Overview | work/ds-curricu (2) - JupyterLab | Data-Science-Lab/finalModel.py | +

vm.wqu.edu/lab/tree/work/ds-curriculum/030-air-quality-in-nairobi/035-assignment.ipynb

File Edit View Run Kernel Tabs Settings Help

Launcher 035-assignment.ipynb 031-data-wrangling-with-mc X Python 3 (ipykernel)

Filter files by name

Name	Last Modified
.ipynb_chec...	27 minutes ago
images	7 minutes ago
031-data-w...	33 minutes ago
032-linear...	a month ago
033-autore...	a month ago
034-arma...	a month ago
035-assign...	seconds ago

Task 3.5.13: Calculate the training residuals for `best_model`, and assign the result to `y_train_resid`. Note that your name of your Series should be "residuals".

```
[91]: y_train_resid = best_model.resid
      y_train_resid.name = "residuals"
      y_train_resid.head()
```

```
[91]: timestamp
2018-01-02 07:00:00+03:00    1.732488
2018-01-02 08:00:00+03:00   -0.381568
2018-01-02 09:00:00+03:00   -0.560971
2018-01-02 10:00:00+03:00   -2.215760
2018-01-02 11:00:00+03:00    0.086468
Freq: H, Name: residuals, dtype: float64
```

That's the right answer. Keep it up!
Score: 1

Task 3.5.14: Create a histogram of `y_train_resid`. Be sure to label the x-axis as "Residuals" and the y-axis as "Frequency". Use the title "Best Model, Training Residuals".

```
[ ]: # Plot histogram of residuals
     # Don't delete the code below
     plt.savefig("images/3-5-14.png", dpi=150)
```

35°C Smoke

Applied Data Science Lab | My Path Assessment Overview | work/ds-curricu (2) - JupyterLab | Data-Science-Lab/030-air-quali... | +

vm.wqu.edu/lab/tree/work/ds-curriculum/030-air-quality-in-nairobi/035-assignment.ipynb

