

1 - Business understanding

Problem Statement

- ▶ We define footfall as the measurement of the number of people entering a shop or shopping mall. We have to predict
 - ▶ peak time in each day of week
 - ▶ footfall in peak time of each day of week

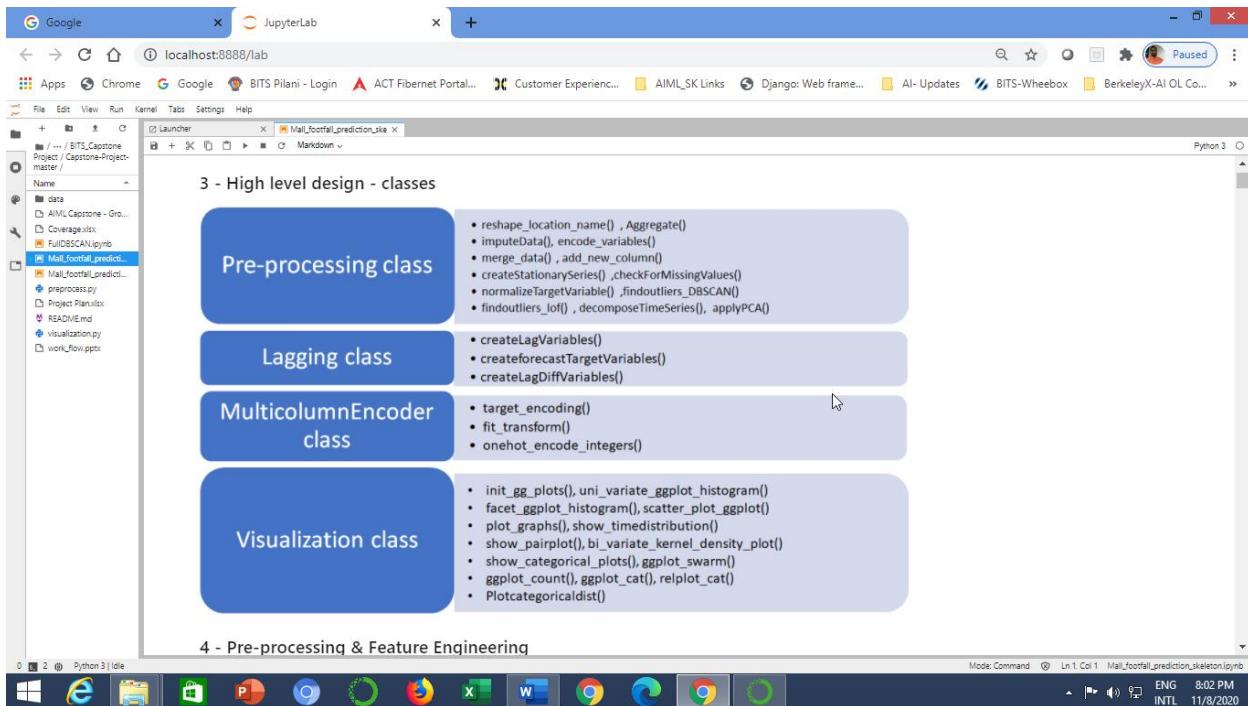
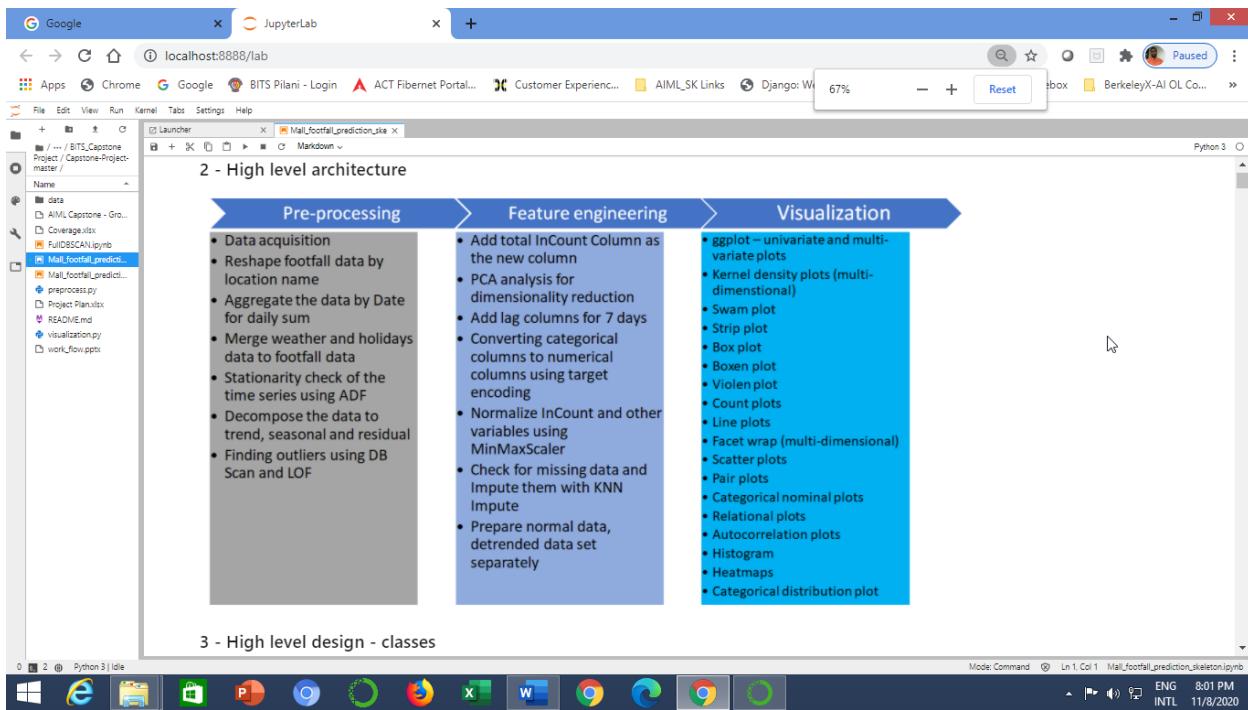
The figure shows a Jupyter Notebook interface with a Python 3 kernel. The notebook title is "Mall_footfall_prediction_skeleton.ipynb". The main content area displays a flowchart for mall footfall prediction:

- Data Sources:** In count, Weather data, Holidays data.
- Pre-processed Feature Engineered Data**
- Process Flow:**
 - Identify and build right machine learning models
 - Select the best model
 - Apply the model
- Validation Loop:** Test, Validate, Train.
- Outcome:** Predict Right, Revenue growth chart.

Revenue growth

Quarter	Revenue
Q1 20-21	60
Q2 20-21	70
Q3 20-21	90
Q4 20-21	104
Q1 21-22	152
Q2 21-22	176

Mode: Command | Line 1, Col 1 | Mall_footfall_prediction_skeleton.ipynb



Google JupyterLab

localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

Launcher Mail_footfall_prediction_ske

Markdown Python 3

4 - Pre-processing & Feature Engineering

4.1 Data Acquisition



```
[1]: pip install --upgrade category_encoders  
pip install pandas  
pip install statsmodels --upgrade  
Collecting category_encoders  
  Downloading https://files.pythontab.org/packages/44/57/fc4f41248797e02e32202099e432adea35555cb370fffcfa22727/category_encoders-2.2.2-py2.py3-none-any.whl (80kB)  
Requirement already satisfied: pandas<1.0.0,>=0.23.4 in /usr/local/lib/python3.6/dist-packages (from category_encoders) (0.23.4)  
Requirement already satisfied: skipping upgrade: pandas<0.23.4,>=0.21.1 in /usr/local/lib/python3.6/dist-packages (from category_encoders) (0.21.1)  
Requirement already satisfied: skipping upgrade: numpy<1.19.0,>=1.17.0 in /usr/local/lib/python3.6/dist-packages (from category_encoders) (1.17.0)  
Requirement already satisfied: skipping upgrade: numpy<1.19.0,>=1.16.0 in /usr/local/lib/python3.6/dist-packages (from category_encoders) (1.16.0)  
Requirement already satisfied: skipping upgrade: numpy<1.19.0,>=1.15.0 in /usr/local/lib/python3.6/dist-packages (from category_encoders) (1.15.0)  
Requirement already satisfied: skipping upgrade: scikit-learn<0.26.8,>=0.23.2 in /usr/local/lib/python3.6/dist-packages (from category_encoders) (0.23.2.post1)  
Requirement already satisfied: skipping upgrade: numpy<1.19.0,>=1.14.0 in /usr/local/lib/python3.6/dist-packages (from category_encoders) (1.14.0)
```

Python 3 Idle

Mode: Command EN1 Ln: Col 1 Mail_footfall_prediction_skeleton.ipynb

INTL 8/2020 11/8/2020

```
Requirement already satisfied, skipping upgrade: pytz<=2017.2 in /usr/local/lib/python3.6/dist-packages (from pandas<0.21.1->category_encoders) (2018.9)
Requirement already satisfied, skipping upgrade: upgrade<0.1.0,!=0.2.3 in /usr/local/lib/python3.6/dist-packages (from pandas<0.21.1->category_encoders) (2.8.1)
Requirement already satisfied, skipping upgrade: six<1.13.0 in /usr/local/lib/python3.6/dist-packages (from statsmodels<0.11.1->category_encoders) (1.15.0)
Requirement already satisfied, skipping upgrade: joblib<0.11 in /usr/local/lib/python3.6/dist-packages (from scikit-learn<0.20.0->category_encoders) (0.17.0)
Installing collected packages: category-encoders
Successfully installed category-encoders-2.2.2
Collecting pmdarima
  Downloading https://files.pythontesting.net/packages/be/62/725b3baebe5de7753adea8139322e7b03ca53f5db8de3b7de87db20/pmdarima-1.7.1-cp36-cp36-manylinux1_x86_64.whl (1.5MB)
Collecting statsmodels<0.12,>=0.11
  Downloading https://files.pythontesting.net/packages/cb/83/549f4d223aa1babec6cc2f0efaf170fc0ac16a1m60b/statsmodels-0.11.1-cp36-cp36-manylinux1_x86_64.whl (8.7MB)
Requirement already satisfied: numpy<1.17.0 in /usr/local/lib/python3.6/dist-packages (from pmdarima) (1.22.2/post1)
Requirement already satisfied: numpy<1.17.0 in /usr/local/lib/python3.6/dist-packages (from pmdarima) (1.18.5)
Requirement already satisfied: scipy<1.1.2 in /usr/local/lib/python3.6/dist-packages (from pmdarima) (1.4.1)
Requirement already satisfied: urllib3<1.25.0 in /usr/local/lib/python3.6/dist-packages (from pmdarima) (1.24.3)
Requirement already satisfied: statsmodels<0.19 in /usr/local/lib/python3.6/dist-packages (from pmdarima) (1.14)
Collecting Cython<0.25.18,>=0.29
  Downloading https://files.pythontesting.net/packages/e7/5b/024f3e9193c7e1cb03578954fe9203e455e857eebe/cython-0.29.17-cp36-cp36-manylinux1_x86_64.whl (8.7MB)
Requirement already satisfied: six<1.18.0 in /usr/local/lib/python3.6/dist-packages (from pmdarima) (0.17.0)
Collecting setuptools<50.0.0
  Downloading https://files.pythontesting.net/packages/c3/0a/5dc324059194812e938a2d214893c2faedff66e057afe7e67d05/setuptools-49.6.0-py3-none-any.whl (80.0kB)
Requirement already satisfied: patso<0.5 in /usr/local/lib/python3.6/dist-packages (from statsmodels<0.11.1->pmdarima) (0.5.1)
Requirement already satisfied: pydateutil<2.7.3 in /usr/local/lib/python3.6/dist-packages (from statsmodels<0.11.1->pmdarima) (2.8.1)
Requirement already satisfied: python-dateutil<2.7.3 in /usr/local/lib/python3.6/dist-packages (from statsmodels<0.11.1->pmdarima) (2.8.1)
Requirement already satisfied: six<1.18.0 in /usr/local/lib/python3.6/dist-packages (from statsmodels<0.11.1->pmdarima) (1.15.0)
ERROR: datscience 0.18.0 has requirement future<0.2.1, but you'll have future 0.18.3 which is incompatible.
Intalling datscience-0.18.0 from https://files.pythontesting.net/packages/0a/0a/datscience-0.18.0.tar.gz
Found existing installation: statsmodels 0.10.2
Uninstalling statsmodels-0.10.2:
  Successfully uninstalled statsmodels-0.10.2
Found existing installation: Cython 0.29.21
Uninstalling Cython-0.29.21:
  Successfully uninstalled Cython-0.29.21
Found existing installation: setuptools 50.3.2
Uninstalling setuptools-50.3.2:
  Successfully uninstalled setuptools-50.3.2
Successfully installed Cython-0.29.17-pmdarima-1.7.1 setuptools-49.6.0 statsmodels-0.11.1
Collecting pmdarima
  Downloading https://files.pythontesting.net/packages/be/4c/9e9435c6a45d8bdffa51581f0a3652b934a2f9e05de67c130926d/statsmodels-0.11.2.1-cp36-cp36-manylinux1_x86_64.whl (9.5MB)
Requirement already satisfied, skipping upgrade: six<1.13.0 in /usr/local/lib/python3.6/dist-packages (from statsmodels) (1.14.4)
Requirement already satisfied, skipping upgrade: scipy<1.1 in /usr/local/lib/python3.6/dist-packages (from statsmodels) (1.4.1)
Requirement already satisfied, skipping upgrade: numpy<1.15 in /usr/local/lib/python3.6/dist-packages (from statsmodels) (1.18.5)
Requirement already satisfied, skipping upgrade: pytz<=2017.2 in /usr/local/lib/python3.6/dist-packages (from statsmodels<0.11.2.1->statsmodels) (2.8.1)
Requirement already satisfied, skipping upgrade: python-dateutil<2.7.3 in /usr/local/lib/python3.6/dist-packages (from statsmodels<0.11.2.1->statsmodels) (2.8.1)
Requirement already satisfied, skipping upgrade: pytz<=2017.2 in /usr/local/lib/python3.6/dist-packages (from pandas<0.21->statsmodels) (2018.9)
Requirement already satisfied, skipping upgrade: six<1.13.0 in /usr/local/lib/python3.6/dist-packages (from statsmodels) (1.16.0)
```



```

[1]: pip install https://h2o-release.s3.amazonaws.com/h2o/rel-zermelo/1/Python/h2o-3.32.0.1-py2.py3-none-any.whl
Collecting h2o>=3.32.0.1
  Downloading https://h2o-release.s3.amazonaws.com/rel-zermelo/1/Python/h2o-3.32.0.1-py2.py3-none-any.whl (164.6MB)
    Requirement already satisfied: requests in /usr/local/lib/python3.6/dist-packages (from h2o>=3.32.0.1) (2.23.0)
    Collecting colorama>=1.3.8
      Downloading https://files.pythonhosted.org/packages/44/b9/588b279fb4e128d5c57249f72141394141e04d450286189448/colorama-0.4.4-py2.py3-none-any.whl
        Requirement already satisfied: future in /usr/local/lib/python3.6/dist-packages (from h2o>=3.32.0.1) (0.16.0)
        Requirement already satisfied: idna>=2.5 in /usr/local/lib/python3.6/dist-packages (from requests>=2.23.0>=3.32.0.1) (2.28.0)
        Requirement already satisfied: certifi>=2018.11.29 in /usr/local/lib/python3.6/dist-packages (from requests>=2.23.0>=3.32.0.1) (2020.6.20)
        Requirement already satisfied: urllib3>=1.25.0,<1.26.0 in /usr/local/lib/python3.6/dist-packages (from requests>=2.23.0>=3.32.0.1) (1.24.3)
        Requirement already satisfied: charset-normalizer>=2.0.0 in /usr/local/lib/python3.6/dist-packages (from requests>=2.23.0>=3.32.0.1) (3.0.4)
    Installing collected packages: colorama, h2o
      Successfully installed colorama-0.4.4 h2o-3.32.0.1

[3]:

```

[4]:

```

import numpy as np # Linear algebra
import pandas as pd # Data processing, CSV file I/O (e.g. pd.read_csv)
import matplotlib.pyplot as plt
import seaborn as sns
sns.set(style="darkgrid")
from sklearn.preprocessing import LabelEncoder
from datetime import datetime
import time
from sklearn.preprocessing import MinMaxScaler
from category_encoders.target_encoder import TargetEncoder
import missingno as nso
from statsmodels.tsa import KNNImputer
import warnings
warnings.filterwarnings('ignore')
from sklearn.impute import SimpleImputer
from sklearn.impute import IterativeImputer
from sklearn.neighbors import NearestNeighbors
from sklearn.cluster import DBSCAN
from sklearn.ensemble import LocalOutlierFactor
from plotnine import *
from scipy import signal, fftpack
from statsmodels.tsa.seasonal import seasonal_decompose
from dateutil.parser import parse
from statsmodels.tsa.seasonal import STL
from dymon import nominal

```

```

[1]: from olympos import numerai
from statsmodels.tsa.stattools import adfuller, acf
from sklearn.decomposition import PCA
from statsmodels.graphics.tsaplots import plot_acf, plot_pacf

# pyplot
from plotly.offline import download_plotlyjs, init_notebook_mode, plot, iplot
from plotly.graph_objs import Scatter, Figure, Layout
import plotly
import plotly.graph_objs as go
from plotly import tools

# yellowbrick visualizations
from yellowbrick.features import ParallelCoordinates
from yellowbrick.model_selection import FeatureImportances
from yellowbrick.regressor import ResidualsPlot
from yellowbrick.regressor import PredictionError

# regression models
from sklearn.linear_model import LinearRegression
from sklearn.linear_model import Lasso
from sklearn.linear_model import Ridge
from sklearn.linear_model import ElasticNet
from sklearn.linear_model import BayesianRidge
from sklearn.linear_model import LassoLars
from sklearn.linear_model import PassiveAggressiveRegressor
from sklearn.linear_model import StochasticRegressor
from sklearn.linear_model import SGDRegressor
from sklearn.tree import DecisionTreeRegressor
from sklearn.tree import ExtraTreeRegressor
from sklearn.svm import SVR
from sklearn.ensemble import AdaBoostRegressor
from sklearn.ensemble import BaggingRegressor
from sklearn.ensemble import RandomForestRegressor
from sklearn.ensemble import GradientBoostingRegressor
from xgboost import XGBRegressor
from lightgbm import LightGBMRegressor
from yellowbrick.regressor import CorrelationMatrix
from sklearn.base import clone
from sklearn.metrics import mean_squared_error, mean_absolute_error, r2_score
from sklearn.model_selection import train_test_split, RandomizedSearchCV, cross_val_predict, cross_val_score, cross_validate
from sklearn.feature_selection import RFECV
from sklearn.metrics import make_scorer
from scipy.stats import loguniform
from Regressor import *
from auto import log

```

```

# neural network model
from sklearn.neural_network import MLPRegressor
# LSTM models
from keras.models import Sequential
from keras.layers import Dense, Add, Concatenate, Reshape
from keras.layers import LSTM, Conv1DND, Bidirectional
from keras.layers import Dropout, Average, RepeatVector, LeakyReLU, BatchNormalization
from keras import optimizers
from keras.optimizers.schedules import ExponentialDecay
from keras.callbacks import ModelCheckpoint
from keras.callbacks import EarlyStopping, ReduceLROnPlateau
from keras.callbacks import LearningRateScheduler
from keras.layers.merge import concatenate

from keras.engine import training
from keras.models import Model, Input

from keras.layers import Flatten
from keras.layers import TimeDistributed
from keras.layers import Dense, Add, Concatenate, Reshape
from keras.layers.convolutional import MaxPooling3D
from keras.callbacks import Callback
from keras.models import load_model
from keras.optimizers import Adam

# H2O models
import h2o
from h2o.estimators.gbm import H2OGradientBoostingEstimator
from h2o.estimators import H2ONaiveBayesEstimator
from h2o.estimators.stackedensemble import H2OStackedEnsembleEstimator
from h2o.grid.grid_search import H2OGridSearchCV
from h2o.model.regression import H2OJ2Score
from h2o.grid.grid_search import H2OGridSearch
from h2o.grid import H2OGridMutator
h2o.init()

# Fb prophet models
from fbprophet import Prophet
from fbprophet.plot import plot_plotly, plot_components_plotly

# ARIMA model
from statsmodels.tsa.arima.model import ARIMA
import statsmodels
from pedarina.pipeline import Pipeline
from pedarina.preprocessing import BoxCoxEndogTransformer

```

```

from pedarina.preprocessing import BoxCoxEndogTransformer
import gdo
import os
from keras import backend as K

Checking whether there is an H2O instance running at http://localhost:54321 ..... not found.
Attempting to start a local H2O server...
Data World: https://www.h2o.ai/h2o/2020-09-20; OpenDK Runtime Environment (build 11.0.9+11-Ubuntu-0ubuntu1.18.04.1); OpenDK 64-Bit Server VM (build 11.0.9+11-Ubuntu-0ubuntu1.18.04.1, mixed mode, sharing)
Starting server from /usr/local/lib/python3.6/dist-packages/h2o/h2o.jar
Ice root: /tmp/h2oblock_wvn
H2O output: /tmp/h2oblock_wvn/h2o_unknowneader_started_from_python.out
H2O status: /tmp/h2oblock_wvn/h2o_unknowneader_started_from_python.err
Server is running at http://127.0.0.1:54321
Server is running at http://127.0.0.1:54321 ... successfull.

H2O cluster uptime: 03 secs
H2O cluster timezone: Etc/UTC
H2O data parsing timezone: UTC
H2O cluster version: 3.32.0.1
H2O cluster version age: 28 days, 20 hours and 55 minutes
H2O cluster name: H2O from python_unknowneader_sp00
H2O cluster total nodes: 1
H2O cluster free memory: 3.180 Gb
H2O cluster total cores: 2
H2O cluster allowed cores: 2
H2O cluster status: accepting new members, healthy
H2O connection url: http://127.0.0.1:54321
H2O connection proxy: ("http": null, "https": null)
H2O internal security: False
H2O API Extensions: Amazon S3, XGBoost, Algos, AutoML, Core V3, TargetEncoder, Core V4
Python version: 3.6.9 final

[5]: mzingpath='https://raw.githubusercontent.com/syechearthii/footfall-prediction-code/master/'
df1 = pd.read_csv(mzingpath+'footfall_Part1.csv')
df2 = pd.read_csv(mzingpath+'footfall_Part2.csv')
df3 = pd.read_csv(mzingpath+'footfall_Part3.csv')

```

Google JupyterLab localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

Launcher Mail_footfall_prediction_ske

/ ... /BITS_Capstone Project Capstone-project-master / Name / data AIML Capstone - Gro... Coverage.xlsx FullDSCAN.ipynb Mail_footfall.predict_ preprocesipy Project Plan.xlsx README.md visualization.py work_flow.pptx

Markdown

```
# Importing required libraries
df_foot_fall = pd.read_csv('df_foot_fall.csv')
df_weather_data = pd.read_csv('mainpath\\overall_weather.csv')
df_easter_sunday_holiday_data = pd.read_csv('mainpath\\easter_sundays.csv')
df_christmas_long_holiday_data = pd.read_csv('mainpath\\christmas_hols_long.csv')
df_us_bank_holidays_data = pd.read_csv('mainpath\\USbankholidays.csv')
df_school_long_holidays_data = pd.read_csv('mainpath\\school_hols_long.csv')
```

Python 3

4.2 Data preparation

```
[6]: df_foot_fall.shape
```

```
[6]: (554368, 15)
```

```
[7]: df_easter_sunday_holiday_data.shape
```

```
[7]: (18, 2)
```

```
[8]: df_foot_fall.info
```

```
[8]: <class 'pandas.core.frame.DataFrame'>
RangeIndex: 554368 entries, 0 to 554367
Data columns (total 15 columns):
 #   Column          Non-Null Count  Dtype  
--- 
 0   Date            554368 non-null  datetime64[ns]
 1   InCount         554368 non-null  int64  
 2   BRCWeek         554368 non-null  int64  
 3   LocationName    554368 non-null  object  
 4   Hour            554368 non-null  int64  
 5   Weekday         554368 non-null  int64  
 6   NewDate         554368 non-null  datetime64[ns]
 7   Month           554368 non-null  int64  
 8   Year            554368 non-null  int64  
 9   Day             554368 non-null  int64  
 10  Day yr          554368 non-null  int64  
 11  Month yr       554368 non-null  int64  
 12  Week yr        554368 non-null  int64  
 13  Hour yr         554368 non-null  int64  
 14  week_no         554368 non-null  int64  
dtypes: datetime64[ns](2), int64(12), object(1)
memory usage: 6.9+ MB
```

```
[8]: [554368 rows x 15 columns]
```

```
[9]: df_foot_fall.head()
```

```
[9]:
```

Date	InCount	BRCWeek	LocationName	Hour	Weekday	NewDate	Month	Year	Day	Day yr	Month yr	Week yr	Hour yr	week_no
0 01-01-2009	268	53	Brigate	0	Thursday	01-Jan-09	1	2009	1	Thursday/01-Jan-2009	Jan-09	00-2009	00,01-Jan-2009	1
1 01-01-2009	366	53	Brigate	1	Thursday	01-Jan-09	1	2009	1	Thursday/01-Jan-2009	Jan-09	00-2009	00,01-Jan-2009	1
2 01-01-2009	273	53	Brigate	2	Thursday	01-Jan-09	1	2009	1	Thursday/01-Jan-2009	Jan-09	00-2009	00,01-Jan-2009	1
3 01-01-2009	219	53	Brigate	3	Thursday	01-Jan-09	1	2009	1	Thursday/01-Jan-2009	Jan-09	00-2009	00,01-Jan-2009	1

Mode: Command Ln: Col 1 Mail_footfall_prediction_skeleton.pyw

0 1 2 Python 3 idle

P O X W C Chrome BerkeleyX-AI OL Co... INT 11/8/2020 ENG 8:08 PM

JupyterLab

localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

Launcher Mail_footfall_prediction_skeleton

4.2.2 MultiColumnEncoder Class

```
[12]: class MultiColumnEncoder:
    def __init__(self, categoricalcolumns = None, includedcolumns=None):
        self._categoricalcolumns = categoricalcolumns # array of column names to encode
        self._included_cols = includedcolumns

    def fit(self,X,y=None):
        return self # not relevant here

    def transform(self,X):
        ...
        Transforms columns of X specified in self.columns using
        LabelEncoder(). If no columns specified, transforms all
        columns in X.

        output = X.copy()
        if self._categoricalcolumns is not None:
            for col in self._categoricalcolumns:
                output[col] = LabelEncoder().fit_transform(output[col])
        else:
            for column,col in output.items():
                output[column] = LabelEncoder().fit_transform(column)
        return output

    def onehot_encode_integers(self,X, df):
        df = df.copy()
        df_encoded = pd.get_dummies(df, columns = self._included_cols)
        return df_encoded

    def fit_transform(self,X,y=None, encoding_type = 'target_encoding'):
        if encoding_type == 'onehot_encoding':
            return self.onehot_encode_integers(X)
        elif encoding_type == 'target_encoding':
            return self.target_encoding(X,y)
        else:
            raise ValueError("encoding_type must be onehot_encoding or target_encoding")

    def target_encoding(self,X, target_column):
        return TargetEncoder(cols=self._included_cols).fit_transform(X,target_column)
```

4.2.3 Lagging Class

```
[13]: class Lagging:
    def __init__(self, location_name_list=None, n_timesteps=7, n_forecast_steps=1):
        self.n_timesteps = n_timesteps
        self.location_name_list = location_name_list # array of Locations
        self.n_forecast_steps = n_forecast_steps

    def createforecasttargetvariables(self, df):
        # copy the InCount to other forecast columns
        for forecast_step in range(1, self.n_forecast_steps+1):
            df['InCountTotal' + '_forecast' + str(forecast_step)] = df['InCountTotal']

        #update multistep output variables
        for i in range(0, df.shape[0]-self.n_forecast_steps):
            for forecast_step in range(1, self.n_forecast_steps+1):
                df['InCountTotal' + '_forecast' + str(forecast_step)][i] = df['InCountTotal'][i+forecast_step]

        #remove the last nforecast rows as these do not have proper values after creating the forecast variables
        df.drop(df.tail(self.n_forecast_steps).index, inplace=True)
        return df
```

Mode Command Ln 1, Col 1 Mail_footfall_prediction_skeleton.ipynb

0 2 Python 3 idle ENGLISH 8:39 PM INTL 11/8/2020

JupyterLab

localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

Launcher Mail_footfall_prediction_skeleton

4.2.3 Lagging Class

```
[13]: class Lagging:
    def __init__(self, location_name_list=None, n_timesteps=7, n_forecast_steps=1):
        self.n_timesteps = n_timesteps
        self.location_name_list = location_name_list # array of Locations
        self.n_forecast_steps = n_forecast_steps

    def createforecasttargetvariables(self, df):
        # copy the InCount to other forecast columns
        for forecast_step in range(1, self.n_forecast_steps+1):
            df['InCountTotal' + '_forecast' + str(forecast_step)] = df['InCountTotal']

        #update multistep output variables
        for i in range(0, df.shape[0]-self.n_forecast_steps):
            for forecast_step in range(1, self.n_forecast_steps+1):
                df['InCountTotal' + '_forecast' + str(forecast_step)][i] = df['InCountTotal'][i+forecast_step]

        #remove the last nforecast rows as these do not have proper values after creating the forecast variables
        df.drop(df.tail(self.n_forecast_steps).index, inplace=True)
        return df

    def createlagvariables(self,df,location_wise_data = True): #Create Lag variables of size=n_timesteps
        if location_wise_data == True:#Location wise data
            for location in self.location_name_list:
                location = 'InCount'+location
                for lag in range(0, self.n_timesteps):
                    df[location+'_'+str(lag+1)] = df[location].shift(periods=lag)
            else:#totalincount only
                for lag in range(0, self.n_timesteps):
                    df['InCountTotal' + '_lag' + str(lag+1)] = df['InCountTotal'].shift(periods=lag)

        df.dropna(inplace = True)
        return df

    def createlagdiffvariables(self,df,location_wise_data = True): #Create lag variables of size=n_timesteps
        if location_wise_data == True:#Location wise data
            for location in self.location_name_list:
                location = 'InCount'+location
                for lag in range(1, self.n_timesteps+1):
                    df[location+'_'+str(lag)] = df[location].diff(periods=lag)
            else:# totalincount only
                for lag in range(1, self.n_timesteps+1):
                    df['InCountTotal' + '_lag' + str(lag)] = df['InCountTotal'].diff(periods=lag)

        df.dropna(inplace = True)
        return df
```

Mode Command Ln 1, Col 1 Mail_footfall_prediction_skeleton.ipynb

0 2 Python 3 idle ENGLISH 8:40 PM INTL 11/8/2020

JupyterLab x +

localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

Launcher x Mail_footfall_prediction_sklearn

Markdown v Python 3

4.2.4 Pre_Processing Class

```
[14]: class pre_processing:
    def __init__(self, data_frame):
        self.data_frame = data_frame.copy() # which data_frame to pre-process
        self.location_name_list = []
        self.date_index = []
        self.scalar_Target = MinMaxScaler(feature_range = (0, 1))
        self.scalar_Input = MinMaxScaler(feature_range = (0, 1))
        self.scalar_Target_each_location = MinMaxScaler(feature_range = (0, 1))
        self.scalar_Input_each_location = MinMaxScaler(feature_range = (0, 1))

    def cleanup_data(self, col_name='Date'):
        # there is a white space in one of the date values in Easter_holidays list. We need to remove that
        if col_name == 'Date':
            self.data_frame[col_name] = self.data_frame[col_name].str.replace(' ', '')
        return self.data_frame

    def aggregate(self, group_by_cols, retain_list=None):
        # aggregate the column values based on conditions
        if retain_list is None:
            retain_list = []
        for location_name in self.location_name_list:
            key = "InCount"+location_name
            retain_list.append(key)
        # now aggregate based on Date
        self.data_frame = self.data_frame.groupby(group_by_cols).agg(retain_list).reset_index()

    def createWeekData(self, data_frame):
        group_by_cols = ['Year', 'Week_no']
        retain_list = { 'InCountTotal':'sum', 'Month':'first', 'mean_temp':'mean', 'rain':'mean', 'wind_speed':'mean' }
        data_frame.groupby(group_by_cols).agg(retain_list).reset_index()
        return df

    def reshape_location_name(self, group_by_cols, target_col='InCount'):# Review this
        data_temp_df = pd.DataFrame(data=self.data_frame[['Date','Weekday','Month']])
        data_temp_df.drop_duplicates(['Date'], inplace=True)

        self.location_name_list = self.data_frame['locationName'].unique()

        self.data_frame = self.data_frame.pivot_table(index=group_by_cols, columns='locationName',values=[target_col]).reset_index()
        self.data_frame.columns = ['+'.join([col] for col in self.data_frame.columns)']

        self.data_frame['Date'] = pd.to_datetime(self.data_frame['Date']).dt.strftime("%d-%b-%Y")
        self.location_name_list = self.data_frame['locationName'].unique()

        self.data_frame = self.data_frame.melt(id_vars=group_by_cols, value_vars=[target_col], reset_index())
        self.data_frame = self.data_frame.pivot_table(index=group_by_cols, columns='locationName',values=[target_col]).reset_index()

    def get_top_n(self, n=10):
        top_n = self.data_frame.groupby('locationName').sum().sort_values('InCount', ascending=False).head(n)
        return top_n
```

JupyterLab

localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

Launcher Mail_football_prediction_ske

```
self.location_name_list = self.data_frame['LocationName'].unique()

self.data_frame = self.data_frame.pivot_table(index='group_by_cols', columns='LocationName', values=[target_col]).reset_index()
self.data_frame.columns = ['_'.join(col) for col in self.data_frame.columns]

return self.data_frame.merge(data_temp_df, on='Date', how='inner')

def encode_variables(self, included_cols, target_variable = 'InCountTotal', encoding_type = 'target_encoding'):
    categorical_target_encoding_of_the_column_value

    self.data_frame = MultiColumnEncoder(includedcolumns = included_cols).fit_transform(self.data_frame,target_variable,encoding_type = 'target_encoding')

    return self.data_frame

def imputevalues(self):
    # fill missing values with mean column values as only missing column i see is InCountAllOnStSouth
    self.data_frame.fillna(self.data_frame.mean(), inplace=True)
    return self.data_frame

def pre_process(self, data_frame_type='FootBall'):
    for different data frame types, do the different pre-processing
    ...

    if(data_frame_type == 'EasterSunday'):
        # cleanup the date columns
        self.data_frame = self.cleanup_data(col_name='Date')

    return self.data_frame

def merge_data(self, data_to_merge, col_by = 'Date'):
    data_to_merge['Date'] = pd.to_datetime(data_to_merge['Date']).dt.strftime('%d-%b-%Y')
    self.data_frame = self.data_frame.merge(data_to_merge, how='inner', left_on=col_by, right_on=col_by)
    return self.data_frame

def add_new_col(self, data_to_add, new_col_name, source_date_format = '%d-%b-%Y', destination_date_format = '%m-%d-%Y'):
    self.data_frame['Tempcolumn'] = pd.to_datetime(self.data_frame['Date'], format = source_date_format).dt.strftime('%d/%m/%Y')

    temp_data = data_to_add.copy()

    temp_data['Tempcolumn'] = pd.to_datetime(temp_data['Date'], format = destination_date_format).dt.strftime('%d/%m/%Y')

    self.data_frame[new_col_name] = (self.data_frame.set_index(['Tempcolumn']).index.isin(temp_data.set_index(['Tempcolumn']).index)).astype(int)
    temp_data.drop(['Tempcolumn'],axis=1)

    self.data_frame.drop(['Tempcolumn'],axis=1, inplace=True)

    return self.data_frame

def createAttributesColumns(self):
    date_time = pd.to_datetime(self.data_frame['Date'], format = '%d-%b-%Y')
```

JupyterLab

localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

Launcher Mail_footfall_prediction_skeleton

```
date_time = pd.to_datetime(self.data_frame['Date'], format='dd-mm-yy')
timestamp_s = date_time.map(datetime.timestamp)

dt_object = datetime.fromtimestamp(timestamp_s[0])
print("dt_object:", dt_object)
print("type(dt_object):", type(dt_object))

day = 24*60*60
week = 7*day
month = 31*day
year = (365.2425)*day

self.data_frame['Day sin'] = np.sin(timestamp_s * (2 * np.pi / day))
self.data_frame['Day cos'] = np.cos(timestamp_s * (2 * np.pi / day))
self.data_frame['Week sin'] = np.sin(timestamp_s * (2 * np.pi / week))
self.data_frame['Week cos'] = np.cos(timestamp_s * (2 * np.pi / week))
self.data_frame['Month sin'] = np.sin(timestamp_s * (2 * np.pi / month))
self.data_frame['Month cos'] = np.cos(timestamp_s * (2 * np.pi / month))
self.data_frame['Year sin'] = np.sin(timestamp_s * (2 * np.pi / year))
self.data_frame['Year cos'] = np.cos(timestamp_s * (2 * np.pi / year))

return self.data_frame
```

def createExtraFeatures(self, n_timesteps=7):
 # Add new columns
 # Additive Decomposition
 result_add = seasonal_decompose(self.data_frame['InCountTotal'], model='additive', freq=30)
 self.data_frame['seasonal_incount'] = result_add.seasonal
 self.data_frame['trend_incount'] = result_add.trend
 # rolling mean column
 self.data_frame['rolling_mean_incount'] = self.data_frame['InCountTotal'].shift(1).rolling(window=n_timesteps).mean()
 # expanding mean column
 self.data_frame['expanding_mean_incount'] = self.data_frame['InCountTotal'].shift(1).expanding(2).mean()
 return self.data_frame

def createTimeSeries(self):
 self.data_frame['Date'] = pd.to_datetime(self.data_frame.Date.astype(str))
 self.data_frame = self.data_frame.set_index('Date').index.month
 self.data_frame['Date'] = self.data_frame.index.year
 self.data_frame['Day'] = self.data_frame.index.day
 self.data_frame['Weekday'] = self.data_frame.index.weekday
 self.data_frame['Week_no'] = self.data_frame.index.week



JupyterLab

localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

Launcher Mail_footfall_prediction_skeleton

```
self.data_frame['Year'] = self.data_frame.index.year
self.data_frame['Day'] = self.data_frame.index.day
self.data_frame['Weekday'] = self.data_frame.index.weekday
self.data_frame['Month'] = self.data_frame.index.month

for i in range(len(self.data_frame)):
    self.data_frame['Weekday'].iloc[i] = calendar.day_name[self.data_frame['Weekday'].iloc[i]]

return self.data_frame
```

def createStationarySeries(self, col_name = 'InCountTotal', lag=4):
 # check for stationarity
 result = adfuller(self.data_frame[col_name].dropna())
 print('ADF test results for ', col_name)
 print('ADF Statistic: %f' % result[0])
 print('p-value: %f' % result[1])
 print('Critical Value: ', result[4])
 print('Full result ', result)
 print('*****')

 for location in self.location_name_col_list:
 result = adfuller(self.data_frame[location].dropna())
 print('ADF test results for ', location)
 print('ADF Statistic: %f' % result[0])
 print('p-value: %f' % result[1])
 print('Critical Value: ', result[4])
 print('Full result ', result)
 print('*****')

 self.data_frame[col_name] = self.data_frame[col_name].diff(lag)
 self.data_frame.dropna(inplace = True)
 return self.data_frame

def checkForMissingValues(self, location_name_list):
 for location in location_name_list:
 self.data_frame[location].loc[self.data_frame['InCount'+location]==0, 'InCount'+location] = np.NaN

def createLog_forecastVariables(self, location_name_list=None, n_timesteps=7, n_forecast_steps = 1, location_wise_data = True):
 nlagClass = Logging(location_name_list,n_timesteps, n_forecast_steps)
 self.data_frame = nlagClass.createLogVariables(self.data_frame,location_wise_data = location_wise_data)
 # create log data for total will include
 totalData = nlagClass.createLogifyVariables(self.data_frame, location_wise_data = location_wise_data)
 self.data_frame = nlagClass.createForecastTargetVariables(self.data_frame)

return self.data_frame

def addTotalInCountColumn(self):
 self.location_name_col_list['InCount']= [x for x in self.location_name_list]
 self.data_frame['InCountTotal'] = self.data_frame[self.location_name_col_list].sum(axis=1)
 self.data_frame['InCountTotal'] = self.data_frame['InCountTotal'].astype(int)



JupyterLab

localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

Launcher Mail_footfall_prediction_ske

Markdown Python 3

```
def addTotalInCountColumn(self):
    self.location_name_col_list[1] = "InCount"+x for x in self.location_name_list
    self.data_frame[["InCountTotal"] = self.data_frame[self.location_name_col_list].sum(axis=1)
    self.data_frame[["InCountTotal"] = self.data_frame[["InCountTotal"]].astype(int)
    return self.data_frame

def normalizeTargetVariable(self):
    #print("normalization")
    self.data_frame[["InCountTotal"] = self.scaler_target.fit_transform(self.data_frame[["InCountTotal"]])

location_name_col_list[1] = "InCount"+x for x in self.location_name_list
for location_col in location_name_col_list:
    self.data_frame[[location_col]] = self.scaler_target_each_location.fit_transform(self.data_frame[[location_col]])
return self.data_frame

def normalizeInputVariables(self, included_cols=None):
    #print("normalize input variables")
    self.data_frame[included_cols] = self.scaler_input.fit_transform(self.data_frame[included_cols])
    return self.data_frame

def visualizeMissingValues(self):
    #print("missing values visualization")
    msno.bar(self.data_frame)

#Input using KNN
def imputeDataUsingKNN(self):
    X = self.data_frame[["InCountCommercialStBarrett", "InCountCommercialStLush"]]
    imputer = KNNImputer(n_neighbors=2)
    imputed_df = imputer.fit_transform(X)
    imputed_df = pd.DataFrame(imputed_df, columns=X.columns)
    #print(imputed_df)
    #print(self.data_frame[["InCountCommercialStBarrett"]])
    self.data_frame[["InCountCommercialStBarrett"]] = imputed_df.values
    return self.data_frame

#Input using IterativeImputer which works similar to MICE
def imputeDataUsingIterativeImputation(self):
    X = self.data_frame[["InCountCommercialStBarrett", "InCountCommercialStLush"]]
    imputer = IterativeImputer(max_iter=10, verbose=0)
    imputed_df = imputer.fit_transform(X)
    imputed_df = pd.DataFrame(imputed_df, columns=X.columns)
    #print(imputed_df)
    #print(self.data_frame[["InCountCommercialStBarrett"]])
    self.data_frame[["InCountCommercialStBarrett"]] = imputed_df.values
    return self.data_frame

def findOutliers_DBSCAN(self,columns,number_neigh):
    
```

JupyterLab

localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

Launcher Mail_footfall_prediction_ske

Markdown Python 3

```
def findoutliers_DBSCAN(self,columns,number_neigh):  
    # This function takes as input, the dataframe name, columns to be analysed for the outlier, nearest neighbours to be considered as input  
    # A new column named outlier will be created with name Outlier_<column name and it will have values 1 If outlier and 0 otherwise  
    # number_neigh is calculated based on instances  
    nco=0  
    fig, ax = plt.subplots(len(columns),figsize=(15,15))  
    for col in columns:  
        # finding the distances based on nearest neighbour to find optimal eps  
        nbrs = NearestNeighbors(n_neighbors=number_neigh).fit(np.array(self.data_frame[col]).reshape(-1,1))  
        distances, indices = nbrs.kneighbors(np.array(self.data_frame[col]).reshape(-1,1))  
        distances = np.sort(distances, axis=0)  
        distances = distances[:,1]  
  
        # If plotted on a graph, a point with maximum curvature . , i.e, a point with maximum slope gives the best eps  
        maxslope=0  
        noslope=0  
        diff=[ ]  
        for i in range(1,distances.shape[1]):  
            slope = (distances[i-1]-distances[i-1])/(indices[i][0]-indices[i-1][0])  
            if (slopel > maxslope):  
                maxslope=slope  
            noslope+=1  
        db=DBSCAN(eps=slope, min_samples=number_neigh).fit(np.array(self.data_frame[col]).reshape(-1,1))  
        labels=db.labels_  
        self.data_frame['Outlier_<'+col+'>'] = [1 if val == -1 else 0 for val in labels]  
        # Visualizing the Outliers marked  
        outliers_index=[]  
        outlier_index=[]  
        dataindex=[]  
        for i in range(len(labels)):  
            if (labels[i]==-1):  
                outlier_index.append(self.data_frame[col][i])  
                outlier_index.append(indices[i][0])  
            else:  
                dataindex.append(self.data_frame[col][i])  
                dataindex.append(indices[i][0])  
        ax[nc].scatter(outlier_index,outlier_index,color='red')  
        ax[nc].scatter(dataindex,outpoint,outpoint,color='blue')  
        ax[nc].set_title(col+' DBSCAN Outliers marked in Red')  
    ncnc=1  
    n_labels = len(set(labels)) - (1 IF 1 in labels ELSE 0)  
    n_outlier = list(labels).count(1)  
    print('col: Analysis is :',ncnc)  
    print('Optimal eps is :',maxslope)
```

JupyterLab

localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

Launcher Mail_footfall_prediction_skeleton

```
n_noise_ = list(label).count(-1)
print("n_noise_ is %d" % n_noise_)
print("Optimal eps is %f" % eps)
print("Estimated number of clusters: %d" % n_clusters_)
print("Estimated number of noise points: %d" % n_noise_)

return self.data_frame
```

def findOutliers_lof(self,columns,number_pear):
This function takes as input, the dataframe name, columns to be analysed for the outlier, nearest neighbours to be considered as input
A new boolean column for outlier will be added with name Outlier+column name and it will have values 1 If outlier and 0 otherwise

```
fig,axs = plt.subplots(len(columns),figsize=(15,15))
fig.tight_layout(pad=0.1)
for col in columns:
    col.set_neighboorFactor(n.neighbors,number_near)
    pred=cif.fit_predict(np.array(self.data_frame[col]).reshape(-1,1))
    #scores = cif.negative_outlier_factor_
    self.data_frame['lof_Outlier'+col]=[1 if val == -1 else 0 for val in pred]
```

Visualizing the Outliers marked

```
outlier=[]
outlierIndex=[]
datapoint=[]
dataindex=[]

for i in range(len(self.data_frame['lof_Outlier'+col])):
    if (self.data_frame['lof_Outlier'+col][i]==-1):
        outlier.append(self.data_frame[col][i])
        outlierIndex.append(i)
    else:
        datapoint.append(self.data_frame[col][i])
        dataindex.append(i)

axs[0].scatter(dataindex,datapoint,color='blue')
axs[0].scatter(outlierIndex,outlier,color='red')
axs[0].set_title(col+" LOF Outliers marked in Red")
ncount+=1
return self.data_frame
```

def detrend_time_series(self,target_col="InCountTotal"):# using signal.detrend()
inCountData = self.data_frame[target_col]
self.data_frame[target_col] = signal.detrend(self.data_frame[target_col])
detrended_InCount = self.data_frame[target_col]

```
from matplotlib import pyplot as plt
plt.rcParams["figure.figsize"] = (16, 8)
plt.plot(detrended_InCount, label="Normal InCount")
plt.plot(de_trended_InCount, label="detrended InCount")
plt.legend(loc="best")
```

Mode Command Ln 1, Col 1 Mail_footfall_prediction_skeleton.ipynb

ENGLISH 8:45 PM INTL 11/8/2020

JupyterLab

localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

Launcher Mail_footfall_prediction_skeleton

```
plt.figure(figsize=(16, 8))
plt.plot(de_trended_InCount, label="detrended InCount")
plt.legend(loc="best")
plt.show()
```

return self.data_frame

def decomposeTimeSeries(self):
Multiplicative Decomposition
result_mul = seasonal_decompose(self.data_frame['InCountTotal'], model="multiplicative", freq=30)
Additive Decomposition
result_add = seasonal_decompose(self.data_frame['InCountTotal'], model="additive", freq=30)

```
detrended = self.data_frame['InCountTotal'].values - result_add.trend - result_add.seasonal

# Plot
plt.rcParams.update({'figure.figsize': (16,8)})
result_add.plot().suptitle('Additive Decompose using seasonal_decompose', fontsize=22)
plt.show()

plt.rcParams.update({'figure.figsize': (16, 8)})
result_add.resid.plot(kde=True, title='residual Density')
plt.show()

self.data_frame['InCountTotal'] = detrended
df_decomposed = pd.concat([result_add.seasonal, result_add.trend, result_add.resid, result_add.observed], axis=1)
df_decomposed.columns = ['seas', 'trend', 'resid', 'actual_incount_total_values']
self.data_frame = pd.concat([self.data_frame, df_decomposed], axis=1)

self.data_frame.dropna(inplace = True)
```

def decomposeTime_series_for_locations(self):
location_name_col_list=[InCount for x in self.location_name_list]
for location_name_col in location_name_col_list:
 # Multiplicative Decomposition
 result_mul = seasonal_decompose(self.data_frame[location_name_col], model="multiplicative", freq=30)
 result_mul = seasonal_decompose(self.data_frame[location_name_col], model="multiplicative", freq=30)
 detrended = self.data_frame[location_name_col].values - result_add.trend - result_add.seasonal
 # Plot
 plt.rcParams.update({'figure.figsize': (16,8)})
 result_add.plot().suptitle('Additive Decompose using seasonal_decompose ' + location_name_col, fontsize=22)
 plt.show()

plt.rcParams(figsize=(16, 8))

Mode Command Ln 1, Col 1 Mail_footfall_prediction_skeleton.ipynb

ENGLISH 8:46 PM INTL 11/8/2020

JupyterLab

localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

Launcher Mail_footfall_prediction_skeleton.ipynb Python 3

```
def decomposeTime_Series_for_locations(self):
    location_name_col_list=[InCount+x for x in self.location_name_list]
    for location_col in location_name_col_list:
        # Add trend to the data frame
        result_add = seasonal_decompose(self.data['frame'][location_col], model='additive', freq=30)
        result_mlt = seasonal_decompose(self.data['frame'][location_col], model='multiplicative', freq=30)
        detrended = self.data['frame'][location_col].values - result_add.trend - result_mlt.seasonal
        # Plot
        plt.rcParams.update({'figure.figsize': (16,8)})
        result_add.plot().suptitle('Additive Decomposition using seasonal_decompose ' + location_col, fontsize=22)
        plt.show()

        plt.figure(figsize=(16, 8))
        plt.subplot(1,2,1).kdeplot(result_mlt.resid, kind='kde', title='Residual Density for ' + location_col )
        plt.show()

        plt.figure(figsize=(16, 8))
        plt.plot(self.data['frame'][location_col]-detrended, label='Detrended ' + location_col)
        plt.plot(self.data['frame'][location_col], label='Normal InCount for ' + location_col)
        plt.show()
        self.data['frame'][location_col] = detrended

    def decompose_STL(self):
        plt.rcParams.update({'figure.figsize': (16,8)})
        res = STL(self.data['frame'][InCountTotal], period=30).fit()

        detrended = self.data['frame'][InCountTotal].values - res.trend - res.seasonal
        res.plot().suptitle('Decompose InCountTotal using STL ', fontsize=22)
        plt.show()

        plt.figure(figsize=(16, 8))
        plt.plot(detrended, label='Detrended data')
        plt.plot(self.data['frame'][InCountTotal], label='Original data')
        plt.title('Detrended and original data using STL')
        plt.show()

    def applyPCA(self,ncomps,varpercent=0):
        # This function takes the datafram on which Principal Component Analysis needs to be done and the % of variance expected to be retained
        # The function then returns the number of optimal principal components and the datafram with those principal components.
        var_cum_variance = np.cumsum(pca.explained_variance_ratio_)
        if((ncomps) & (ncomps>len(pca.components_))):
            pca = PCA(ncomps)
            principalComponents = pca.fit_transform(pcadf)
            var_cum_variance = np.cumsum(pca.explained_variance_ratio_)
```

Mode Command Ln 1, Col 1 Mail_footfall_prediction_skeleton.ipynb

0 2 Python 3 idle ENGLISH 8:50 PM INTL 11/8/2020

JupyterLab

localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

Launcher Mail_footfall_prediction_skeleton.ipynb Python 3

```
def findEuclideanData(self, number_nearest_neighbours='euclidean'):
    odf=pd.DataFrame(self.data['frame'])
    ncol=len(odf.columns)

    # finding the distances based on nearest neighbour to find optimal eps
    nbs = np.zeros((ncol,ncol))
    for i in range(ncol):
        for j in range(ncol):
            nbs[i][j]=np.sqrt(np.sum(np.square(odf[i]-odf[j])))

    nbrs = neigh.kneighbors(nbs, n_neighbors=number_nearest_neighbours, metric='euclidean')
    distances, indices = nbrs.kneighbors(np.array(odf).reshape(-1,ncol))

    maxslope=0
    maxslope0=0
    diff=[0] # A list to hold the differences in the distances
    for i in range(len(distances)):
        distance.append(np.sum(distances[i]))
        distance = np.sort(distance, axis=0)
        diff.append(np.sum(distances[i])-distance[0])
    # If plotted on a graph, a point with maximum curvature, i.e, a point with maximum slope gives the best eps
    maxslope0=0
    for i in range(1, len(diff)-1):
        diff[i] = abs(diff[i]-diff[i-1])
        if (diff[i]>maxslope0):
            maxslope0=diff[i]
            maxslope=indices[i][0]
            maxslopeslope = np.sum(np.diff(indices[i]))/(np.sum(np.diff(distance[i])))
```

pca = PCA(ncomps)
principalComponents = pca.fit_transform(pcadf)
principalPdF = pd.DataFrame(data = principalComponents, columns = ['D1','D2'])

Mode Command Ln 1, Col 1 Mail_footfall_prediction_skeleton.ipynb

0 2 Python 3 idle ENGLISH 8:51 PM INTL 11/8/2020

JupyterLab

localhost:8888/lab

```

File Edit View Run Kernel Tabs Settings Help
+ - + X
Code
Project / Capstone-Project-master / Name
data AML Capstone - Gro... Coverage.xls FullDSCAN.ipynb Mail_footfall.predict_...
preprocess.py Project Plan.xlsx README.md visualization.py work_flow.pptx
Mail_footfall_prediction_skeleton.ipynb
pca = PCA(n_components=2)
principalComponents = pca.fit_transform(odata)
principalDF = pd.DataFrame(data = principalComponents, columns = ['01','02'])

db=DSCAN(eps=neps, min_samples=number_neigh,metric=metrics).fit(np.array(odata).reshape(-1,ncol))
labels=db.labels_

clf = LocalOutlierFactor(n_neighbors=n_neighbors)
predclf.fit_predict(np.array(odata).reshape(1,ncol))
odf['Outliers'] = [1 if val == -1 else 0 for val in labels]

n_clusters_ = len(set(labels)) - (1 if -1 in labels else 0)
n_noise_ = len([l for l in labels if l == -1])
print("DBSCAN Analysis on Full Data")
print("Optimal eps is: %s" % neps)
print("Estimated number of clusters: %d" % n_clusters_)
print("Estimated number of noise points: %d" % n_noise_)

n_clusters_ = len(set(pred))
n_noise_ = len([l for l in pred if l == -1])
print("LOF Analysis on Full Data")
print("Estimated number of clusters: %d" % n_clusters_)
print("Estimated number of noise points: %d" % n_noise_)

#Visualizing the Outliers marked
outlier=[]
outlierIndex=[]
datapoint=[]
dataindex=[]

for j in range(len(labels)):
    if (labels[j]==-1):
        outlier.append((np.sum(np.array(principalDF),loc[j])))
        outlierIndex.append((principalDF.loc[j][['01','02']]))

    else:
        datapoint.append(np.sum(np.array(principalDF),loc[j]))
        datapoint.append((principalDF.loc[j][['01','02']]))

        dataindex.append((principalDF.loc[j][['01','02']]))

fig = plt.figure(figsize=(20,10))
plt.scatter(datapoint,dataindex,color="blue")
plt.scatter(outlier,outlierIndex,color="red")
plt.title("DBSCAN Outliers on Full Data marked in Red")
plt.show()

score = clf.negative_outlier_factor_
print(len(pred),len(principalDF))
for i in range(len(pred)):
    odf['LoF_Outlier']=1 if val == -1 else 0 for val in pred]

```

Mode: Edit | Line: 432, Col: 57 | Mail_footfall_prediction_skeleton.ipynb | ENGLISH | 8:52 PM | INTL | 11/8/2020

JupyterLab

localhost:8888/lab

```

File Edit View Run Kernel Tabs Settings Help
+ - + X
Code
Project / Capstone-Project-master / Name
data AML Capstone - Gro... Coverage.xls FullDSCAN.ipynb Mail_footfall.predict_...
preprocess.py Project Plan.xlsx README.md visualization.py work_flow.pptx
Mail_footfall_prediction_skeleton.ipynb
else:
    datapoint.append(np.sum(np.array(principalDF),loc[j]))
    datapoint.append((principalDF.loc[j][['01','02']]))

    dataindex.append((principalDF.loc[j][['01','02']]))

fig = plt.figure(figsize=(20,10))
plt.scatter(datapoint,dataindex,color="blue")
plt.scatter(outlier,outlierIndex,color="red")
plt.title("DBSCAN Outliers on Full Data marked in Red")
plt.show()

score = clf.negative_outlier_factor_
print(len(pred),len(principalDF))
for i in range(len(pred)):
    # print(i, pred[i])
    odf['LoF_Outlier']=1 if val == -1 else 0 for val in pred]

#Visualizing the Outliers marked
outlier=[]
outlierIndex=[]
datapoint=[]
dataindex=[]

for j in range(len(odf['LoF_Outlier'])):
    if (odf['LoF_Outlier'][j]==-1):
        outlier.append((np.sum(np.array(principalDF),loc[j])))
        outlierIndex.append((principalDF.loc[j][['01','02']]))

    else:
        datapoint.append(np.sum(np.array(principalDF),loc[j]))
        datapoint.append((principalDF.loc[j][['01','02']]))

        dataindex.append((principalDF.loc[j][['01','02']]))

fig = plt.figure(figsize=(20,10))
plt.scatter(datapoint,dataindex,color="blue")
plt.scatter(outlier,outlierIndex,color="red")
plt.title("LOF Outliers on Full Data marked in Red")
plt.show()

return odf

```

4.2.5 Do_All_Pre_Processing Function

```

[15]: def do_all_pre_processing(for_visualization_only = False):
    global df_easter_sunday_holiday_data, df_university_long_holidays_data, df_school_long_holidays_data, df_UK_bank_holidays_data
    principalDF = []
    countries = []
    countries.append('US')

```

Mode: Edit | Line: 432, Col: 57 | Mail_footfall_prediction_skeleton.ipynb | ENGLISH | 8:52 PM | INTL | 11/8/2020

JupyterLab

localhost:8888/lab

File Edit View Run Kernel Help

Launcher Mail_footfall_prediction_ske

Python 3

4.2.5 Do_All_Pre_Processing Function

```
[15]: def do_all_pre_processing(for_visualization_only = False):
    global df_easter_sunday_holiday_data, df_university_long_holidays_data, df_school_long_holidays_data, df_UK_bank_holidays_data
    print(df_easter_sunday_holiday_data)
    print(df_university_long_holidays_data)
    print(df_school_long_holidays_data)
    print(df_UK_bank_holidays_data)

    var_covariance = []
    # process the main data set
    pre_process_footfall = pre_processing(df_footfall)

    pre_process_footfall.reshape_location_name(['Date','Hour'])
    pre_process_footfall.data_frame = pre_process_footfall.inputData.read()

    pre_process_footfall.aggregate('Date')
    # add the InCountTotal column for the total mull's Incount
    pre_process_footfall.addTotalIncountColumn()

    data_prophet = pre_process_footfall.data_frame[['Date','InCountTotal']].copy()

    pre_process_footfall.createStationarySeries()

    if for_visualization_only == False:
        # Normalize the InCount columns
        pre_process_footfall.normalizeTargetVariable()

        # Merge weather data and foot fall data based on date
        data_weather = pre_process_footfall.merge_data(df_weather_data, col_by='Date')

        # Add new column for Easter Sunday
        pre_process_easter_sunday_data = pre_processing(df_easter_sunday_holiday_data)
        df_easter_Sunday_holiday_data = pre_process_easter_sunday_data.pre_process(data.frame_type='EasterSunday')
        pre_process_footfall.add_new_column(df_easter_sunday_holiday_data,'EasterSundayHoliday', destination_date_format = '%d-%b-%Y' )

        # add a new column for university holidays
        pre_process_footfall.add_new_column(df_university_long_holidays_data,'University_Holidays', destination_date_format = '%d-%b-%Y' )

        # add a new column for School holidays
        pre_process_footfall.add_new_column(df_school_long_holidays_data,'School_Holidays', destination_date_format = '%d-%b-%Y' )

        # add a new column for UK Bank Holidays
        pre_process_footfall.add_new_column(df_UK_bank_holidays_data,'UKBankHoliday', destination_date_format = '%d-%b-%Y' )

        #create extra columns from the date
        pre_process_footfall.createExtraDateColumns()
        # Create time series
        pre_process_footfall.createTimeSeries()
```

JupyterLab

localhost:8888/lab

File Edit View Run Kernel Table Settings Help

Launcher Code Python 3

Mail_footfall_prediction_skeleton.ipynb

```
print('Number of Principal components are', n_pca_comp)
#Add the target variable to PCA datframe
prinadf['InCountTotal'] = data_final['InCountTotal_forecast1'].values
pre_process_footfall.decomposeTrend()
pre_process_footfall.decomposeSeasonal()
pre_process_footfall.decomposeTrend()
data_final_detrended = pre_process_footfall.data_frame_copy()

data_outliersremoved = data_final[data_final.columns[[data_final['Outliers_InCountTotal']==0]]]
data_outliersremoved.to_csv("outliers_removed.csv")
data_outliers.to_csv("outlier_data.csv")
data_final_detrended.to_csv("detrended_data.csv")
prinadf.to_csv("pca_features.csv")

data_weekly = pre_process_footfall.createWeeklyData(data_final)
data_weekly.to_csv("weekly_data.csv")
```

return data_final, data_final_detrended, data_outliersremoved, data_prophet, data_weekly, location_name_list, pre_process_footfall.scaler_target, pre_process_footfall.scaler_input, pre_process_footfall.scaler_target_each_location, scaled_prophet, prinadf, compdf, var_cum_variance = do_all_pre_processing()

[16]:

```
data_final, data_final_detrended, data_outliersremoved, data_prophet, data_weekly, location_name_list, scaler_target, scaler_input, scaler_target_each_location, scaled_prophet, prinadf, compdf, var_cum_variance = do_all_pre_processing()

ADF test results for InCountTotal
ADF Statistic: -7.638551
p-value: 1.921495340734742e-11
full result: (-7.638551478845224, 1.92149534300742e-11, 28, 2869, {'1%': -3.4326313386809552, '5%': -2.8625479402718623, '10%': -2.5673865559730883), 70840.73050705886)

ADF test results for InCountBriggate
ADF Statistic: -7.780988
p-value: 8.160138000516e-11
full result: (-7.78098845691125, 1.3401466510758263e-11, 27, 2870, {'1%': -3.43263095438048927, '5%': -2.8625475888928142, '10%': -2.567386368909217), 61610.27124697735)

ADF test results for InCountBriggateAtMCDs
ADF Statistic: -6.974079
p-value: 8.5170159705347e-10
full result: (-6.974079038567212, 8.517015973653547e-10, 28, 2869, {'1%': -3.4326313386809552, '5%': -2.8625479402718623, '10%': -2.5673865559730883), 56063.748851183336)

ADF test results for InCountAlbionStNorth
ADF Statistic: -8.167221
p-value: 8.8167513000516e-11
full result: (-8.16722147111441, 8.8167513000516e-11, 28, 2869, {'1%': -3.4326313386809552, '5%': -2.8625479402718623, '10%': -2.5673865559730883), 57609.079992860838)

ADF test results for InCountAlbionStSouth
ADF Statistic: -8.534749
```

Mode Command L 432 C 57 Mail_footfall_prediction_skeleton.ipnb EN 8:55 PM INTL 11/8/2020

JupyterLab

localhost:8888/lab

File Edit View Run Kernel Table Settings Help

Launcher Code Python 3

Mail_footfall_prediction_skeleton.ipynb

```
full result: (-8.46720991711141, 8.8167513000516e-13, 28, 2869, {'1%': -3.4326313386809552, '5%': -2.8625479402718623, '10%': -2.5673865559730883), 57609.079992860838)

ADF test results for InCountAlbionStSouth
ADF Statistic: -8.534749
p-value: 1.01510459792997e-13
full result: (-8.53474903633883, 1.01510459792997e-13, 27, 2870, {'1%': -3.43263095438048927, '5%': -2.8625475888928142, '10%': -2.567386368909217), 61804.4250888043)

ADF test results for InCountCommercialStLush
ADF Statistic: -8.598819
p-value: 8.458490722e-14
full result: (-8.59881900040518, 8.458490722e-14, 27, 2870, {'1%': -3.43263095438048927, '5%': -2.8625475888928142, '10%': -2.567386368909217), 60241.799302121464)

ADF test results for InCountMeadow
ADF Statistic: -7.933177
p-value: 3.469806056456294e-12
full result: (-7.93317700879057, 3.469806056456294e-12, 27, 2870, {'1%': -3.43263095438048927, '5%': -2.8625475888928142, '10%': -2.567386368909217), 56877.64679189137)

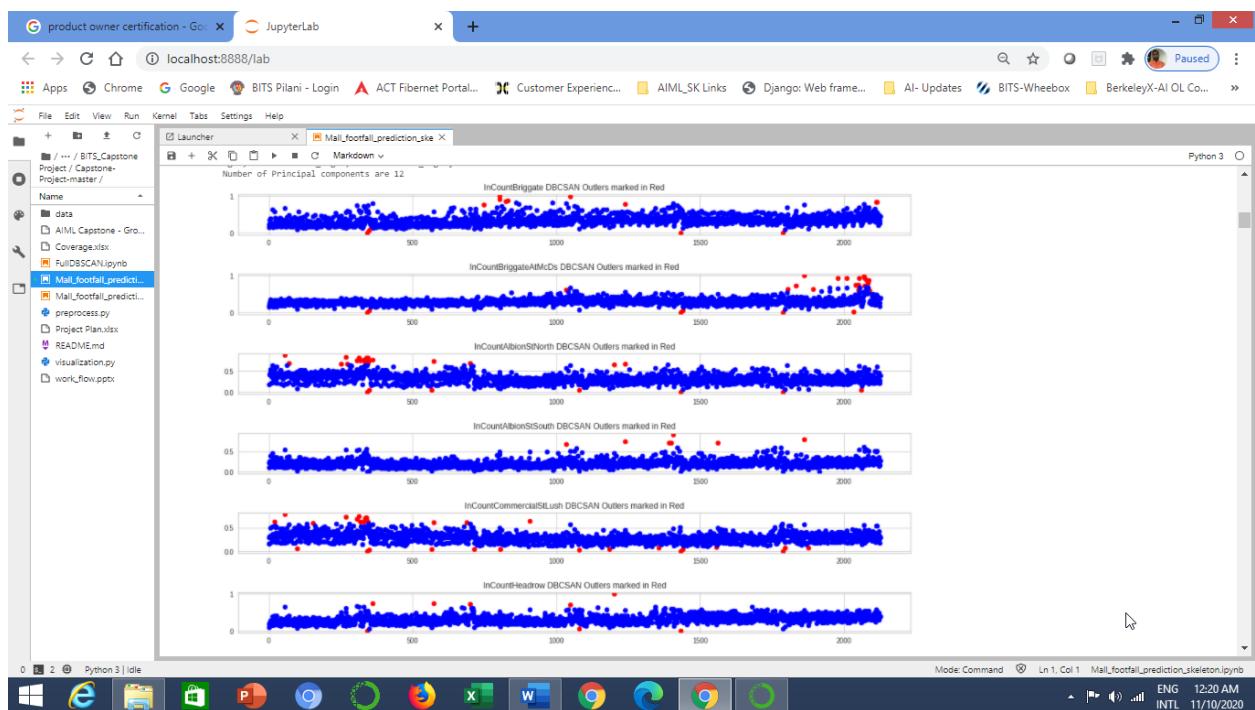
ADF test results for InCountOrdnance
ADF Statistic: -7.714764
p-value: 1.2366744186966634e-11
full result: (-7.71476425929241, 1.2366744186966634e-11, 26, 2871, {'1%': -3.43263095438048927, '5%': -2.86254723775867, '10%': -2.567386181959156), 59389.320373311576)

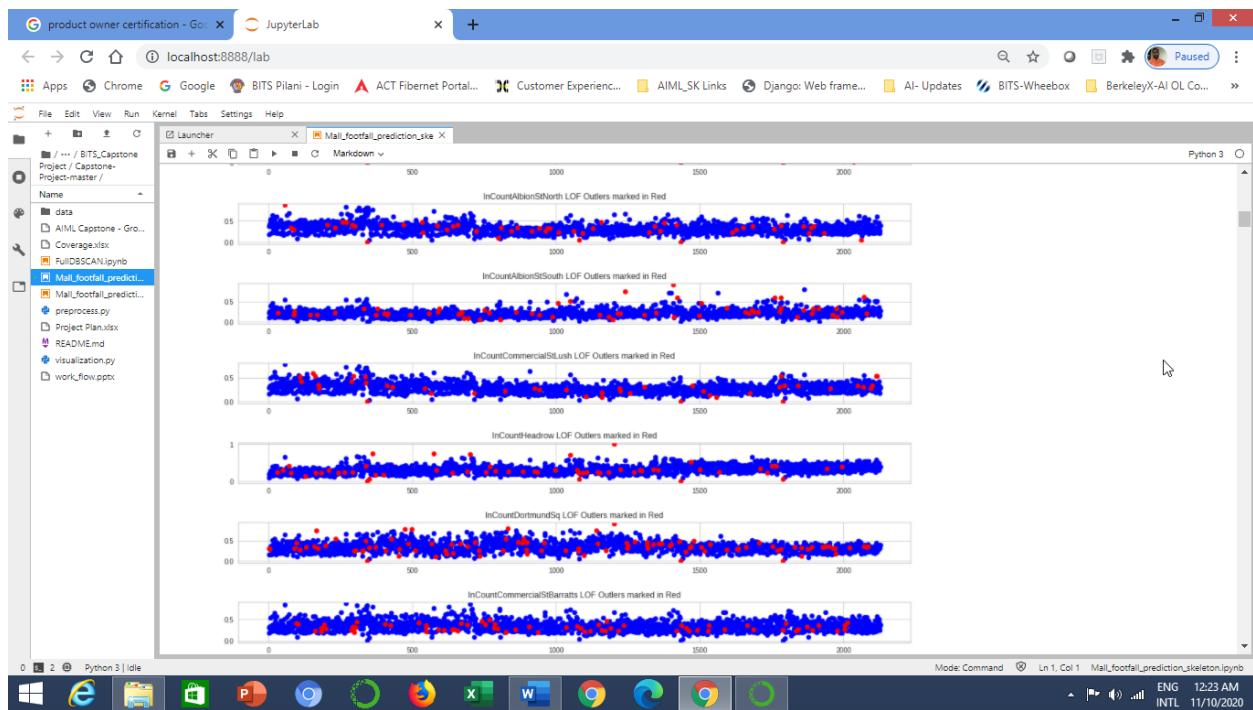
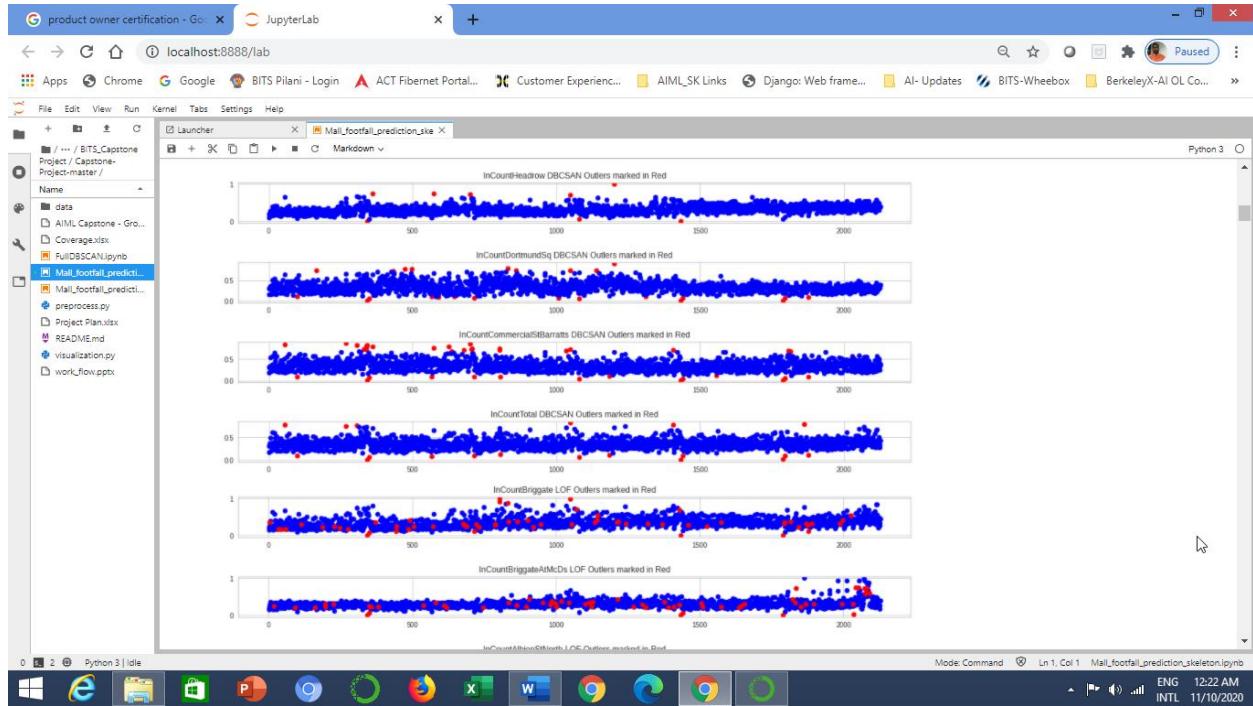
ADF test results for InCountCommercialStBarrett
ADF Statistic: -7.981598
p-value: 1.213514220981595e-13
full result: (-8.981598402192032, 1.213514220981595e-13, 28, 2869, {'1%': -3.4326313386809552, '5%': -2.8625479402718623, '10%': -2.5673865559730883), 58255.70110518673)

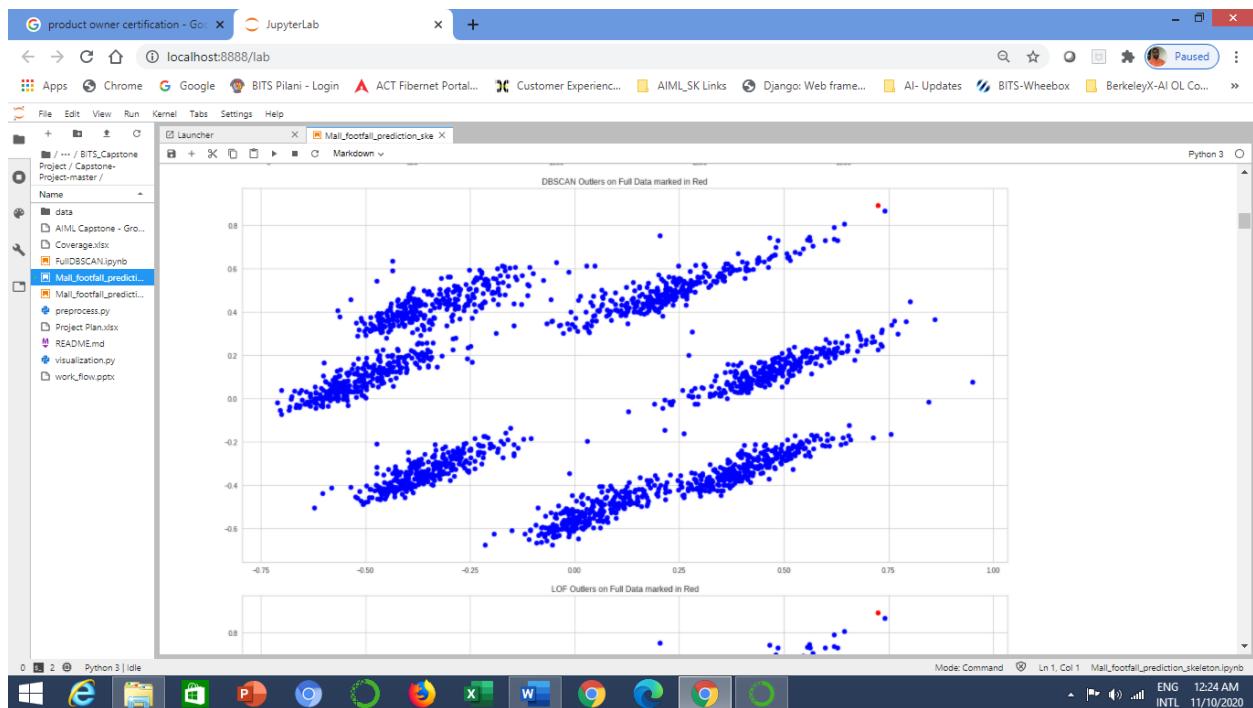
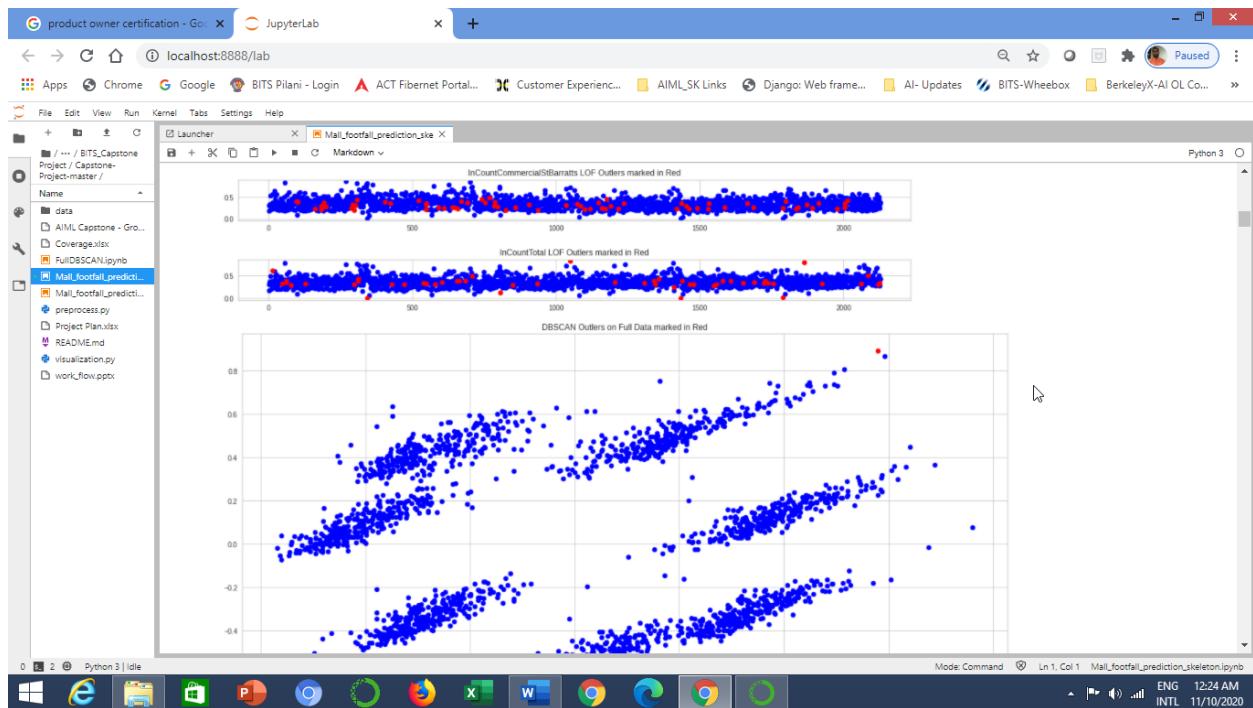
dt_object: 2018-04-01 00:00:00
type(dt_object): <class 'datetime.datetime'>
InCountBriggettedAnalysis
Optimal eps is 0.04996958795940593
Estimated number of clusters: 3
InCountBriggettedAnalysis
Optimal eps is 0.0371521849459006
Estimated number of clusters: 2
Estimated number of noise points: 27
InCountCommercialStLush
Optimal eps is 0.0083709522554644
Estimated number of clusters: 4
Estimated number of noise points: 28
InCountAlbionStSouth
Optimal eps is 0.04996958795940593
Estimated number of clusters: 1
Estimated number of noise points: 7
InCountCommercialStLushAnalysis
Optimal eps is 0.01803696810462136
Estimated number of clusters: 1
Estimated number of noise points: 33
InCountMeadowAnalysis
```

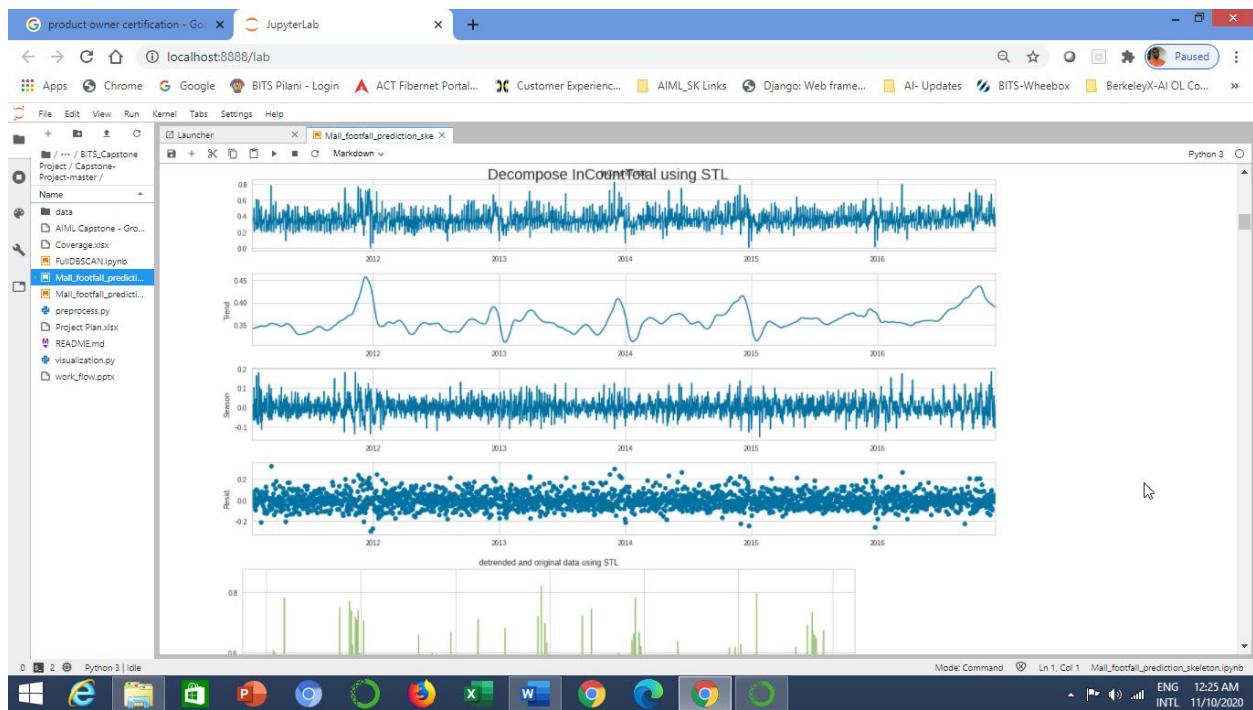
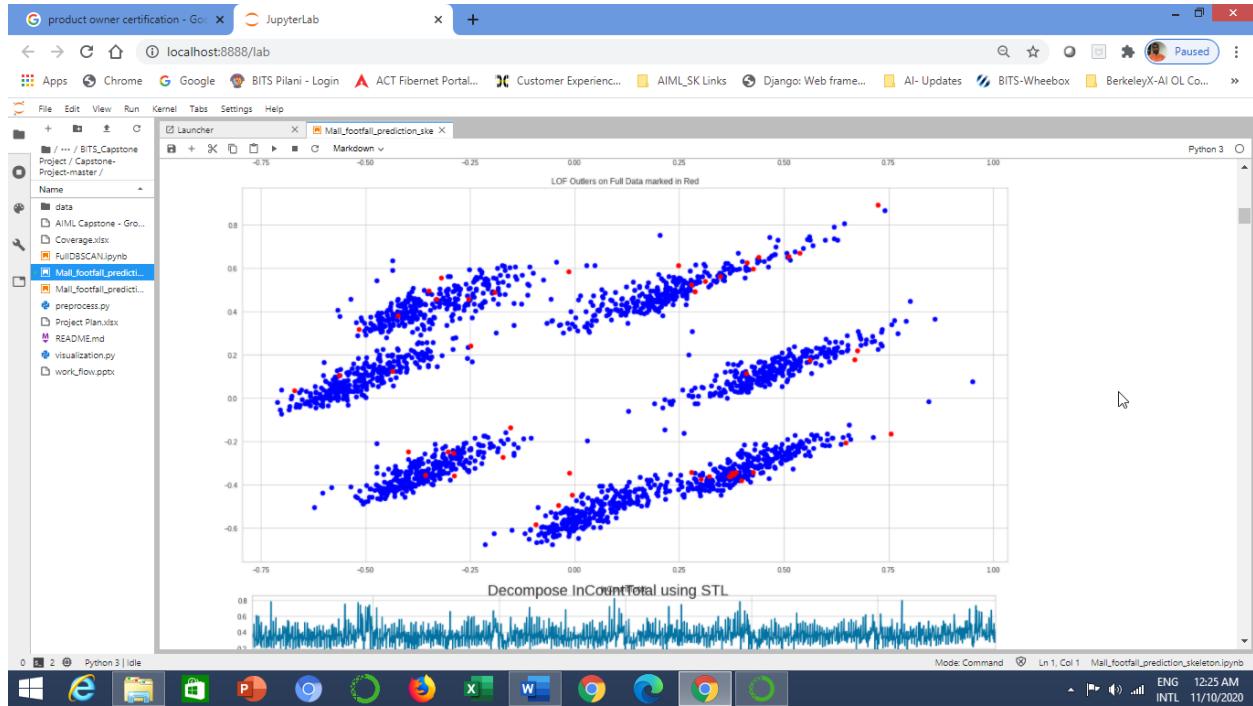
Mode Command L 432 C 57 Mail_footfall_prediction_skeleton.ipnb EN 8:57 PM INTL 11/8/2020

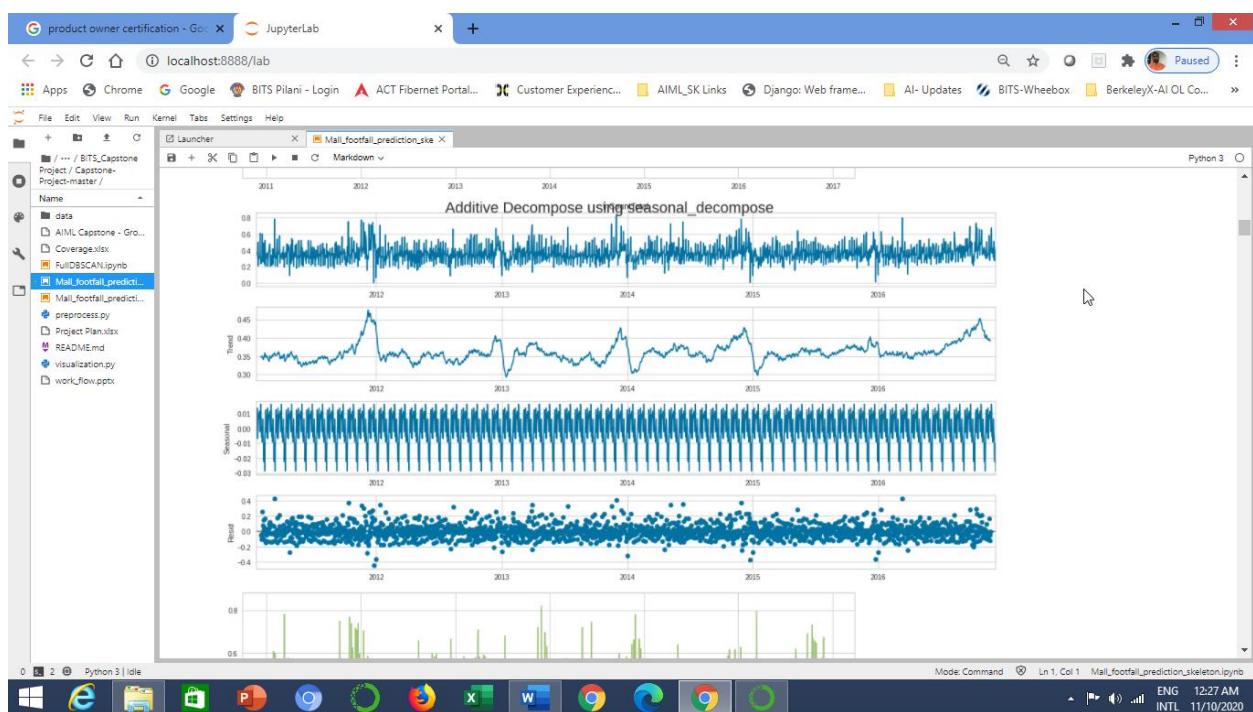
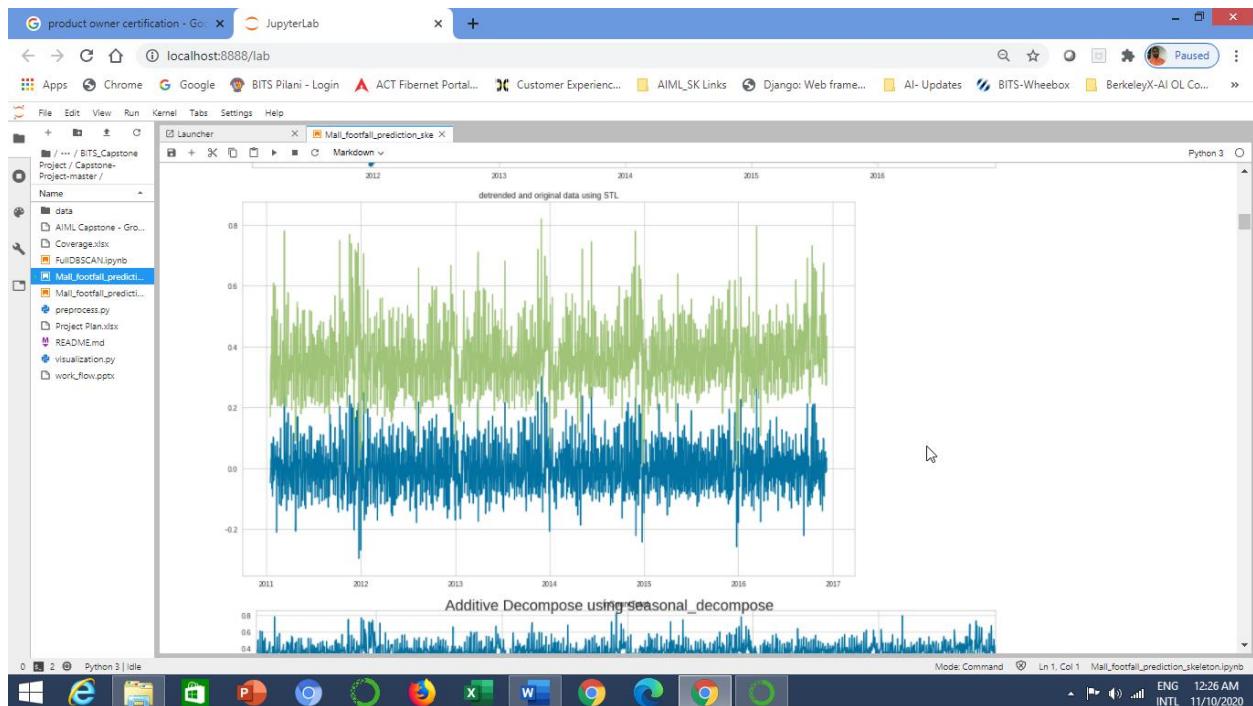
The screenshot shows a Jupyter Notebook interface with two plots. The top plot is titled "InCountBriagte DBSCAN Outliers marked in Red" and the bottom plot is titled "InCountBriagteAIMDs DBSCAN Outliers marked in Red". Both plots show a scatter of blue points with several red points highlighted as outliers. The x-axis for both plots ranges from 0 to 3000. The y-axis ranges from 0 to 1. The plots are overlaid on a background of a large number of blue points.