File Edit View Run Kernel Tabs Settings Help

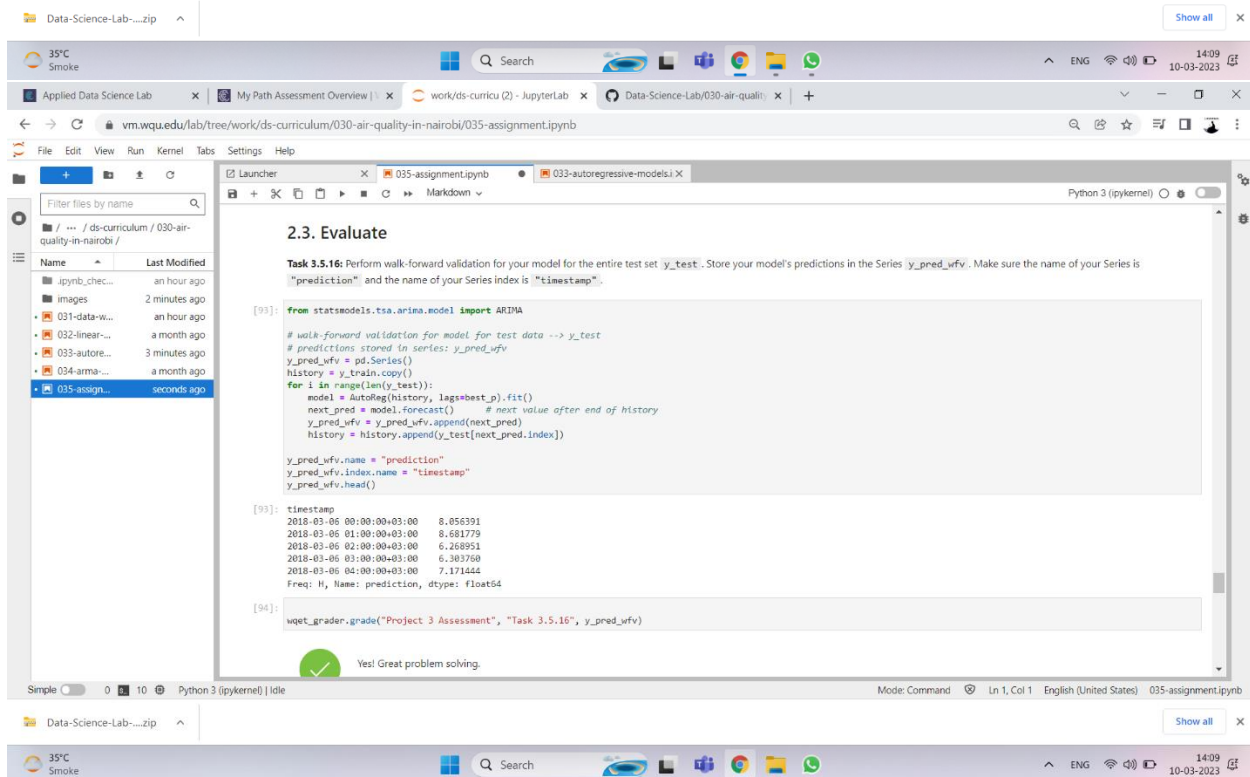
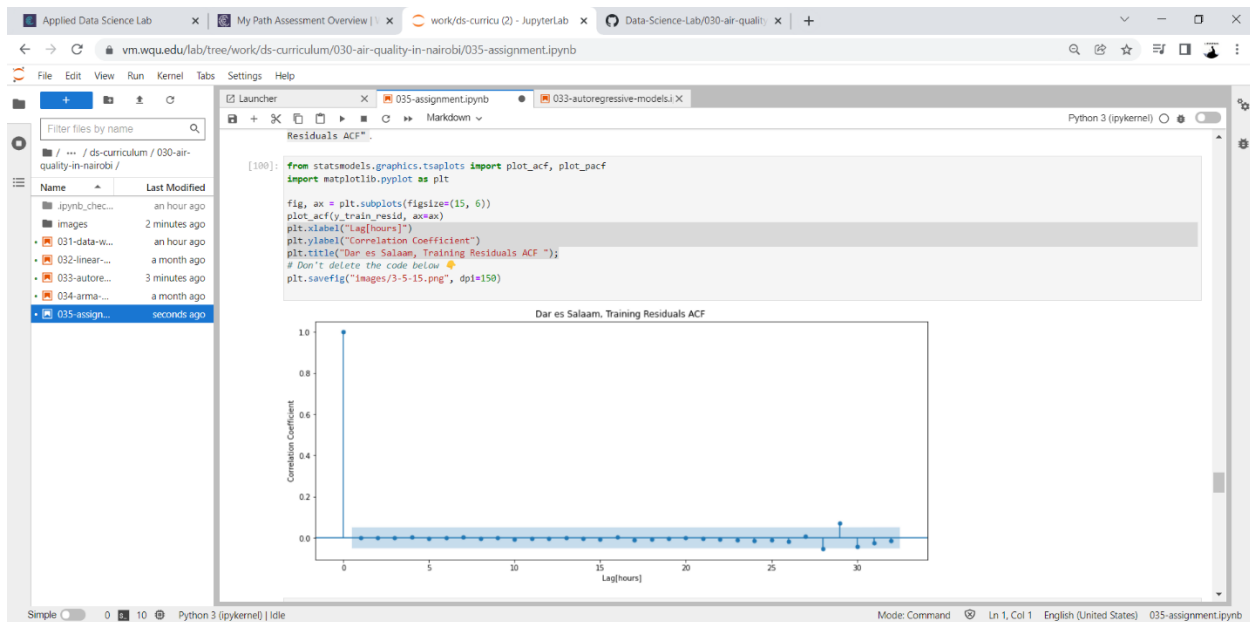
Launcher 035-assignment.ipynb 033-autoregressive-models.i X Python 3 (ipykernel)