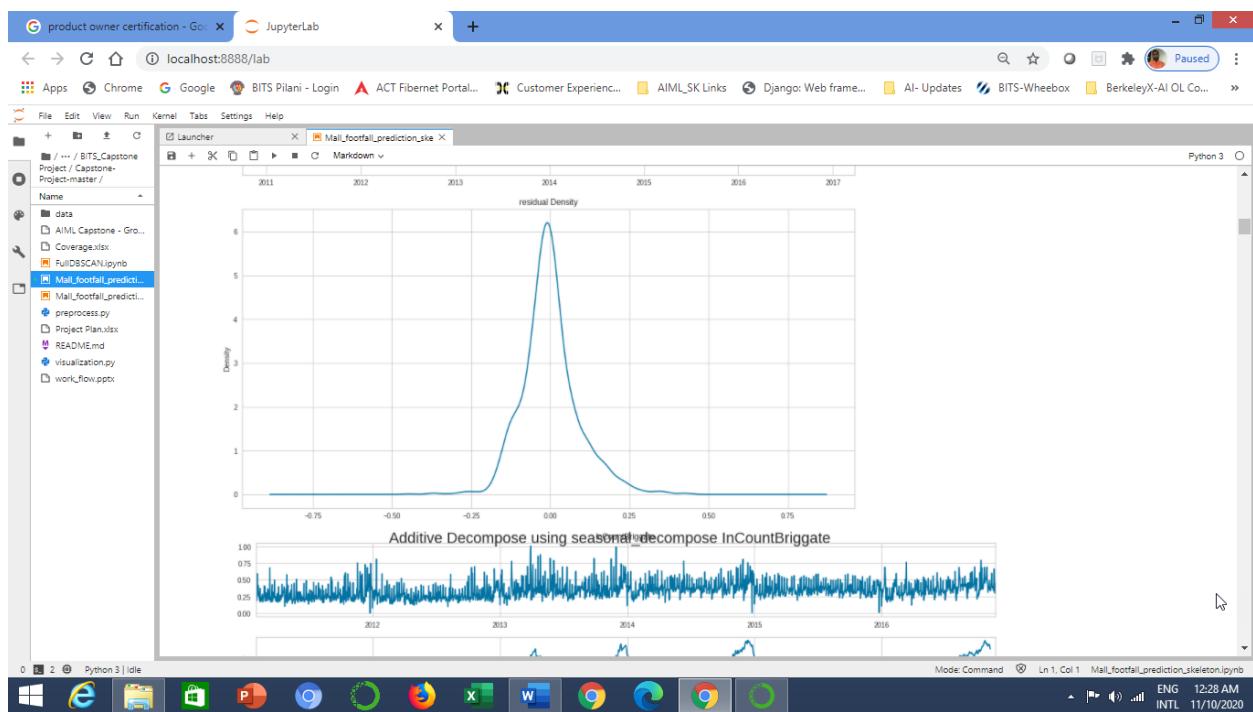
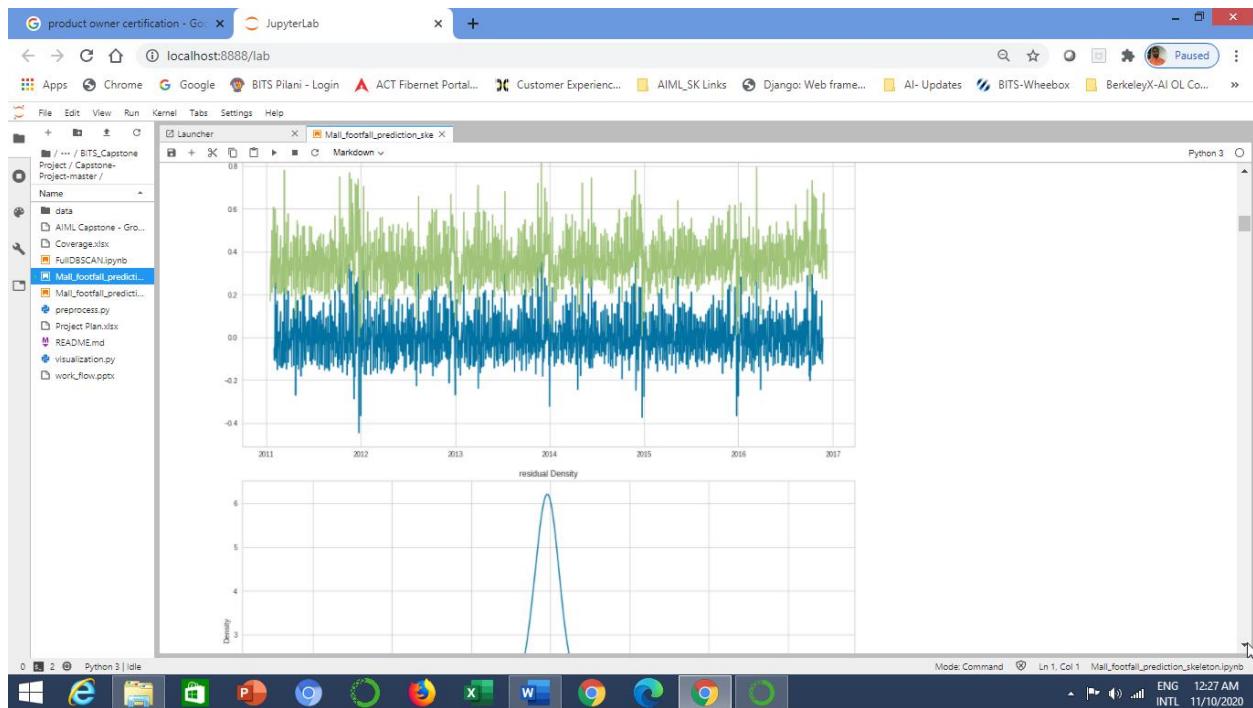


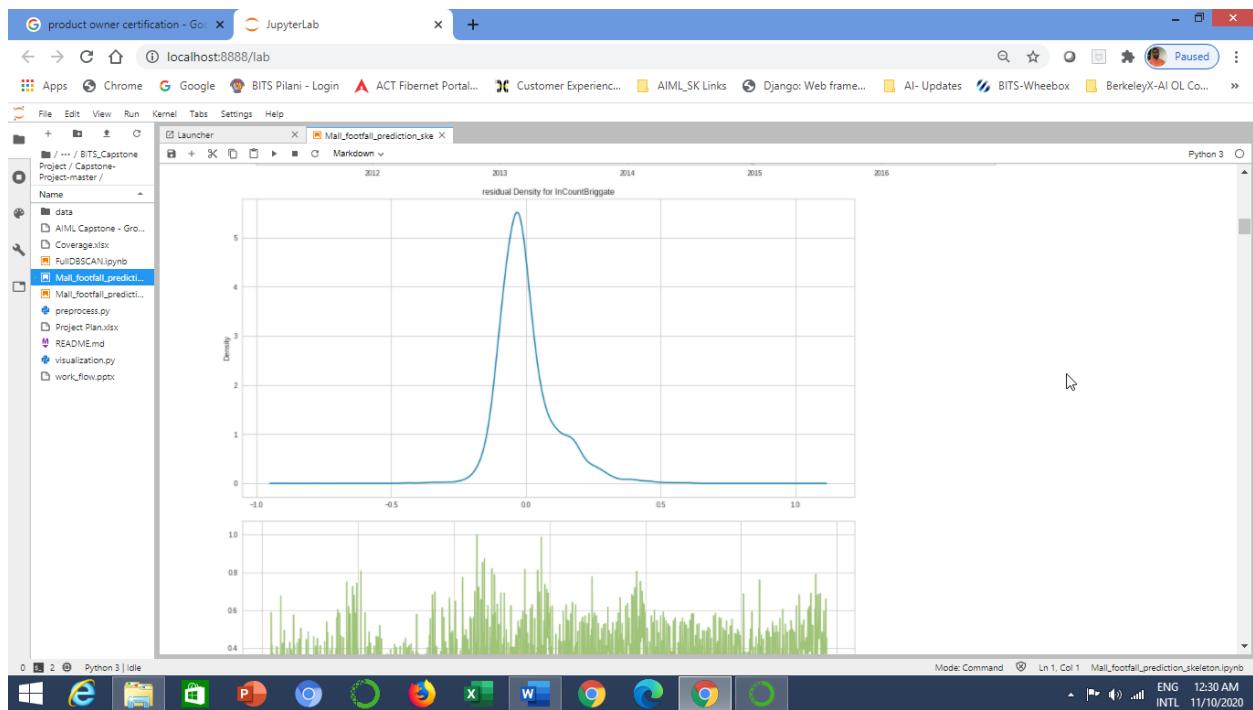
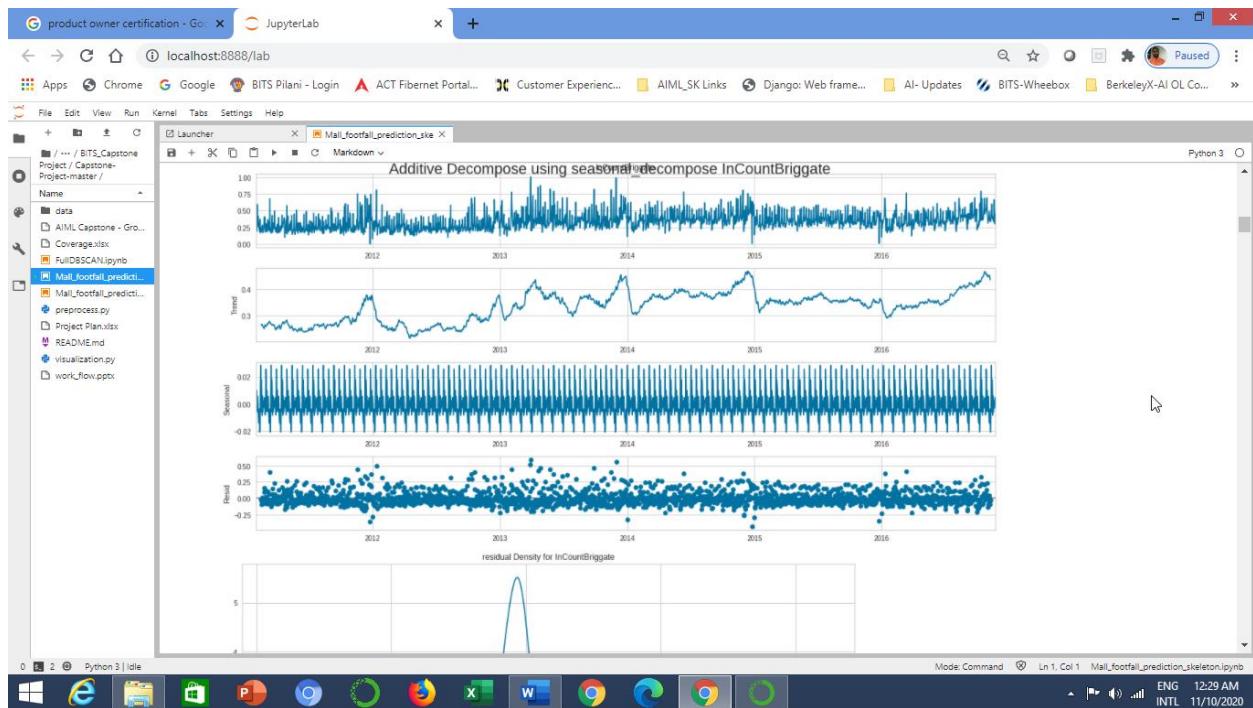


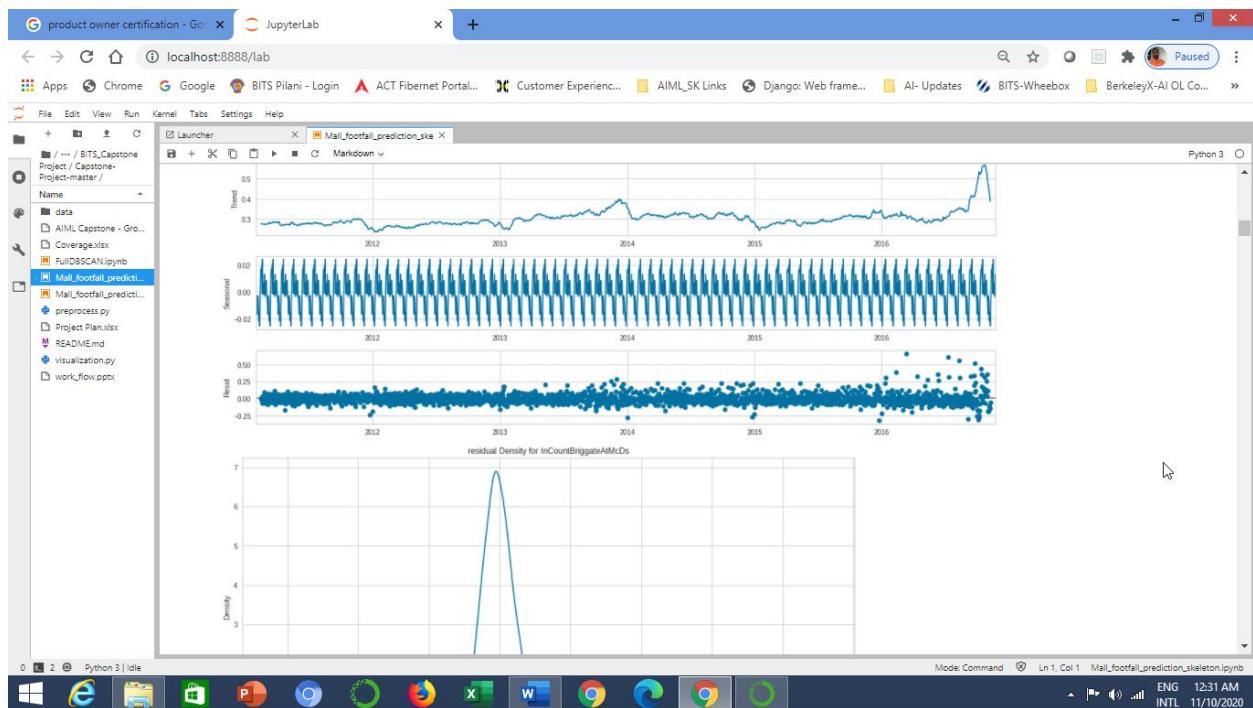
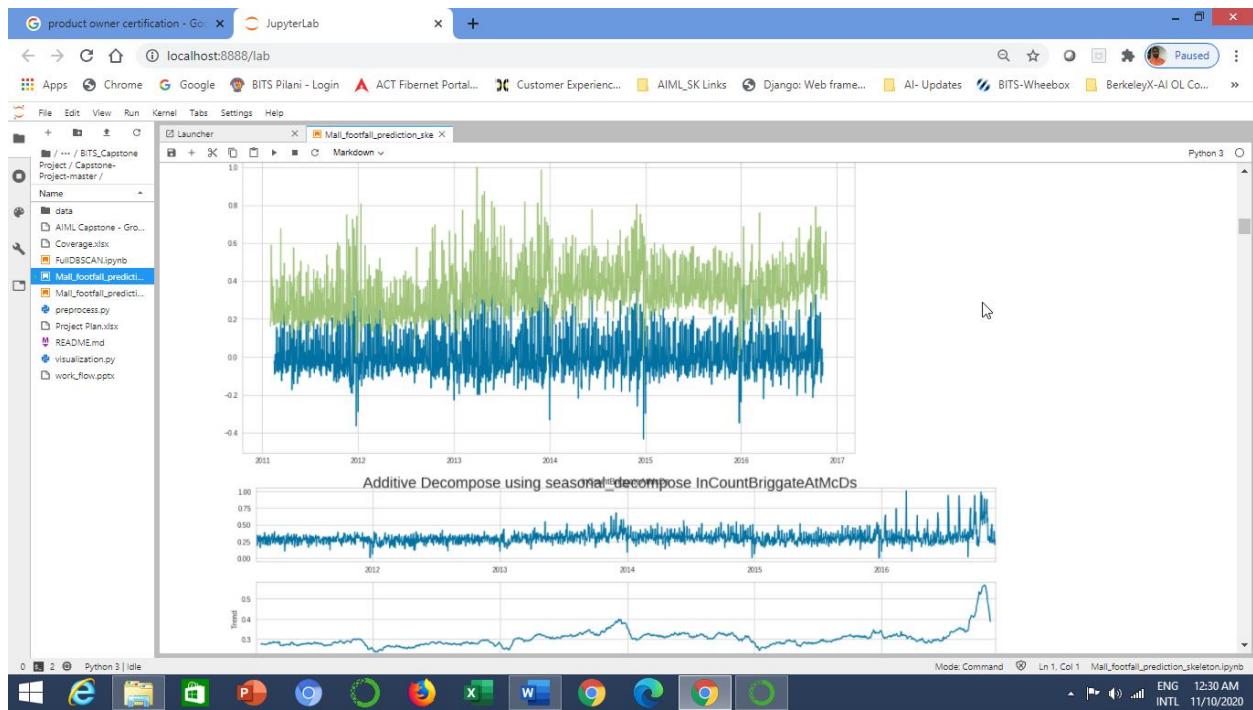


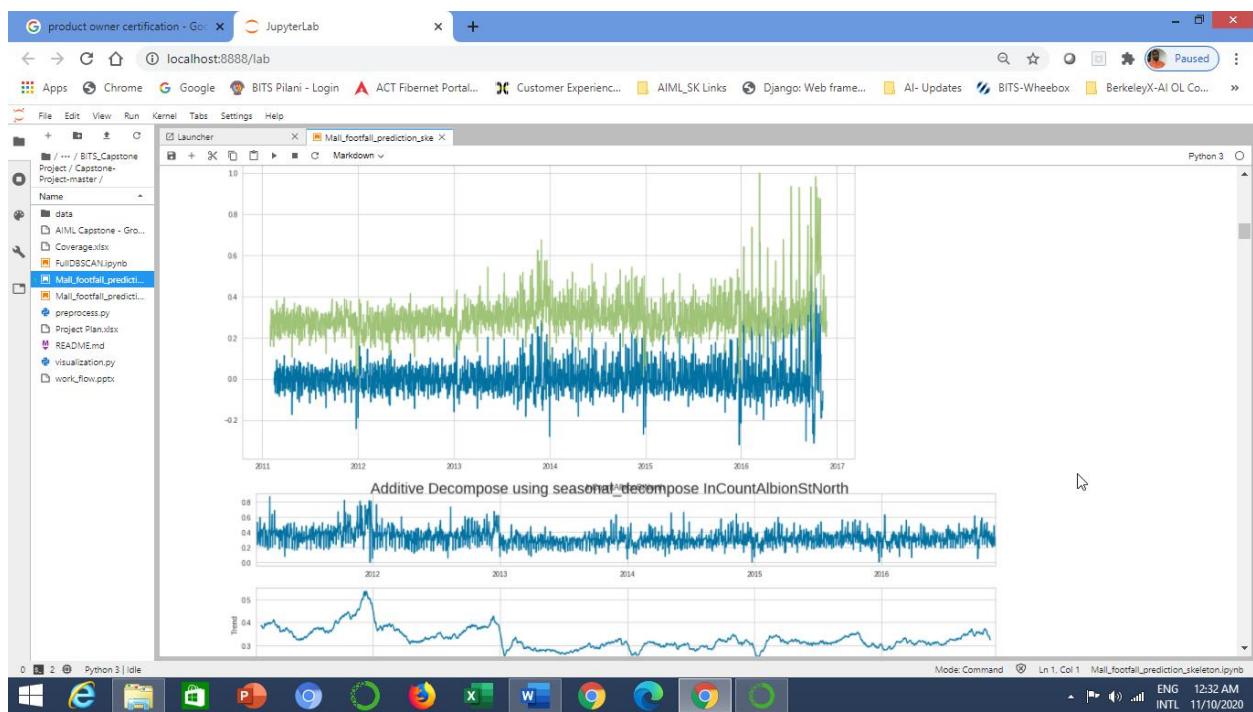
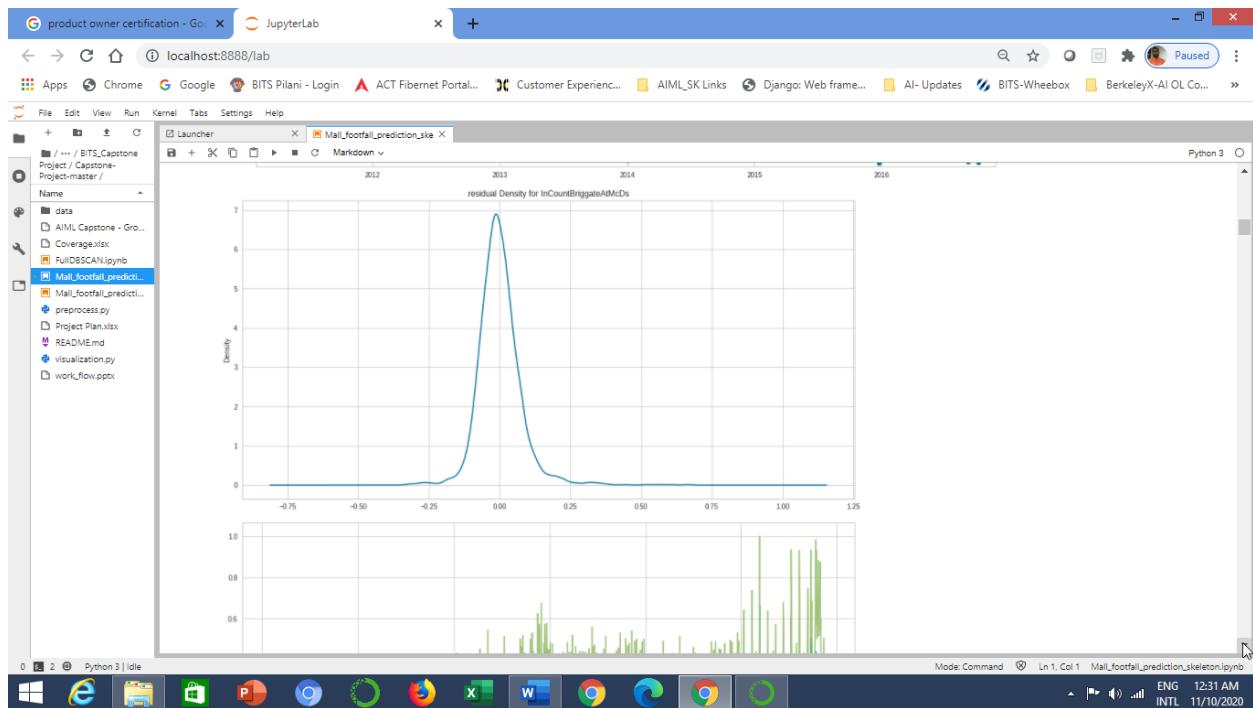


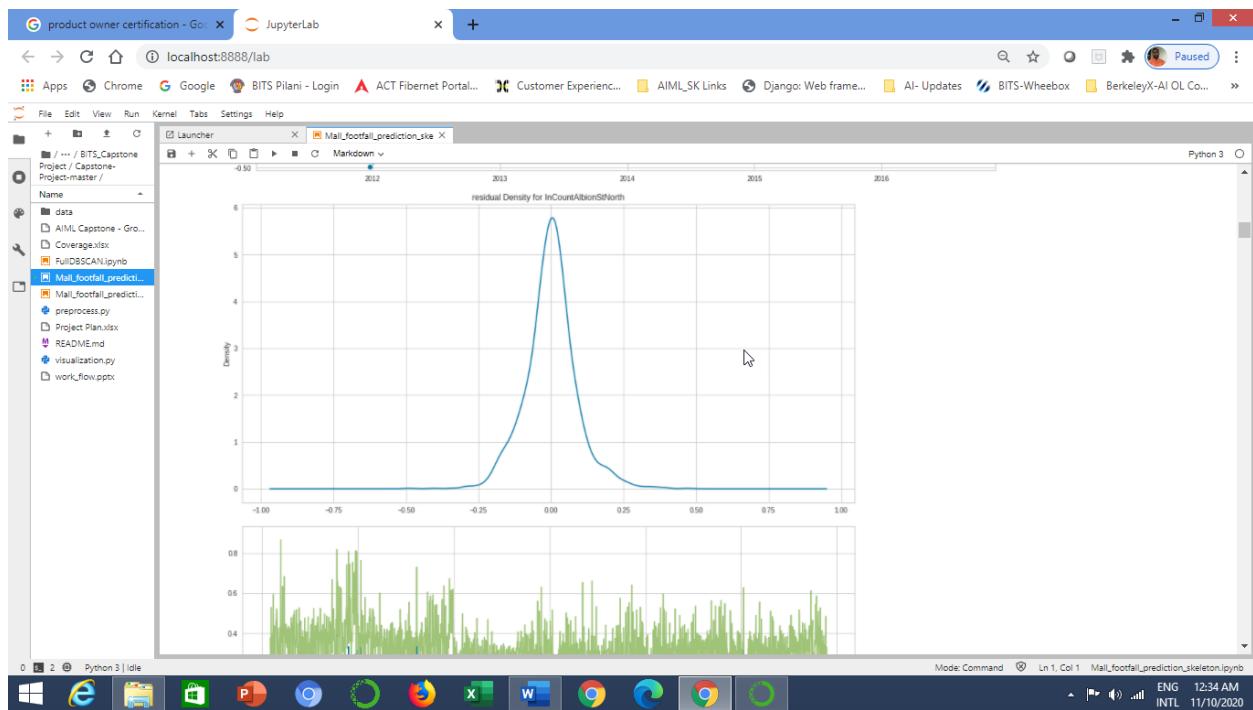
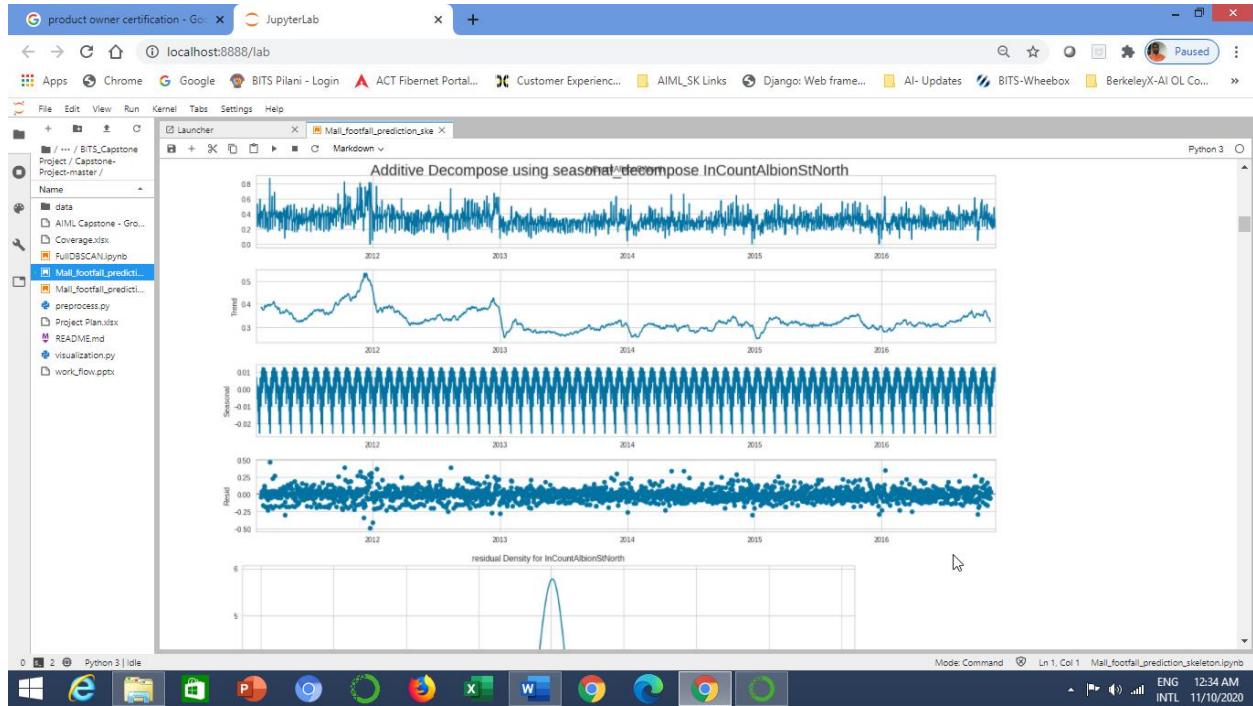


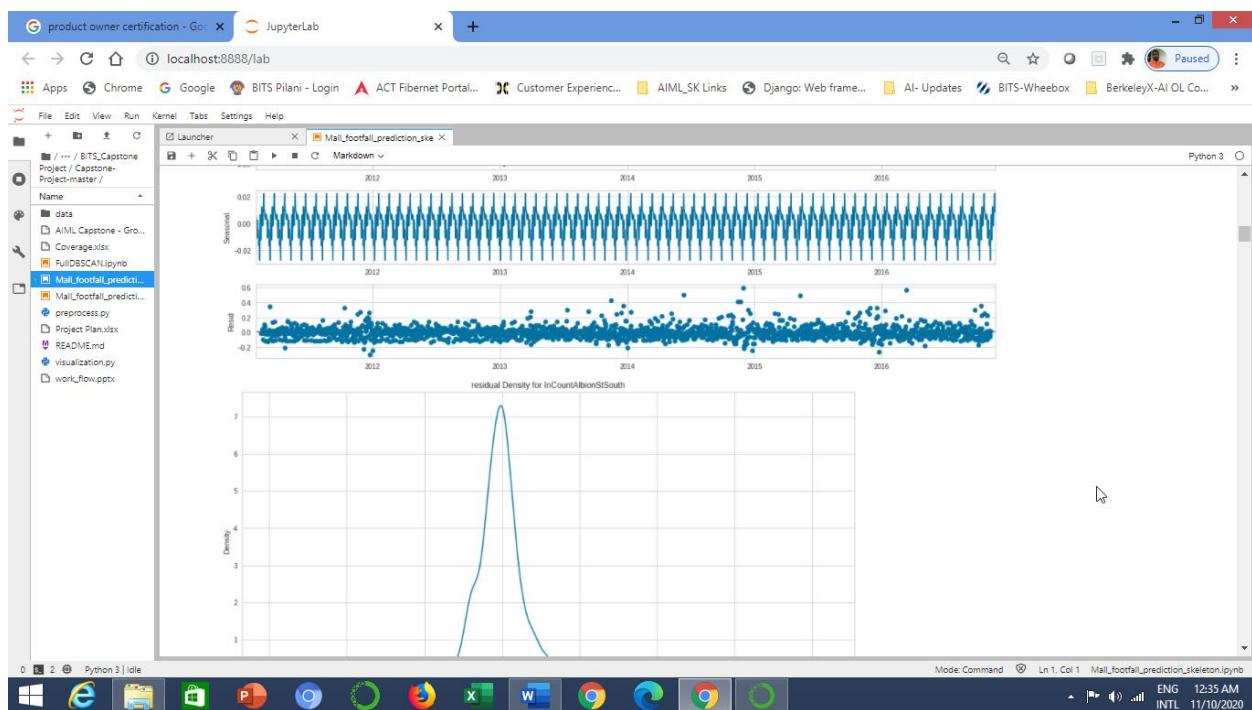
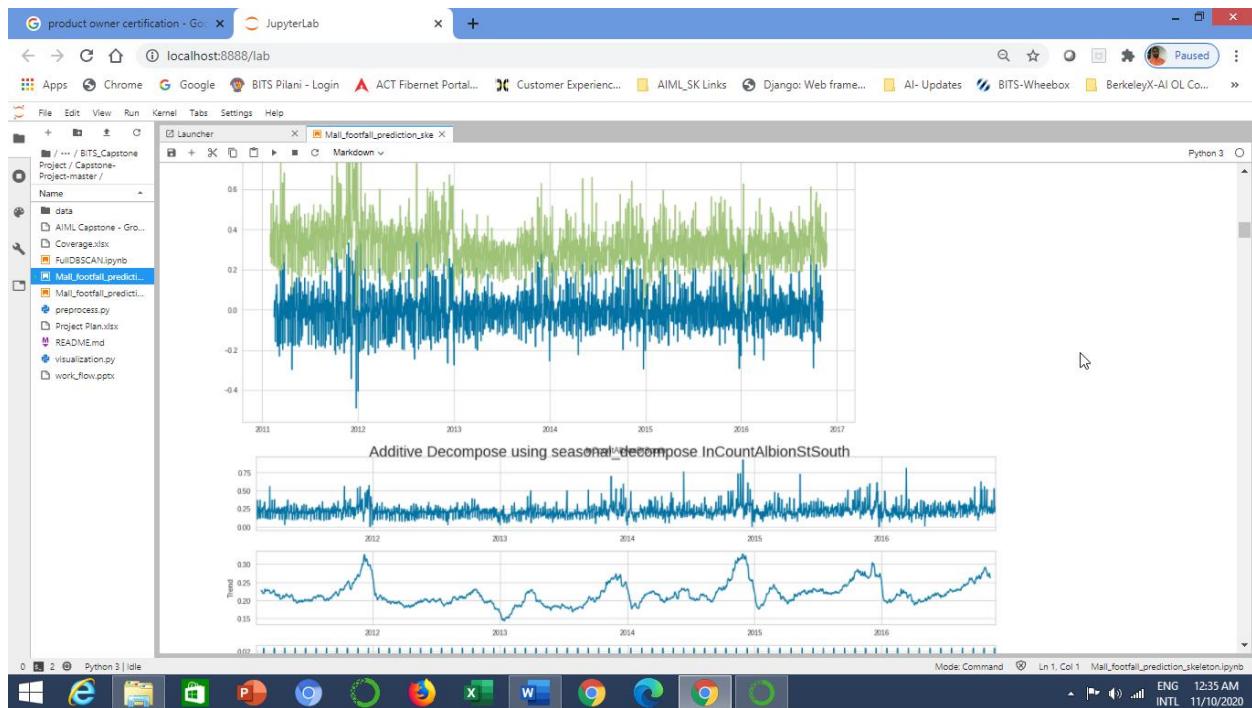


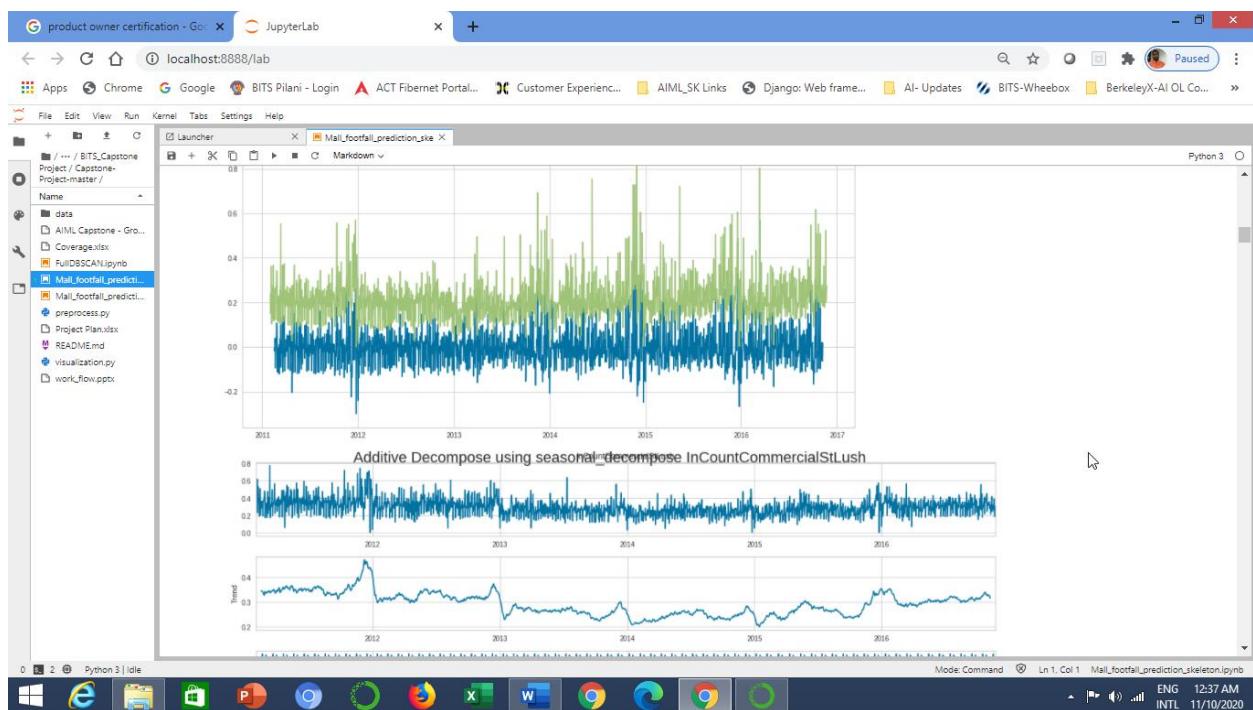
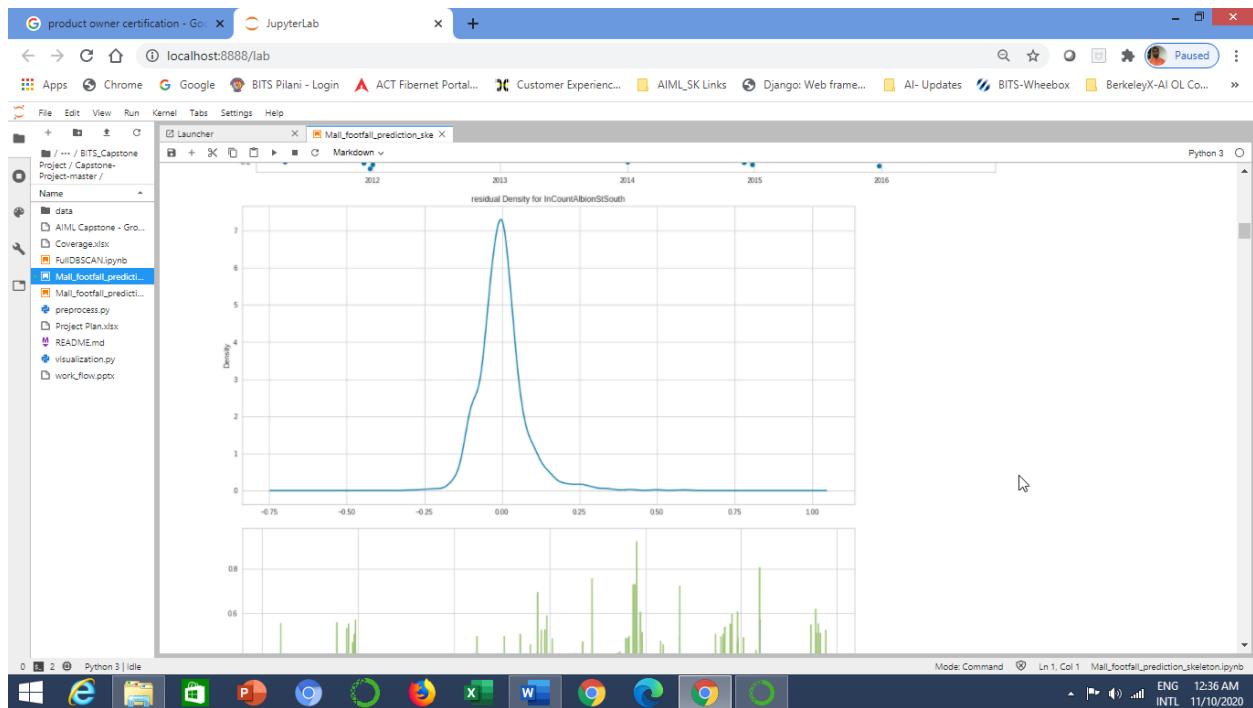


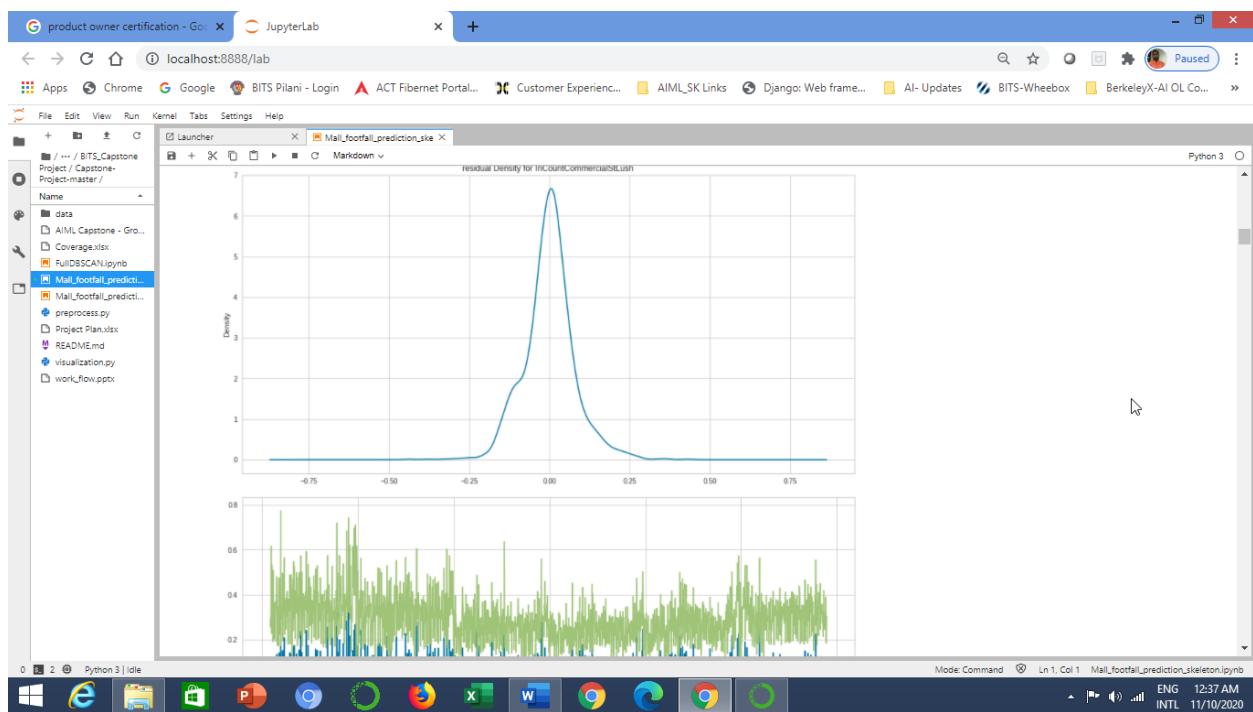
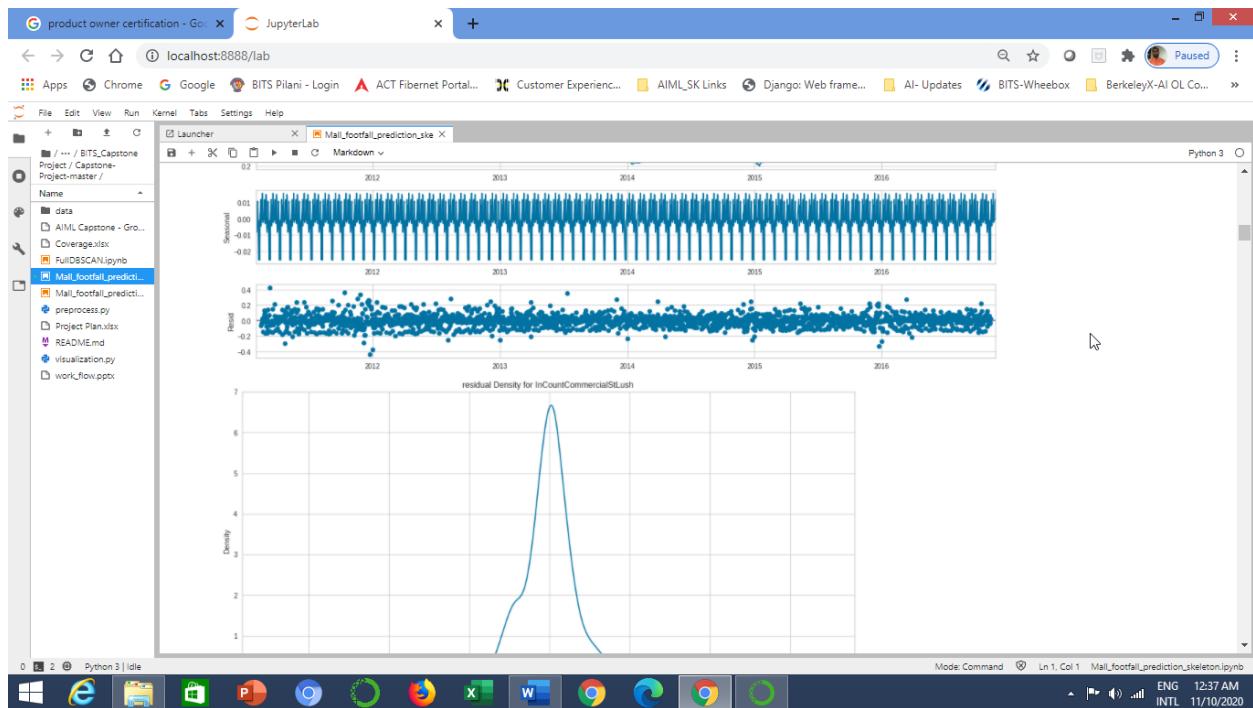


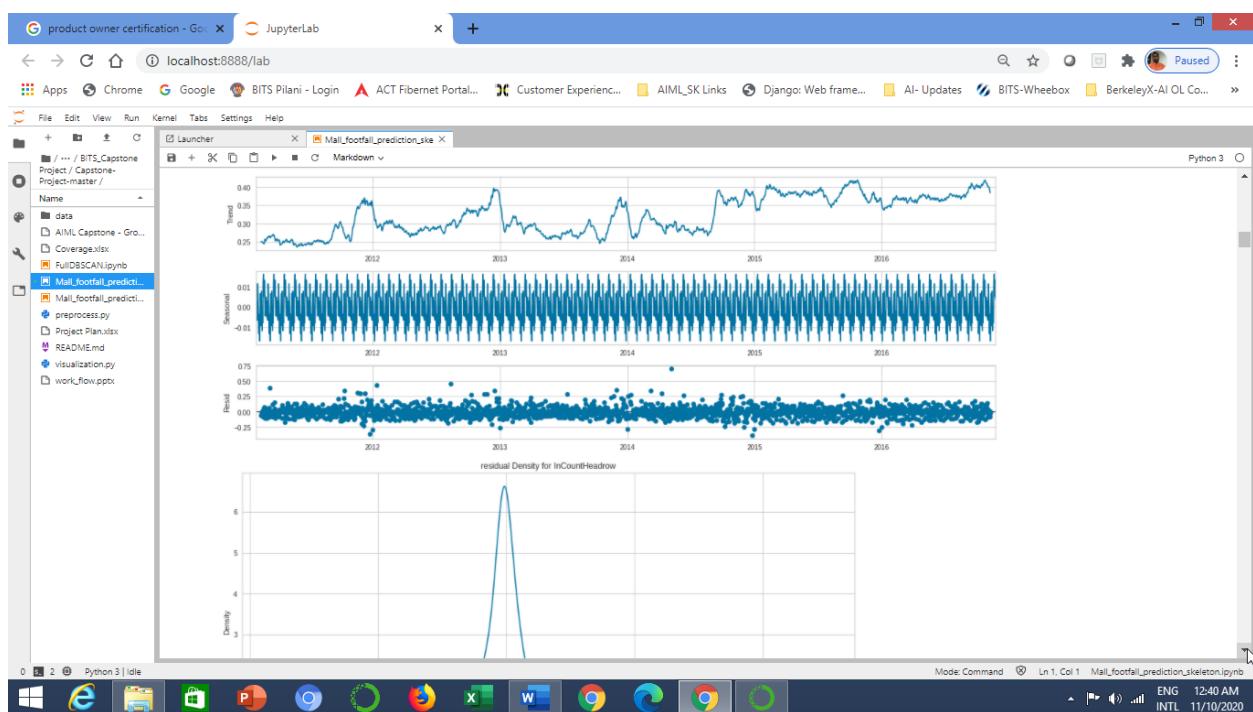
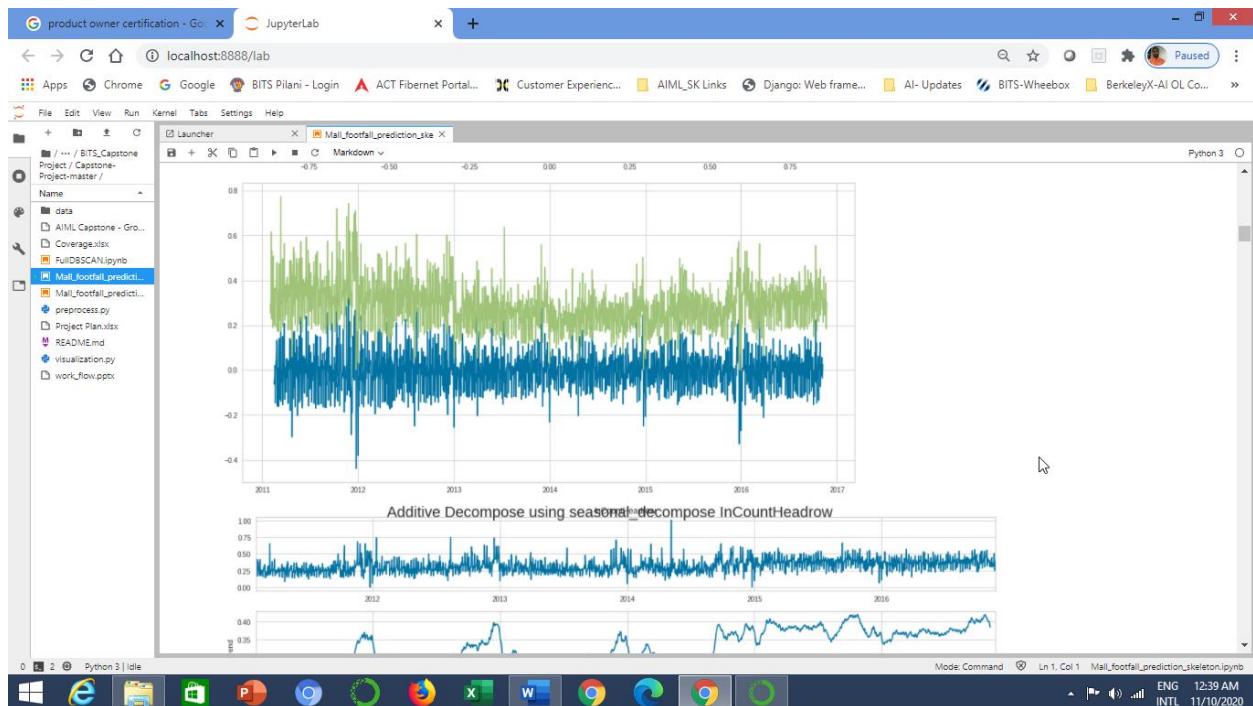


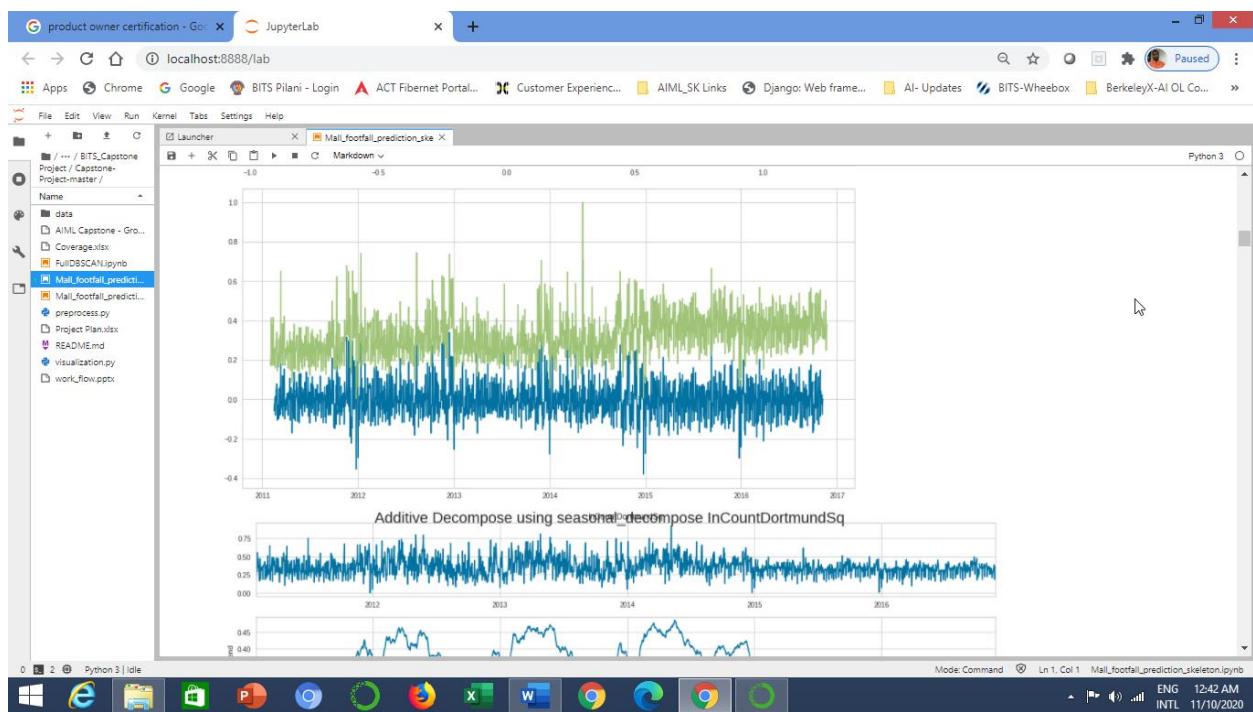
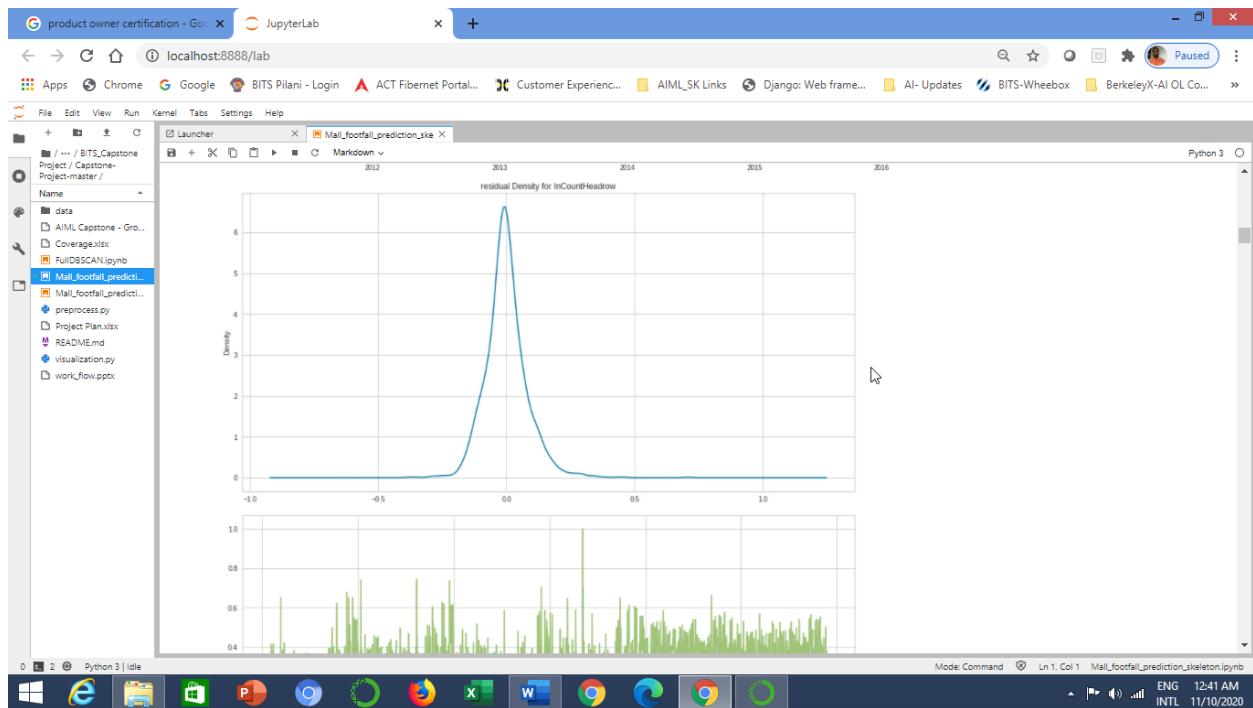


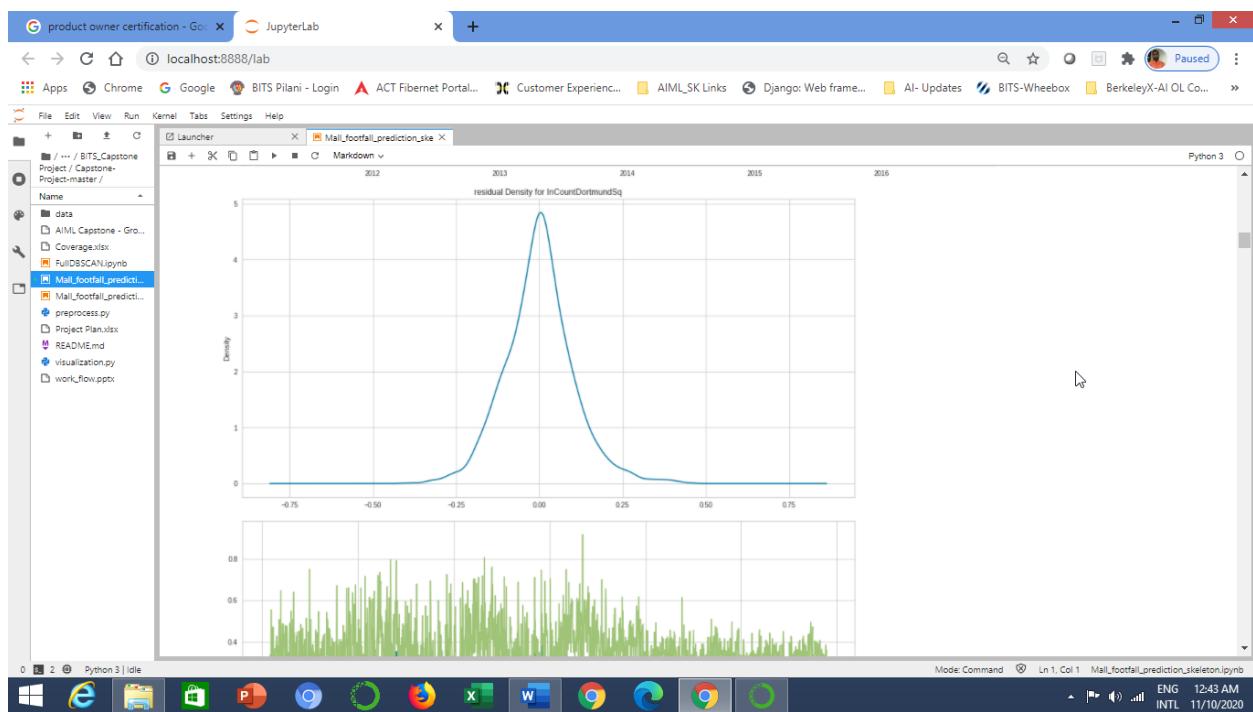
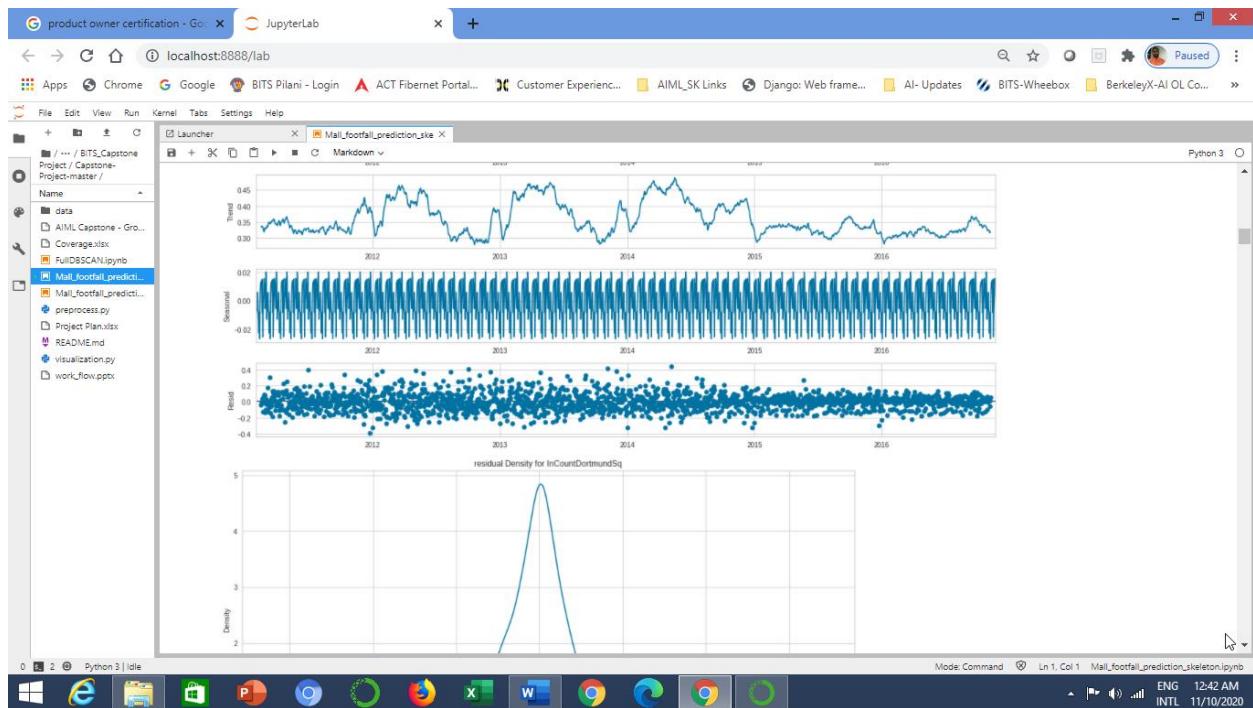


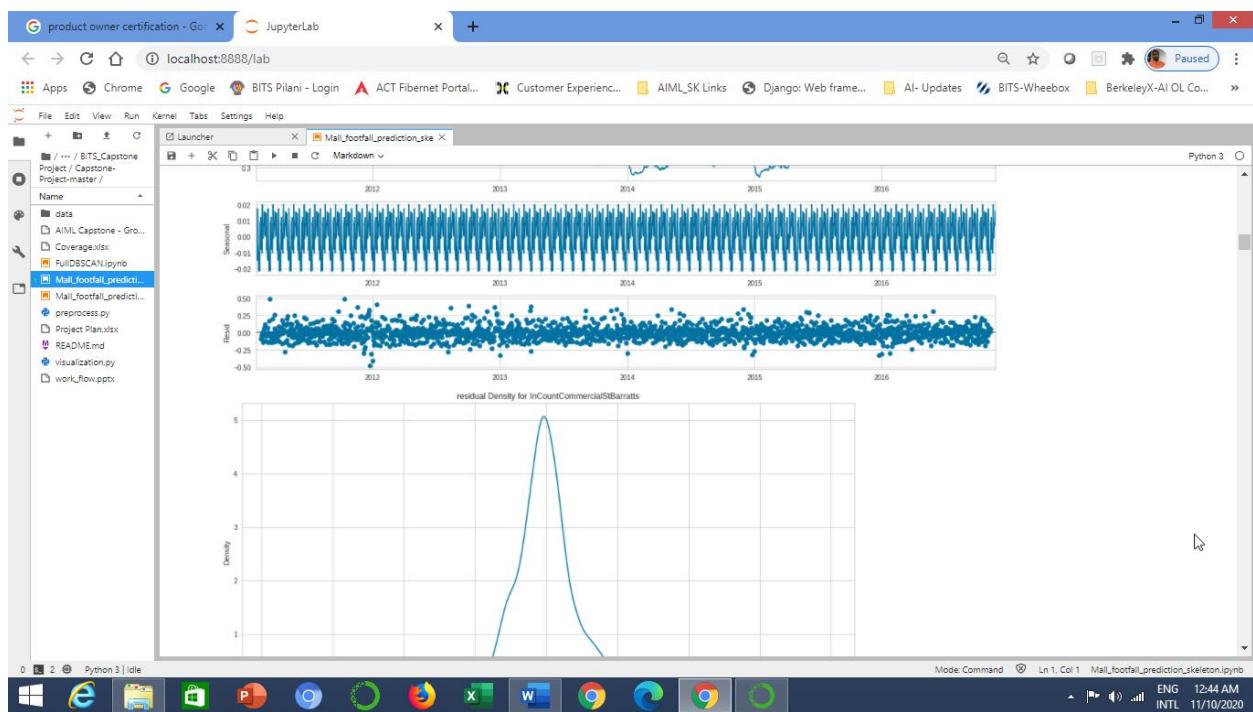
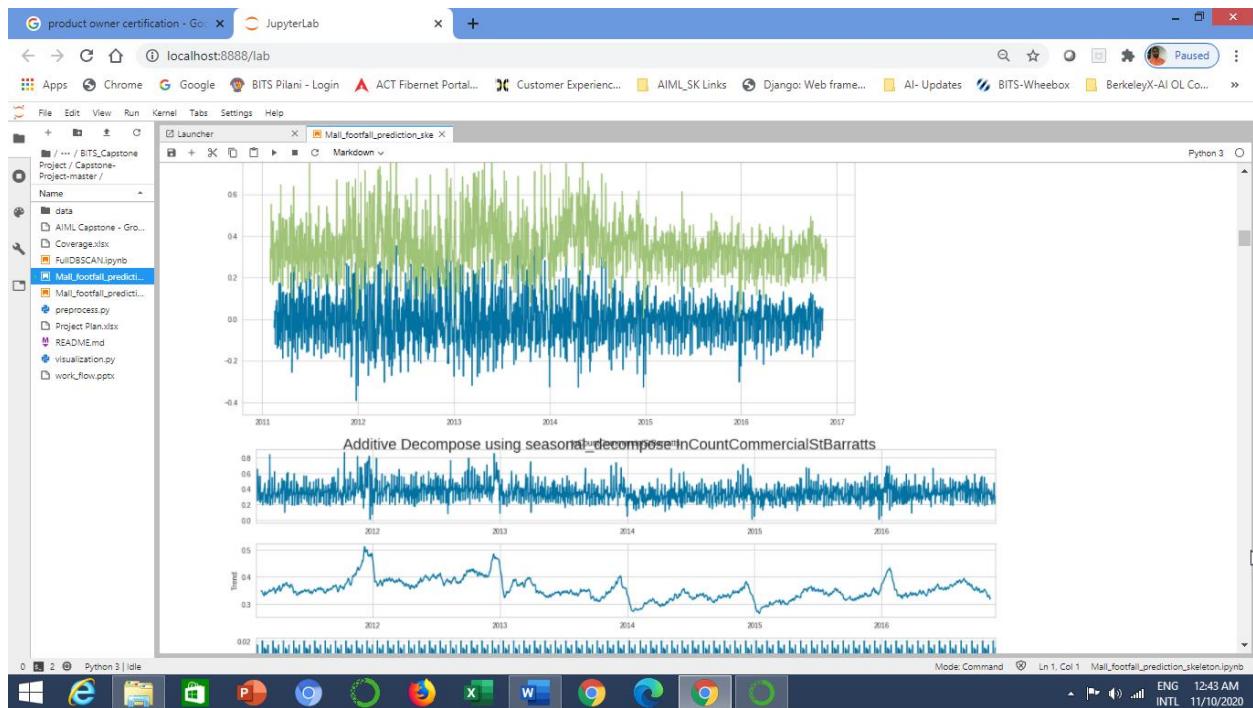


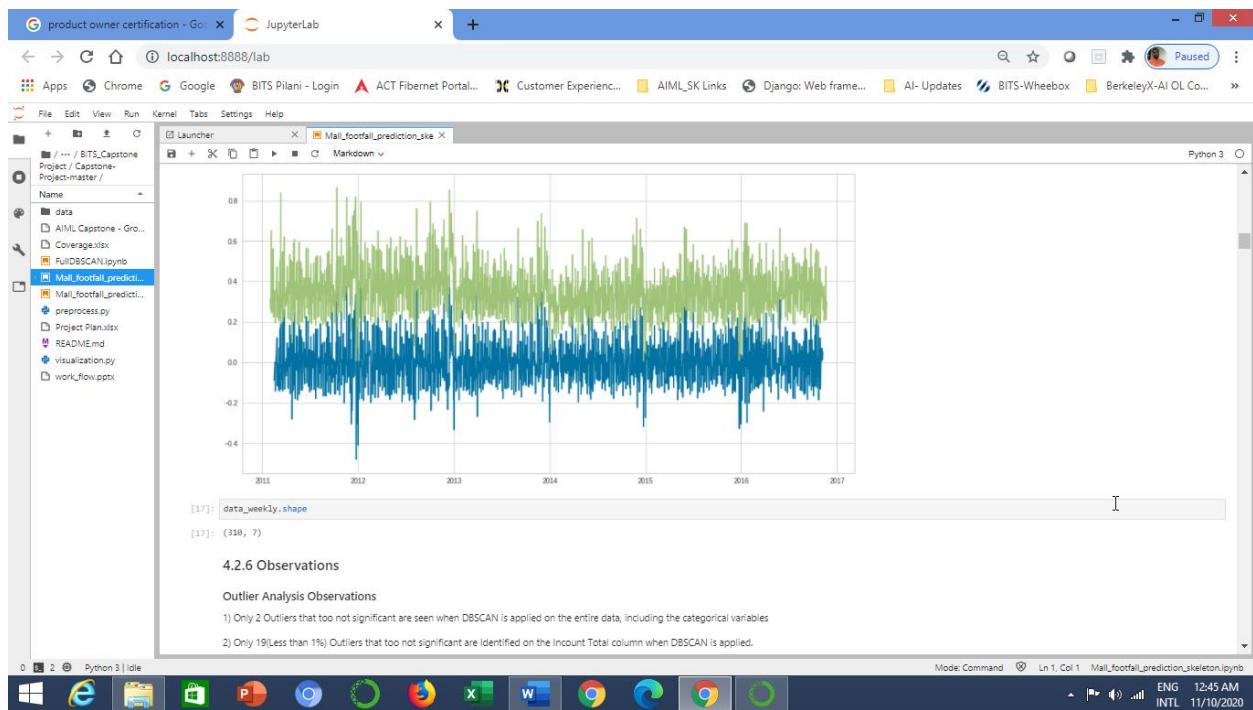
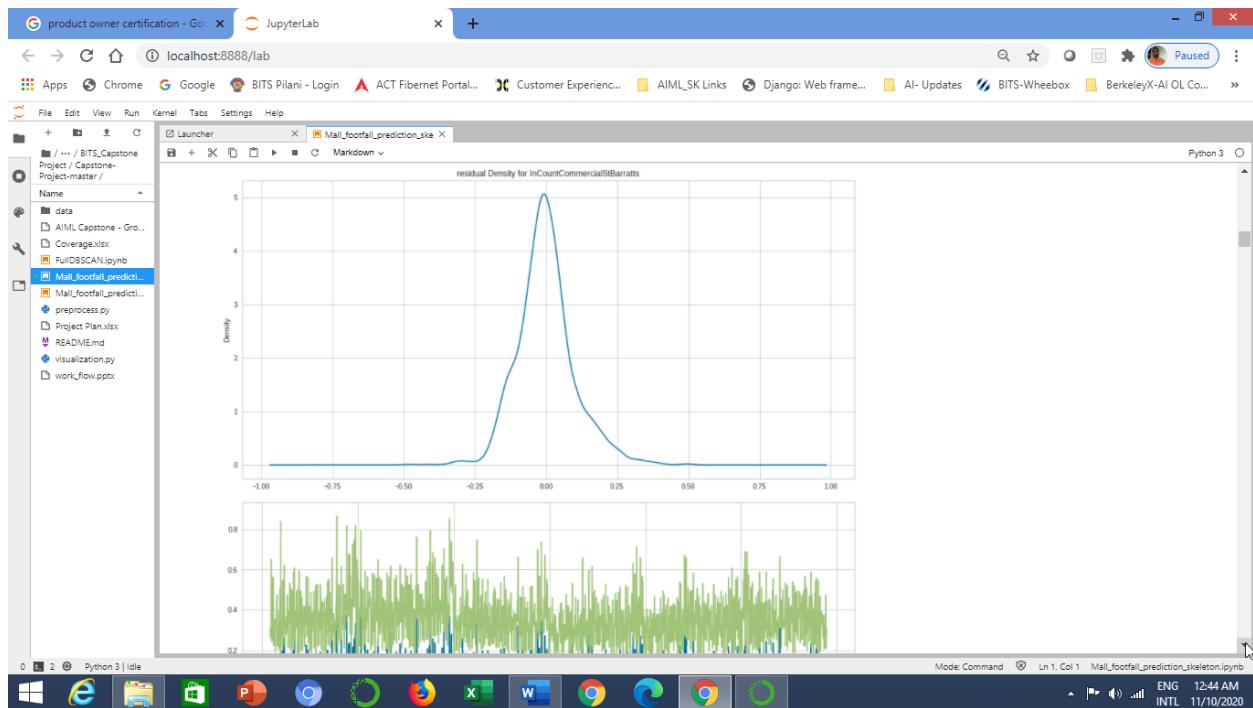












The screenshot shows a JupyterLab interface with the following details:

- Header:** product owner certification - Go... > JupyterLab
- Toolbar:** File, Edit, View, Run, Kernel, Tabs, Settings, Help
- Launcher:** Mail_footfall_prediction_ske
- Section:** Outlier Analysis Observations
- Text:** Only 2 Outliers that too not significant are seen when DBSCAN is applied on the entire data, including the categorical variables.
- Text:** Only 19(Less than 1%) Outliers that too not significant are identified on the Incount Total column when DBSCAN is applied.
- Text:** LOF is giving far more outliers compared to DBSCAN, but does not seem very relevant on the current problem statement
- Section:** Decomposing Data to Trend, Seasonality, noise, residue
- Text:** There is no visible trend seen in the Incount data using seasonal_decompose
- Text:** STL (Seasonal-Trend decomposition using LOESS - Locally Estimated Scatterplot Smoothing) shows a smoother trend
- Text:** The residual density data for all the gates and IncountTotal almost follows a normal distribution
- Section:** PCA Observations
- Text:** Just 11 components are able to explain 90% of the variance of the entire data(including Categorical variables that are encoded)
- Data Preview:** A preview of the data frame 'data_final.describe()' showing statistical summary statistics for various columns.

product owner certification - Go X JupyterLab x + localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

Launcher x Mail_footfall_prediction_ske x

data_final.head()

PCA Output

```
[95]: princompdf
```

	0	1	2	3	4	5	6	7	8	9	InCountTotal
0	0.501186	-0.073384	-0.308815	0.007210	-0.073542	0.033058	0.159510	0.042412	0.083294	0.066633	0.306589
1	0.473293	-0.160209	0.374247	-0.102031	-0.212009	-0.044546	0.141883	0.035981	0.022237	-0.020510	0.333570
2	0.450700	-0.253826	0.089324	-0.202355	-0.177654	-0.071719	0.105374	-0.061570	-0.069490	-0.116211	0.360520
3	0.370065	-0.355813	-0.193561	-0.195929	-0.150807	-0.111659	0.014388	-0.069540	-0.050136	0.142080	0.370871
4	0.298813	-0.405949	-0.477639	-0.290647	-0.221103	-0.096917	-0.026453	-0.065144	0.122334	0.018815	0.356540
...
212	-0.274768	-0.411468	-0.115133	-0.089755	-0.076417	0.087092	0.054679	0.095781	0.097241	0.158022	0.393165
213	-0.359710	-0.346921	-0.398152	0.084086	0.171626	0.075083	-0.016854	-0.013887	0.059490	-0.055338	0.272368
214	-0.422223	-0.263957	0.278268	0.151538	-0.135616	0.064667	0.107704	0.104022	-0.019381	-0.016178	0.348800
215	-0.449588	-0.174788	-0.007275	0.167731	-0.037955	0.040813	0.058860	-0.008645	-0.004975	0.083027	0.433450
216	-0.497112	-0.097951	-0.292263	0.250000	0.077615	0.056800	0.019170	-0.150700	-0.001552	-0.021012	0.379947

2127 rows x 11 columns

The contribution of each Feature for the principal components

```
[96]: compdf
```

	Month	Weekday	Day	Week_no	mean_temp	rain	wind_speed	abnormal_rain	high_temp	low_temp	high_wind	EasterSundayHoliday	University_holidays	School_holidays	UKBankHoliday	Day sin	Month sin	Month cos
0	0.000143	-0.002021	0.004861	0.007279	-0.014059	-0.007177	-0.024392	0.000234	0.000204	-0.000099	0.000178	0.000168	0.000004	-0.000123	0.000763	0.006648	-0.160914	-0.985591
1	0.003005	0.00149	0.002103	0.001417	-0.006693	0.001153	0.013375	0.000330	0.000167	-0.002211	-0.000303	0.000469	0.000025	-0.000326	0.001441	-0.013345	0.988276	-0.161999
2	0.000493	-0.00724	0.045852	0.00149	0.004766	-0.000284	-0.000102	-0.00390	0.000957	0.000045	-0.000203	-0.000214	0.000019	0.000005	0.000042	-0.000259	-0.014242	0.002430
3	0.003200	0.00747	0.000330	0.001454	0.000160	0.000044	0.003300	0.000460	0.000205	0.000443	0.000180	0.000047	0.000137	0.000044	0.000074	0.000254	0.002354	0.001603

Google product owner certification - Go... JupyterLab

localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

Launcher Mail_footfall_prediction_skeleton

Python 3

```
[96]: compdf
```

	Month	Weekday	Day	Week_no	mean_temp	rain	wind_speed	abnormal_rain	high_temp	low_temp	high_wind	EasterSundayHoliday	University_holiday	School_holidays	UKBankHoliday	Day sin	Month sin	Month cos
0	0.000143	-0.000201	0.004861	0.007279	-0.014059	-0.007177	-0.024392	0.000234	0.000204	-0.000099	0.000178	0.000168	0.000004	-0.000123	0.000763	0.006665	-0.160914	-0.985591
1	0.000095	0.000149	0.002103	0.001417	-0.006632	0.001153	0.013375	0.000330	-0.000167	-0.000221	-0.000030	0.000469	0.000025	-0.000326	0.001441	-0.013345	0.985276	-0.161999
2	0.000492	-0.000724	0.004582	0.001049	0.004785	-0.000284	-0.00102	-0.000350	0.000057	0.000045	-0.000203	-0.000214	0.000019	0.000005	0.000042	-0.000262	-0.014242	-0.004390
3	-0.003089	-0.000715	-0.002358	-0.001626	0.998510	-0.000811	-0.022386	0.000650	-0.007903	0.008412	0.001319	-0.000045	0.000177	0.001414	0.000871	0.004868	0.003854	-0.014602
4	0.020332	-0.000998	-0.004040	0.001481	0.003544	0.193777	0.965184	-0.010307	0.001585	0.001363	-0.011864	-0.000972	0.000202	0.000899	-0.002978	0.000443	-0.013477	-0.019482
5	0.081877	-0.015730	0.006713	0.112470	-0.010910	-0.172346	-0.088966	0.008718	0.000297	0.000940	0.000422	0.000249	-0.000267	0.001099	-0.000662	0.002111	0.026341	0.039903
6	-0.009988	-0.148497	-0.047704	0.004189	0.034174	0.165891	-0.097310	-0.012342	-0.000864	-0.000167	0.001315	-0.000972	-0.000074	0.000887	-0.011650	-0.005353	0.015080	-0.009263
7	-0.007176	-0.134108	-0.000490	-0.018102	0.245692	-0.106679	-0.019833	-0.001217	-0.000618	0.001201	-0.006844	0.000035	-0.000662	-0.003525	-0.003535	-0.001383	0.001930	
8	0.002127	-0.009142	-0.013567	0.005591	0.003781	0.119543	-0.038245	-0.008967	0.003371	0.000170	0.000423	-0.005415	0.000096	0.000367	-0.002419	0.004137	0.002331	-0.000226
9	-0.004048	-0.066067	-0.009104	0.005025	-0.015090	-0.067714	0.005955	0.000167	0.000558	0.000256	-0.002074	-0.000199	-0.001343	0.005811	0.000916	-0.000247	0.001273	
10	0.008006	-0.053799	-0.007368	0.004452	-0.001175	0.393137	-0.071588	-0.028317	0.000669	0.0001022	0.000379	-0.002824	0.000134	0.000366	0.007288	-0.000934	0.001979	0.000853

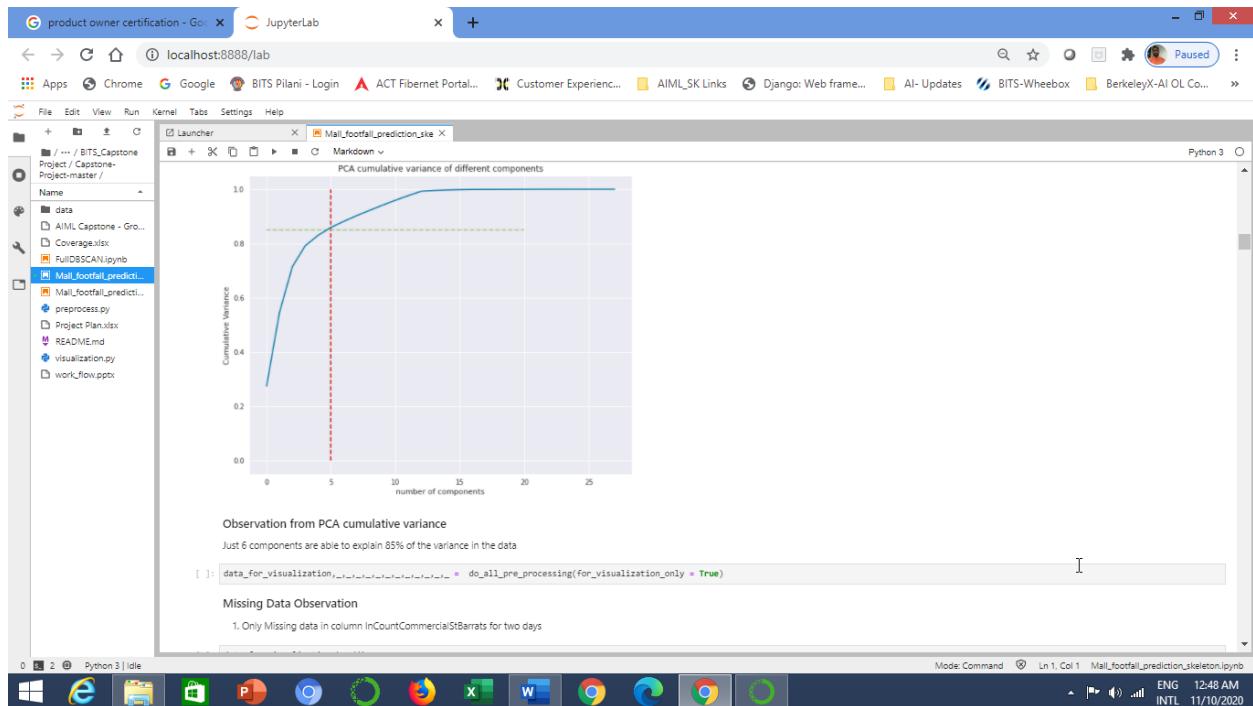
```
[98]: print(var_cum_variance[0:8])
# Plotting the scree plot with a mark on 85 percentile of variance being captured
fig = plt.figure(figsize=[10,8])
plt.vlines([5, 18], ymin=0, ymax=1, colors="r", linestyles="--")
plt.vlines([5, 18], y1=0.85, y2=1, minv=0, maxv=1, colors="g", linestyles="--")
plt.plot(var_cum_variance)
plt.xlabel("Cumulative Variance")
plt.ylabel("Number of components")
plt.title("PCA cumulative variance of different components")
plt.show()
```

0.27410866 0.54519811 0.71412959 0.79147247 0.83003109 0.85913497
0.8816237 0.90256981

PCa cumulative variance of different components

Mode Command Ln 1, Col 1 Mail_footfall_prediction_skeleton.ipynb

ENGLISH 12:47 AM INTL 11/10/2020



product owner certification - Google

JupyterLab

localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

Project / Capstone-Project-master /

Name

data

AIML Capstone - Gro...

Coverage.xlsx

FullDBSCAN.ipynb

Mail_footfall_prediction.ipynb

Mail_footfall_prediction...

preprocess.py

Project Plan.xlsx

README.md

visualization.py

work_flow.pptx

Launcher

Mail_footfall_prediction_skeleton.ipynb

Python 3

```
[1]: data_for_visualization.head()
[2]: data_final[['InCountTotal_lag7','InCountTotal_lag6','InCountTotal_lag5','InCountTotal_lag4','InCountTotal_lag3','InCountTotal_lag2','InCountTotal_lag1','InCountTotal','InCountTotal_forecast1']].head(15)
[3]: data_final[data_final.isnull().any(1)]
[4]: data_for_visualization[data_for_visualization.isnull().any(1)]
[5]: data_for_visualization.loc['2014-08-13':'2014-08-13']
[6]: print(location_name_list)
```

5 - Data Visualization

5.1 Visualisation Class

```
[1]: class Visualisation:
    def __init__(self,data_frame = None,location_name_list=None):
        self.data_frame = data_frame # array of column names to encode
        self.location_name_list = location_name_list
        self.x_limit_list = []
        self.y_limit_list = []
        self.distrapping = {'W': 'Weekly', 'M': 'Monthly', 'Y': 'Yearly'}
        self.cols_to_keep = ['mean_temp','rain','wind_speed','abnormal_rain','high_temp','low_temp','high_wind','EasterSundayHoliday','UniversityHolidays','School_holidays','UKBankHoliday']
        self.modified_location_name_list=[InCount+x for x in self.location_name_list]
        self.init_gg_plots()

    def init_gg_plots(self):
        theme_set(
            theme_S3R()
        )
        theme(
            figure_size = (8, 4),
            text = element_text(
                size = 8,
                color = "black",
                family = "DejaVu Sans"
            ),
            plot_title = element_text(
                color = "black",
                family = "DejaVu Sans",
                weight = "bold",
                size = 12
            ),
            axis_title = element_text(
                color = "black",
                family = "DejaVu Sans",
                weight = "bold",
                size = 6
            ),
        )
```

product owner certification - Google

JupyterLab

localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

Project / Capstone-Project-master /

Name

data

AIML Capstone - Gro...

Coverage.xlsx

FullDBSCAN.ipynb

Mail_footfall_prediction.ipynb

Mail_footfall_prediction...

preprocess.py

Project Plan.xlsx

README.md

visualization.py

work_flow.pptx

Launcher

Mail_footfall_prediction_skeleton.ipynb

Python 3

```
[1]: theme(
    figure_size = (8, 4),
    text = element_text(
        size = 8,
        color = "black",
        family = "DejaVu Sans"
    ),
    plot_title = element_text(
        color = "black",
        family = "DejaVu Sans",
        weight = "bold",
        size = 12
    ),
    axis_title = element_text(
        color = "black",
        family = "DejaVu Sans",
        weight = "bold",
        size = 6
    ),
)

def uni_variate_ggplot_histogram(self, x_list,x_limit_list, y_limit_list ):
    index=0
    for col in x_list:
        fig = ggplot(self.data_frame) + geom_histogram(aes(x = col),fill = 'blue', color = 'red') + labs(title=col+' distribution - (median = dashed line; mean = solid line)', x = col, y = 'Frequency', subtitle='Histogram')
        + scale_x_continuous(
            limits = (0, y_limit_list[index]),
            labels = self.labels(0, x_limit_list[index], int(x_limit_list[index]/10)),
            breaks = self.breaks(0, x_limit_list[index], int(x_limit_list[index]/10))
        )+ scale_y_continuous(
            limits = (0, y_limit_list[index]),
            labels = self.labels(0, y_limit_list[index], int(y_limit_list[index]/10)),
            breaks = self.breaks(0, y_limit_list[index], int(y_limit_list[index]/10))
        )+ geom_vline(aes(xintercept = self.data_frame[col].mean()), color = "gray") + geom_vline(aes(xintercept = self.data_frame[col].median()), linetype = 'dashed', color = 'gray')
        fig.draw()
        index += 1

def facet_ggplot_histogram(self, x, Facet_wrap_cols, x_limit, y_limit):
    fig = ggplot(self.data_frame) + geom_histogram(aes(x = x),fill = 'blue', color = 'red') + labs(title=x+' distribution', x = x, y = 'Frequency', subtitle='Histogram')
    + scale_x_continuous(
        limits = (0, x_limit),
        labels = self.labels(0, x_limit, int(x_limit/10)),
        breaks = self.breaks(0, x_limit, int(x_limit/10))
    )+ scale_y_continuous()
```

```
Google product owner certification - Go... JupyterLab
localhost:8888/lab
File Edit View Run Kernel Tabs Settings Help
Project / Capstone-Project-master /
Name
data
AIML Capstone - Gro...
Coverage.xlsx
fullDBSCAN.ipynb
Mail_footfall_predict...
preprocess.py
Project Plan.xlsx
README.md
visualisation.py
work_flow.pptx
Python 3
label = self.labels[0, x_limit, int(x_limit/10)],  
breaks = self.breaks[0, x_limit, int(x_limit/10)]  
)  
> scale_y_continuous()  
limits = (0, y_limit),  
labels = self.labels[0, y_limit, int(y_limit/10)],  
breaks = self.breaks[0, y_limit, int(y_limit/10)]  
) theme((figure_size = (16, 24)) + facet_wrap(facet_wrap_cols, ncol = len(facet_wrap_cols), labeller="label_both"))  
fig.draw()  
  
def scatter(self, hue="rain", **kwargs):  
    fig = ggplot((self.data_frame) + geom_point(aes(x=x,y=y, fill=hue),alpha = 0.5, color = color) + labs(title="Scatter plot of 'x' vs 'y', x = x, y = y,"  
)  
> scale_x_continuous()  
limits = (0, x_limit),  
labels = self.labels[0, x_limit, int(x_limit/10)],  
breaks = self.breaks[0, x_limit, int(x_limit/10)]  
)  
> scale_y_continuous()  
limits = (0, y_limit),  
labels = self.labels[0, y_limit, int(y_limit/10)],  
breaks = self.breaks[0, y_limit, int(y_limit/10)]  
)  
> theme((figure_size = (12, 8)))  
fig.draw()  
  
def label(self, from_, to_, step_):  
    ret = pd.Series(np.arange(from_, to_ + step_, step_)).apply(lambda x: '{:.0f}'.format(x).tolist()  
def break(self, from_, to_, step_):  
    return pd.Series(np.arange(from_, to_ + step_, step_)).tolist()  
  
def plot_graph(self, plot_type = 'boxplot', col_name = 'InCount', nplotrows=4, nplotcols=2):  
    location_area_index = 0  
    for i in range(nplotrows):# 4 rows of plots  
        fig = plt.figure(figsize=(16,8))  
        for j in range(nplotcols):# 2 columns of sub plots  
            ax = fig.add_subplot(1, 2, j+1)  
            if(col_name == 'InCount'):  
                target_variable = self.modified_location_name_list[location_area_index]  
            else:  
                target_variable = col_name  
  
            if(plot_type == 'boxplot'):  
                ax.boxplot(self.data_frame[target_variable])  
            elif(plot_type == 'hist'): # histogram  
                n, bins, patches = ax.hist(self.data_frame[target_variable], bins=200, color=self.colors[j], density=True)  
            elif(plot_type == 'heatmap'):  
  
def show_tiledistribution(self, distribution_type='W', col_name='InCount'): #  
    for location_name in self.location_name_list:  
        if(col_name == 'InCount'):  
            target_variable = col_name + location_name  
        else:  
            target_variable = col_name  
        self.data_frame[target_variable].resample(distribution_type).mean().plot(figsize = (20,8))  
  
        xlabel = self.distrimapping[distribution_type]  
        plt.xlabel(label)  
        plt.ylabel('Mail foot fall')  
        plt.legend(self.location_name_list)  
        plt.title("Mail Foot Fall for different locations - " + xlabel)  
def show_pairplot(self, one_location_only = True, kind='scatter'): #  
    sns.set()  
    cols = self.cols_to_keep  
    cols.append('InCount'+self.location_name_list[0])  
    sns.pairplot(self.data_frame[cols], kind=kind)  
  
def comparegates(self):  
    gatecorrdatas(data_final)if(modified_location_name_list)  
    sns.pairplot(gatecorrdatas)  
    plt.show()  
    sns.heatmap(gatecorrdatas.corr(), annot = True)  
    plt.show()  
  
def bi_variate_kernel_density_plot(self, x, y_list, hue_list=None):  
    nplotrows = len(hue_list)  
    nplotcols = len(y_list)  
    for i in range(nplotrows):# 4 rows of plots  
        fig = plt.figure(figsize=(16,8))  
        for j in range(nplotcols):# 2 columns of sub plots
```

```
Google product owner certification - Go... JupyterLab
localhost:8888/lab
File Edit View Run Kernel Tabs Settings Help
Project / Capstone-Project-master /
Name
data
AIML Capstone - Gro...
Coverage.xlsx
fullDBSCAN.ipynb
Mail_footfall_predict...
preprocess.py
Project Plan.xlsx
README.md
visualisation.py
work_flow.pptx
Python 3
sns.heatmap(self.data_frame[[col]].corr(), annot=True, linewidths=.2)  
        ax.set_xlabel("Mail foot fall - " + col_name + " for " + self.location_name_list[location_area_index] )  
        ax.set_ylabel("frequency")  
        ax.set_title("Mail foot fall ' + plot_type + ' for " + col_name + " for " + self.location_name_list[location_area_index] )  
        location_area_index+=1  
  
def show_tiledistribution(self, distribution_type='W', col_name='InCount'): #  
    for location_name in self.location_name_list:  
        if(col_name == 'InCount'):  
            target_variable = col_name + location_name  
        else:  
            target_variable = col_name  
        self.data_frame[target_variable].resample(distribution_type).mean().plot(figsize = (20,8))  
  
        xlabel = self.distrimapping[distribution_type]  
        plt.xlabel(label)  
        plt.ylabel('Mail foot fall')  
        plt.legend(self.location_name_list)  
        plt.title("Mail Foot Fall for different locations - " + xlabel)  
def show_pairplot(self, one_location_only = True, kind='scatter'): #  
    sns.set()  
    cols = self.cols_to_keep  
    cols.append('InCount'+self.location_name_list[0])  
    sns.pairplot(self.data_frame[cols], kind=kind)  
  
def comparegates(self):  
    gatecorrdatas(data_final)if(modified_location_name_list)  
    sns.pairplot(gatecorrdatas)  
    plt.show()  
    sns.heatmap(gatecorrdatas.corr(), annot = True)  
    plt.show()  
  
def bi_variate_kernel_density_plot(self, x, y_list, hue_list=None):  
    nplotrows = len(hue_list)  
    nplotcols = len(y_list)  
    for i in range(nplotrows):# 4 rows of plots  
        fig = plt.figure(figsize=(16,8))  
        for j in range(nplotcols):# 2 columns of sub plots
```

```
for i in range(nplotrows):# 4 rows of plots
    fig = plt.figure(figsize=(16,8))
    for j in range(nplotcols):# 2 columns of sub plots
        ax1 = fig.add_subplot(1, 2, j+1)

        sns.kdeplot(data = self.data_frame, x=x, y=y_list[j], hue=hue_list[i])

        ax1.set_xlabel('Plot ' + x)
        ax1.set_ylabel(y_list[i])
        ax1.set_title(KDE plot x' + x + ' and ' + y_list[j])
```

```
def plot_graphs_for_list(self, plot_type = 'boxplot', col_list = ['rain'], target_col='InCountTotal'):

    col_index = 0
    nplotrows = 1
    nplotcols = 1
    if len(col_list)>1:
        nplotrows = int(len(col_list)/2)
        nplotcols = 2
    for i in range(nplotrows):# 4 rows of plots
        fig = plt.figure(figsize=(16,8))
        for j in range(nplotcols):# 2 columns of sub plots
            ax1 = fig.add_subplot(1, 2, j+1)

            col_name = col_list[col_index]
            if(plot_type == "boxplot"):
                ax1 = sns.boxplot(x=x, y=y, data=self.data_frame[col_name])
            elif(plot_type == "histogram"):
                n, bins, patches = ax1.hist(self.data_frame[col_name], bins=10, color=self.colors[j], density=True)
            elif(plot_type == "kdeplot"):
                sns.kdeplot(data = self.data_frame, x=x, target_col, shade=True, hue=col_name)
            ax1.set_xlabel(col_name + " for " + self.location_name_list[0])
            ax1.set_ylabel("Frequency")
            ax1.set_title(plot_type + " - " + col_name)
            col_index +=1
```

```
def show_categorical_plots(self, style="swarm", x="InCountTotal", y="Weekday", hue="Month", order=["Sunday", "Monday", "Tuesday", "Wednesday", "Thursday", "Friday", "Saturday"]):
    if style == "swarm":
        sns.swarmplot(x=x, y=y, data=self.data_frame, order=order)
    elif style == "swarm_with_hue":
        sns.swarmplot(x=x, y=y, hue=hue, data=self.data_frame, order=order)
    elif style == "stripplot":
        sns.stripplot(x=x, y=y, hue=hue, data=self.data_frame, order=order)
```

```
elif style == "swarmplot":
    sns.swarmplot(x=x, y=y, hue=hue, data=self.data_frame, order=order)
elif style == "stripplot_yhue":
    sns.stripplot(x=x, y=y, data=self.data_frame, order=order)
elif style == "boxplot_yhue":
    sns.boxplot(x=x, y=y, hue=hue, data=self.data_frame, whis=np.inf, order=order)
elif style == "violin":
    sns.violinplot(x=x, y=y, data=self.data_frame, order=order)
elif style == "box_with_hue":
    sns.boxplot(x=x, y=y, hue=hue, data=self.data_frame, whis=np.inf, order=order)
elif style == "boxen":
    sns.boxenplot(x=x, y=y, data=self.data_frame, order=order)
elif style == "box_and_whisker":
    sns.boxplot(x=x, y=y, data=self.data_frame, order=order)
else:
    print("No style is selected")
```

```
def ggpplot_swarm(self, dataFrame, x_column, y_column, hue=None):
    ax = sns.swarmplot(x=x_column, y=y_column, hue=hue, data=dataFrame, size=6)
    ax.set_title(x_column + " and " + y_column + " analysis.")
```

```
def ggpplot_count(self, dataFrame, x_column, y_column, hue=None):
    ax = sns.countplot(x=x_column, y=y_column, data=dataframe)
    ax.set_title(x_column + " and " + y_column + " analysis.")
```

```
def ggpplot_cat(self, dataFrame, x_column, y_column, col, hue=None):
    g = sns.catplot(x=x_column, y=y_column,
                    hue=hue, col=col,
                    data=dataframe, kind="swarm",
                    height=6, aspect=.7);
```

```
def relplot_cat(self, dataFrame, x_column, y_column, col, row, size, hue=None):
    g = sns.relplot(x=x_column, y=y_column,
                    hue=hue, col=col,
                    row=row, size=size,
                    data=dataframe, palette=["b", "r"]);
```

```
def plotcategoryalidt(self, catcollist,featlist):
    dd=pd.DataFrame(self.data_frame)
    rnum=len(dd)
    fig, axs = plt.subplots(len(catcollist),figsize=(15,40))
    fig.tight_layout(pad=2.0)
    rc=0
```

The screenshot shows a JupyterLab interface with a Python script titled "Mail_footfall_prediction_skeleton.ipynb". The script contains code for generating bar charts to show the percentage distribution of categorical features in the dataset. It uses the `pandas` library to count unique values and the `matplotlib` library to create the charts.

```
sig_tight_layout(pad=2.0)
rc=0

for col in(catcollist):
    uniq=dd[col].unique()
    pcnt=[]
    for val in uniq:
        pcnt.append(100*(dd[col][dd[col]==val].count())/numb)

    hts=xs[nc].bar(uniq,pcnt)

    for ht in (hts):
        ht.set_height()
        axis[nc].text(ht.get_x(), 1.05*ht, str(round(ht,2)))
        axis[nc].set_title(catcol+'Percentage Distribution in the Data')

    nc+=1
plt.show()

fig, axs = plt.subplots(len(catcollist), figsize=(15,40))
sig_tight_layout(pad=2.0)
rc=0

for catcol in catcollist:
    temp[catcol].unique()
    rn = np.arange(len(uniq))
    bwidht=.05

    for feature in fealist:
        temp[feature]
        pcnt=[]
        for colval in uniq:
            temp.append(np.sum(dd[feature][dd[catcol]==colval]))
            pcnt.append(100*(dd[catcol][dd[catcol]==colval].count())/numb)
        axis[nc].bar(rn,temp, width=bwidht, label=feature)

    rn = [x + bwidht for x in rn]
    axis[nc].set_title(catcol+'Distribution in the Data')
    axis[nc].set_xticks([r + (bwidht*len(uniq)) for r in range(len(uniq))], uniq)
    nc+=1

plt.show()
```

The screenshot shows a JupyterLab interface with a Python script titled "Mail_footfall_prediction_skeleton.ipynb". The script includes imports for visualization, defines a variable "graphs", and calls "graphs.show_timedistribution('W')". Below this, sections for "5.2 Observations" and "5.3 Plots" are shown. The "5.3 Plots" section contains code for "Line Plot" and "KDE plots". The "KDE plots" section includes a call to "graphs.plot_graphs_for_list(plot_type = 'kdeplot', target_col = 'InCountTotal', col_list= ['Weekday','Month','high_temp','low_temp','abnormal_rain','high_wind','EasterSundayHoliday','University_holidays']". The "Observations" section lists three points: 1. weekday : Saturday seems to be having more incount as expected 2. Month : December is seen as the busy month with more footfall 3. high_temp, low_temp, abnormal_rain : Incount is high in the absence of these

```
[ 1]: #from visualization import Visualization
graphs = Visualization(data_for_visualization,location_name_list)
graphs.show_timedistribution("W")



## 5.2 Observations



1. Incount is higher during the last weeks of December  
2. 'Brigate' : There seems to be downward trend in 2015; overall, this seems to be the main gate  
3. 'AlbionStNorth' : There is a downward trend throughout



## 5.3 Plots



### Line Plot



```
[1]: graphs.show_timedistribution('W')
[1]: graphs.show_timedistribution('Y')
```



### KDE plots



#### KDE plots for 'InCountTotal'



```
[1]: graphs.plot_graphs_for_list(plot_type = 'kdeplot', target_col = 'InCountTotal', col_list= ['Weekday','Month','high_temp','low_temp','abnormal_rain','high_wind','EasterSundayHoliday','University_holidays'])
```



Observations:



1. weekday : Saturday seems to be having more incount as expected
2. Month : December is seen as the busy month with more footfall
3. high_temp, low_temp, abnormal_rain : Incount is high in the absence of these

```

Google product owner certification - Go... JupyterLab

localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

Launcher Mail_footfall_prediction_ske

Python 3

```
[ ]: graphs.bi_variate_kernel_density_plot('InCountTotal',['rain','wind_speed','mean_temp'], hue_list=['Weekday','Month','high_temp','low_temp','abnormal_rain','high_wind','EasterSundayHoliday','University_holiday'])
```

Observations:

1. high_temp, low_temp, abnormal_rain : InCount is high in the absence of these
2. EasterSunday/Holiday : InCount is very low
3. Other holidays : We do not see any significant impact on the InCount.

Bi-variate KDE plots

```
[ ]: graphs.bi_variate_kernel_density_plot('InCountTotal',['rain','wind_speed','mean_temp'], hue_list=['Weekday','Month','high_temp','low_temp','abnormal_rain','high_wind','EasterSundayHoliday','University_holiday'])
```

Observations:

1. weekday : when the rain is low Saturday seems to be having more incount as expected
2. Month : December is seen as the busy month with more footfall regardless of rain/temp situation
3. high_temp, low_temp, abnormal_rain : InCount is high in the absence of these
4. EasterSunday/Holiday : InCount is very low
5. Other holidays : We do not see any significant impact on the InCount.