Filter files by name

Name	Last Modified
.ipynb_chec...	an hour ago
images	a minute ago
031-data-w...	an hour ago
032-linear...	a month ago
033-autore...	3 minutes ago
034-arma...	a month ago
035-assign...	seconds ago

Task 3.5.14: Create a histogram of `y_train_resid`. Be sure to label the x-axis as "Residuals" and the y-axis as "Frequency". Use the title "Best Model, Training Residuals".

```
[106]: fig, ax = plt.subplots(figsize=(15, 6))
        # Plot histogram of residuals
        y_train_resid.hist()
        plt.xlabel("Residuals")
        plt.ylabel("Frequency")
        plt.title("Best Model, Training Residuals");
        # Don't delete the code below
        plt.savefig("images/3-5-14.png", dpi=150)
```

Applied Data Science Lab | My Path Assessment Overview | work/ds-curricu (2) - JupyterLab | Data-Science-Lab/030-air-quality- | +

vm.wqu.edu/lab/tree/work/ds-curriculum/030-air-quality-in-nairobi/035-assignment.ipynb

File Edit View Run Kernel Tabs Settings Help

Launcher 035-assignment.ipynb 033-autoregressive-models.ipynb

Filter files by name

Name	Last Modified
031-data-w...	an hour ago
032-linear...	a month ago
033-autore...	3 minutes ago
034-arma...	a month ago
035-assign...	seconds ago

3. Communicate Results

Task 3.5.18: Put the values for `y_test` and `y_pred_wfv` into the DataFrame `df_pred_test` (don't forget the index). Then plot `df_pred_test` using `plotly` express. Be sure to label the x-axis as "Date" and the y-axis as "PM2.5 Level". Use the title "Dar es Salaam, WfV Predictions".

```
[109]: import plotly.express as px
import pandas as pd

# Put test and walk-forward validation values
# in a dataframe and plot df
df_pred_test = pd.DataFrame(
    {"y_test": y_test, "y_pred_wfv": y_pred_wfv}
)

fig = px.line(df_pred_test, labels={"value": "PM2.5"})
fig.update_layout(
    title="Dar es Salaam, WfV Predictions",
    xaxis_title="Date",
    yaxis_title="PM2.5 Level",
)

# Don't delete the code below
fig.write_image("images/3-5-18.png", scale=1, height=500, width=700)
fig.show()
```

Dar es Salaam, WfV Predictions

variable
— y_test
— y_pred_wfv

Simple 0 10 Python 3 (ipykernel) | Idle Mode: Command Ln 1, Col 1 English (United States) 035-assignment.ipynb

35°C Smoke

Applied Data Science Lab | My Path Assessment Overview | work/ds-curricu (2) - JupyterLab | Data-Science-Lab/030-air-quality- | +

vm.wqu.edu/lab/tree/work/ds-curriculum/030-air-quality-in-nairobi/035-assignment.ipynb

File Edit View Run Kernel Tabs Settings Help

Launcher 035-assignment.ipynb 033-autoregressive-models.ipynb

Filter files by name

Name	Last Modified
031-data-w...	an hour ago
032-linear...	a month ago
033-autore...	3 minutes ago
034-arma...	a month ago
035-assign...	a minute ago

```
[109]: {"y_test": y_test, "y_pred_wfv": y_pred_wfv}
fig = px.line(df_pred_test, labels={"value": "PM2.5"})
fig.update_layout(
    title="Dar es Salaam, WfV Predictions",
    xaxis_title="Date",
    yaxis_title="PM2.5 Level",
)

# Don't delete the code below
fig.write_image("images/3-5-18.png", scale=1, height=500, width=700)
fig.show()
```

Dar es Salaam, WfV Predictions

variable
— y_test
— y_pred_wfv

```
[110]: with open("images/3-5-18.png", "rb") as file:
wget_grader.grade("Project 3 Assessment", "Task 3.5.18", file)
```

Simple 0 10 Python 3 (ipykernel) | Idle Mode: Command Ln 1, Col 1 English (United States) 035-assignment.ipynb

35°C Smoke