Univariate Histogram ggplot

```
[ ]: graphs.uni_variate_ggplot_histogram(['InCountTotal', 'rain','mean_temp','wind_speed'],[200000,10,30,10],[250,250,250,250])
```

Observations:

1. mean of InCount : 100000
2. mean of rain < 1
3. mean of temp : 10.5
4. mean of wind_speed : ~1.75

Facet Histogram ggplot

```
[ ]: graphs.facet_ggplot_histogram('InCountTotal',('Weekday','high_temp','abnormal_rain'),200000,50)
```

Observations:

- 1) The Incount is significantly low during the days with Abnormal rain and also High temperature

Mode Command Ln 1, Col 1 Mail_footfall_prediction_skeleton.ipynb

ENGLISH 12:52 AM INTL 11/10/2020

Google product owner certification - Go... JupyterLab

localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

Launcher Mail_footfall_prediction_ske

Python 3

```
[ ]: graphs.scatter_plot_ggplot(x = 'InCountTotal',y='rain', hue='Weekday', x_limit=200000, y_limit=15)
```

```
[ ]: graphs.scatter_plot_ggplot(x = 'InCountTotal',y='wind_speed', hue='Weekday', color = 'blue', x_limit=200000, y_limit=10)
```

```
[ ]: graphs.scatter_plot_ggplot(x = 'InCountTotal',y='mean_temp', hue='Weekday', color = 'green', x_limit=200000, y_limit=30)
```

```
[ ]: graphs.scatter_plot_ggplot(x = 'InCountTotal',y='mean_temp', hue='Month', color = 'green', x_limit=200000, y_limit=30)
```

```
[ ]: graphs.scatter_plot_ggplot(x = 'InCountTotal',y='rain', hue='Month', color = 'purple', x_limit=200000, y_limit=15)
```

Scatter ggplot

```
[ ]: graphs.show_categorical_plots(style="swarm")
```

Swarm plot

```
[ ]: graphs.show_categorical_plots(style="swarm_with_hue")
```

```
[ ]: graphs.show_categorical_plots(style="swarm_with_hue", x="Month",hue="Weekday",order=[1,2,3,4,5,6,7,8,9,10,11,12])
```

Swarm plot with Hue

```
[ ]: graphs.show_categorical_plots(style="stripplot")
```

Strip plot

```
[ ]: graphs.show_categorical_plots(style="stripplot")
```

Mode Command Ln 1, Col 1 Mail_footfall_prediction_skeleton.ipynb

ENGLISH 12:53 AM INTL 11/10/2020

product owner certification - Go x JupyterLab x

localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

Launcher Mail_footfall_prediction_ske

Kernel: Python 3

File Edit View Run Kernel Tabs Settings Help

Launcher Mail_footfall_prediction_ske

[1]: graphs.show_categorical_plots(style="stripplot")

Strip plot with Hue

[1]: graphs.show_categorical_plots(style="stripplot_with_hue")

Observations from Scatter Plots (Swarm and Strip with Hue):

1. Weekdays : There is a gradual increase in Incount from Sunday to Saturday. Saturday being highest Incount
2. Months : we can see a seasonal trend in months - peaks in the months of Dec, Nov and March. Lower in the summer months(Jul and August) and also in Jan

[1]: graphs.show_categorical_plots(style="box")

Observations for Box Plot:

1. The Total variation of Incount is seen similar on Sunday, Wednesday, Thursday and Saturday.
2. The variation within the middle 50% is higher during Saturday followed by Sunday.
3. However the median is highest on Saturday, which is more than the 75% of overall InCount on the other days.
4. We see a Left skew on Saturday, Right Skew on Sunday, while the other days we see Symmetric distribution.

[1]: graphs.show_categorical_plots(style="violin")

Violin plot

[1]: graphs.show_categorical_plots(style="boxen")

Boxen plot

Mode: Command Ln 1, Col 1 Mail_footfall_prediction_skeleton.ipynb ENG 125.4 AM INT 11/10/2020

product owner certification - Go X JupyterLab x +

localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

Launcher Mail_footfall_prediction_ske

2. Monday to Friday the distribution is similar

3. On some of the days the distribution is also minimum except monday and saturdays.

4. Saturday we can see that there are some data which has count of ,400000 count compared to other days where data is normally at 350000

Autocorrelation plots

```
[1]: plot_acf(data_final['InCountTotal'], lags=100)
plt.show()
```

Observations:

- every 7th day seems to have some correlations -> inline with the trend on different days
- every 4 weeks (say, month - 30 days), there is a correlation

Histogram plot

```
[1]: graphs.plot_graphs(plot_type = 'histogram')
```

Observations

The Incount for the gates - Briggate, CommercialStLush, Headrow and CommercialStBarrats are slightly skewed towards higher incount. Rest of the gates Incount is almost normally distributed.

```
[1]: graphs.plot_graphs_for_list(plot_type = 'histogram', col_list =[ 'rain','mean_temp','wind_speed','abnormal_rain'])
```

Observations

- Incount is lower as windSpeed is higher
- Incount is lower, as the rain is higher
- Incount is significantly low during Abnormal rain
- Incount is almost normally distributed over mean temperature. The Incount is lower when the mean temperature is low or high.

Python 3 | idle Mode: Command Ln 1, Col 1 Mail_footfall_prediction_skeleton.ipynb

The screenshot shows a JupyterLab environment with several tabs open. The active tab is titled 'Mail_footfall_prediction_skeleton.ipynb'. The left sidebar displays a file tree with various notebooks and Python files. The main area contains code cells and their outputs.

Box plot (outlier) for InCount

```
[1]: graphs.plot_graphs(plot_type = 'boxplot')  
[2]: for location in graphs.modified_location_name_list:  
    print(data_for_visualization[location].nlargest(2))  
  
Observations : There are a few days where incount is much higher than other days
```

Boxplot for other attributes

```
[1]: graphs.plot_graphs_for_list(plot_type = 'boxplot', col_list = ['rain','mean_temp','wind_speed','abnormal_rain'])  
  
Observations :  
1. mean_temp , abnormal_rain : no outliers  
2. rain : There is one day with rain = 13.8  
3. wind speed is high on a few days
```

Pair plot & Heatmap

```
[1]: graphs.show_pairplot()  
  
Observations:  
1. rain : when there is no rain or less rain, Incount is higher  
2. wind-speed : When the wind-speed is lower, the Incount is higher  
3. abnormal_rain : when there is no abnormal_rain, Incount is higher, linked with rain attribute  
4. high_temp : when there is no high_temp, Incount is higher  
5. low_temp : There is no impact on Incount if low_temp or not  
6. high_wind : If there is no high_wind, Incount is higher  
7. EasterSundayHoliday : If it's EasterSunday, the InCount is very low  
8. School Holiday/University Holiday : we do not see any major changes in the in-flow  
9. UK Bank holiday - Incount is slightly lower on bank holidays
```

product owner certification - Go x JupyterLab x

localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help Python 3

+ C Mail_footfall_prediction_skeleton

[1]: graphs.comparegates()

Observation: The gates CommercialStLush and CommercialStBarratts are closely correlated

Categorical Columns Correlation

[1]: nominal.associations(data_for_visualization[['Month', 'Year', 'Day', 'Weekday', 'Week_no', 'abnormal_rain', 'high_temp', 'low_temp', 'high_wind', 'EasterSundayHoliday', 'University_holidays', 'School_holiday']])

Observations for Categorical Correlation:

1. School and University Holidays has correlation with Month, Week No and Day
2. School and University Holidays are correlated

Swarm plot

[1]: df_weekday_Hour = df_foot_fall.groupby(['Weekday', 'Hour'], as_index=False)[['InCount']].sum()
graphs.ggpolt_swarm(df_weekday_hour, 'Hour', 'InCount', 'Weekday')

Observations from Weekday wise Hour and Incount:

1. On Friday and Saturday mall has more people between 11:00 to 16:00 hours
2. On Monday, Tuesday, Wednesday and Thursday mall has more people between 12:00 to 14: hours.
3. Sunday has more people between 12:00 to 16:00 hours.

We can also see that hours 13:00-17:00 are busy hours.

Cat plot

Weather impact on Incount - Multivariate

[1]: graphs.ggpolt_cat(data_for_visualization, "high_wind", "InCountTotal", "abnormal_rain", "high_temp")
[1]: graphs.ggpolt_cat(data_for_visualization, "high_wind", "InCountTotal", "abnormal_rain", "low_temp")

Mode: Command Ln 1, Col 1 Mail_footfall_prediction_skeleton.ipynb

0 1 2 3 Python 3 | idle ENG 12:56 AM INTL 11/10/2020

product owner certification - Go x JupyterLab x

localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

Launcher Mail_footfall_prediction_ske

Observations from weather impact:

1. Incount is high during no abnormal rain, no high wind and normal temperature (not high nor low).
2. Incount is less during abnormal rain, high wind, low & high temperature.

Holidays impact on Incount - Multivariate

```
[1]: graphs.gplot_cat(data_for_visualization, "EasterSundayHoliday", "InCountTotal", "University_holidays", "UKBankHoliday")  
[1]: graphs.gplot_cat(data_for_visualization, "EasterSundayHoliday", "InCountTotal", "School_holidays", "UKBankHoliday")  
[1]: graphs.gplot_cat(data_for_visualization, "UKBankHoliday", "InCountTotal", "School_holidays", "University_holidays")
```

Observations from holidays impact:

1. Incount is high during no easter sunday holiday and no UK bank holiday.
2. Incount is very low when there is easter sunday holiday or UK bank holiday.

Rel plot

```
[1]: graphs.relplot_cat(data_for_visualization, "Weekday", "InCountTotal", "high_temp", "Month", "high_wind", "University_holidays")  
[1]: graphs.relplot_cat(data_for_visualization, "Weekday", "InCountTotal", "high_temp", "Month", "high_wind", "School_holidays")  
[1]: graphs.relplot_cat(data_for_visualization, "Weekday", "InCountTotal", "high_temp", "Month", "high_wind", "UKBankHoliday")  
[1]: graphs.relplot_cat(data_for_visualization, "Weekday", "InCountTotal", "high_temp", "Month", "high_wind", "EasterSundayHoliday")
```

Observations for MultiDimensional Categorical Columns:

University Holidays: 1, For Month 7, Month 8 and Month 9 there are more university holidays, where the incount is between 100000 to 250000 when there is high temp or low temp

2, For Month 10 and Month 11 there is more incount even though there is no university holidays, it sometimes got incount more than 350000

3, For Month 12 if there are university holidays the incount is more when the temperature is less compared to temperature which is high

product owner certification - Go x JupyterLab x

localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

Launcher Mail_footfall_prediction_skeleton.ipynb

Markdown Python 3

2, For Month 10 and Month 11 there is more incount even though there is no university holidays, it sometimes got incount more than 350000
3, For Month 12 if there are university holidays the incount is more when the temperature is less compared to temperature which is high
School Holidays: 1, For Month 8 there are more school holidays so the incount is around 250000 even if the there is high or low temperature
2, For Month 10 and Month 11 there is more incount even though there is no university holidays, it sometimes got incount more than 350000
3, For Month 12 if there are university holidays the incount is more when the temperature is less compared to temperature which is high
UKBankHolidays and EasterSundayHolidays: 1, The incount is not more when there is a UK Bank holiday and EasterSunday Holidays, it seems Bank holidays/EasterSunday Holidays are not impacting much to the incount

[1]: data_final.shape

Bar plot

Plotting the distribution of each Categorical variable

[1]: catcollist=['Month', 'Year', 'Day', 'Weekday', 'WeekNo', 'abnormal_rain', 'high_temp', 'low_temp', 'high_wind', 'EasterSundayHoliday', 'University_holidays', 'School_holidays', 'UKBankHoliday']
featlist=[['InCountBriggate', 'InCountBriggateAtMCNs', 'InCountAlbionStNorth', 'InCountAlbionStSouth', 'InCountCommercialStLush', 'InCountHeadrow', 'InCountBorntmundSq', 'InCountCommercialStBarnatts']]
graphs.plotcategoricaldist(catcollist,featlist)

Observations from the Categorical variables Percentage Distribution:

Incount on dates 9th and 11th are slightly lower compared to other dates.
9.78 % of days are with Abnormal Rain
19.1% of the Total days are with High Temperature
16.1% of the Total days are with Low Temperature
14.76% of the Total days are with High Wind
43.96% of the Total days are University Holidays
26.11% of the Total days are School Holidays
9.39% of the Total days are with Easter Sunday

Mode: Command Ln 1, Col 1 Mail_footfall_prediction_skeleton.ipynb

product owner certification - Google

JupyterLab

localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

Launcher

Mail_footfall_prediction_skeleton.ipynb

Python 3

45.98% of the Total days are University Holidays
 26.11% of the Total days are School Holidays
 2.26% of the Total days are Bank Holidays
 Biggate has more Incount compared to other gates
 Saturday has more Incount compared to other days
 Incount is low during Abnormal Rains, High Wind, High Temperature, Low Temperature.
 Incount is very low on Easter Sundays and UK Bank Holidays
 Incount is slightly higher on school and university holidays

5.4 Graphs for detrended data

```
[ ]: graphs = Visualization(data_Final_detrended,location_name_list)
graphs.show_timedistribution('W')

[ ]: graphs = Visualization(data_Final_detrended,location_name_list)
graphs.show_timedistribution('M')

[ ]: graphs.show_pairplot()
```

5.5 Summary of top observations

- Missing data is seen with for 2 days only with one specific gate. This is imputed with KNN Impute
- Stationarity check on the data : Augmented Dickey-Fuller (ADF) test shows p-value < 0.05, so the series is already stationary
- Decomposing Data to Trend, Seasonality, noise, residuals ==> a) There is no visible trend seen in the Incount data, b) there is seasonality visible c)The residual density data for all the gates and IncountTotal almost follows a normal distribution
- PCA - 85% variance achievable with 6 principal components and 90% variance preserved with 11 components
- Outlier detection - a) Only 2 Outliers that are not significant are seen when DBSCAN is applied on the entire data, including the categorical variables b) LOF is giving far more outliers compared to DBSCAN, but does not seem very relevant on the current problem statement
- Autocorrelation plot - every 7th day or multiples of 7th seem to have correlation

Mode: Command | Ln 1, Col 1 | Mail_footfall_prediction_skeleton.ipynb | ENG | 1:00 AM | INTL | 11/10/2020

product owner certification - Google

JupyterLab

localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

Launcher

Mail_footfall_prediction_skeleton.ipynb

Python 3

4. PCA - 85% variance achievable with 6 principal components and 90% variance preserved with 11 components
 5. Outlier detection - a) Only 2 Outliers that are not significant are seen when DBSCAN is applied on the entire data, including the categorical variables b) LOF is giving far more outliers compared to DBSCAN, but does not seem very relevant on the current problem statement
 6. Autocorrelation plot - every 7th day or multiples of 7th seem to have correlation
 7. rain : when there is no rain or less rain, Incount is higher
 8. wind-speed : When the wind-speed is lower, the Incount is higher
 9. Incount is low during Abnormal Rains, High Wind, High Temperature, Low Temperature.
 10. Abnormal rain seems to have more impact on the Incount than the High temperature. Days with only Abnormal rain has slightly low footfall than the days with only High Temperature
 11. EasterSunday/Holiday : if it's EasterSunday, the InCount is very low
 12. School Holiday/University Holiday - we do not see any major changes in the In-flow
 13. UK Bank Holiday - Incount is slightly lower on bank holidays
 14. Weekdays : There is a gradual increase in Incount from Sunday to Saturday, Saturday having highest Incount
 15. Month : we can see a seasonal trend in months - peaks in the months of Dec, Nov and March. Lower in the summer months(Jul and August) and also in Jan
 16. Hourly analysis: busy hours are 13:00-17:00 a) On Friday and Saturday mall has more people between 11:00 to 16:00 hours, b) On Monday, Tuesday, Wednesday and Thursday mall has more people between 12:00 to 14:00 hours. c) Sunday has more people between 12:00 to 16:00 hours.
 17. School holidays and University Holidays are highly correlated to each other

6 - Model preparation

6.1 train_test_validation_data_set Class

```
[28]: class train_test_validation_data_set:
    def __init__(self,data_frame,train_split=global_parameters["train_split"], validation_split=global_parameters["validation_split"], test_split= global_parameters["test_split"]):
        self.data_frame = data_frame.copy() # which data frame to pre-process
        self.train_split = train_split
        self.validation_split = validation_split
        self.test_split = test_split

    def prepare_pca_datasets(self,n_timesteps=global_parameters["n_timesteps"],n_forecast_steps = global_parameters["n_forecast_steps"], shuffle=False):
        X_columns = self.data_frame.columns.values
        X_columns = X_columns[:-1]
        X_data = self.data_frame.iloc[:, :-1]
        y_col = "Incount"
        y_data = self.data_frame[y_col].values
```

Mode: Command | Ln 1, Col 1 | Mail_footfall_prediction_skeleton.ipynb | ENG | 1:02 AM | INTL | 11/10/2020

The screenshot shows a Jupyter Notebook interface with a Python script titled "Mail_footfall_prediction_skeleton.ipynb". The script defines a function `prepare_regression_datasets` which reads a CSV file, performs data processing, and saves the result to another CSV file. The code includes imports for pandas, numpy, and other modules.

```
def prepare_regression_datasets(self,n_timesteps=global_parameters["n_timesteps"],n_forecast_steps = global_parameters["n_forecast_steps"], input_cols=None,shuffle=False):  
    # input variables  
    df_temp = self.data_frame[input_cols]  
  
    # remove the first n_timesteps rows from this dataset.  
    df_temp.drop(df_temp.index[0:n_timesteps],inplace=True)  
    # remove the last n_timesteps rows from this dataset.  
    df_temp.drop(df_temp.tail(n_forecast_steps).index,inplace=True)  
    df_temp = df_temp.values  
  
    target_data = self.data_frame['InCountTotal'].values  
    target_data = target_data.reshape(-1,1)  
  
    X_dataset = []  
    y_dataset = []  
  
    for i in range(n_timesteps, target_data.shape[0]-n_forecast_steps):  
        X_dataset.append(target_data[i-n_timesteps:i])  
        y_dataset.append(np.concatenate((y_dataset,target_data[i:i+n_forecast_steps,0].flatten()),axis=0))  
  
    X_columns = list(['InCountTotal_lag'+str(i) for i in range(n_timesteps,0,-1)])  
    X_columns = X_columns + list(input_cols)  
  
    X_dataset = pd.DataFrame(X_dataset)  
    df_temp = pd.DataFrame(df_temp)  
  
    df = pd.concat([X_dataset,df_temp],axis=1)  
    df.columns = X_columns  
    df.to_csv('regression_dataset_X.csv')  
    X_dataset = df.values  
    data_set_len = X_dataset.shape[0]  
  
    df_y = pd.DataFrame(y_dataset,columns=['InCountTotal'])  
    df_y.to_csv('regression_dataset_y.csv')  
  
    if shuffle == True:  
        df_x = pd.concat([pd.read_csv('regression_dataset_X.csv'),pd.read_csv('regression_dataset_y.csv')], axis=1)  
        df_x = df_x.sample(n=len(df_x),random_state=42)  
        X_dataset = df_x.iloc[:,1:len(X_columns)].values  
        y_dataset = df_x.iloc[:,0].values
```

The screenshot shows a Jupyter Notebook interface with a Python script titled "Mail_footfall_prediction_skeleton.ipynb". The script defines a function `createTimeSteps_LSTM` which shuffles data and creates time steps for an LSTM model. It also includes a function `shuffleDataset` for deterministic data.

```
def shuffleDataset(self,X_data,y_data,deterministic_data,n_timesteps,n_features,n_forecast_steps,len_deterministic_data_col):  
    len_X_data_cols = n_timesteps*n_features  
    len_y_data_cols = n_forecast_steps  
  
    X_data = X_data.reshape(X_data.shape[0],len_X_data_cols)  
    y_data = y_data.reshape(y_data.shape[0],len_y_data_cols)  
    deterministic_data = deterministic_data.reshape(deterministic_data.shape[0],len_deterministic_data_col)  
  
    df_data = pd.concat((pd.DataFrame(X_data),pd.DataFrame(y_data),pd.DataFrame(deterministic_data)), axis=1)  
  
    df_data = df_data.sample(n=len(df_data), random_state=42)  
  
    X_data_shuffled = df_data.iloc[:,len_X_data_cols:values  
    y_data_shuffled = df_data.iloc[:,len_X_data_cols:len_X_data_cols+len_y_data_cols].values  
    deterministic_data_shuffled = deterministic_data_shuffled.reshape(deterministic_data.shape[0],n_forecast_steps, len_deterministic_data_col)  
  
    X_data_shuffled,y_data_shuffled,deterministic_data_shuffled  
  
    return X_data_shuffled,y_data_shuffled,deterministic_data_shuffled  
  
def createTimeSteps_LSTM(self,n_timesteps=n_forecast_steps+1, input_cols=['InCountTotal'],deterministic_features=None, shuffle=False):  
  
    input_data = self.data_frame[input_cols].values  
    n_features = len(input_cols)  
  
    X_dataset = []  
    y_dataset = []  
  
    for i in range(n_timesteps, input_data.shape[0]-n_forecast_steps):  
        X_dataset.append(input_data[i-n_timesteps:i])  
        y_dataset.append(np.concatenate((y_dataset,input_data[i:i+n_forecast_steps,1:].flatten()),axis=0))  
  
    deterministic_dataset = []  
    if deterministic_features != None:
```

The screenshot shows a JupyterLab interface with a Python notebook open. The notebook title is "Mail_footfall_prediction_ske". The code in the notebook is as follows:

```
if deterministic_features != None:
    deterministic_feature_data = self.data_frame[deterministic_features].values
    for i in range(n_timesteps, deterministic_feature_data.shape[0]-n_forecast_steps):
        deterministic_dataset = np.concatenate((deterministic_dataset,deterministic_feature_data[i:i+n_forecast_steps,:].flatten()),axis=0)
    ...
LSTM input shape must be in the following order.
Samples - One sequence is one sample. A batch is comprised of one or more samples.
Time Steps - One time step is one point of observation in the sample.
Features - One feature is one observation at a time step.
...
X_dataset = np.array(X_dataset)
lendataset = X_dataset.shape[0]

X_dataset = np.reshape(X_dataset, (X_dataset.shape[0], n_timesteps, len(input_cols)))
y_dataset = np.reshape(y_dataset, (X_dataset.shape[0], n_forecast_steps))
if deterministic_features != None:
    deterministic_dataset = np.reshape(deterministic_dataset,(lendataset,n_forecast_steps, len(deterministic_features)))

return X_dataset,y_dataset,deterministic_dataset,n_features

def prepare_LSTM_datasets(self,n_timesteps=7,n_forecast_steps = 7, input_cols=[InCountTotal], deterministic_features=None, validation=True, shuffle=False):

    # input variables
    X_dataset, y_dataset,deterministic_dataset, n_features = self.createTimeSteps_LSTM(n_timesteps=n_timesteps,n_forecast_steps = n_forecast_steps, input_cols = input_cols,deterministic_features=deterministic_features)

    data_set_len = X_dataset.shape[0]

    if validation == True:
        # We do the split of train /validation/test
        X_train, X_validation, X_test = np.split(X_dataset, [int(self.train_split*data_set_len),int((self.train_split+self.validation_split)*data_set_len)])
        y_train, y_validation, y_test = np.split(y_dataset, [int(self.train_split*data_set_len),int((self.train_split+self.validation_split)*data_set_len)])
        deterministic_train, deterministic_validation, deterministic_test = np.split(deterministic_dataset, [int(self.train_split*data_set_len),int((self.train_split+self.validation_split)*data_set_len)])
    else:
        # If no validation, then train/test split
        X_train, X_test = np.split(X_dataset, [int((1-self.test_split)*data_set_len)])
        y_train, y_test = np.split(y_dataset, [int((1-self.test_split)*data_set_len)])
        deterministic_train, deterministic_test = np.split(deterministic_dataset, [int((1-self.test_split)*data_set_len)])
    X_validation = X_test
    y_validation = y_test
    deterministic_validation = deterministic_test
```

The screenshot shows a Jupyter Notebook environment with the following details:

- Title Bar:** product owner certification - Go X JupyterLab X
- URL Bar:** localhost:8888/lab
- File Menu:** File Edit View Run Kernel Tabs Settings Help
- Launcher:** Mail_footfall_prediction_ske
- Code Cell:** Python 3
- Code Content:**

```
if(shuffle == True):
    X_train, y_train , deterministic_train = self.shuffleDataset(X_train, y_train, deterministic_train,n_timesteps,n_features,n_forecast_steps,len(deterministic_features))

    return np.array(X_train), np.array(y_train), np.array(X_validation), np.array(y_validation), np.array(X_test), np.array(y_test) ,deterministic_train, deterministic_validation, deterministic_test, n_features
```

6.2 model_utility class

```
[21]:
```

```
class model_utility:
    def __init__(self):
        dummy = []
    def calculate_metrics(self, actual, forecast,n_sample_size=100,n_input_variables=1, model1_name=None):
        eps = 1e-10
        mape = np.mean(np.abs(forecast - actual))/np.abs(actual+eps)*100 # MAPE
        me = np.mean(np.abs(forecast - actual)) # ME
        mae = np.mean(np.abs(forecast - actual)) # MAE
        mpe = np.mean((forecast - actual)/actual+eps) # MPE
        rmse = np.mean((forecast - actual)**2)**.5 # RMSE
        mse = np.mean((forecast - actual)**2) # MSE

        n = n_sample_size
        p = n_input_variables
        AIC = n * log(mse) + 2 * p # AIC
        BIC = n * log(mse) + p * log(n) # BIC

        r2 = r2_score(actual,forecast) # R2 score
        adjR2_score = 1 - (1-r2)*(n-1)/(n-p-1) # Adj R2 score = 1-(1-R2)*(n-1)/(n-p-1)

        return {'mape':mape, 'mse':mse, 'me':me, 'mae':mae,
                'mpe':mpe, 'rmse':rmse, 'R2 score':r2, 'Adj R2 score':adjR2_score, 'AIC':AIC, 'BIC': BIC}

    def convert_dict_dataframe(self, loss_values):
        # appending columns
        columns = []
        columns.append('Model_type')
        first_dict_item = list(loss_values.values())[0]
        for loss_type, _ in first_dict_item.items():
            columns.append(loss_type)
```
- Bottom Status Bar:** Mode: Command Ln 1 Col 1 Mail_footfall_prediction_ske.ipynb ENG 10:04 AM

product owner certification - Go... JupyterLab

localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

Launcher Mail_footfall_prediction_skeleton.ipynb

```

columns.append('Model_type')
first_dict_itm = list(loss_values.values())[0]
for loss_type, _ in first_dict_item.items():
    columns.append(loss_type)

# rows [list initialization
rows = []
# appending rows
for name, value in loss_values.items():
    data_row = []
    data_row.append(name)
    for _, loss_value in value.items():
        data_row.append(loss_value)
    rows.append(data_row)

# using data frame
df = pd.DataFrame(rows, columns=columns)
return df

```

6.2 regression_model Class

```

class regression_model:
    def __init__(self, dataset, dataset_type='normal'):
        self.dataset = dataset.copy()
        self.dataset.dropna(inplace=True, how='normal')

```

product owner certification - Go... JupyterLab

localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

Launcher Mail_footfall_prediction_skeleton.ipynb

6.2 regression_model Class

```

class regression_model:
    def __init__(self, dataset, dataset_type='normal'):
        self.dataset = dataset.copy()
        self.dataset_type = dataset_type # normal or detrend
        self.trend = []
        self.seasonal = []
        self.grid_params = {}
        self.grid_params.update({'linear_regression': {'fit_intercept': [True, False], "normalize": [True, False]}})
        self.grid_params.update({'lassoRegression': { 'alpha': [1.0, 0.1, 0.001, 0.00001], "fit_intercept": [True, False], "normalize": [True, False], "precompute": [True, False], "max_iter": [10, 100, 1000], "warm_start": [True, False], "positive": [True, False], "selection": ['cyclic', 'random'] }})
        self.grid_params.update({'ridgeRegressor': { 'alpha': [1.0, 0.1, 0.001, 0.00001], "fit_intercept": [True, False], "normalize": [True, False], "max_iter": [10, 100, 1000], "solver": ['auto', 'svd', 'cholesky', 'lsqr', 'sparse_cg', 'sag', 'saga'] }})
        self.grid_params.update({'elastic_net': { 'alpha': [1.0, 0.1, 0.001, 0.00001], "fit_intercept": [True, False], "normalize": [True, False], "precompute": [True, False], "max_iter": [10, 100, 1000], "warm_start": [True, False], "positive": [True, False], "selection": ['cyclic', 'random'] }})
        self.grid_params.update({'huberRegressor': { 'epsilon': [1.1, 1.1, 1.2, 1.35, 2], "alpha": [1.0, 0.1, 0.001, 0.00001], "fit_intercept": [True, False], "max_iter": [10, 100, 1000], "warm_start": [True, False], "tol": [1e-04, 1e-05, 1e-06] }})
        self.grid_params.update({'lasso_lars': { "alpha": [1.0, 0.1, 0.001, 0.00001], "fit_intercept": [True, False], "max_iter": [10, 100, 1000], "precompute": [True, False], "fit_path": [True, False], "positive": [True, False] }})
        self.grid_params.update({'passive_aggressive': { "C": [1.0, 0.1, 0.001, 0.00001], "fit_intercept": [True, False], "max_iter": [10, 100, 1000], "tol": [1e-03, 1e-04, 1e-05], "early_stopping": [True, False], "shuffle": [True, False], "warm_start": [True, False], "average": [True, False] }})
        self.grid_params.update({'sgd': { "penalty": ['l1', 'l1', 'elasticnet'], "alpha": [1.0, 0.1, 0.001, 0.00001, 0.000001], "l1_ratio": [0.15, 0.2, 0.25], "fit_intercept": [True, False], "max_iter": [10, 100, 1000], "tol": [1e-03, 1e-04, 1e-05], "shuffle": [True, False], "learning_rate": ['invscaling', 'constant', 'optimal', 'adaptive'], "eta0": [0.01, 0.001], "early_stopping": [True, False], "warm_start": [True, False], "average": [True, False] }})
        self.grid_params.update({'knn_regressor': { "n_neighbors": [4, 5, 6, 7], "weights": ['uniform', 'distance'], "algorithm": ['auto', 'ball_tree', 'kd_tree', 'brute'], "n_jobs": [-1, 1] }})

```

```
Google product owner certification - Go... JupyterLab
localhost:8888/lab
File Edit View Run Kernel Tabs Settings Help
Project / Capstone-Project-master/
Name
data
AIML Capstone - Gro...
Coverage.xlsx
FullDBSCAN.ipynb
Mail_footfall_predict...
preprocess.py
Project Plan.xlsx
README.md
visualization.py
work_flow.ppt
Launcher Mail_footfall_prediction_skeleton.ipynb Python 3
self.grid_params.update({ "knn_regressor": { "n_neighbors": [4,5,6,7], "weights": ['uniform', 'distance'], "algorithm": ['auto', 'ball_tree','kd_tree','brute'], "leaf_size": [20,30,40], "p": [1,2], "metric": ['chebyshev','minkowski'] } })
self.grid_params.update({ "decisiontree_regressor": { "criterion": ['mse', 'friedman_mse','mae'], "splitter": ['best','random'], "max_depth": [2,3,4,5,6], "min_samples_split": [0.05, 0.01, 0.05, 0.1], "min_samples_leaf": [0.05, 0.01, 0.05, 0.1], "max_features": ['auto','sqrt','log2'] } })
self.grid_params.update({ "extra_trees_regression": { "criterion": ['mse', 'friedman_mse', 'mae'], "max_depth": [2,3,4,5,6], "max_features": ['auto','sqrt','log2'] } })
self.grid_params.update({ "svm_regression": { "kernel": ['linear', 'poly','rbf','sigmoid'], "gamma": ['scale','auto'], "tol": [1e-03,1e-04,1e-05], "C": [1.0, 0.1, 0.001,0.0001], "shrinking": [True,False] } })
self.grid_params.update({ "adaBoost": { "base_estimator": [DecisionTreeRegressor(max_depth = x) for x in range(1,6)], "learning_rate": [0.001, 0.01, 0.05, 0.1, 0.25, 0.5, 0.75, 1.0], "n_estimators": [100,200,300] } })
self.grid_params.update({ "bagging": { "base_estimator": [DecisionTreeRegressor(max_depth = x) for x in range(1,6)], "n_estimators": [100,200,300], "max_features": [0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1.0] } })
self.grid_params.update({ "random_forest": { "n_estimators": [100,200,300], "max_features": ["auto", "sqrt", "log2"], "max_depth": [3, 4, 5, 6, 7, 8], "min_samples_split": [0.05, 0.01, 0.05, 0.1], "min_samples_leaf": [0.05, 0.01, 0.05, 0.1] } })
self.grid_params.update({ "extra_trees_regression": { "n_estimators": [100,200,300], "criterion": ['mse', 'friedman_mse', 'mae'], "min_samples_split": [2, 3], "max_features": ['auto','sqrt','log2'] } })
self.grid_params.update({ "gradient_boosting_regression": {"learning_rate":[0.15,0.1,0.05,0.01,0.005,0.001], "n_estimators": [100,200,300], "max_depth": [2,3,4,5,6], "min_child_weight":range(1,6,2), "gamma": [1/10.0 for i in range(0,5)], "estimators": [100,200,300], "subsample": [i/10.0 for i in range(6,10)], "colsample_bytree": [i/10.0 for i in range(6,10)], "reg_alpha": [1e-5, 1e-2, 0.1, 1, 100], "reg_lambda": [1e-4, 1e-2, 0.1, 1, 100], "learning_rate": [0.001, 0.01, 0.1, 0.2, 0.3] } })
self.grid_params.update({ "LightGBM": {"max_depth":range(3,10,2), "min_child_weight":range(1,6,2), "gamma": [1/10.0 for i in range(0,5)], "n_estimators": [100,200,300], "max_depth": [2,3,4,5,6], "min_child_weight":range(1,6,2), "gamma": [1/10.0 for i in range(0,5)], "n_estimators": [100,200,300], "subsample": [i/10.0 for i in range(6,10)], "colsample_bytree": [i/10.0 for i in range(6,10)], "reg_alpha": [1e-5, 1e-2, 0.1, 1, 100], "reg_lambda": [1e-4, 1e-2, 0.1, 1, 100], "learning_rate": [0.001, 0.01, 0.1, 0.2, 0.3] } })
self.grid_params.update({ "categoricalGradientBoosting": {"max_depth":range(3,10,2), "n_estimators": [100,200,300], "subsample": [i/10.0 for i in range(6,10)], "learning_rate": [0.001, 0.01, 0.1, 0.2, 0.3] } })
self.grid_params.update({ "NeuralNet_MLP": {"hidden_layer_sizes":[(100, ), (100,200 )], "activation": ["tanh",'relu','logistic'], "solver":['sgd','adam','lbfgs'], "alpha": [0.001, 0.0001], "batch_size": [100,200,300], "learning_rate":['constant','invscaling','adaptive'], "max_iter": [100,200] } })
Mode Command L1 Col 1 Mail_footfall_prediction_skeleton.ipynb
ENGLISH 11:07 AM INTL 11/10/2020
```

```
Google product owner certification - Go... JupyterLab
localhost:8888/lab
File Edit View Run Kernel Tabs Settings Help
Project / Capstone-Project-master/
Name
data
AIML Capstone - Gro...
Coverage.xlsx
FullDBSCAN.ipynb
Mail_footfall_predict...
preprocess.py
Project Plan.xlsx
README.md
visualization.py
work_flow.ppt
Launcher Mail_footfall_prediction_skeleton.ipynb Python 3
self.grid_params.update({ "gradient_boosting_regression": {"learning_rate": [0.15, 0.1, 0.05, 0.01, 0.005, 0.001], "n_estimators": [100, 200, 300], "max_depth": [2, 3, 4, 5, 6], "min_child_weight": range(1, 6, 2), "gamma": [1/10.0 for i in range(0, 5)], "estimators": [100, 200, 300], "subsample": [i/10.0 for i in range(6, 10)], "colsample_bytree": [i/10.0 for i in range(6, 10)], "reg_alpha": [1e-5, 1e-2, 0.1, 1, 100], "reg_lambda": [1e-4, 1e-2, 0.1, 1, 100], "learning_rate": [0.001, 0.01, 0.1, 0.2, 0.3] } })
self.grid_params.update({ "XGBoost": {"max_depth": range(3, 10, 2), "min_child_weight": range(1, 6, 2), "gamma": [1/10.0 for i in range(0, 5)], "n_estimators": [100, 200, 300], "subsample": [i/10.0 for i in range(6, 10)], "colsample_bytree": [i/10.0 for i in range(6, 10)], "reg_alpha": [1e-5, 1e-2, 0.1, 1, 100], "reg_lambda": [1e-4, 1e-2, 0.1, 1, 100], "learning_rate": [0.001, 0.01, 0.1, 0.2, 0.3] } })
self.grid_params.update({ "LightGBM": {"max_depth": range(3, 10, 2), "min_child_weight": range(1, 6, 2), "gamma": [1/10.0 for i in range(0, 5)], "n_estimators": [100, 200, 300], "subsample": [i/10.0 for i in range(6, 10)], "colsample_bytree": [i/10.0 for i in range(6, 10)], "reg_alpha": [1e-5, 1e-2, 0.1, 1, 100], "reg_lambda": [1e-4, 1e-2, 0.1, 1, 100], "learning_rate": [0.001, 0.01, 0.1, 0.2, 0.3] } })
self.grid_params.update({ "categoricalGradientBoosting": {"max_depth": range(3, 10, 2), "n_estimators": [100, 200, 300], "subsample": [i/10.0 for i in range(6, 10)], "learning_rate": [0.001, 0.01, 0.1, 0.2, 0.3] } })
self.grid_params.update({ "NeuralNet_MLP": {"hidden_layer_sizes": [(100, ), (100, 200 )], "activation": ["tanh", "relu", "logistic"], "solver": ['sgd', 'adam', 'lbfgs'], "alpha": [0.001, 0.0001], "batch_size": [100, 200, 300], "learning_rate": ['constant', 'invscaling', 'adaptive'], "max_iter": [100, 200] } })
Mode Command L1 Col 1 Mail_footfall_prediction_skeleton.ipynb
ENGLISH 11:07 AM INTL 11/10/2020
```

```
def prepare_datasets(self, shuffle=False, extra_input_features=False, pca_features=False, weekly_data = False):
    # create dictionary to store different datasets
    # keys : LSTM, Regression, LSTM_detrend, Regression_detrend
    dataset_dict = {}

    if pca_features == False:
        if extra_input_features == True:
            input_cols = global_parameters['input_cols_reg_with_extra_features']
        else:
            input_cols = global_parameters['input_cols_regression']

    if weekly_data == True:
        input_cols = global_parameters['input_cols_regression_weekly']
        n_timesteps = 4 weeks Log
        n_forecast_steps = 1 # 1 week forecast

    # prepare normal regression dataset for train dataset
    model_dataset = train_test_validation_data_set(self.dataset[:global_parameters['test_data_set_weeks']])
    X_train, y_train, X_columns = model_dataset.prepare_regression_datasets(n_timesteps=n_timesteps, n_forecast_steps=n_forecast_steps, input_cols = input_cols, shuffle=shuffle)

    # prepare normal regression dataset for test dataset
    model_dataset = train_test_validation_data_set(global_parameters['test_data_set_weeks'])
    X_test, y_test, _ = model_dataset.prepare_regression_datasets(n_timesteps=n_timesteps, n_forecast_steps=n_forecast_steps, input_cols = input_cols)

    # prepare normal regression dataset for train dataset
    model_dataset = train_test_validation_data_set(global_parameters['data_set_days'])
    X_train, y_train, X_columns = model_dataset.prepare_regression_datasets(input_cols = input_cols, shuffle=shuffle)

    # prepare normal regression dataset for test dataset
    model_dataset = train_test_validation_data_set(global_parameters['test_data_set_days'])
    X_test, y_test, _ = model_dataset.prepare_regression_datasets(input_cols = input_cols)

    # prepare PCA dataset for train dataset
    model_dataset = train_test_validation_data_set(self.dataset[:global_parameters['data_set_days']]:global_parameters['test_data_set_days'])

    X_train, y_train, X_columns = model_dataset.prepare_pca_datasets(shuffle=shuffle)

    # prepare PCA dataset for test dataset
    model_dataset = train_test_validation_data_set(global_parameters['test_data_set_days'])

    X_test, y_test, _ = model_dataset.prepare_pca_datasets(input_cols = input_cols)

    if self.dataset_type == 'detrend':
        # prepare PCA dataset for train dataset
        model_dataset = train_test_validation_data_set(self.dataset[:global_parameters['data_set_days']]:global_parameters['test_data_set_days'])

        X_train, y_train, X_columns = model_dataset.prepare_pca_datasets(shuffle=shuffle)

        # prepare PCA dataset for test dataset
        model_dataset = train_test_validation_data_set(global_parameters['test_data_set_days'])

        X_test, y_test, _ = model_dataset.prepare_pca_datasets(input_cols = input_cols)
```

```
if self.dataset_type == 'detrend':
    self.trend = self.dataset['trend'][-len(y_test):].values
    self.seasonal = self.dataset['seas'][-len(y_test):].values

dataset_dict['regression'+self.dataset_type] = {"X_train": X_train, "y_train": y_train, "X_test": X_test, "y_test":y_test, "X_columns": X_columns}

return dataset_dict

# prepare a list of ml models
def get_models(self, models_dict):
    # linear models
    models["linear_regression"] = LinearRegression()
    models["lassoRegressor"] = Lasso()
    models["ridgeRegressor"] = Ridge()
    models["elastic_net"] = ElasticNet()
    models["passive_aggressive"] = PassiveAggressiveRegressor(max_iter=1000, tol=1e-3)
    models["sgd"] = SGDRegressor(max_iter=1000, tol=1e-3)
    # non-linear models
    models["knn_regressor"] = KNeighborsRegressor(n_neighbors=7)
    models["decisiontree_regressor"] = DecisionTreeRegressor()
    models["extra_tree_regressor"] = ExtraTreeRegressor()
    models["extra_random_forest_regressor"] = ExtraRandomForestRegressor()
    models["svr_regressor"] = SVR()

    # ensemble models
    n_trees = 200
    models["adaboost"] = AdaBoostRegressor(n_estimators=n_trees)
    models["bagging"] = BaggingRegressor(n_estimators=n_trees)
    models["random_forest"] = RandomForestRegressor(n_estimators=n_trees)
    models["extra_trees_regressor"] = ExtraTreesRegressor(n_estimators=n_trees)
    models["gradientboosting"] = GradientBoostingRegressor(n_estimators=n_trees)
    models["xgbboost"] = XGBRegressor(objective="reg:squarederror")
    models["lightgbmboosting"] = LightGBMRegressor()
    models["categoricalgradientboosting"] = CatBoostRegressor(n_estimators=n_trees, verbose=0)
    # neural network models
    models["neuralnet_tf"] = MLPRegressor()

    print("Defined %d models" % len(models))
    return models

def randomizedSearch(self, paramgrid, model, X_train, y_train, X_test, y_test, name):
    # train and predict each model with randomized search CV
    kFold = KFold(shuffle=True, random_state=50)
    rscv = RandomizedSearchCV(model, param_distributions=paramGrid, n_jobs=-1, scoring="neg_mean_squared_error", cv=kFold, n_iter=10, verbose=0, random_state=42)
```

product owner certification - Google

JupyterLab

localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

Python 3

```
#!/usr/bin/env python3
# coding: utf-8
# KFold = KFold(shuffle=True, random_state=50)
rscv = RandomizedSearchCV(model, param_distributions=paramGrid, n_jobs=-1, scoring='neg_mean_squared_error', cv=kFold,
    n_iter=10, verbose=0, random_state=42)
# Fit rscv
print("Now tuning (model).")
rscv.fit(X_train, y_train)

prediction = rscv.best_estimator_.predict(X_test)
test_loss = mean_squared_error(y_test, prediction)
return rscv.best_estimator_, rscv.best_params_, prediction, test_loss

def mape(self,actual,forecast):
    return np.mean(np.abs( forecast - actual))/np.abs(actual)*100 # MAPE

def rmse(self,actual,forecast):
    return np.sqrt(np.mean((forecast - actual)**2)) # RMSE

def mae(self,actual,forecast):
    return np.mean(np.abs( forecast - actual)) # MAE

def mpe(self,actual,forecast):
    return np.mean((forecast - actual)/actual)*100 # MPE
def mase(self,actual,forecast):
    return np.mean(np.abs( forecast - actual)) # MASE
def mspe(self,actual,forecast):
    return np.mean((forecast - actual)**2)*.5 # RMSE
def mse(self,actual,forecast):
    return np.mean((forecast - actual)**2) # MSE

def crossVal(self,model,X_test,y_test):
    # predict with cross_val_predict with CV
    scoring = ("mape": make_scoring(self.mape),
               "mse": make_scoring(self.mse), "mae": make_scoring(self.mae),
               "mpe": make_scoring(self.mpe), "rmse":make_scoring(self.rmse), "R2 score":r2_score
              )

    kFold = KFold(shuffle=True, random_state=50)

    cvPredictions = cross_val_predict(model, X_test, y_test, cv=kFold)
    score = cross_validate(model, X_test, y_test, cv=kFold, scoring=scoring)

    loss_values = {'mape':np.mean(score['test_mape']), 'mse':np.mean(score['test_mse']), 'mae': np.mean(score['test_mae']),
                  'mpe': np.mean(score['test_mpe']), 'rmse':np.mean(score['test_rmse']), 'R2 score':np.mean(score['test_R2 score'])}
    return loss_values

def rfECV(self, model,X_train,y_train,X_test,y_test,X_columns, n_params,name):
    rfcv = RFECV(estimator=model, step=1,cv= KFold(), scoring= 'r2')
    rfcv.fit(X_train,y_train)
    prediction_rfcv = rfcv.predict(X_test)
    if(self.dataset_type == 'time series' or self.invert_trend):
        # inverse trend for the predictions
        prediction_rfcv = self.inverse_detrend(prediction_rfcv,self.trend,self.seasonal)

    loss_values = model.utility().calculate_metrics(y_test,prediction_rfcv,len(y_test),n_params,name)
    print("RFECV - optimal number of features ", rfcv.n_features_)

    optimal_features = []
    for i in range(len(X_columns)):
        if rfcv.ranking_[i]==1:
            optimal_features.append(X_columns[i])
    print("## CV - optional features ", optimal_features)

    plt.figure(figsize=(16, 8))
    plt.title('Recursive Feature Elimination with CV - ' + name, fontsize=18, fontweight='bold', pad=20)
    plt.xlabel('Number of features selected', fontsize=14, labelpad=20)
    plt.ylabel('R^2 score', fontsize=14, labelpad=20)
    plt.plot(range(1, len(rfcv.grid_scores_) + 1), rfcv.grid_scores_, color='#303F9F', linewidth=3)
    plt.show()

    return prediction_rfcv,loss_values

# fit a single model
def fit_model(self, model, X, y):
    # clone the model configuration
    local_model = clone(model)
    # Fit the model
    local_model.fit(X, y)
    return local_model

def visualizeFeatureImportance(self,model,X_train,y_train,X_columns,name):
    if name not in ['knn_regressor','svm_regressor','bagging','LightGradientBoosting','CategoricalGradientBoosting','NeuralNet_MLP']:
        fig = plt.figure(figsize=(16, 8))
        visualize = FeatureImportances(model_fit,ax=plt.axes(),is_fitted=True, labels=X_columns)
        visualize.fit(X_train, y_train)
        visualize.show()

def predictionErrorPlot(self,model,X_train,y_train,X_test,y_test):
    fig = plt.figure(figsize=(16, 8))
    visualizer = PredictionError(model, ax=plt.axes(), is_fitted=True)
    visualizer.fit(X_train, y_train)
    visualizer.show()
```

product owner certification - Google

JupyterLab

localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

Python 3

```
def rfECV(self, model,X_train,y_train,X_test,y_test,X_columns, n_params,name):
    rfcv = RFECV(estimator=model, step=1,cv= KFold(), scoring= 'r2')
    rfcv.fit(X_train,y_train)
    prediction_rfcv = rfcv.predict(X_test)
    if(self.dataset_type == 'time series' or self.invert_trend):
        # inverse trend for the predictions
        prediction_rfcv = self.inverse_detrend(prediction_rfcv,self.trend,self.seasonal)

    loss_values = model.utility().calculate_metrics(y_test,prediction_rfcv,len(y_test),n_params,name)
    print("RFECV - optimal number of features ", rfcv.n_features_)

    optimal_features = []
    for i in range(len(X_columns)):
        if rfcv.ranking_[i]==1:
            optimal_features.append(X_columns[i])
    print("## CV - optional features ", optimal_features)

    plt.figure(figsize=(16, 8))
    plt.title('Recursive Feature Elimination with CV - ' + name, fontsize=18, fontweight='bold', pad=20)
    plt.xlabel('Number of features selected', fontsize=14, labelpad=20)
    plt.ylabel('R^2 score', fontsize=14, labelpad=20)
    plt.plot(range(1, len(rfcv.grid_scores_) + 1), rfcv.grid_scores_, color='#303F9F', linewidth=3)
    plt.show()

    return prediction_rfcv,loss_values

# fit a single model
def fit_model(self, model, X, y):
    # clone the model configuration
    local_model = clone(model)
    # Fit the model
    local_model.fit(X, y)
    return local_model

def visualizeFeatureImportance(self,model,X_train,y_train,X_columns,name):
    if name not in ['knn_regressor','svm_regressor','bagging','LightGradientBoosting','CategoricalGradientBoosting','NeuralNet_MLP']:
        fig = plt.figure(figsize=(16, 8))
        visualize = FeatureImportances(model_fit,ax=plt.axes(),is_fitted=True, labels=X_columns)
        visualize.fit(X_train, y_train)
        visualize.show()

def predictionErrorPlot(self,model,X_train,y_train,X_test,y_test):
    fig = plt.figure(figsize=(16, 8))
    visualizer = PredictionError(model, ax=plt.axes(), is_fitted=True)
    visualizer.fit(X_train, y_train)
    visualizer.show()
```

```
# In -> %%pylab inline
# In -> %load Mail_footfall_prediction_skeleton.ipynb
# In -> # visualize = PredictionError(model, ax=plt.axes(), is_Fitted=True)
# In -> visualizer.fit(X_train, np.ravel(y_train))
# In -> y_test = np.ravel(y_test)
# In -> visualizer.score(X_test, y_test)
# In -> visualizer.show()

# evaluate_a_suite_of_models
def train_and_evaluate_models(self, RFECross_validation = True, shuffle=False, extra_input_features=False, pcaFeatures=False, weekly_data=False):
    models = self.get_models()
    trained_models = []
    dataset_dict = self.prepare_datasets(shuffle, extra_input_features, pcaFeatures, weekly_data)

    X_train = dataset_dict['regression'] + self.dataset_type[[ '_X_train' ]]
    y_train = dataset_dict['regression'] + self.dataset_type[[ '_y_train' ]]
    X_test = dataset_dict['regression'] + self.dataset_type[[ '_X_test' ]]
    y_test = dataset_dict['regression'] + self.dataset_type[[ '_y_test' ]]
    X_columns = dataset_dict['regression'] + self.dataset_type[[ '_X_columns' ]]

    # Initialize the prediction data frame
    modelPredictions = np.zeros((X_test.shape[0], len(models)))
    modelPredictions = pd.DataFrame(modelPredictions)

    cvPredictions = np.zeros((X_test.shape[0], len(models)))
    cvPredictions = pd.DataFrame(cvPredictions)

    rscvPredictions = np.zeros((X_test.shape[0], len(models)))
    rscvPredictions = pd.DataFrame(rscvPredictions)

    rfcvPredictions = np.zeros((X_test.shape[0], len(models))) # not applicable for 'knn', 'svm', 'naive bayes', 'mlp-nn', 'qda', 'bagging'
    rfcvPredictions = pd.DataFrame(rfcvPredictions)

    n_forecast_steps = global_parameters['n_forecast_steps']
    loss_values = {}

    n_params = X_train.shape[1] + 1
    for i, (name, model) in enumerate(models.items()):

        # 1. fit default models
        model_fit = self.fit_model(model, X_train, y_train)
        # make predictions
        predictions = model_fit.predict(X_test)
        print("*****")
        print(name)

        n_params = X_train.shape[1] + 1
        for i, (name, model) in enumerate(models.items()):

            # 1. fit default model fitting...
            if self.dataset_type == 'detrend':# inverse the detrend by adding trend and seasonal info
                # inverse detrend for the predictions
                predictions = self.inverse_detrend(predictions, self.trend, self.seasonal)
                y_test_real = self.inverse_detrend(y_test, self.trend, self.seasonal)
            else:
                y_test_real = y_test

            loss_values[name+'_'+self.dataset_type+'_'+model._default_model'] = model_utility().calculate_metrics(y_test_real,predictions,len(y_test_real),n_params,name )

            #inverse transform
            predictions = np.reshape(predictions, (predictions.shape[0], n_forecast_steps))
            predictions = scaler_target.inverse_transform(predictions)

            modelPredictions.loc[:, i] = predictions
            modelPredictions = modelPredictions.rename(columns={i: name})

            trained_models[name+self.dataset_type+'standard'] = model_fit

        # 2. fit randomized CV models
        model_fit_rscv, best_params, predicton_rscv, loss_rscv = self.randomizedSearch(self.grid_params[name],model,X_train, y_train,X_test, y_test, name)
        # visualize the feature importance for the predicton_rscv
        # self.grid_params[name].feature_importance(model_fit_rscv,X_train,y_train,X_columns, name)
        # self.grid_params[name].feature_importance(model_fit_rscv,X_train,y_train,X_test,y_test)

        if self.dataset_type == 'detrend':# inverse the detrend by adding trend and seasonal info
            # inverse detrend for the predictions
            predicton_rscv = self.inverse_detrend(predicton_rscv, self.trend, self.seasonal)
            y_test_real = self.inverse_detrend(y_test, self.trend, self.seasonal)

        else:
            y_test_real = y_test

        loss_values[name+'_'+self.dataset_type+'_'+rscv_model'] = model_utility().calculate_metrics(y_test_real,predicton_rscv,len(y_test_real),n_params,name )

        #inverse transform
        predicton_rscv = np.reshape(predicton_rscv, (predicton_rscv.shape[0], n_forecast_steps))
        predicton_rscv = scaler_target.inverse_transform(predicton_rscv)

        trained_models[name+self.dataset_type+'rscv'] = model_fit_rscv
        rscvPredictions.loc[:, i] = predicton_rscv
        rscvPredictions = rscvPredictions.rename(columns={i: name})
```

```
# In -> %%pylab inline
# In -> %load Mail_footfall_prediction_skeleton.ipynb
# In -> # visualize = PredictionError(model, ax=plt.axes(), is_Fitted=True)
# In -> visualizer.fit(X_train, np.ravel(y_train))
# In -> y_test = np.ravel(y_test)
# In -> visualizer.score(X_test, y_test)
# In -> visualizer.show()

# evaluate_a_suite_of_models
def train_and_evaluate_models(self, RFECross_validation = True, shuffle=False, extra_input_features=False, pcaFeatures=False, weekly_data=False):
    models = self.get_models()
    trained_models = []
    dataset_dict = self.prepare_datasets(shuffle, extra_input_features, pcaFeatures, weekly_data)

    X_train = dataset_dict['regression'] + self.dataset_type[[ '_X_train' ]]
    y_train = dataset_dict['regression'] + self.dataset_type[[ '_y_train' ]]
    X_test = dataset_dict['regression'] + self.dataset_type[[ '_X_test' ]]
    y_test = dataset_dict['regression'] + self.dataset_type[[ '_y_test' ]]
    X_columns = dataset_dict['regression'] + self.dataset_type[[ '_X_columns' ]]

    # Initialize the prediction data frame
    modelPredictions = np.zeros((X_test.shape[0], len(models)))
    modelPredictions = pd.DataFrame(modelPredictions)

    cvPredictions = np.zeros((X_test.shape[0], len(models)))
    cvPredictions = pd.DataFrame(cvPredictions)

    rscvPredictions = np.zeros((X_test.shape[0], len(models)))
    rscvPredictions = pd.DataFrame(rscvPredictions)

    rfcvPredictions = np.zeros((X_test.shape[0], len(models))) # not applicable for 'knn', 'svm', 'naive bayes', 'mlp-nn', 'qda', 'bagging'
    rfcvPredictions = pd.DataFrame(rfcvPredictions)

    n_forecast_steps = global_parameters['n_forecast_steps']
    loss_values = {}

    n_params = X_train.shape[1] + 1
    for i, (name, model) in enumerate(models.items()):

        # 1. fit default model fitting...
        if self.dataset_type == 'detrend':# inverse the detrend by adding trend and seasonal info
            # inverse detrend for the predictions
            predictions = self.inverse_detrend(predictions, self.trend, self.seasonal)
            y_test_real = self.inverse_detrend(y_test, self.trend, self.seasonal)
        else:
            y_test_real = y_test

        loss_values[name+'_'+self.dataset_type+'_'+model._default_model'] = model_utility().calculate_metrics(y_test_real,predictions,len(y_test_real),n_params,name )

        #inverse transform
        predictions = np.reshape(predictions, (predictions.shape[0], n_forecast_steps))
        predictions = scaler_target.inverse_transform(predictions)

        modelPredictions.loc[:, i] = predictions
        modelPredictions = modelPredictions.rename(columns={i: name})

        trained_models[name+self.dataset_type+'standard'] = model_fit

    # 2. fit randomized CV models
    model_fit_rscv, best_params, predicton_rscv, loss_rscv = self.randomizedSearch(self.grid_params[name],model,X_train, y_train,X_test, y_test, name)
    # visualize the feature importance for the predicton_rscv
    # self.grid_params[name].feature_importance(model_fit_rscv,X_train,y_train,X_columns, name)
    # self.grid_params[name].feature_importance(model_fit_rscv,X_train,y_train,X_test,y_test)

    if self.dataset_type == 'detrend':# inverse the detrend by adding trend and seasonal info
        # inverse detrend for the predictions
        predicton_rscv = self.inverse_detrend(predicton_rscv, self.trend, self.seasonal)
        y_test_real = self.inverse_detrend(y_test, self.trend, self.seasonal)

    else:
        y_test_real = y_test

    loss_values[name+'_'+self.dataset_type+'_'+rscv_model'] = model_utility().calculate_metrics(y_test_real,predicton_rscv,len(y_test_real),n_params,name )

    #inverse transform
    predicton_rscv = np.reshape(predicton_rscv, (predicton_rscv.shape[0], n_forecast_steps))
    predicton_rscv = scaler_target.inverse_transform(predicton_rscv)

    trained_models[name+self.dataset_type+'rscv'] = model_fit_rscv
    rscvPredictions.loc[:, i] = predicton_rscv
    rscvPredictions = rscvPredictions.rename(columns={i: name})
```

product owner certification - Go x JupyterLab x

localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

Launcher Mail_footfall_prediction_ske

Markdown v

```
trained_models[name==self.dataset_type+'rscv'] = model_fit_rscv
rscvPredictions.iloc[:, i] = predictor_rscv
rscvPredictions = rscvPredictions.rename(columns={i: name})

# 3. cross CV models # 7000
print("Cross CV validation...")
predictor_crossCV, loss_cross_CV_values = self.crossCV(model_fit_rscv,X_test,y_test_real)

cvPredictions.iloc[:, i] = predictor_crossCV
cvPredictions = cvPredictions.rename(columns={i: name})

loss_values[name+'_'+self.dataset_type+'_crosscv_model'] = loss_cross_CV_values

if RFECv_Cross_validation == True:
    # 4. RFE CV # 7000
    if name not in ['knn_regressor','svm_regressor','bagging','LightGradientBoosting','CategoricalGradientBoosting','NeuralNet_MLP']:
        print("RFE CV tuning...")
        prediction_rfecv, loss_rfecv = self.RFECv(model_fit_rscv,X_train,y_train,X_test,y_test_real,X_columns,n_params,name)

        loss_values[name+'_'+self.dataset_type+'_rfecv_model'] = loss_rfecv

        rfecvPredictions.iloc[:, i] = prediction_rfecv
        rfecvPredictions = rfecvPredictions.rename(columns={i: name})

y_test_real = np.reshape(y_test_real,(y_test_real.shape[0], n_forecast_steps))
y_test_real = scaler_target.inverse_transform(y_test_real)

self.plot_predictions_models(model1Predictions,y_test_real)
self.plot_predictions_models(rscvPredictions,y_test_real,model_type='Randomized Search CV')

loss_values = model1Utility().convert_dict_to_dataframe(loss_values)
return trained_models, model1Predictions,rscvPredictions, rfcvPredictions, loss_values

def plot(self,data,model_type):
    fig = plt.figure(figsize=(24,8))
    data.plot()
    plt.xlabel('Time')
    plt.ylabel('Daily Mall Footfall')
    plt.title('Mall Foot Fall prediction with ' + model_type)
    plt.legend()
    plt.show()
```

The screenshot shows a Jupyter Notebook interface with the following details:

- Title Bar:** product owner certification - Go X JupyterLab X
- Header:** localhost:8888/lab
- Toolbar:** Apps, Chrome, Google, BITS Pilani - Login, ACT Fibernet Portal..., Customer Experience..., AIIML_SK Links, Django: Web frame..., AI- Updates, BITS-Wheebox, BerkeleyX-AI OL Co... (with a Paused icon)
- File Menu:** File, Edit, View, Run, Kernel, Tabs, Settings, Help
- Launcher:** Mail_footfall_prediction_skeleton.ipynb
- Code Cell:** The cell contains Python code for a 'Mall_footfall_prediction_skeleton' notebook. It includes imports, plotting functions for different models, an inverse_detrend function, and a main plot_predictions function that handles different model types (Normal Regression, ARIMA, SARIMA, Exponential Smoothing) and forecast steps.
- Bottom Status Bar:** Mode: Command, Ln 1, Col 1, Mail_footfall_prediction_skeleton.ipynb, ENG, 1:11 AM, INTL, 11/10/2020

product owner certification - Google

JupyterLab

localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

Launcher

Project / Capstone - Project-master /

Name

data

AIML Capstone - Gro...

Coverage.xlsx

FullDBSCAN.ipynb

Mail_footfall_prediction.ipynb

Mail_footfall_predict...

preprocess.py

Project Plan.xlsx

README.md

visualization.py

work_flow.pptx

```

predicts = scaler_target.inverse_transform(predicts)
y_test_real = scaler_target.inverse_transform(y_test_real)

fig = plt.figure(figsize=(24,8))
plt.plot(y_test_real, color = 'red', label = 'Real mall footfall')
plt.plot(predicts, color = 'blue', label = 'Predicted mall footfall')
plt.title('Validation with ' + name)
plt.xlabel('Time')
plt.ylabel('Daily Mall Footfall')
plt.legend()
plt.show()

```

6.3 LSTM model classes

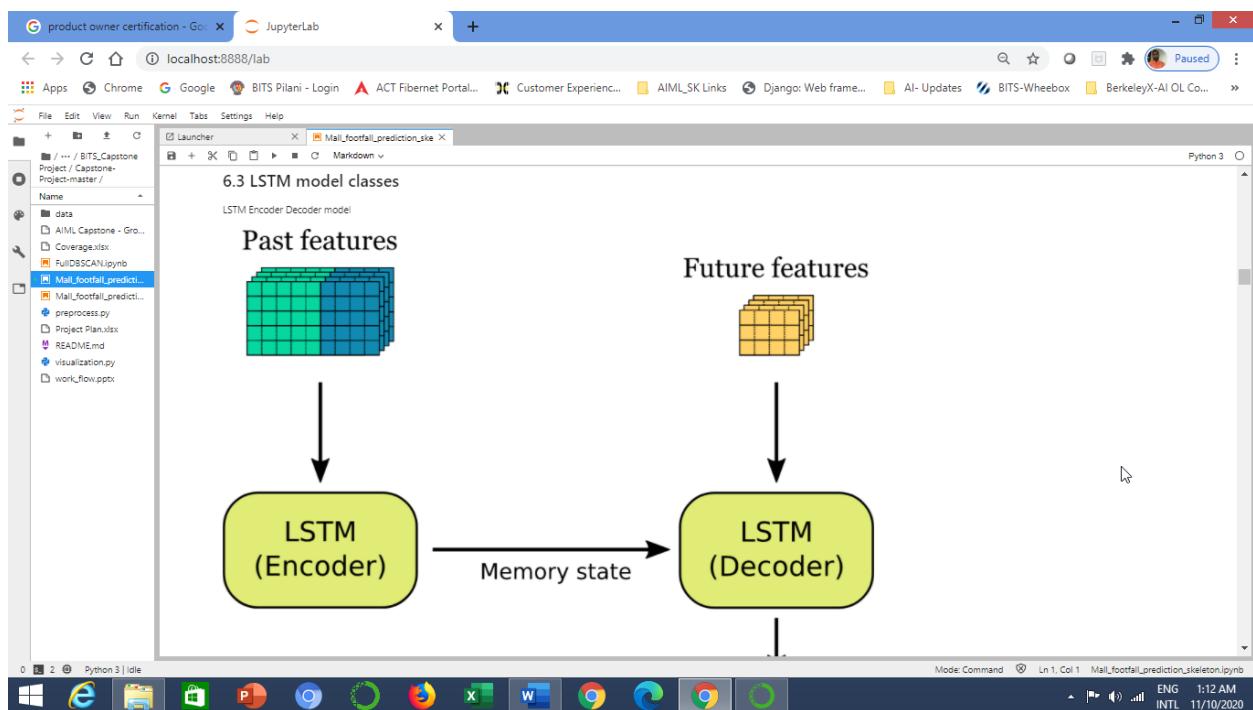
LSTM Encoder Decoder model

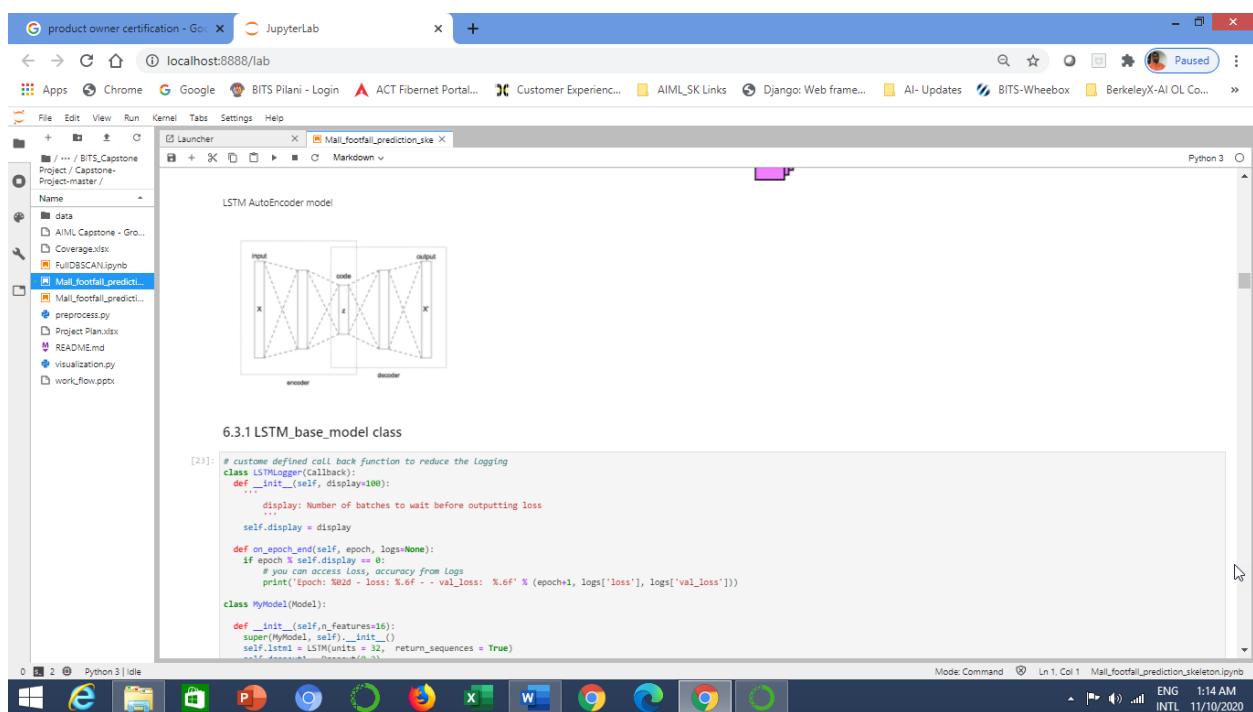
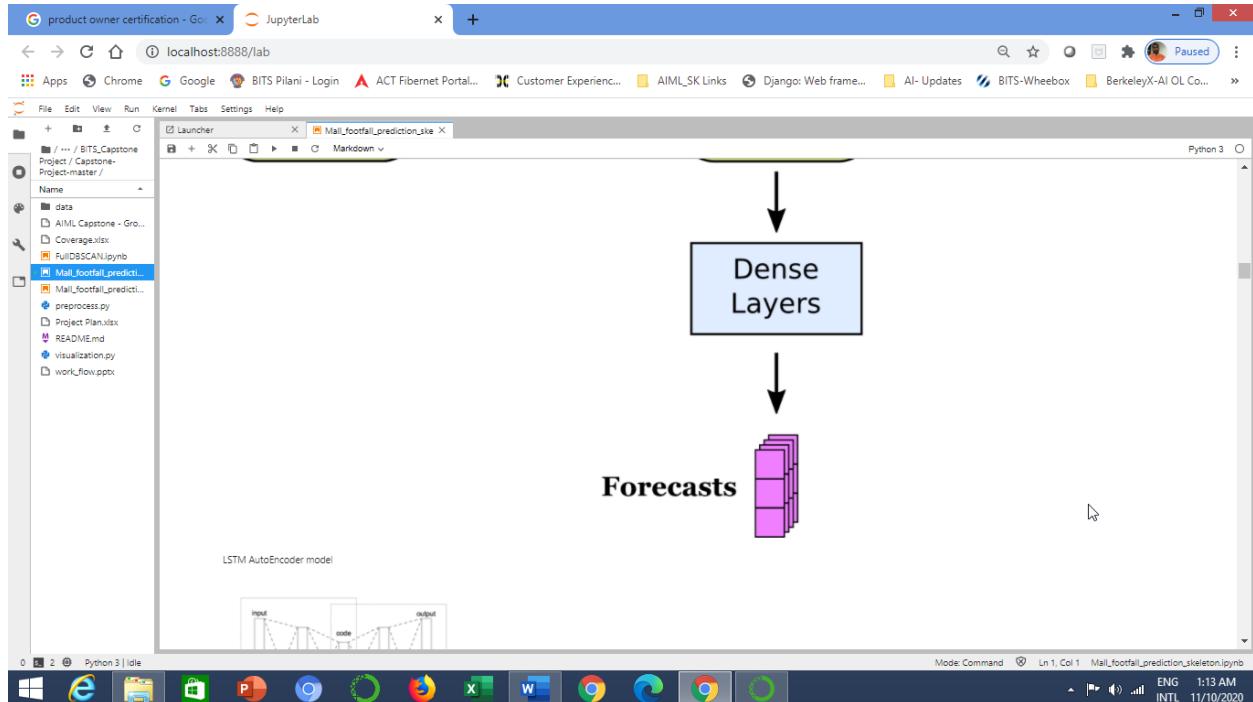
Past features

Future features

Mode Command Ln 1, Col 1 Mail_footfall_prediction_skeleton.ipynb

ENGLISH 11:11 AM INTL 11/10/2020





product owner certification - Go... JupyterLab

localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

Project / Capstone-Project-master /

Name

data

AIML Capstone - Gro...

Coverage.xlsx

FullDBSCAN.ipynb

Mail_footfall_predict...

Mail_footfall_prediction_skeleton.ipynb

preprocess.py

Project Plan.xlsx

READMEEnd

visualization.py

work_flow.ppt

```
your can access loss, accuracy true logs
print('Epoch: %d - loss: %.6f' % (epoch+1, logs['loss'], logs['val_loss']))
```

```
class MyModel(Model):
    def __init__(self,n_features=16):
        super().__init__()
        self.lstm1 = LSTM(units = 32, return_sequences = True)
        self.dropout1 = Dropout(0.2)
        self.lstm2 = LSTM(units = 16, return_sequences = False)
        self.dense1 = Dense(16, activation='relu')
        self.dense2 = Dense(16, activation='relu')
        self.output1 = Dense(1,activation='relu')

    def call(self, inputs):
        x = self.lstm1(inputs)
        x = self.dropout1(x)
        x = self.lstm2(x)
        x = self.dense1(x)
        x = self.dense2(x)
        x = self.output1(x)
        print(x.shape)
        print(inputs.shape)

        return inputs*x

class ResidualWrapper(Model):
    def __init__(self, model):
        super().__init__()
        self.model = model

    def call(self, inputs, *args, **kwargs):
        delta = self.model(inputs, *args, **kwargs)

        # The prediction for each timestep is the input
        # from the previous time step plus the delta
        # calculated by the model.
        return inputs + delta
```

```
[24]: class LSTM_base_model:
    def __init__(self, n_forecast=1):
        self.n_forecast = n_forecast

    def LSTM_simple_model_detrended_data(self,layer_parameters=dict(), time_steps = 7,n_forecast = 7, n_features=1, model_name=None):
```

Mode Command ↵ Ln 1, Col 1 Mail_footfall_prediction_skeleton.ipynb

ENGLISH 1:16 AM INTL 11/10/2020

product owner certification - Go... JupyterLab

localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

Project / Capstone-Project-master /

Name

data

AIML Capstone - Gro...

Coverage.xlsx

FullDBSCAN.ipynb

Mail_footfall_predict...

Mail_footfall_prediction_skeleton.ipynb

preprocess.py

Project Plan.xlsx

READMEEnd

visualization.py

work_flow.ppt

```
[24]: class LSTM_base_model:
    def __init__(self, n_forecast=1):
        self.n_forecast = n_forecast

    def LSTM_simple_model_detrended_data(self,layer_parameters=dict(), time_steps = 7,n_forecast = 7, n_features=1, model_name=None):

        model_input = Input(shape=(time_steps, n_features))

        x = LSTM(units = layer_parameters["LSTM_layer1"], return_sequences = True)(model_input)
        x = Dropout(0.2)(x)
        x = LSTM(units = layer_parameters["LSTM_layer2"], return_sequences = False)(x)
        x = Dropout(0.2)(x)

        x = Dense(units = n_forecast)(x)
        model = Model(model_input, x, name=model_name)
        return model

    def LSTM_simple_model(self,layer_parameters=dict(), time_steps = 7,n_forecast = 7, n_features=1,model_name=None):

        model_input = Input(shape=(time_steps, n_features))

        x = LSTM(units = layer_parameters["LSTM_layer1"], return_sequences = True)(model_input)
        x = Dropout(0.2)(x)
        x = LSTM(units = layer_parameters["LSTM_layer2"], return_sequences = False)(x)
        x = Dropout(0.2)(x)
        x = Dense(layer_parameters["Dense_layer1"],activation='relu')(x)
        x = Dense(layer_parameters["Dense_layer2"],activation='relu')(x)

        x = Dense(units = n_forecast,activation='relu')(x)
        model = Model(model_input, x, name=model_name)
        return model
```

```
def CNN_LSTM_encoder_decoder(self,layer_parameters=dict(), time_steps = 7,n_forecast = 7, n_features=1, n_deterministic_features=4, n_subtime_steps=1,model_name=None):
    # (samples, time, rows, cols, channels)
    # CNN LSTM part input shape - (sub_timesteps, rows=1, col=time_steps/n_subtime_steps, channels)
    past_inputs = Input(shape=(n_subtime_steps, 1, int(time_steps/n_subtime_steps), n_features))
    x = Conv1D(n_deterministic_features, "ConvLSTM2D", filters=layer_parameters["ConvLSTM2D_layer1"], kernel_size=layer_parameters["kernel_size_layer1"],return_sequences=True)(past_inputs)
    x = LeakyReLU()(x)
    x = BatchNormalization()(x)
    x = Dropout(0.2)(x)
    x = Flatten()(x)
    x = RepeatVector(n_forecast)(x)

    # First branch of the net is an Lstm which finds an embedding for the past
    # Encoding the past
```

Mode Command ↵ Ln 1, Col 1 Mail_footfall_prediction_skeleton.ipynb

ENGLISH 1:16 AM INTL 11/10/2020

```

# First branch of the net is an lstm which finds an embedding for the past
# Encoding the past
encoder = LSTM(layer_parameters["LSTM_layer1"], return_state=True)
encoder_outputs, state_h, state_c = encoder(x)

future_inputs = Input(shape=(n_forecast, n_deterministic_features), name='future_inputs')
# Combining future inputs with recurrent branch output
decoder_lstm = LSTM(layer_parameters["LSTM_layer1"], return_sequences=True)
x = decoder_lstm(future_inputs, initial_state=[state_h, state_c])
x = Dense(layer_parameters["Dense_layer1"], activation='relu')(x)
x = Dense(layer_parameters["Dense_layer2"], activation='relu')(x)
output = Dense(units = n_forecast, activation='relu')(x)

model = Model(inputs=[past_inputs, future_inputs], outputs=output, name=model_name)
return model

def LSTM_model_with_extra_features(self,layer_parameters=dict(), time_steps = 7,n_forecast = 7, n_features=1, n_deterministic_features=4, n_subtime_steps=1, model_name=None):
    md1_input1 = Input(shape=(time_steps, n_features))

    x = LSTM(units = layer_parameters["LSTM_layer1"], return_sequences = True)(md1_input1)
    x = Dropout(0.2)(x)
    x = LSTM(units = layer_parameters["LSTM_layer2"])(x)
    x = Dropout(0.2)(x)

    md1_input2 = Input(shape=(n_deterministic_features,))
    #add extra features in the middle
    concat = Concatenate()([x, md1_input2])

    x = Dense(layer_parameters["Dense_layer1"], activation='relu')(concat)
    x = Dense(layer_parameters["Dense_layer2"], activation='relu')(x)
    output = Dense(units = n_forecast, activation='relu')(x)

    model = Model(inputs=[md1_input1, md1_input2], outputs=output, name=model_name)
    return model

def LSTM_encoder_decoder(self,layer_parameters=dict(),time_steps = 7,n_forecast = 7, n_features=1, n_deterministic_features=4,model_name=None):
    # First branch of the net is an lstm which finds an embedding for the past
    past_inputs = Input(shape=(time_steps, n_features), name='past_inputs')
    # Encoding the past
    encoder = LSTM(layer_parameters["LSTM_layer1"], return_state=True)
    encoder_outputs, state_h, state_c = encoder(past_inputs)

    future_inputs = Input(shape=(n_forecast, n_deterministic_features), name='future_input_1')
    # Combining future inputs with recurrent branch output
    decoder_lstm = LSTM(layer_parameters["LSTM_layer1"], return_sequences=True)
    x = decoder_lstm(future_inputs, initial_state=[state_h, state_c])
    x = Dense(layer_parameters["Dense_layer1"], activation='relu')(x)
    x = Dense(layer_parameters["Dense_layer2"], activation='relu')(x)

    output = Dense(units = n_forecast, activation='relu')(x)

    model = Model(inputs=[past_inputs, future_inputs], outputs=output, name=model_name)
    return model

def LSTM_encoder_decoder_with_renet(self,layer_parameters=dict(),time_steps = 7,n_forecast = 7, n_features=1, n_deterministic_features=4,model_name=None):
    # First branch of the net is an lstm which finds an embedding for the past
    past_inputs = Input(shape=(time_steps, n_features), name='past_inputs')
    # Encoding the past
    encoder = LSTM(layer_parameters["LSTM_layer1"], return_state=True)
    encoder_outputs, state_h, state_c = encoder(past_inputs)

    future_inputs = Input(shape=(n_forecast, n_deterministic_features), name='future_input_1')
    # Combining future inputs with recurrent branch output
    decoder_lstm = LSTM(layer_parameters["LSTM_layer1"], return_sequences=True)
    x = decoder_lstm(future_inputs, initial_state=[state_h, state_c])
    # Final step: Add shortcut value to main path, and pass it through a RELU activation
    x = Add()([x, past_inputs])

    x = Dense(layer_parameters["Dense_layer1"], activation='relu')(x)
    x = Dense(layer_parameters["Dense_layer2"], activation='relu')(x)

    output = Dense(units = n_forecast, activation='relu')(x)

    model = Model(inputs=[past_inputs, future_inputs], outputs=output, name=model_name)
    return model

```

```

def LSTM_auto_encoder_MLP(self, layer_parameters=dict(),n_encoded_features=5, n_timesteps=7, n_forecast=7, n_deterministic_features=4, model_name=None):
    # First branch of the net is an lstm which finds an embedding for the past
    past_inputs = Input(shape=(time_steps, n_features), name='past_inputs')
    # Encoding the past
    encoder = LSTM(layer_parameters["LSTM_layer1"], return_state=True)
    encoder_outputs, state_h, state_c = encoder(past_inputs)

    future_inputs = Input(shape=(n_forecast, n_deterministic_features), name='future_input_1')
    # Combining future inputs with recurrent branch output
    decoder_lstm = LSTM(layer_parameters["LSTM_layer1"], return_sequences=True)
    x = decoder_lstm(future_inputs, initial_state=[state_h, state_c])
    # Final step: Add shortcut value to main path, and pass it through a RELU activation
    x = Add()([x, past_inputs])

    x = Dense(layer_parameters["Dense_layer1"], activation='relu')(x)
    x = Dense(layer_parameters["Dense_layer2"], activation='relu')(x)

    output = Dense(units = n_forecast, activation='relu')(x)

    model = Model(inputs=[past_inputs, future_inputs], outputs=output, name=model_name)
    return model

```

```
def LSTM_auto_encoder_NLP(self, layer_parameters=dict(), n_encoded_features=5, n_timesteps=7, n_forecast=7, n_deterministic_features=4, model_name=None):
    past_inputs = Input(shape=(n_encoded_features), name='encoded_inputs')
    future_inputs = Input(shape=(n_deterministic_features), name='future_input_1')
    x = Concatenate([past_inputs, future_inputs])
    x = Dense(64, activation='relu')(x)
    x = Dense(64, activation='relu')(x)
    output = Dense(units=n_forecast, activation='relu')(x)

    model = Model(inputs=[past_inputs, future_inputs], outputs=output, name=model_name)
    return model

def prepare_auto_encoder_output(self, model_name, X_train, X_validation, X_test):
    train_encoded, validation_encoded, test_encoded = [], [], []

    fileexists, model_file = model_utility().CheckPointModelExisting(model_name)
    if fileexists == True:
        print('loaded %s' % model_file)
        auto_encoder_model = load_model(model_file)
        encoder = Model(inputs=auto_encoder_model.inputs, outputs=auto_encoder_model.layers[5].output, name=model_name+'encoder')
        train_encoded = encoder.predict(X_train)
        validation_encoded = encoder.predict(X_validation)
        test_encoded = encoder.predict(X_test)

    return train_encoded, validation_encoded, test_encoded

def LSTM_encoder_decoder_detrended_data(self, layer_parameters=dict(), time_steps=7, n_forecast=7, n_features=1, n_deterministic_features=4, model_name=None):

    # First branch of the net is an lstm which finds an embedding for the past
    past_inputs = Input(shape=(time_steps, n_features), name='past_inputs')
    # Encoding the past
    encoder = LSTM(layer_parameters['LSTM_layer1'], return_state=True)
    encoder_outputs, state_h, state_c = encoder(past_inputs)

    future_inputs = Input(shape=(n_forecast, n_deterministic_features), name='future_input_2')
    # Combining future inputs with recurrent branch output
    decoder_lstm = LSTM(layer_parameters['LSTM_layer1'], return_sequences=True)
    x = Decoder_lstm(future_inputs, initial_state=[state_h, state_c])
    x = Dense(layer_parameters['Dense_layer1'], activation='tanh')(x)
    v = Dense(layer_parameters['Dense_layer2'], activation='tanh')(x)
```

```
initial_state=[state_h, state_c]
x = Dense(layer_parameters['Dense_layer1'], activation='tanh')(x)
x = Dense(layer_parameters['Dense_layer2'], activation='tanh')(x)
output = Dense(units=n_forecast)(x)

model = Model(inputs=[past_inputs, future_inputs], outputs=output, name=model_name)
return model

def LSTM_auto_encoder(self, layer_parameters=dict(), time_steps=7, n_forecast=7, n_features=1, n_out_features=1, model_name=None):

    model_input = Input(shape=(time_steps, n_features))
    # encoder
    x = LSTM(layer_parameters['LSTM_layer1'], activation='relu', return_sequences=True)(model_input)
    x = Dropout(0.2)(x)
    x = LSTM(layer_parameters['LSTM_layer2'], activation='relu', return_sequences=True)(x)
    x = Dropout(0.2)(x)
    output = LSTM(n_out_features, activation='relu')(x)
    x = RepeatVector(time_steps)(output)
    # decoder
    x = LSTM(layer_parameters['LSTM_layer2'], activation='relu', return_sequences=True)(x)
    x = Dropout(0.2)(x)
    x = LSTM(layer_parameters['LSTM_layer1'], activation='relu', return_sequences=True)(x)
    x = Dropout(0.2)(x)
    x = TimeDistributed(Dense(n_out_features))(x)

    model = Model(model_input, x, name=model_name)
    return model

def ensemble(self, models, model_name=None):
    outputs = [model.outputs[0] for model in models]
    y = Average()(outputs)

    for i in range(len(models)):
        model = models[i]
        for layer in model.layers:
            # rename to avoid 'unique layer name' issue
            layer._name = 'default_ensemble_' + str(i+1) + '_' + layer.name

    model_inputs = [model.input for model in models]
    print(model_inputs)
    model = Model(model_inputs, y, name=model_name)
    return model
```

```
model = Model(model_inputs, y, name=model_name)
return model

# Load models from file
def load_all_LSTM_models(self, model_files):
    all_models = []
    for filename in model_files:
        if os.path.isfile(filename):
            model = load_model(filename)
            all_models.append(model)
            print('Loaded %s' % filename)
    return all_models

# Define stacked model from multiple member input models
def define_stacked_LSTM_ensemble_model(self, members):
    # Update old members in all models to not be trainable
    for i in range(len(members)):
        model = members[i]
        for layer in model.layers:
            if not layer.trainable:
                layer.trainable = True
                if len(layer.name) > 1 and 'unique' in layer.name:
                    layer.name = 'ensemble_' + str(i+1) + '_' + layer.name
    # Define multi-headed input
    ensemble_visible = [model.input for model in members]
    # Concatenate member output from each model
    model_output_list = []
    for model in members:
        if len(model.output.shape) == 2:
            model.output = K.reshape(model.output, (-1, 1))(model.output)
            print(model.output.shape)
            model.output_list.append(model.output)
        else:
            model.output_list.append(model.output)
    ensemble_outputs = [model.output for model_output in model_output_list]
    merge = concatenate(ensemble_outputs)
    hidden = Dense(32, activation='relu')(merge)
    hidden = Dense(16, activation='relu')(hidden)
    output = Dense(1, activation='linear')(hidden)
    model = Model(inputs=ensemble_visible, outputs=output)
    # Compile
    model.compile(optimizer='adam', loss='mean_squared_error')
    return model
```

```
# Fit a stacked model
def fit_stacked_LSTM_model(self, model, X_train, y_train,X_validation, y_validation, name):
    # Prepare input data
    checkpoint_path = 'saved_model_' + name +'.hdf5'
    fileexists, model_file = model_utility().checkPointModelExisting(name)
    if fileexists == False:
        # Train the model from scratch
        cp_callbacks = [ModelCheckpoint(filepath=checkpoint_path,
                                        save_weights_only=False,
                                        monitor='val_loss',
                                        mode='min',
                                        save_best_only=True,
                                        verbose=0)]
        # Fit model
        model.fit(X_train, y_train, epochs=300, verbose=1, validation_data=(X_validation, y_validation), callbacks=cp_callbacks)

    model_object = load_model(checkpoint_path)
    return model_object

def train_stacked_LSTM_ensemble_model(self,model_files,X_train,y_train,X_validation, y_validation, name):
    members = self.load_all_LSTM_models(model_files)
    print('Loaded %d models' % len(members))
    y_validation = np.array(y_validation)
    stacked_model = self.define_stacked_LSTM_ensemble_model(members)
    # Fit stacked model on test dataset
    stacked_model = self.fit_stacked_LSTM_model(stacked_model,X_train, y_train,X_validation, y_validation, name)
    #stacked_model.fit(X_test, y_test, epochs=300, verbose=1)
    return stacked_model

def scheduler(self, epoch,1):
    if epoch < 50:
        return 0.001
    elif epoch < 400:
        return 0.0001
    else:
        return 0.00001

def compile_and_train(self,model, data_set_dict, num_epochs = 10, batch_size=16, learning_rate = 0.0001,compile=True):
    lr_schedule = LearningRateScheduler(self.scheduler, verbose=0)
```

product owner certification - Go... JupyterLab localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

Launcher Mail_footfall_prediction_ske

def compile_and_train(self, model, data_set_dict, num_epochs = 10, batch_size=16, learning_rate = 0.0001,compile=True):
 lr_scheduler = LearningRateScheduler(self.scheduler, verbose=0)

 if compile == True:
 opt = optimizers.Adam(learning_rate=learning_rate, beta_1=0.9, beta_2=0.999, epsilon=1e-07, amsgrad=True)
 model.compile(optimizer = opt, loss = 'mean_squared_error')
 else:
 # update the learning rate
 print("INFO] old learning rate: (" + K.get_value(model.optimizer.lr))
 K.set_value(model.optimizer.lr, learning_rate)
 print("INFO] new learning rate: (" + K.get_value(model.optimizer.lr))

 # Output Loss every 100 epochs
 logger = LSTMLogger(display=100)

 checkpoint_path = "saved_model_ " + model.name +str(batch_size)+str(learning_rate)+".hdfs"
 cp_callbacks = [ModelCheckpoint(filepath=checkpoint_path,
 save_weights_only=False,
 mode="val_loss",
 mode="min",
 save_best_only=True,
 verbose=0),
 logger
]

 history = model.fit(data_set_dict["X_train"], data_set_dict["y_train"], epochs = num_epochs,
 validation_data=(data_set_dict["X_validation"], data_set_dict["y_validation"]),
 batch_size = batch_size, callbacks=cp_callbacks, verbose=0)

 return history.history['loss'], history.history['val_loss'], checkpoint_path

def plot_history(self, train_loss, val_loss):
 from matplotlib import pyplot
 fig = pyplot.figure(figsize=(16,8))
 pyplot.plot(train_loss, label="train")
 pyplot.plot(val_loss, label="validation")
 pyplot.plot(np.argmin(val_loss), np.min(val_loss), marker="x", color="r", label="best model")
 pyplot.legend()
 pyplot.show()

Python 3 | idle Mode: Command Ln 1 Col 1 Mail_footfall_prediction_skeleton.ipynb

product owner certification - Go X JupyterLab X

localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

Launcher Mail_footfall_prediction_skeleton.ipynb

6.3.2 LSTM_model_train_and_evaluation class

```
[25]: class LSTM_model_train_and_evaluation:  
    def __init__(self):  
        self.lstm_model_parameters = { "time_steps": global_parameters['n_timesteps'], "forecast_steps":global_parameters['n_forecast_steps'],  
        "features":len(global_parameters['input_is_neural_networks']), "deterministic_features":len(global_parameters['deterministic_features']),  
        "sub_timesteps":1, "batch_size":50, "out_features":5, "learning_rate":0.001 }  
        self.lstm_layer_parameters = { "cnn_lstm_encoder_decoder": { "Conv1TMD_layer1": 4, "kernel_size_layer1": (1,3), "LSTM_layer1": 8, "Dense_layer1": 8, "Dense_layer2":8 },  
        "lstm_encoder_decoder": { "LSTM_layer1": 8, "Dense_layer1": 8, "Dense_layer2":8 },  
        "lstm_encoder_decoder_determ": { "LSTM_layer1": 32, "Dense_layer1": 16, "Dense_layer2":16 },  
        "lstm_encoder_decoder_outlier": { "LSTM_layer1": 32, "Dense_layer1": 16, "Dense_layer2":16 },  
        "lstm_encoder_decoder_determ_outlier": { "LSTM_layer1": 32, "Dense_layer1": 16, "Dense_layer2":16 },  
        "lstm_auto_encoder": { "LSTM_layer1": 16, "LSTM_layer2": 16 },  
        "cnn_lstm_encoder_decoder_outlier_removed": { "Conv1TMD_layer1": 4, "kernel_size_layer1": (1,3), "LSTM_layer1": 8, "Dense_layer1": 8, "Dense_layer2":8 },  
        "lstm_encoder_decoder_outlier_removed": { "LSTM_layer1": 8, "Dense_layer1": 8, "Dense_layer2":8 },  
        "lstm_encoder_decoder_outlier_removed_determ": { "LSTM_layer1": 32, "Dense_layer1": 16, "Dense_layer2":16 },  
        "lstm_encoder_decoder_outlier_removed_determ_outlier": { "LSTM_layer1": 32, "Dense_layer1": 8, "Dense_layer2":8 },  
        "lstm_simple_model_reset": { "LSTM_layer1": 64, "LSTM_layer2": 32, "LSTM_layer3": 16, "LSTM_layer4": 16, "Dense_layer1": 16, "Dense_layer2":16 },  
        "lstm_model_with_extra_features":{ "LSTM_layer1": 16, "LSTM_layer2": 16, "Dense_layer1": 8, "Dense_layer2":8 }  
    },  
    "ensemble_lstm" : ["lstm_encoder_decoder",  
        "cnn_lstm_encoder_decoder"],  
    "ensemble_lstm_outlier_removed" : ["lstm_encoder_decoder_outlier_removed",  
        "cnn_lstm_encoder_decoder_outlier_removed"],  
    "ensemble_lstm_with_auto_encoder" : ["lstm_encoder_decoder_outlier_removed",  
        "cnn_lstm_encoder_decoder_outlier_removed"],  
    "ensemble_lstm_with_auto_encoder_outlier_removed" : ["lstm_encoder_decoder_outlier_removed",  
        "cnn_lstm_encoder_decoder_outlier_removed",  
        "lstm_auto_encoder_MLP"],  
    "stacked_ensemble_lstm" : ["lstm_encoder_decoder", "cnn_lstm_encoder_decoder",  
        "lstm_encoder_decoder_outlier_removed",  
        "cnn_lstm_encoder_decoder_outlier_removed"],  
    "stacked_ensemble_lstm2" : ["lstm_encoder_decoder_outlier_removed",  
        "cnn_lstm_encoder_decoder_outlier_removed"],  
    "stacked_ensemble_lstm3" : ["lstm_encoder_decoder_outlier_removed",  
        "cnn_lstm_encoder_decoder_outlier_removed"],  
    "stacked_ensemble_lstm4" : ["lstm_encoder_decoder_outlier_removed",  
        "cnn_lstm_encoder_decoder_outlier_removed"],  
    "stacked_ensemble_lstm5" : ["lstm_encoder_decoder",  
        "cnn_lstm_encoder_decoder"],  
    "ensemble_of_ensembles_lstm" : [ "stacked_ensemble_lstm",  
        "stacked_ensemble_lstm2",  
        "stacked_ensemble_lstm3",  
        "stacked_ensemble_lstm4"],  
    "ensemble_of_ensembles_lstm5" : [ "stacked_ensemble_lstm",  
        "stacked_ensemble_lstm2",  
        "stacked_ensemble_lstm3",  
        "stacked_ensemble_lstm4",  
        "stacked_ensemble_lstm5"]  
}
```

```

    "lstm_auto_encoder_MLP",
    "ensemble_of_ensembles_lstm1": ["stacked_ensemble_lstm1",
                                    "stacked_ensemble_lstm2",
                                    "stacked_ensemble_lstm3",
                                    "stacked_ensemble_lstm4"],
    "ensemble_of_ensembles_lstm2": ["stacked_ensemble_lstm2",
                                    "stacked_ensemble_lstm3",
                                    "stacked_ensemble_lstm4"],
    "ensemble_of_ensembles_lstm3": ["stacked_ensemble_lstm1",
                                    "stacked_ensemble_lstm2",
                                    "stacked_ensemble_lstm4"]
]

# prepare a list of ml models
def get_ml_models(self):
    models = dict()
    n_timesteps = self.list_ml_model_params['time_steps']
    n_forecast_steps = self.list_ml_model_params['forecast_steps']
    n_features = self.list_ml_model_params['features']
    n_deterministic_features = self.list_ml_model_params['deterministic_features']
    sub_timesteps = self.list_ml_model_params['sub_timesteps']
    n_out_auto_encoder_features = self.list_ml_model_params['n_out_features']

    layer_params = self.list_ml_model_params['list_layer_parameters']

    models['lstm_encoder_decoder'] = LSTM_base_model().LSTM_encoder_decoder(layer_params['lstm_encoder_decoder'], n_timesteps, n_forecast_steps, n_features, n_deterministic_features, 'lstm_encoder_decoder')
    models['lstm_simple_model'] = LSTM_base_model(layer_params['lstm_simple_model'], n_timesteps, n_forecast_steps, n_features, 'lstm_simple_model')
    models['lstm_auto_encoder'] = LSTM_base_model().LSTM_auto_encoder(layer_params['lstm_auto_encoder'], n_timesteps, n_forecast_steps, n_features, 'lstm_auto_encoder')
    models['cnn_lstm_encoder_decoder'] = LSTM_base_model().CNN_LSTM_encoder_decoder(layer_params['cnn_lstm_encoder_decoder'], n_timesteps, n_forecast_steps, n_features, n_deterministic_features, 'cnn_lstm_encoder_decoder')
    models['lstm_auto_encoder_MLP'] = LSTM_base_model().LSTM_auto_encoder_MLP(layer_params['lstm_auto_encoder'], n_timesteps, n_forecast_steps, n_features, n_deterministic_features, 'lstm_auto_encoder_MLP')

    print('Defined %d models' % len(models))
    return models

def create_LSTM_train_dataset(self):
    # create train/test/validation datasets
    # keep all LSTM, LSTM_detrend
    dataset_dict = {}
    # prepare LSTM datasets
    input_cols = global_parameters['input_cols_neural_networks']
    deterministic_features = global_parameters['deterministic_features']

    # create train/test/validation split with detrend data
    LSTM_model.dataset = train_test_validation_data_set(data_final_detrend)
    X_train_lstm_detrend, y_train_lstm_detrend, X_validation_lstm_detrend, y_validation_lstm_detrend, X_test_lstm_detrend, y_test_lstm_detrend, deterministic_train_detrend, deterministic_validation_detrend, test_detrend_values = data_final_detrend['train'][len(y_train_lstm_detrend):]
    test_detrend_values = data_final_detrend['test'][len(y_train_lstm_detrend):]

    # create the LSTM data set
    LSTM_model.dataset = train_test_validation_data_set(data_final)
    X_train_lstm, y_train_lstm, X_validation_lstm, y_validation_lstm, X_test_lstm, y_test_lstm, deterministic_train, deterministic_validation, deterministic_test, n_features = LSTM_model.dataset.prepare_LSTM()
    input_cols = input_cols, deterministic_features=deterministic_features, shuffle=True

    # create the LSTM data set with outlier removed data
    LSTM_model.dataset = train_test_validation_data_set(data_outliersremoved)
    X_train_lstm_outlier, y_train_lstm_outlier, X_validation_lstm_outlier, y_validation_lstm_outlier, X_test_lstm_outlier, y_test_lstm_outlier, deterministic_train_outlier, deterministic_validation_outlier, deterministic_test_outlier, n_features = LSTM_model.dataset.prepare_LSTM()
    input_cols = input_cols, deterministic_features=deterministic_features, validation=False, shuffle=True

    sub_timesteps = self.list_ml_model_params['sub_timesteps']
    n_timesteps = self.list_ml_model_params['time_steps']

    # samples, sub_timesteps, rowsx1, colstime_steps/n_subtime_steps, channels]
    X_train_CNN_LSTM = np.reshape(X_train_lstm, (X_train_lstm.shape[0], sub_timesteps, 1, int(n_timesteps/sub_timesteps), n_features))
    X_validation_CNN_LSTM = np.reshape(X_validation_lstm, (X_validation_lstm.shape[0], sub_timesteps, 1, int(n_timesteps/sub_timesteps), n_features))
    X_test_CNN_LSTM = np.reshape(X_test_lstm, (X_test_lstm.shape[0], sub_timesteps, 1, int(n_timesteps/sub_timesteps), n_features))

    # outlier removed data
    X_train_CNN_LSTM_outlier = np.reshape(X_train_lstm_outlier, (X_train_lstm_outlier.shape[0], sub_timesteps, 1, int(n_timesteps/sub_timesteps), n_features))
    X_validation_CNN_LSTM_outlier = np.reshape(X_validation_lstm_outlier, (X_validation_lstm_outlier.shape[0], sub_timesteps, 1, int(n_timesteps/sub_timesteps), n_features))
    X_test_CNN_LSTM_outlier = np.reshape(X_test_lstm_outlier, (X_test_lstm_outlier.shape[0], sub_timesteps, 1, int(n_timesteps/sub_timesteps), n_features))

    # extra deterministic feature shape
    deterministic_train_extra_features = np.reshape(deterministic_train_outlier, (X_train_lstm_outlier.shape[0], n_deterministic_features))
    deterministic_validation_extra_features = np.reshape(deterministic_validation_outlier, (X_validation_lstm_outlier.shape[0], n_deterministic_features))
    deterministic_test_extra_features = np.reshape(deterministic_test_outlier, (X_test_lstm_outlier.shape[0], n_deterministic_features))

    # normalize encoded data from auto encoder

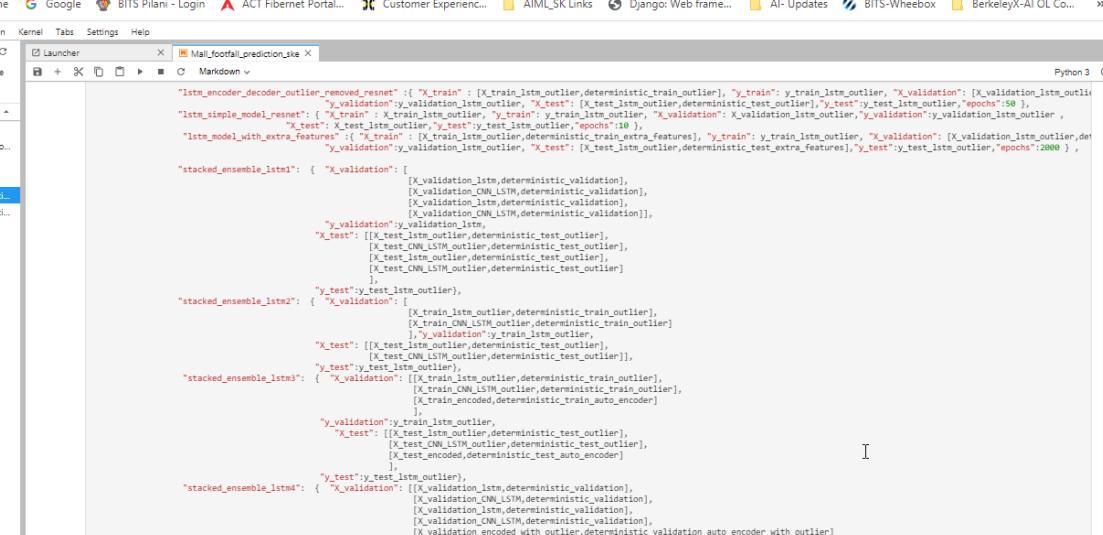
```

```

    deterministic_train_extra_features = np.reshape(deterministic_train_outlier, (X_train_lstm_outlier.shape[0], n_deterministic_features))
    deterministic_validation_extra_features = np.reshape(deterministic_validation_outlier, (X_validation_lstm_outlier.shape[0], n_deterministic_features))
    deterministic_test_extra_features = np.reshape(deterministic_test_outlier, (X_test_lstm_outlier.shape[0], n_deterministic_features))

    # normalize encoded data from auto encoder

```



The screenshot shows a browser window with multiple tabs open. The active tab is titled 'Mail_footfall_prediction_skeleton.ipynb' and contains a large amount of Python code. The code defines various models and their validation and testing configurations. It includes sections for 'lstm_encoder_decoder_outlier_removed_resnet', 'lstm_simple_model_resnet', 'lstm_model_with_extra_features', 'stacked_ensemble_lstm1', 'stacked_ensemble_lstm2', 'stacked_ensemble_lstm3', 'stacked_ensemble_lstm4', and 'stacked_ensemble_lstm5'. Each section specifies training and validation datasets, as well as test datasets for different models like LSTM, CNN-LSTM, and Auto-Encoder.

```
File Edit View Run Kernel Tabs Settings Help
+ Mail_footfall_prediction_skeleton.ipynb Python 3
[...]
"lstm_encoder_decoder_outlier_removed_resnet": { "X_train": [X_train_lstm_outlier,deterministic_train_outlier], "y_train": y_train_lstm_outlier, "X_validation": [X_validation_lstm_outlier,deterministic_validation], "y_validation": y_validation_lstm_outlier, "X_test": [X_test_lstm_outlier,deterministic_test_outlier], "y_test": y_test_lstm_outlier, "epochs": 50 },
"lstm_simple_model_resnet": { "X_train": X_train_lstm_outlier, "y_train": y_train_lstm_outlier, "X_validation": X_validation_lstm_outlier, "y_validation": y_validation_lstm_outlier, "X_test": X_test_lstm_outlier, "y_test": y_test_lstm_outlier, "epochs": 10 },
"lstm_model_with_extra_features": { "X_train": [X_train_lstm_outlier,deterministic_train_extra_features], "y_train": y_train_lstm_outlier, "X_validation": X_validation_lstm_outlier, "y_validation": y_validation_lstm_outlier, "X_test": [X_test_lstm_outlier,deterministic_test_extra_features], "y_test": y_test_lstm_outlier, "epochs": 2000 },
"stacked_ensemble_lstm1": { "X_validation": [
    [X_validation_lstm,deterministic_validation],
    [X_validation_CNN_LSTM,deterministic_validation],
    [X_validation_lstm,deterministic_validation],
    [X_validation_CNN_LSTM,deterministic_validation]
], "y_validation": y_validation_lstm,
"X_test": [[X_test_lstm_outlier,deterministic_test_outlier],
    [X_test_CNN_LSTM_outlier,deterministic_test_outlier],
    [X_test_lstm_outlier,deterministic_test_outlier],
    [X_test_CNN_LSTM_outlier,deterministic_test_outlier]
],
"y_test": y_test_lstm_outlier},
"stacked_ensemble_lstm2": { "X_validation": [
    [X_train_lstm_outlier,deterministic_train_outlier],
    [X_train_CNN_LSTM_outlier,deterministic_train_outlier]
], "y_validation": y_train_lstm_outlier,
"X_test": [[X_test_lstm_outlier,deterministic_test_outlier],
    [X_test_CNN_LSTM_outlier,deterministic_test_outlier],
    "y_test": y_test_lstm_outlier},
"stacked_ensemble_lstm3": { "X_validation": [[X_train_lstm_outlier,deterministic_train_outlier],
    [X_train_CNN_LSTM_outlier,deterministic_train_outlier],
    [X_train_encoded,deterministic_train_auto_encoder]
], "y_validation": y_train_lstm_outlier,
"X_test": [[X_test_lstm_outlier,deterministic_test_outlier],
    [X_test_CNN_LSTM_outlier,deterministic_test_outlier],
    [X_test_encoded,deterministic_test_auto_encoder]
],
"y_test": y_test_lstm_outlier},
"stacked_ensemble_lstm4": { "X_validation": [[X_validation_lstm,deterministic_validation],
    [X_validation_CNN_LSTM,deterministic_validation],
    [X_validation_lstm,deterministic_validation],
    [X_validation_CNN_LSTM,deterministic.validation],
    [X.validation_encoded_with_outlier,deterministic_validation_auto_encoder]
],
"y_validation": y_validation_lstm,
"X_test": [[X_test_lstm_outlier,deterministic_test_outlier],
    [...]
```


The screenshot shows a JupyterLab environment with a Python 3 kernel. The code in the notebook cell is as follows:

```
[...]
    ],
    "y_validation": "y_validation_lstm",
    "X_test": [[X_lstm_outlier,deterministic_test_outlier],
               [X_lstm_CNN_lstm_outlier,deterministic_test_outlier],
               [X_lstm_lstm_outlier,deterministic_test_outlier],
               [X_lstm_CNN_LSTM_outlier,deterministic_test_outlier],
               [...],
               [X_lstm_lstm_outlier,deterministic_test_outlier],
               [X_lstm_CNN_LSTM_outlier,deterministic_test_outlier]
              ],
    "y_test": "y_test_lstm_outlier"
  }
  return lstm_data_set_dict
}

def train_LSTM_models(self,models=dict(), epochs=10):
  lstm_data_set_dict = self.create_LSTM_train_datasets()
  model_training_results = dict()
  for name, model in models.items():
    for batch_size in self.lstm_model_parameters["batch_size"]:
      for learning_rate in self.lstm_model_parameters["learning_rate"]:
        print("Now training (%s)..." % name)
        num_epochs = lstm_data_set_dict[name]['epochs']

        fileExists, model_file = model_utility().CheckPointModelExisting(name)
        if fileExists == True:
          #Load the model with the checkpoint file
          print(model_file)
          model = self.loadModel(model_file)
          compile = False;
          learning_rate = 0.0001
          num_epochs = 1000
          batch_size = 16
        else:
          compile = True
        print("number of epochs - (%s)" % (str(num_epochs)))
        train_loss, val_loss, checkpoint_path = LSTM_base_model().compile_and_train(model,lstm_data_set_dict[name], num_epochs = num_epochs, batch_size=batch_size,learning_rate=learning_rate,compile=compile)

        train_loss = pd.DataFrame(train_loss, columns=['train_loss'])
        val_loss = pd.DataFrame(val_loss, columns=['val_loss'])

        [...]
```

```
File Edit View Run Kernel Tabs Settings Help
+ - + localhost:8888/lab
File Edit View Run Kernel Tabs Settings Help
+ - + Mail_football_prediction_ske
Code
val_loss = pd.DataFrame([val_loss], columns=['val_loss'])
plot_data_set = pd.concat([train_loss,val_loss], axis=1)
model_utility().plotly_graphs(plot_data_set, name)
model_training_result[name, batch_size, learning_rate] = {"model_weights": checkpoint_path, "val_loss": val_loss, "train_loss":train_loss}

def evaluate_LSTM_models(self, model_training_results, data_set_dict):
    models = self.get_LSTM_models()
    ensemble_models = []
    ensemble_models['ensemble_lstm'] = []
    ensemble_models['ensemble_lstm_outlier_removed'] = []
    ensemble_models['ensemble_lstm_with_auto_encoder'] = []

    loss_values = {}
    for name, model_object in models.items():
        fileExists, model_file = model_utility().CheckPointModelExisting(name)
        if fileExists == True:
            load(model_file)
            model_object = load(model_file)
            predicts,y_actual, loss_values[name] = self.predict_model(model_object,data_set_dict[name][X_test]),data_set_dict[name][y_test])
            print("best loss value", name, loss_values[name])
            predicts = pd.DataFrame(predicts, columns=['Predicted'])
            y_test_df = pd.DataFrame(y_actual, columns=['Actual'])
            plot_data_set = pd.concat([predicts,y_test_df], axis=1)
            model_utility().plotly_graphs(plot_data_set, name)

            if name in self.list_model_parameters['ensemble_lstm']:
                ensemble_models['ensemble_lstm'].append(model_object)
            if name in self.list_model_parameters['ensemble_lstm_outlier_removed']:
                ensemble_models['ensemble_lstm_outlier_removed'].append(model_object)
            if name in self.list_model_parameters['ensemble_lstm_with_auto_encoder']:
                ensemble_models['ensemble_lstm_with_auto_encoder'].append(model_object)

            for ensemble_name in ensemble_lstm['ensemble_lstm_outlier_removed','ensemble_lstm_with_auto_encoder']:
                ensemble_name = np.array(ensemble_models[ensemble_name])
                ensemble_model = LSTM_base_model().ensemble(ensemble_name,ensemble_name)
                predicts,y_actual, loss_values[ensemble_name] = self.predict_model(ensemble_model, data_set_dict[ensemble_name][X_test]),data_set_dict[ensemble_name][y_test])

                predicts = pd.DataFrame(predicts, columns=['Predicted'])
                y_test_df = pd.DataFrame(y_actual, columns=['Actual'])
                loss_values[ensemble_name] = np.count_nonzero(predicts != y_test_df['Actual'])
```

```
predicts = pd.DataFrame(predicts, columns=['predicted'])
y_test_df = pd.DataFrame(y_actual, columns=['actual'])
plot_data_set = pd.concat([predicts,y_test_df], axis=1)
model_utility().plotly_graph(plot_data_set,ensemble_name)

# stacked ensemble models
for ensemble_name in [stacked_ensemble_list[i]+str(i+1) for i in range(4)] + ['ensemble_of_ensembles_lstm'+str(j+1) for j in range(3)]:
    stacked_model_files = []
    for model_name in self.lstn_model_parameters[ensemble_name]:
        fileExist(model_file = model_utility().CheckPointModellisting(model_name))
        if fileExist(model_file):
            stacked_model_file.append(model_file)
    stacked_lstm_model = LSTMbase(model)
    train_stacked_lstm_model(stacked_model_files,data_set_dict[ensemble_name][y_validation"],data_set_dict[ensemble_name][y_validation"])
    predicts, y_actual, loss_values[ensemble_name] = self.predict_model(stacked_lstm_model,data_set_dict[ensemble_name][y_test"],data_set_dict[ensemble_name][y_validation"])

    predicts = pd.DataFrame(predicts, columns=['predicted'])
    y_test_df = pd.DataFrame(y_actual, columns=['actual'])
    plot_data_set = pd.concat([predicts,y_test_df], axis=1)
    model_utility().plotly_graphs(plot_data_set,ensemble_name)

    loss_values_df = model_utility().convert_dict_dataframe(loss_values)

return loss_values_df

def inverse_detrend(self,predicts,trend,seasonal):
    predicts = predicts + trend
    return predicts

def predict_model(self,model,input,y_test,model_type=None,trend=None,seasonal=None):
    # show the validation trend for a few cell i
    predicts = model.predict(input)
    if(model.name in ['lstm_auto_encoder']):
        predicts = np.reshape(predicts,[1])
    print('predicts shape:',predicts.shape)
    predicts = np.reshape(predicts,[predicts.shape[0],self.lstn_model_parameters['forecast_steps']])
    if(model_type=='trend'): inverse the detrend by adding trend and seasonal info
    trend = np.reshape(trend,[trend.shape[0], 1])
    seasonal = np.reshape(seasonal,[seasonal.shape[0], 1])
    predicts = inverse_detrend(predicts,trend,seasonal)
    y_test = inverse_detrend(y_test,trend,seasonal)

    loss_values = model_utility().calculate_metrics(y_test, predicts,n_sample_size=y_test.shape[0],n_input_variables=int(len(global_parameters['input_cols_neural_networks']))*global_parameters['n_timestep']:
```

```
loss_values = model_utility().calculate_metrics(y_test, predicts,n_sample_size=y_test.shape[0],n_input_variables=int(len(global_parameters['input_cols_neural_networks']))*global_parameters['n_timestep'])

predicts = scaler_.target_inverse_transform(predicts)
y_actual = scaler_.target_inverse_transform(y_test)
return predicts, y_actual, loss_values
```

```
[26]: class H2o_model:
    def __init__(self, panda_df,dataset_type):
        h2o_no_error()
        self.panda_df = panda_df.copy()
        self.dataset_type = dataset_type
        self.grid_params = {}
        self.grid_params.update({'H2OGridSearchConfig': {'learn_rate': [1 * 0.01 for i in range(1, 11)],
                                                       'max_depth': [11, 10, 9, 8, 7, 6, 5, 4, 3, 2, 1],
                                                       'lambda': [1 * 0.1 for i in range(5, 11)],
                                                       'col_sample_rate': [1 * 0.1 for i in range(1, 11)]}})

        self.grid_params.update({'H2OAutoMLConfig': [
            'max_runtime_secs': [1000, 10000, 100000, 1000000],
            'sample_rate': [1 * 0.1 for i in range(5, 11)],
            'col_sample_rate_change_per_level': [1 * 0.1 for i in range(1, 11)], 'mtries':[-1,-2,2]
        ]})

    def autoML(self,max_models=50,max_runtime_secs=1000, explain=True, full_lag_data=False):
        if full_lag_data == True:
            dataset_dict = self.prepare_auto_ML_full_lag_dataset()
        else:
            dataset_dict = self.prepare_datasets()

        train_frame = dataset_dict['h2o'].self.dataset_type]['train_frame']
        validation_frame = dataset_dict['h2o'].self.dataset_type]['Validation_frame']
        test = dataset_dict['h2o'].self.dataset_type]['test_frame']
        y_col = dataset_dict['h2o'].self.dataset_type]['target_col']
        X_columns = dataset_dict['h2o'].self.dataset_type]['input_columns']

        aml = H2OAutoML(max_models=max_models,max_runtime_secs=max_runtime_secs,
                        seed = 1)
```

```
aml = H2OAutoML(max_models=max_models,max_runtime_secs=max_runtime_secs,
                 seed=1)
aml.train(x=X.columns,
          y=y_col,
          training_frame=train_frame,
          validation_frame=validation_frame,
          leaderboard_frame=test)

lb = aml.leaderboard
print(lb.head(rows=lb.nrows)) # Print all rows instead of default (10 rows)
# The Leader model is stored here
print(aml.leader)

if explain == True:
    ro_plot = aml.explain(test)

# get all the models from AutoML
models = aml.get_models()['model_id'].as_data_frame().iloc[:,0]
loss_values = []
for m_id in model_ids:
    print(m_id)
    md = h2o.get_model(m_id)
    predictions = md.predict(test)
    predictions = predictions.as_data_frame().values
    predictions = predictions.reshape(-1)

    y_test_real = test[y_col].as_data_frame().values
    y_test_real = y_test_real.reshape(-1)

    loss_values[m_id+self.dataset_type] = model_utility().calculate_metrics(y_test_real,predictions,len(y_test_real),len(X.columns),'AutoML')

loss_values_df = model_utility().convert_dict_dataframe(loss_values)

return loss_values_df
```

def get_models(self, models=dict()):
 # ensemble models
 models["H2OGradientBoosting"] = H2OGradientBoostingEstimator(distribution="gaussian",ntrees=100,min_rows=2,nfolds=5,keep_cross_validation_predictions=True,seed=1)
 models["H2ORandomForest"] = H2ORandomForestEstimator(ntrees=100,nfolds=5,keep_cross_validation_predictions=True,seed=1)
 print("Defined %d models" % len(models))

```
models["H2ORandomForest"] = H2ORandomForestEstimator(ntrees=100,nfolds=5,keep_cross_validation_predictions=True,seed=1)

print("Defined %d models" % len(models))
return models

def prepare_auto_ml_full_lag_dataset(self):
    dataset_dict = {}
    n_timesteps = global_parameters["n_timesteps"]
    n_forecast_steps = global_parameters["n_forecast_steps"]
    df2 = self.panda_data_df[global_parameters["input_cols_neural_networks"]].copy()

    for column in df2:
        for lag in range(1,n_timesteps+1):
            df2[column+"_"+str(lag)] = df2[column].shift(lag-(n_forecast_steps-1))

    df2.dropna(inplace=True)

    mask = (df2.columns.str.contains("InCountTotal") | df2.columns.str.contains("lag"))
    df2_processed = df2[df2.columns[mask]]

    df2_processed.reset_index(drop=True,inplace=True)

    # the columns in the processed data frame

    df2_processed.to_csv("auto_ml_full_lag_dataset.csv")

    data_frame_slice = df2_processed[0:-global_parameters["test_data_set_days"]].copy()

    h2o_data_frame = h2o.H2OFrame(data_frame_slice)
    train, valid = h2o_data_frame.split_frame(ratios=[.8], seed = 1)

    test_frame_slice = df2_processed[-global_parameters["test_data_set_days"]:].copy()
    h2o_test_frame = h2o.H2OFrame(test_frame_slice)

    y_col = "InCountTotal"
    input_columns = train.columns
    input_columns.remove(y_col)

    dataset_dict["h2o_"+self.dataset_type] = {"train_frame": train, "validation_frame": valid, "test_frame": h2o_test_frame,
                                              "input_columns":input_columns, "target_col": y_col}

    return dataset_dict

def prepare_datasets(self):
```

This screenshot shows a JupyterLab interface with two tabs open: 'product owner certification - Google' and 'localhost:8888/lab'. The code editor displays the 'Mail_footfall_prediction_skeleton.ipynb' notebook. The code defines a class with methods for preparing datasets, fitting models, and performing randomized search.

```
def prepare_datasets(self):
    # create dictionary to store different datasets
    dataset_dict = {}
    # prepare regression dataset for train dataset
    n_timesteps = global_parameters['n_timesteps']
    input_columns = [f'InCountTotal_lag_{str(i)}' for i in range(n_timesteps, 0, -1)] + list(global_parameters['input_cols_regression'])
    train_columns = input_columns + list(global_parameters['target_col_h2o'])

    data_frame_slice = self.panda_data_df[0:-global_parameters['test_data_set_days']].copy()
    data_frame_slice = data_frame_slice[train_columns]

    h2o_data_frame = H2OFrame(data_frame_slice)
    train, valid = h2o_data_frame.split(ratios=[.8], seed=1)

    test_frame_slice = self.panda_data_df[-global_parameters['test_data_set_days']:].copy()
    h2o_test_frame = H2OFrame(test_frame_slice[train_columns])

    dataset_dict["h2o"] = self.data_type_map = {"train": train, "validation": valid, "test": h2o_test_frame,
                                                "input_columns": input_columns, "target_col": global_parameters['target_col_h2o'][0]}

    return dataset_dict

# fit a single model
def fit_model(self, model, train, validation, X_col, y_col):
    # click on the model configuration
    model.train(x=X_col, y=y_col, training_frame=train, validation_frame=validation)

def randomizedSearch(self, paramGrid, model, train, validation, test, X_columns, y_col):
    # train and predict each model with grid search CV

    # Search criteria
    search_criteria = {'strategy': 'RandomDiscrete', 'max_models': 20, 'seed': 1}
    # Train and validate a random grid of GBMs
    model_grid = H2OGridSearch(paramGrid,
                               grid_idname,
                               hyper_params=paramGrid,
                               search_criteria=search_criteria
    )
    print("Now tuning [model]")
    model_grid.train(x=X_columns, y=y_col)
```

This screenshot continues the JupyterLab session with the same notebook. The code editor now includes additional logic for training a grid of GBMs, selecting the top model based on validation AUC, and evaluating its performance on a test set. It also includes code for generating a 2-model ensemble (GBM + RF) and comparing it to base learners.

```
model_grid.train(x=X_columns, y=y_col,
                 training_frame=train,
                 validation_frame=validation,
                 ntrees=100,
                 seed=1)

# Get the grid results, sorted by validation AUC
model_gridperf2 = model_grid.get_grid(sort_by='mse', decreasing=False)
print(model_gridperf2)

# Grab the top GBM model, chosen by validation AUC
best_model = model_gridperf2.models[0]

# Now let's evaluate the model performance on a test set
# so we get an honest estimate of top model performance
best_model_perf2 = best_model.model_performance(test)
print(best_model_perf2)

predictions = best_model.predict(test)
predictions = predictions.as_data_frame().values
predictions = predictions.reshape(-1)
y_test_real = test[y_col].as_data_frame().values
y_test_real = y_test_real.reshape(-1)

loss_values = model.utility().calculate_metrics(y_test_real, predictions, len(X_columns), name)
return best_model, predictions, loss_values

def ensemble_stack_models(self, fitted_models, train, validation, test, X_columns, y_col, name):
    # 1. Generate a 2-model ensemble (GBM + RF)
    # Train a stacked ensemble using the GBM and RF above
    ensemble = H2OStackedEnsembleEstimator(model_id=name,
                                             base_models=fitted_models)
    ensemble.train(x=X_columns, y=y_col, training_frame=train, validation_frame=validation)

    # Evaluate ensemble performance on the test data
    perf_stack_test = ensemble.model_performance(test)

    # Compare to base Learner performance on the test set
    perf_gbm_test = fitted_models[0].model_performance(test)
    perf_rf_test = fitted_models[1].model_performance(test)
    baselinelearner_best_mse_test = baselinelearner_best_mse.gbm_test.mse(), perf_rf_test.mse()
    stack_mse_test = perf_stack_test.mse()
    print("Best Base-learner Test MSE: ({})".format(baselinelearner_best_mse_test))
    print("Ensemble Test MSE: ({})".format(stack_mse_test))
```

```
stack_mse_test = perf_stack_test.mse()
print("Best Base-Learner Test MSE: ({})".format(baselearner_best_mse_test))
print("Ensemble Test MSE: ({})".format(stack_mse_test))

# Generate predictions on a test set (if necessary)
predictions = ensemble.predict(test)

predictions = prediction_as_data_frame().values
predictions = predictions.reshape(-1, 1)
y_test_real = test[y_col].as_data_frame().values
y_test_real = y_test_real.reshape(-1)

loss_values = model_utility().calculate_metrics(y_test_real,predictions,len(y_test_real),len(X_columns),name )
return ensemble, predictions, loss_values

def ensemble_random_grid_models(self, paramgrid, model,train, validation,test, X_columns, y_col, name):

# 2. Generate a random grid of models and stack them together
search_criteria = {"strategy": "RandomDiscrete", "max_models": 20, "seed": 1}

# Train the grid
grid = H20GridSearch(model=model,
                     hyper_params=paramGrid,
                     search_criteria=search_criteria,
                     grid_id=name)
grid_id_name = grid.grid_id_name
grid.train(x=X_columns, y=y_col, training_frame=train, validation_frame = validation)

# Train a stacked ensemble using the GBM grid
ensemble = H2OStackedEnsembleEstimator(model_id=name+'random_ensemble',
                                         base_models=grid.model_ids)
ensemble.train(x=X_columns, y=y_col, training_frame=train, validation_frame = validation)

# Evaluate ensemble performance on the test data
perf_stack_test = ensemble.model_performance(test)

# Compare to base Learner performance on the test set
baselearner_best_mse_test = min([h2o.get_model(model).model_performance(test_data=test).mse() for model in grid.model_ids])
stack_mse_test = perf_stack_test.mse()
print("Best Base-Learner Test MSE: ({})".format(baselearner_best_mse_test))
print("Ensemble Test MSE: ({})".format(stack_mse_test))

# Generate predictions on a test set (if necessary)
```

The screenshot shows a JupyterLab environment with a Python 3 notebook open. The code in the notebook is as follows:

```
# Generate predictions on a test set (if necessary)
predictions = ensemble.predict(test)

predictions = predictions.astype(np.int64)
predictions = predictions.reshape(-1)
y_test_real = test[y_col].values
y_test_real = y_test_real.reshape(-1)

loss_values = model_utility().calculate_metrics(y_test_real,predictions,len(y_test_real),len(X_columns),name )

return ensemble, predictions, loss_values

def train_and_evaluate_models(self):

    models = self.get_models()
    trained_models = []
    dataset_dict = self.prepare_datasets()
    train_frame = dataset_dict['h2o_' + self.dataset_type]['train_frame']
    validation_frame = dataset_dict['h2o_' + self.dataset_type]['validation_frame']
    test_frame = dataset_dict['h2o_' + self.dataset_type]['test_frame']
    y_col = dataset_dict['h2o_' + self.dataset_type]['target_col']
    X_columns = dataset_dict['h2o_' + self.dataset_type]['input_columns']

    # initialize the prediction data frames
    modelPredictions = np.zeros((test_frame.shape[0], len(models)))
    modelPredictions = pd.DataFrame(modelPredictions)

    cvPredictions = np.zeros((test_frame.shape[0], len(models)))
    cvPredictions = pd.DataFrame(cvPredictions)

    rscvPredictions = np.zeros((test_frame.shape[0], len(models)))
    rscvPredictions = pd.DataFrame(rscvPredictions)

    #ensemble stack models
    ensemblePredictions = np.zeros((test_frame.shape[0], len(models)+1))
    ensemblePredictions = pd.DataFrame(ensemblePredictions)

    n_forecast_steps = global_parameters['n_forecast_steps']
    loss_values = {}

    for i, (name, model) in enumerate(models.items()):
```

```
File Edit View Run Kernel Tabs Settings Help
+ Mail_Footfall_prediction_skeleton.ipynb
Python 3
# 1. fit default models
for i, (name, model) in enumerate(models.items()):
    # 1. fit default models
    self.fit_model(model, train_frame, validation_frame, X_columns, y_col)
    print(model.model._performance(test_f_frame))
    # make predictions
    predictions = model.predict(test_f_frame)
    print("*****")
    print(name)
    print("Default fit model fitting...")
    
    predictions = predictions.as_data_frame().values
    predictions = predictions.reshape(1, -1)
    y_test_real = test_f_frame[y_col].as_data_frame().values
    y_test_real = y_test_real.reshape(1, -1)

    loss_values[name] = self.dataset_type['default_model'] = model_utility().calculate_metrics(y_test_real, predictions, len(y_test_real), len(X_columns), name)

# inverse transform
predictions = np.reshape(predictions, (predictions.shape[0], n_forecast_steps))
predictions = scaler_target.inverse_transform(predictions)

modelPredictions.iloc[:, i] = predictions
modelPredictions = modelPredictions.rename(columns={i: name})

trained_models[name+self.dataset_type+'standard'] = model

# 2. fit randomized CV models
model_fit_rscv, predicton_rscv, loss_values[name] = self.randomizedSearch(self.grid_params[name], model, train_frame, validation_frame, test_frame, X_columns, y_col, name)

# inverse transform
predicton_rscv = np.reshape(predicton_rscv, (predicton_rscv.shape[0], n_forecast_steps))
predicton_rscv = scaler_target.inverse_transform(predicton_rscv)

#trained_models[name+self.dataset_type+'rscv'] = model_fit_rscv
rscvPredictions.iloc[:, i] = predicton_rscv
rscvPredictions = rscvPredictions.rename(columns={i: name})

# 3. Fit random grid ensembles
model_fit_ensemble, predicton_ensemble, loss_values[name] = self.ensemble_random_grid_models(self.grid_params[name], model, train_frame, validation_frame, X_columns, y_col, name)

# inverse transform
```

```
model_fit_ensemble, predicton_ensemble, loss_values[name+'_'+self.dataset_type+'_random_grid_ensemble'] = self.ensemble_random_grid_models(self.grid_params[name],model, train_frame, validation_frame)

#inverse transform
predicton_ensemble = np.reshape(predicton_ensemble, (predicton_ensemble.shape[0], n_forecast_steps))
predicton_ensemble = scalar_target.inverse_transform(predicton_ensemble)

#trained_models[name+self.dataset_type+'_random_grid_ensemble'] = model_fit_ensemble
ensemblePredictions.iloc[:, i] = predicton_ensemble
ensemblePredictions = ensemblePredictions.rename(columns=(i: name+'_random_grid_ensemble'))

#ensemble stack models
name = "H2O_Stacked_Ensemble"
models = [ model for _, model in trained_models.items()]
model_fit_ensemble, predicton_ensemble, loss_values[name+'_'+self.dataset_type+'_ensemble'] = self.ensemble_stack_models(models, train_frame, validation_frame, test_frame,X_columns, y_col_name)

#inverse transform
predicton_ensemble = np.reshape(predicton_ensemble, (predicton_ensemble.shape[0], n_forecast_steps))
predicton_ensemble = scalar_target.inverse_transform(predicton_ensemble)

trained_models[name+self.dataset_type+'_ensemble'] = model_fit_ensemble
ensemblePreditcions.iloc[:, len(models)] = predicton_ensemble
ensemblePreditcions = ensemblePreditcions.rename(columns=[len(models): name])

y_test_real = np.reshape(y_test_real, (y_test_real.shape[0], n_forecast_steps))
y_test_real = scalar_target.inverse_transform(y_test_real)

y_test_df = pd.DataFrame(y_test_real, columns=['Actual'])
plot_data_set = pd.concat([modelPredictions,y_test_df], axis=1)
model_utility().plotly_graphs(plot_data_set,"H2O default models")

plot_data_set = pd.concat([rsrcvPredictions,y_test_df], axis=1)
model_utility().plotly_graphs(plot_data_set,"H2O Randomized Search CV models")

plot_data_set = pd.concat([ensemblePredictions,y_test_df], axis=1)
model_utility().plotly_graphs(plot_data_set,"H2O Ensemble models")

loss_values_df = model_utility().convert_dict_dataframe(loss_values)

return loss_values_df
```

6.5 FB Prophet class

```
[27]: class fb_prophet:
    def __init__(self, data):
        self.data = data.copy()

    def createFB_prophet_dataset(self,test_data_days=365):
        self.data['Date'] = pd.to_datetime(self.data.Date.astype(str))
        self.data = self.data.set_index('Date')
        self.data = self.data.sort_index()
        self.data['Date'] = self.data.index

        # Rename the features: These names are NEEDED for the model fitting
        self.data = self.data.rename(columns = {"Date":'ds',"inCountTotal":'y'}) #renaming the columns of the dataset

        data_train = self.data[0:-test_data_days]
        data_test = self.data[-test_data_days:]

        return data_train, data_test

    def train(self):
        test_data_days = global_parameters['test_data_set_days']
        data_train, data_test = self.createFB_prophet_dataset(test_data_days)
        model_train = Prophet(daily_seasonality = True) # the Prophet class (model)
        model_train.fit(data_train)

        future = model_train.make_future_dataframe(periods=test_data_days) # we need to specify the number of days in future
        prediction = model_train.predict(future)

        yhat = np.array(prediction['yhat']).tail(test_data_days).reshape(-1,1)
        y_test = data_test['y'].values.reshape(-1,1)

        y_test_df = pd.DataFrame(y_test, columns=['Actual'])
        predictions = pd.DataFrame(yhat, columns=['Predicted'])

        plot_data_set = pd.concat([predictions,y_test_df], axis=1)
        model_utility().plotly_graphs(plot_data_set,"FB Prophet Model")

        fig = plotly.offline.plot(model_train, prediction)
        fig.show()
```

Mode: Command | Ln 343. Col 25 | Mail_footfall_prediction_skeleton.ipynb | ENG 1:48 AM INTL 11/10/2020

7 - Model training and evaluation

7.1 Model Training with Regression and Ensemble models

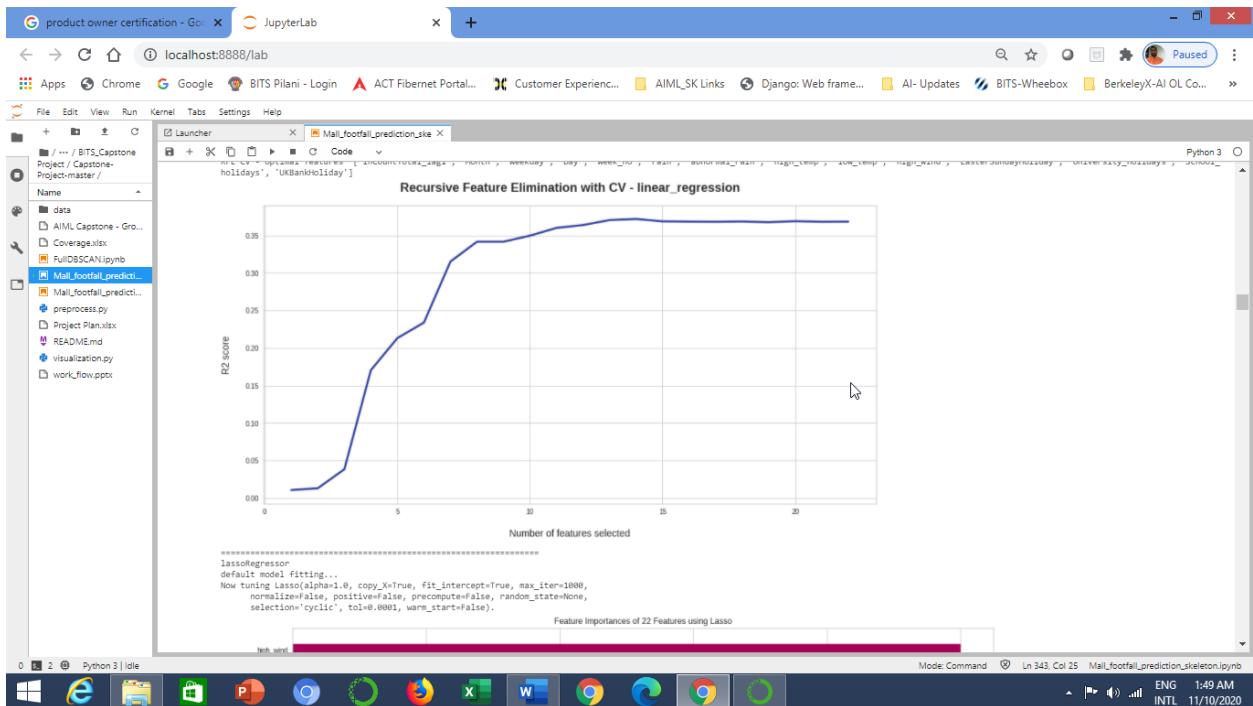
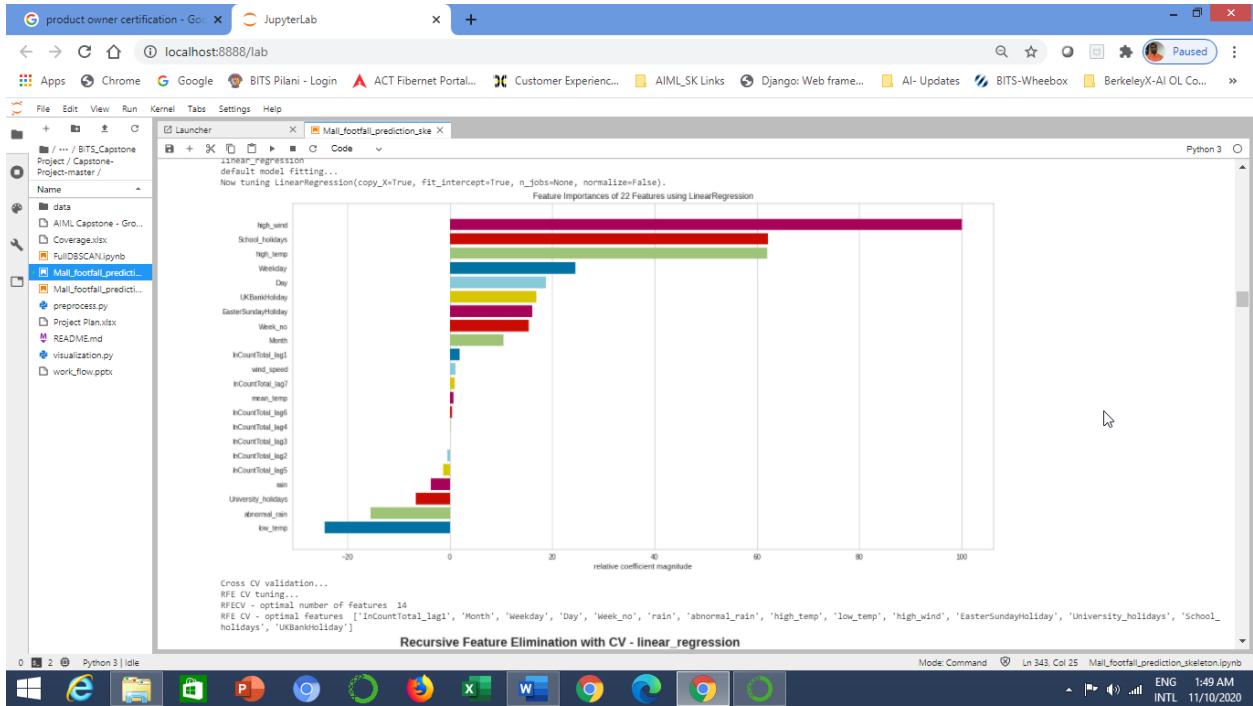
7.1.1 Model training with normal data

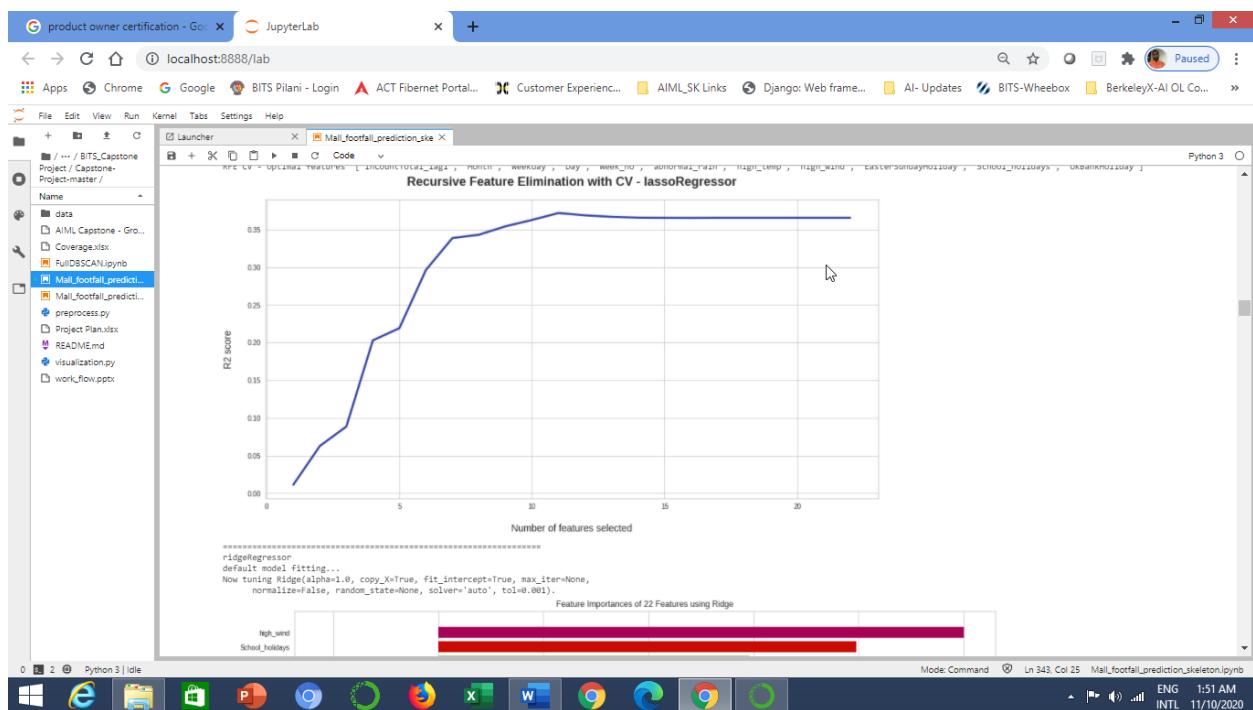
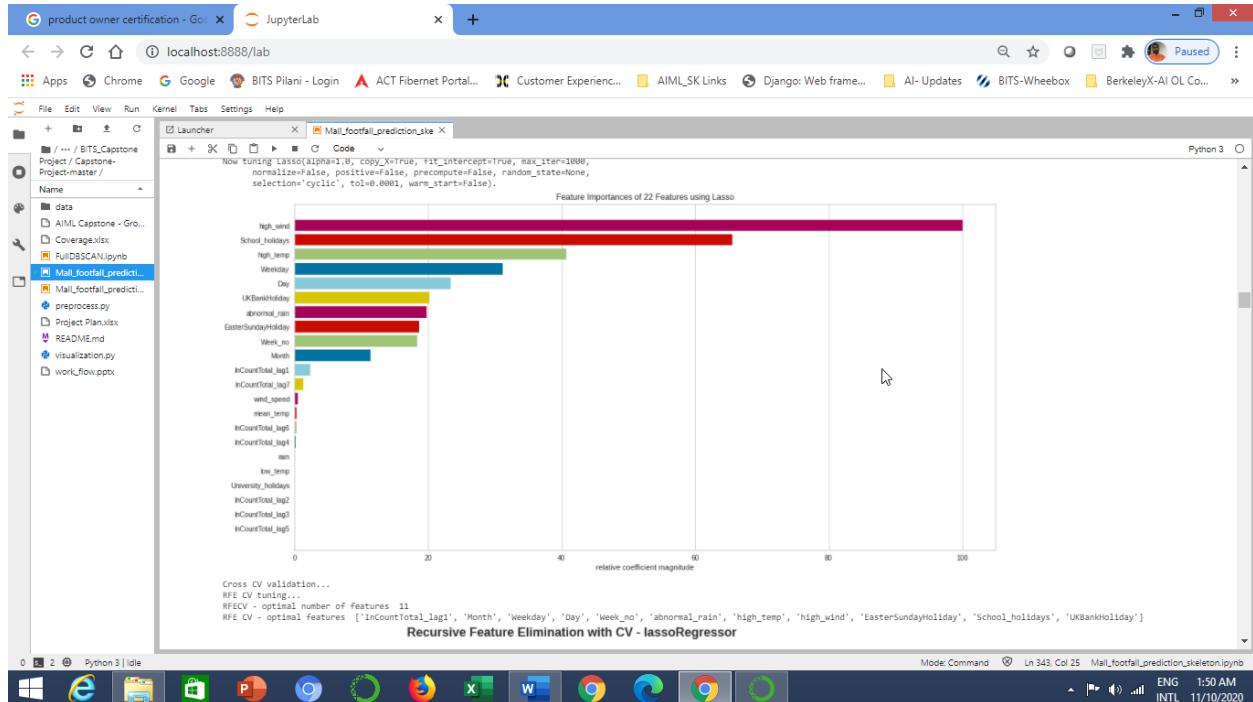
```
[28]: loss_values_all_models = []

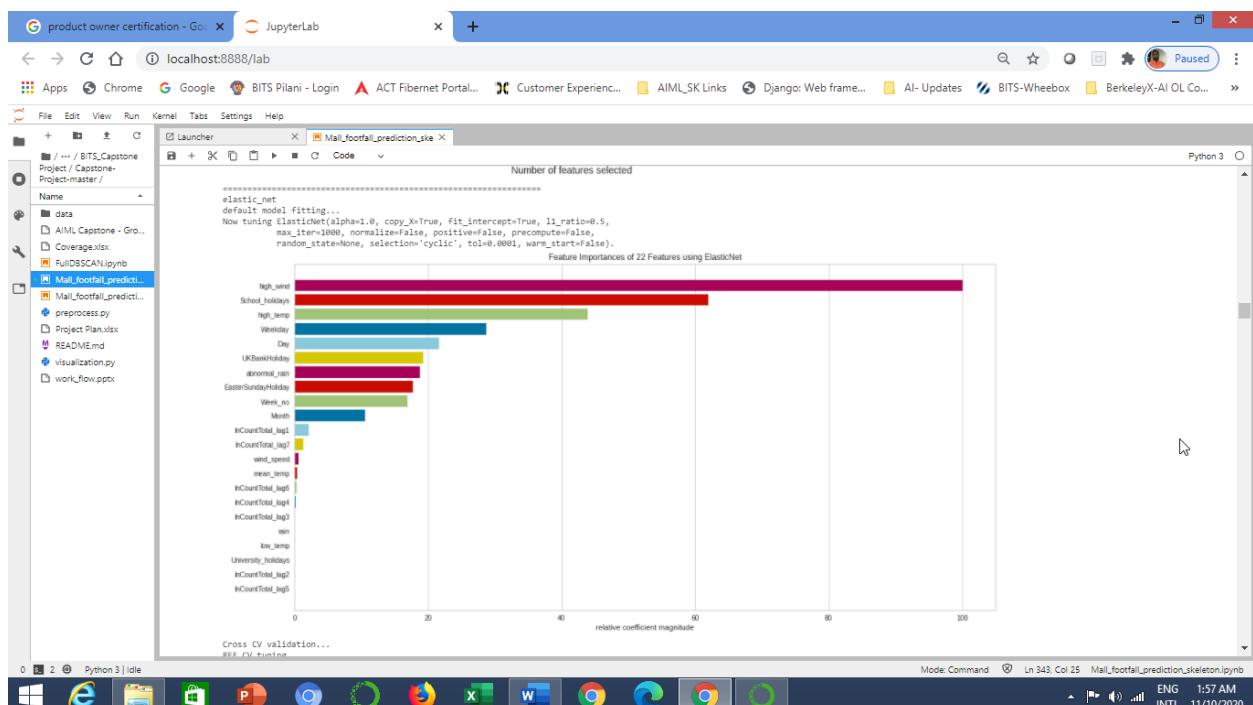
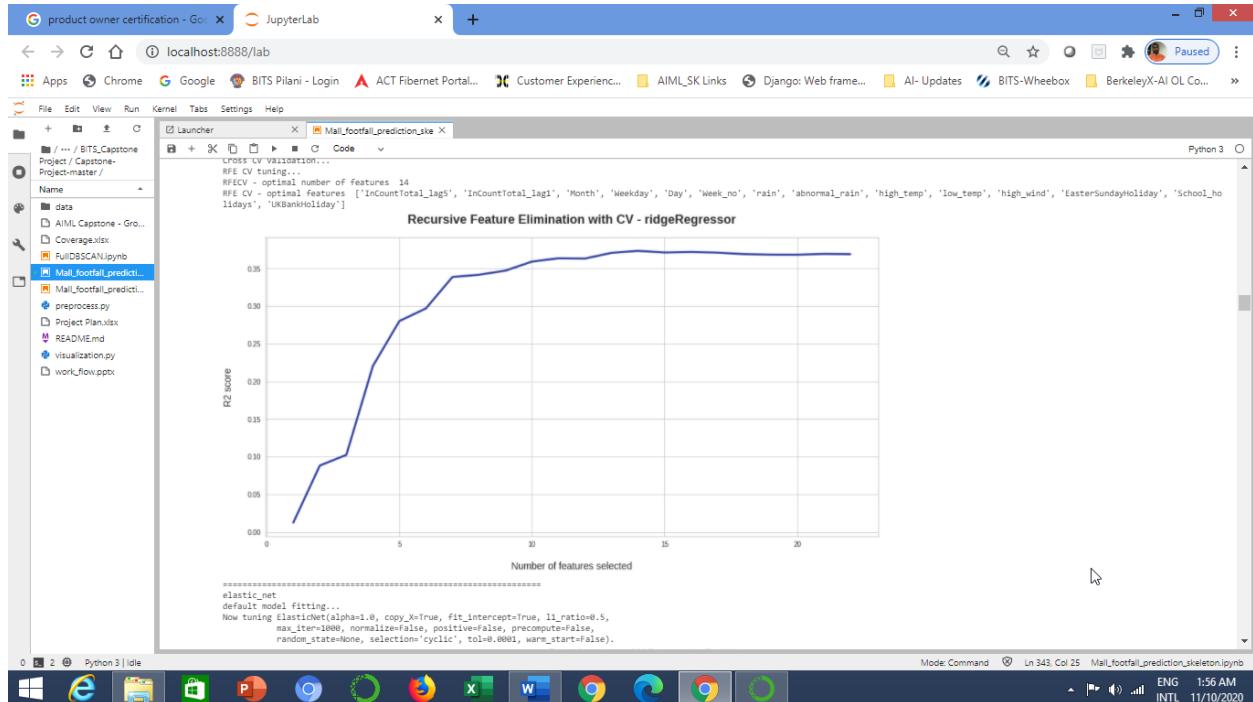
[29]: reg_models = regression_model(data_final,dataset_type='normal_shuffled')
for _, loss_normal_df in reg_models.train_and_evaluate_models(shuffle=True):
    loss_values_all_models.append(loss_normal_df)

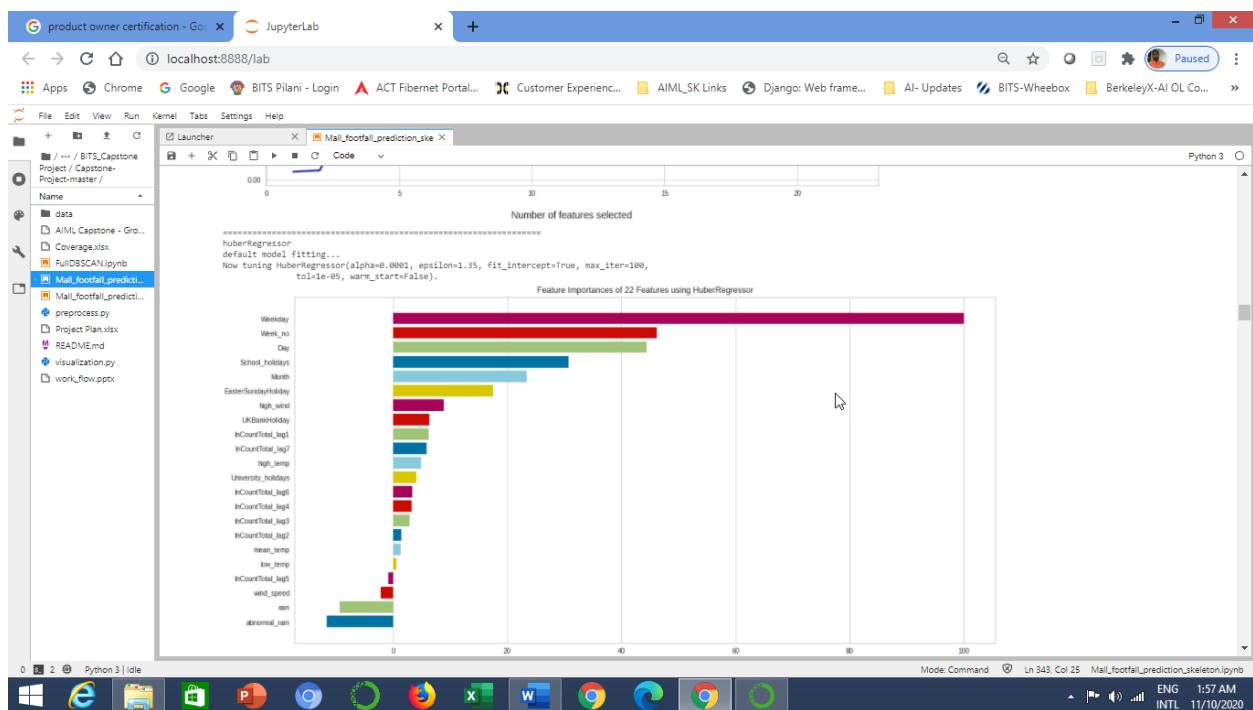
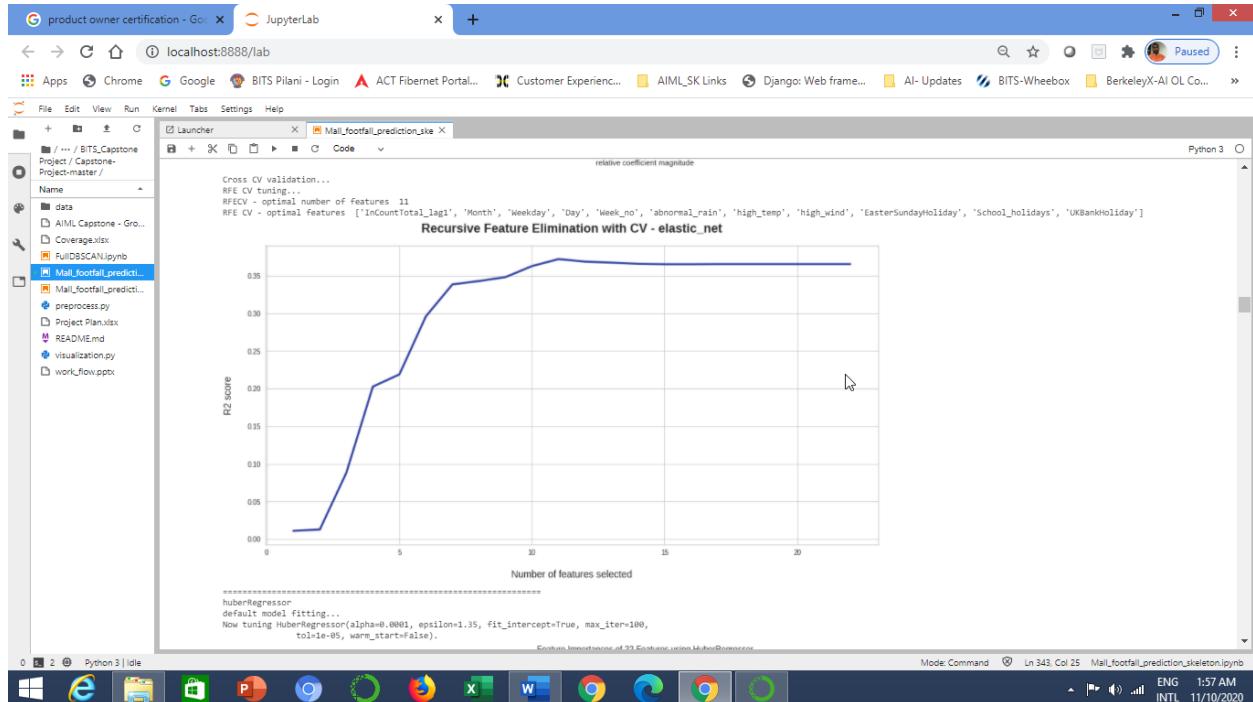
Defined 20 models
=====
linear_regression
Default model fitting...
Now tuning LinearRegression(copy_X=True, fit_intercept=True, n_jobs=None, normalize=False).
Feature Importances of 22 Features using LinearRegression
```

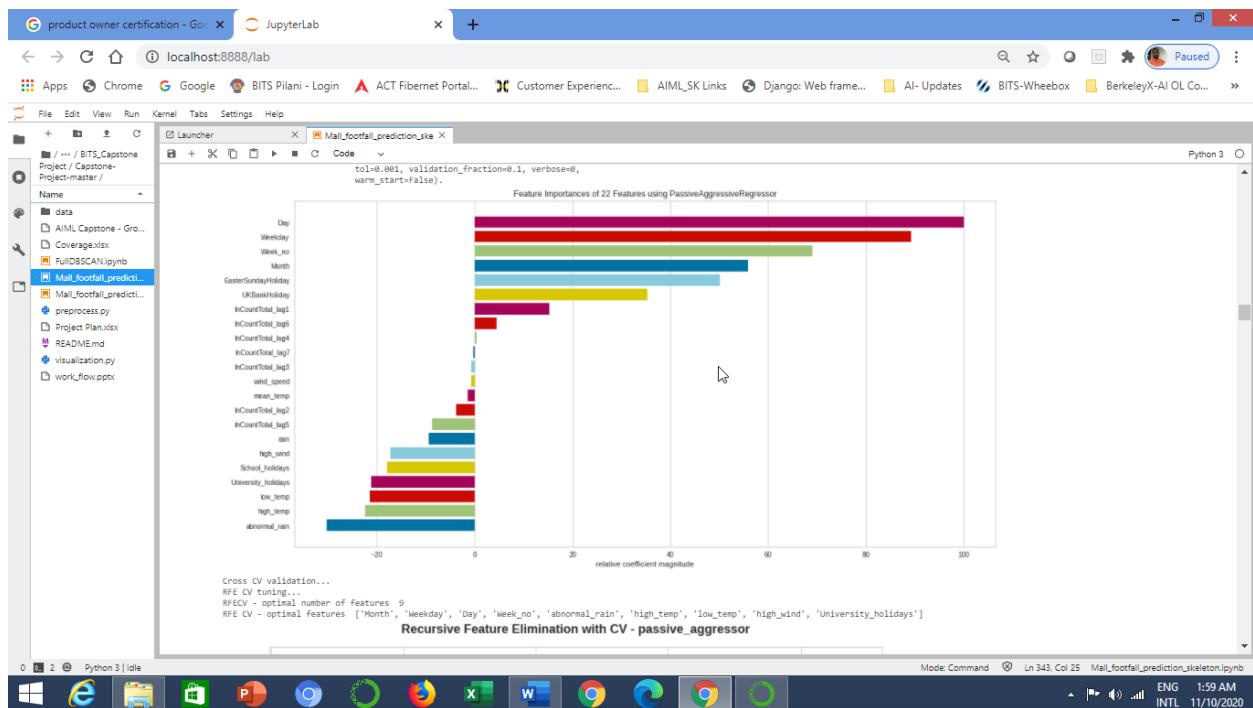
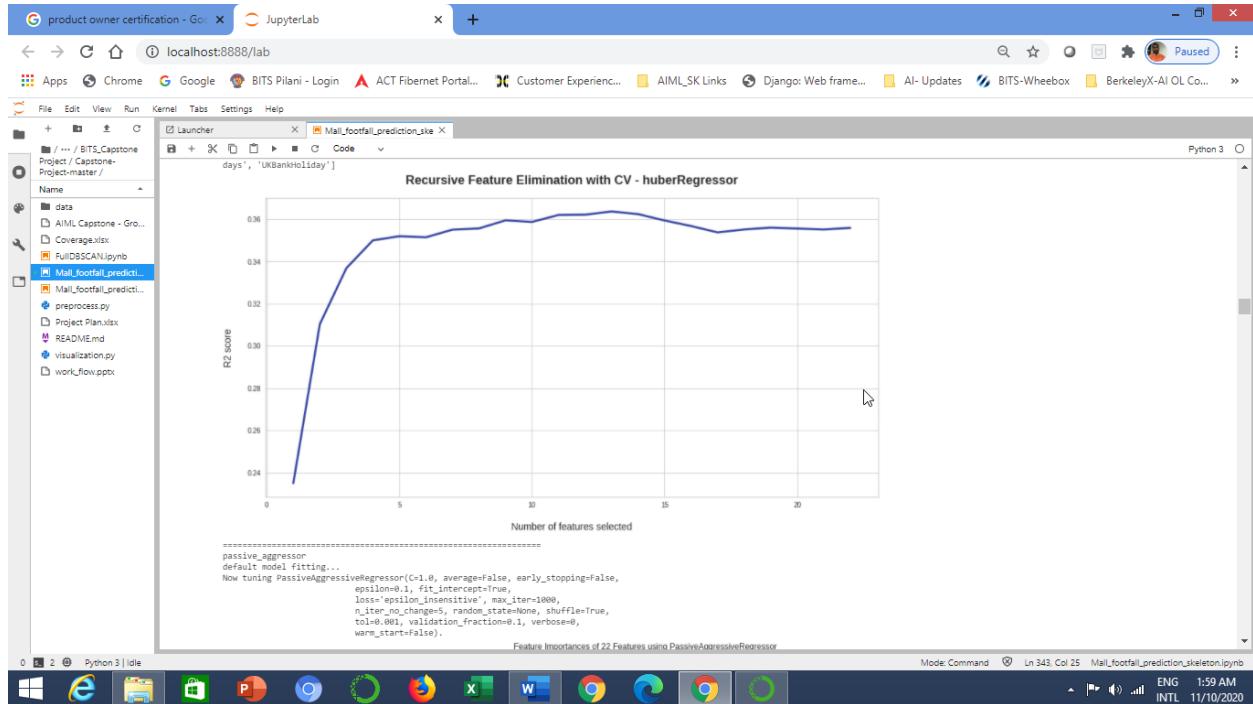
Mode: Command | Ln 343. Col 25 | Mail_footfall_prediction_skeleton.ipynb | ENG 1:49 AM INTL 11/10/2020

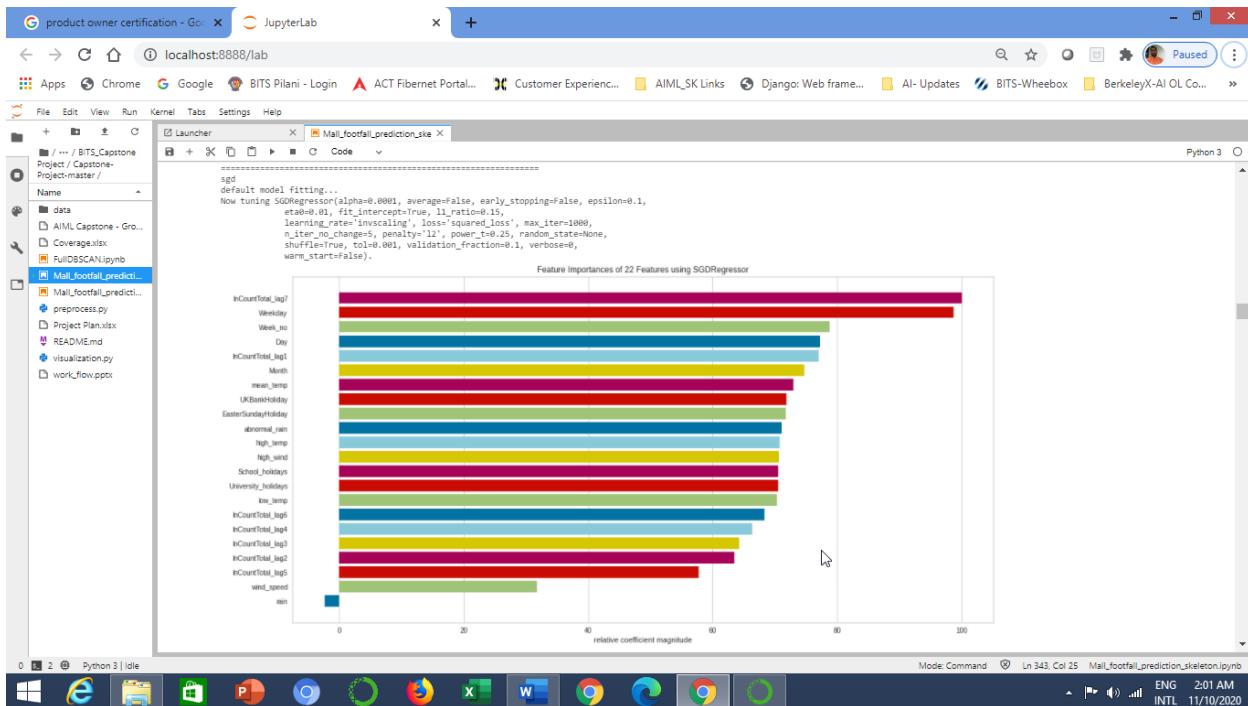
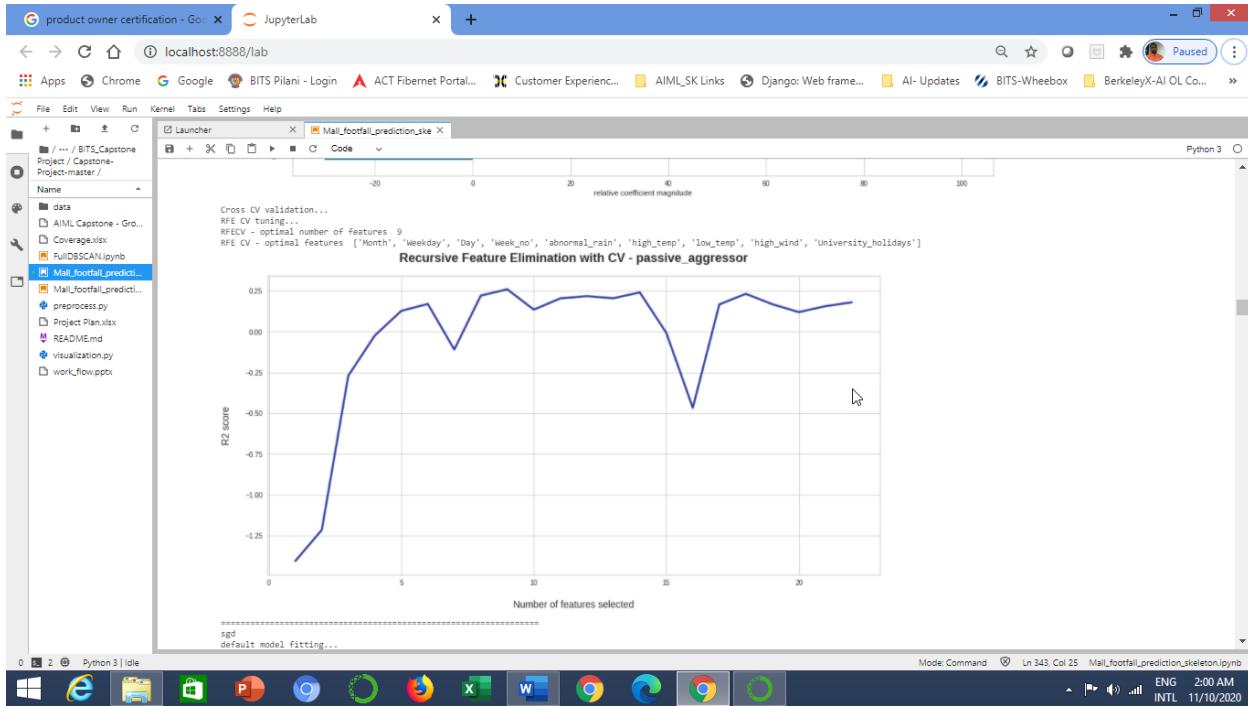


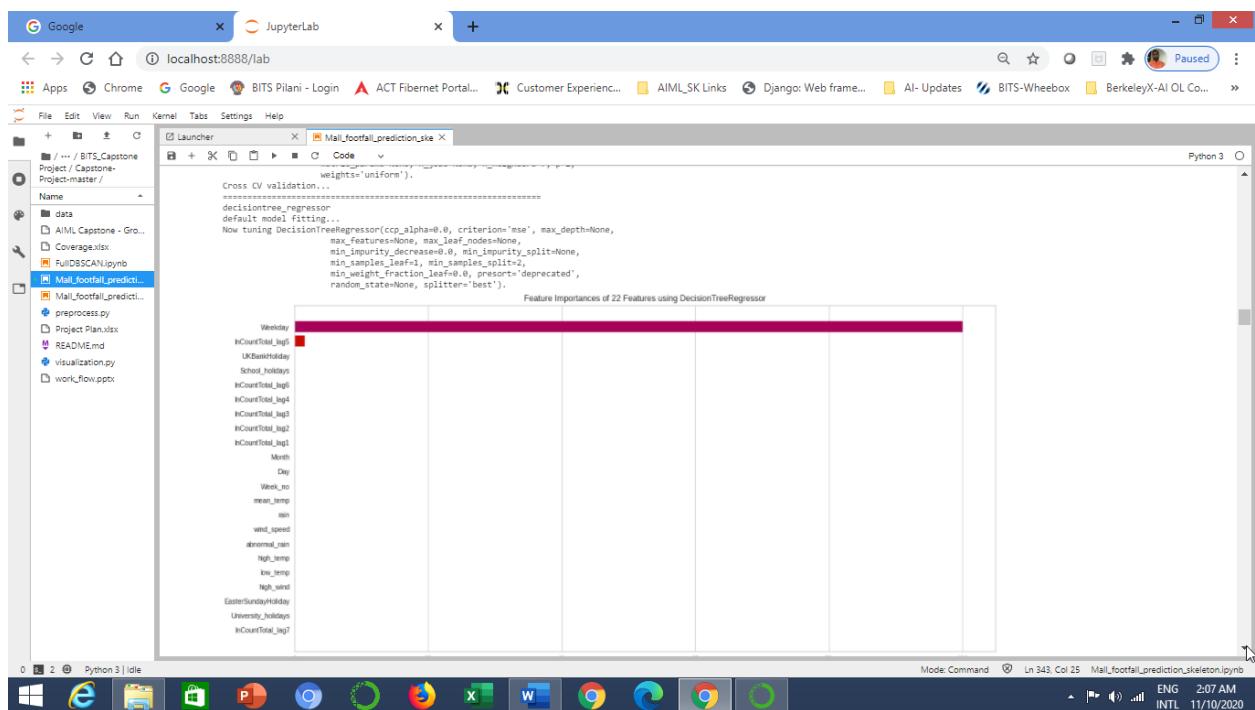
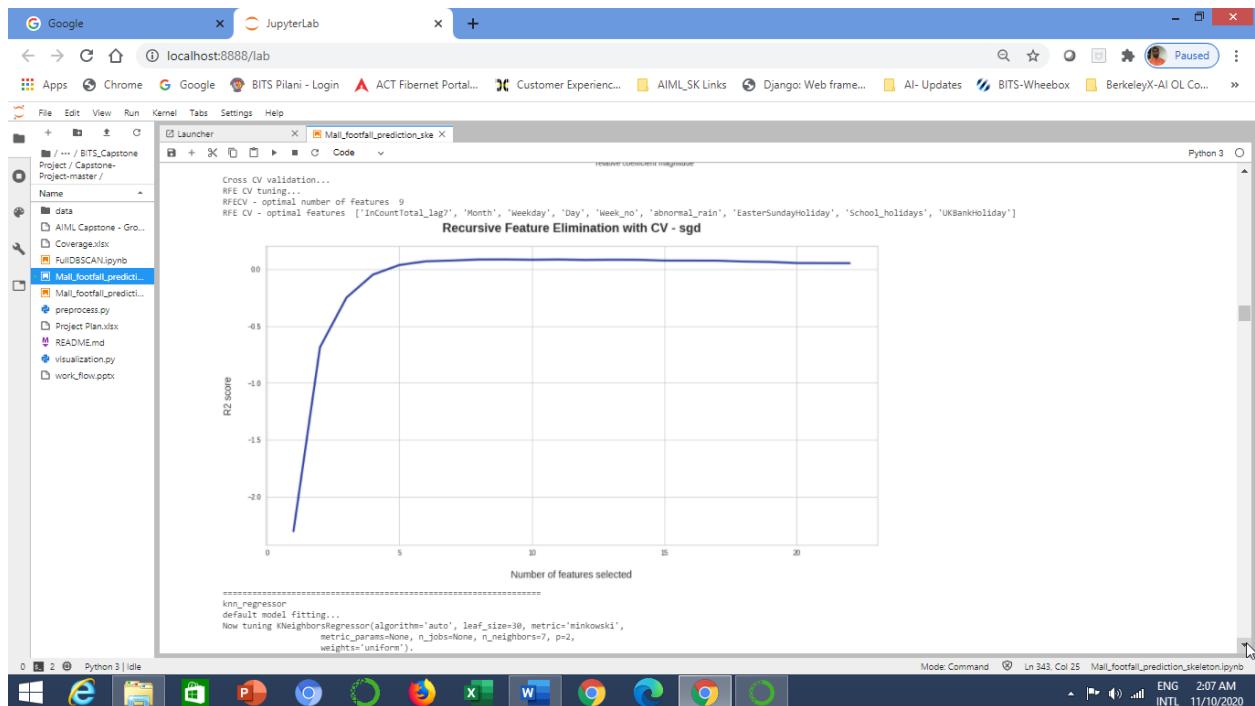


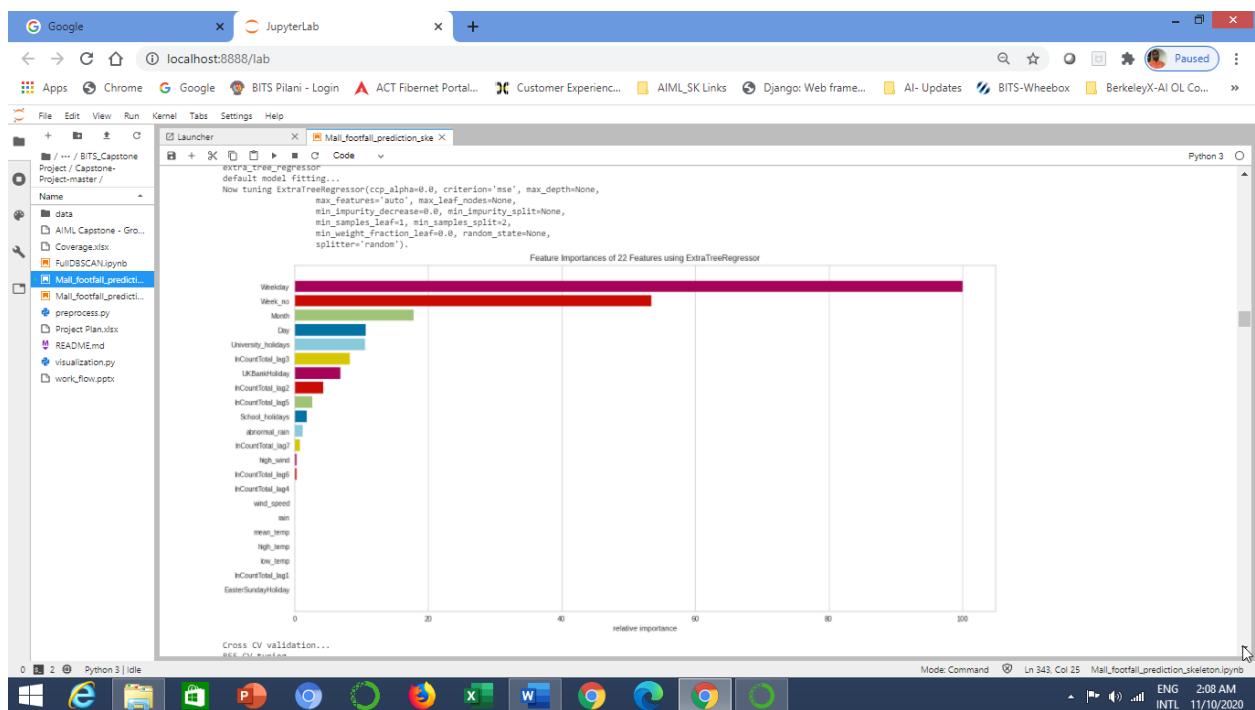
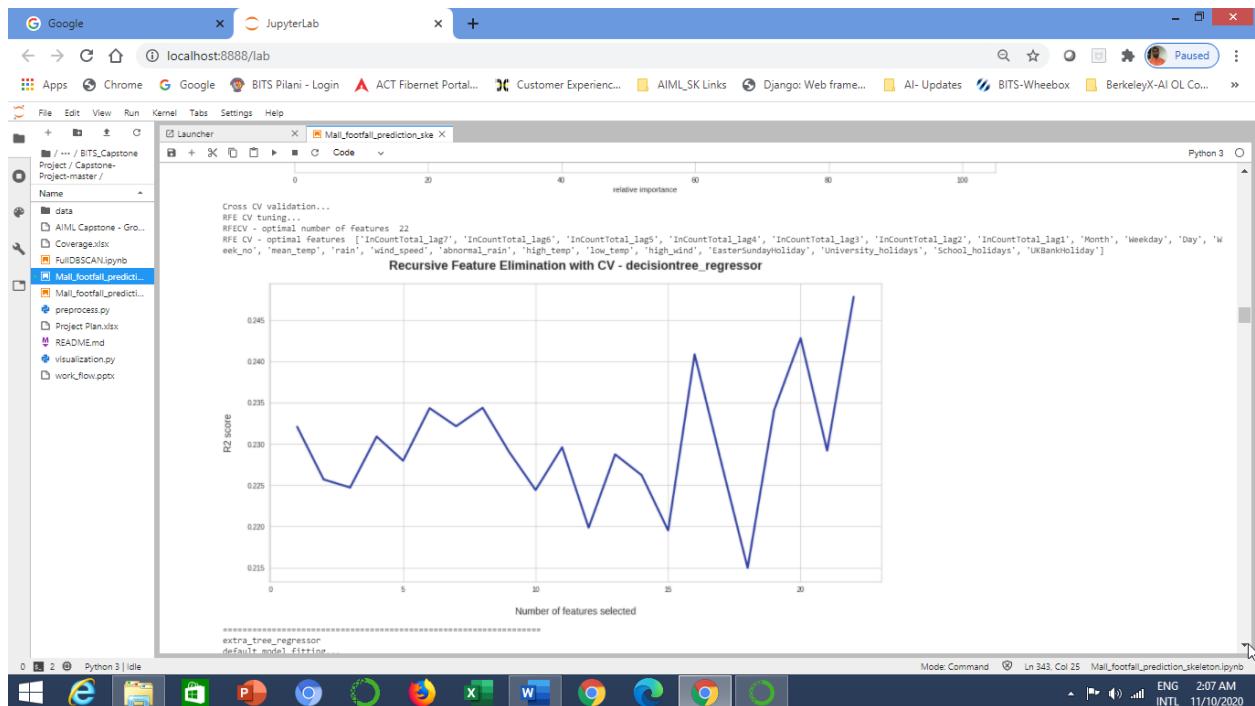












Google JupyterLab

localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

Launcher Mail_footfall_prediction_skeleton.ipynb

```
Cross CV validation...
RFE CV tuning...
RFECV - optimal number of Features 7
RFE CV - optimal features ['InCountTotal_lag1', 'Weekday', 'Day', 'Week_no', 'wind_speed', 'University_holidays', 'School_holidays']

Recursive Feature Elimination with CV - extra_tree_regressor
```

R2 score

Number of features selected

```
=====
SVM regression
default model fitting...
Now tuning SVR(C=1.0, cache_size=200, coef0=0.0, degree=3, epsilon=0.1, gamma='scale',
kernel='rbf', max_iter=-1, shrinking=True, tol=0.001, verbose=False).
Cross CV validation...
```

Mode: Command Ln 343. Col 25 Mail_footfall_prediction_skeleton.ipynb

Google JupyterLab

localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

Launcher Mail_footfall_prediction_skeleton.ipynb

```
=====
kernel='pyspark', max_iter=-1, shrinking=True, tol=0.001, verbose=False).
Cross CV validation...

adaBoost
default model fitting...
Now tuning AdaBoostRegressor(base_estimator=None, learning_rate=1.0, loss='linear',
n_estimators=200, random_state=None).
```

Feature Importances of 22 Features using AdaBoostRegressor

relative importance

```
Cross CV validation...
RFE CV tuning...
RFECV - optimal number of Features 7
```

Mode: Command Ln 343. Col 25 Mail_footfall_prediction_skeleton.ipynb

Google JupyterLab

localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

RFE CV tuning...

```
RFE CV - optimal number of features 18
RFE CV - optimal features ['InCountTotal_lag6', 'InCountTotal_lag5', 'InCountTotal_lag4', 'InCountTotal_lag3', 'InCountTotal_lag2', 'InCountTotal_lag1', 'Month', 'Weekday', 'Week_no', 'mean_temp', 'rain', 'wind_speed', 'EasterSundayHoliday', 'University_holidays', 'School_holidays', 'UKBankHoliday']
```

Recursive Feature Elimination with CV - adaBoost

R2 score

Number of features selected

```
=====
bagging
default model fitting...
Now tuning BaggingRegressor(base_estimator=None, bootstrap=True, bootstrap_features=False,
max_features=1.0, max_samples=1.0, n_estimators=200,
n_jobs=None, oob_score=False, random_state=None, verbose=0,
warm_start=False).
```

Python 3 | idle

Mode: Command Ln 343. Col 25 Mail_footfall_prediction_skeleton.ipynb

ENGLISH 2:09 AM INTL 11/10/2020

Google JupyterLab

localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

Cross CV validation...

```
=====
random_forest
default model fitting...
Now tuning RandomForestRegressor(bootstrap=True, ccp_alpha=0.0, criterion='mse',
max_depth=None, max_features='auto', max_leaf_nodes=None,
max_samples=None, min_impurity_decrease=0.0,
min_impurity_split=None, min_samples_leaf=1,
min_samples_split=2, min_weight_fraction_leaf=0.0,
n_estimators=200, n_jobs=None, oob_score=False,
random_state=None, verbose=0, warm_start=False).
```

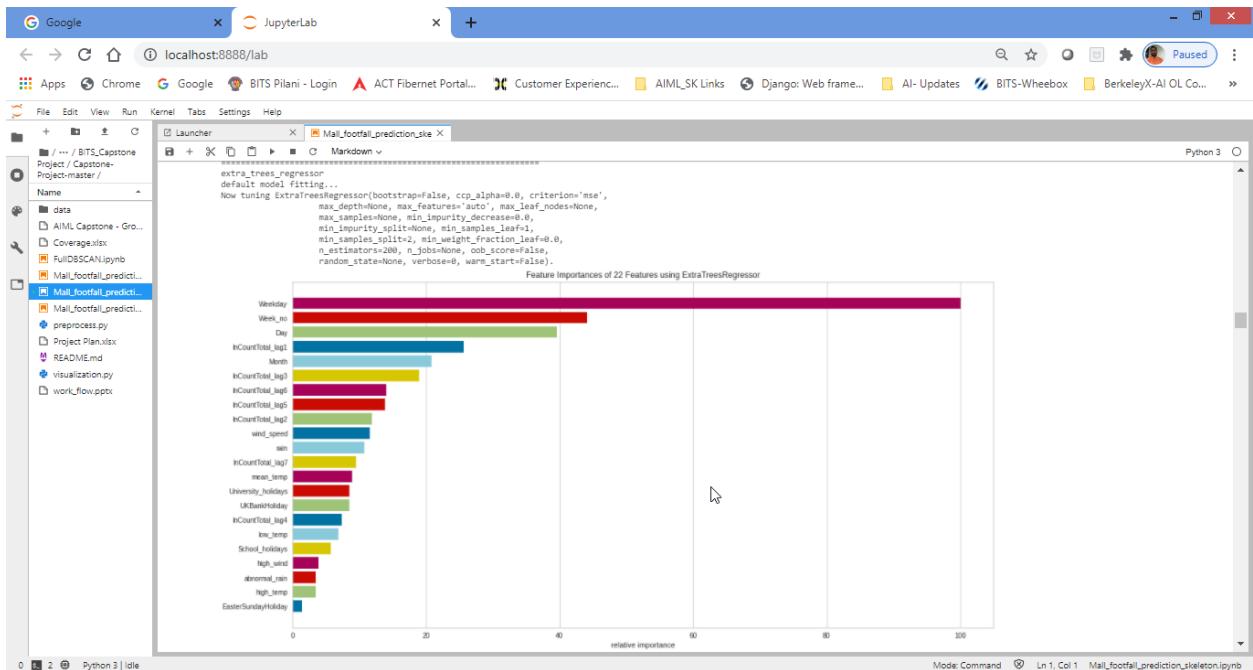
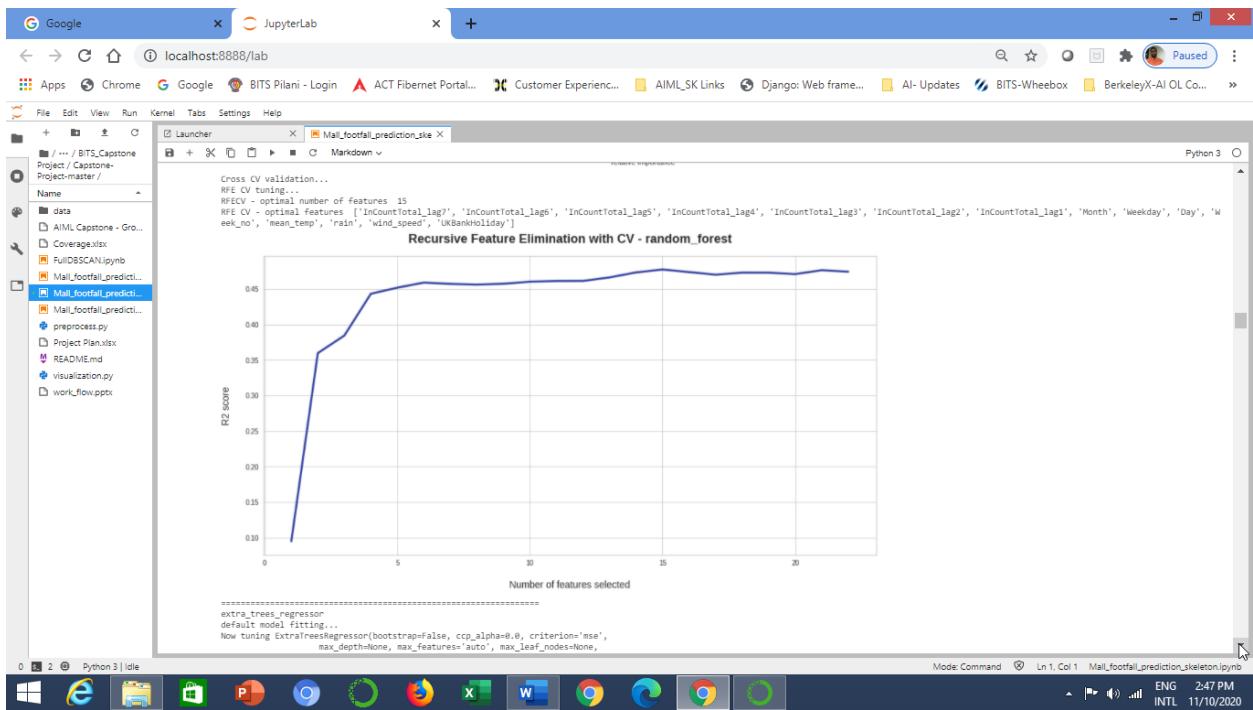
Feature Importances of 22 Features using RandomForestRegressor

Weekday
Week_no
Day
InCountTotal_lag1
InCountTotal_lag2
InCountTotal_lag3
InCountTotal_lag4
Month
InCountTotal_lag5
mean_temp
InCountTotal_lag6
rain
wind_speed
InCountTotal_lag7
UKBankHoliday
School_holidays
InCountTotal_lag8
InCountTotal_lag9
InCountTotal_lag10
InCountTotal_lag11
InCountTotal_lag12
InCountTotal_lag13
InCountTotal_lag14
InCountTotal_lag15
InCountTotal_lag16
InCountTotal_lag17
InCountTotal_lag18
InCountTotal_lag19
InCountTotal_lag20
InCountTotal_lag21
InCountTotal_lag22

Python 3 | idle

Mode: Command Ln 343. Col 25 Mail_footfall_prediction_skeleton.ipynb

ENGLISH 2:10 AM INTL 11/10/2020



Google JupyterLab

localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

Launcher Mail_footfall_prediction_ske

```
gradient_boosting_regressor
default model fitting...
Now tuning GradientBoostingRegressor(alpha=0.9, ccp_alpha=0.0, criterion='friedman_mse',
init=None, learning_rate=0.1, loss='ls', max_depth=3,
max_features=None, max_leaf_nodes=None, min_impurity_decrease=0.0,
min_impurity_split=None, min_samples_leaf=1, min_samples_split=2,
min_weight_fraction_leaf=0.0, n_estimators=200,
n_iter_no_change=None, presort='deprecated',
random_state=None, subsample=1.0, validation_fraction=0.1,
verbose=0, warm_start=False).
```

Feature Importances of 22 Features using GradientBoostingRegressor

0 2 Python 3 | idle Mode Command Ln 1, Col 1 Mail_footfall_prediction_skeleton.ipynb

Google JupyterLab

localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

Launcher Mail_footfall_prediction_ske

```
Cross CV validation...
RFE CV tuning...
RFE CV - optimal number of features: 13
RFE CV - optimal features: ['InCountTotal_lag7', 'InCountTotal_lag6', 'InCountTotal_lag5', 'InCountTotal_lag3', 'InCountTotal_lag2', 'InCountTotal_lag1', 'Month', 'Weekday', 'Day', 'Week_no', 'mean_temp', 'rain', 'wind_speed']
```

Recursive Feature Elimination with CV - extra_trees_regressor

relative importance

R2 score

Number of features selected

```
gradient_boosting_regressor
default model fitting...
Now tuning GradientBoostingRegressor(alpha=0.9, ccp_alpha=0.0, criterion='friedman_mse',
```

0 2 Python 3 | idle Mode Command Ln 1, Col 1 Mail_footfall_prediction_skeleton.ipynb

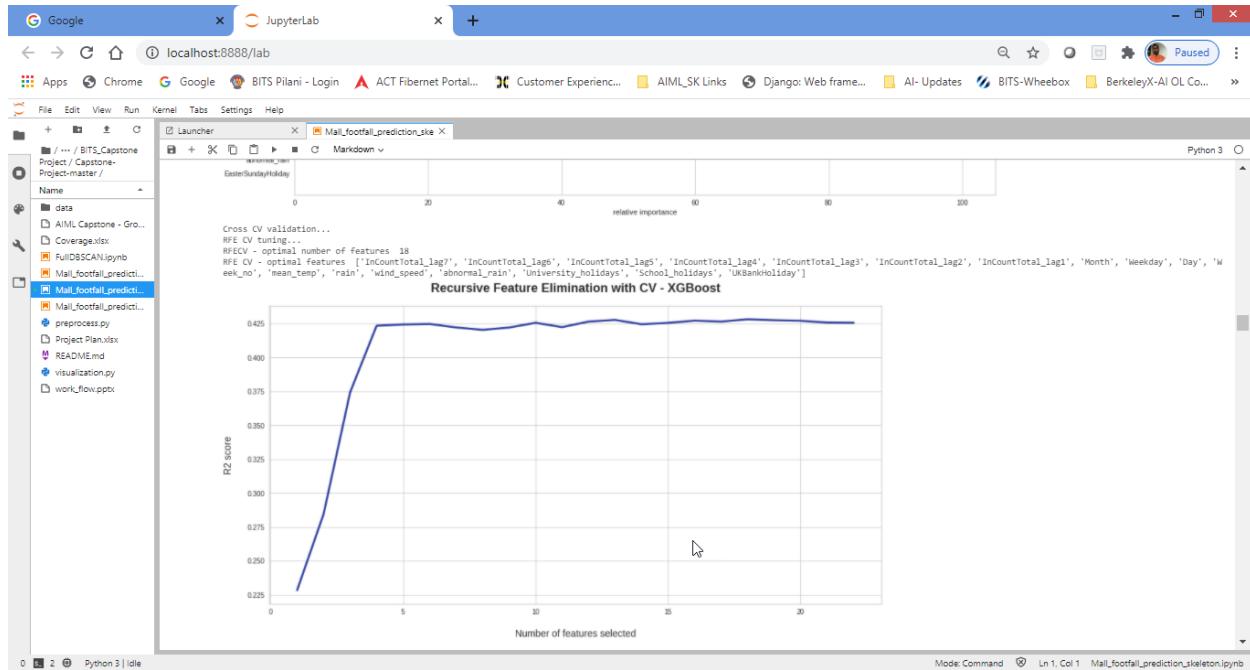
The screenshot shows a JupyterLab interface with a Python 3 kernel. The left sidebar displays a file tree for a 'BITS_Capstone' project, including files like 'data', 'preprocess.py', and 'Mail_footfall_prediction.ipynb'. The main area shows the code for fitting an XGBoost model and a corresponding feature importance plot.

```
File Edit View Run Kernel Tabs Settings Help
+ - C
localhost:8888/lab
File Launcher Mail_footfall_prediction_skeleton.ipynb
File Edit View Run Kernel Tabs Settings Help
+ - C
BITS Pilani - Login ACT Fibernet Portal... Customer Experience... AIML_SK Links Django: Web frame... AI- Updates BITS-Wheebox BerkeleyX-AI OL Co...
File Edit View Run Kernel Tabs Settings Help
+ - C
Python 3
Number of features selected
XGBoost
default model fitted.
Now tuning regression (base_score=0.5, booster='gbtree', colsample_bytree=1,
colsample_bynode=1, colsample_bytree=1, gamma=1, importance_type='gain', learning_rate=0.1, max_delta_step=0,
max_depths=3, min_child_weight=1, missing=None, n_estimators=100,
n_jobs=1, nthread=None, objective='reg:squarederror',
random_state=0, reg_alpha=0, reg_lambda=1, scale_pos_weight=1,
seed=None, silent=None, subsample=1, verbosity=1)
Feature Importances of 22 Features using XGBRegressor
Weekday
Week_no
InCountTotal_1stP
Day
InCountTotal_2ndP
Month
UKBankholiday
ian
InCountTotal_3rdP
InCountTotal_4thP
InCountTotal_5thP
InCountTotal_6thP
InCountTotal_7thP
University_holidays
wind_speed
mean_temp
high_temp
School_holidays
low_temp
high_rain
abnormal_rain
Catering_outside

```

```
File Edit View Run Kernel Tabs Settings Help
+ C Mail_footfall_prediction_skeleton.ipynb
File Edit View Run Kernel Tabs Settings Help
+ Mail_footfall_prediction_skeleton.ipynb
Python 3
Kernel: Mail_footfall_prediction_skeleton.ipynb
Name: Mail_footfall_prediction_skeleton.ipynb
data
Project Capstone - Group 1
Coverage.xlsx
FullDBSCAN.ipynb
Mail_footfall_prediction_skeleton.ipynb
Mail_footfall_prediction.ipynb
preprocess.py
Project Plan.xlsx
README.md
visualization.py
work_flow.pptx

=====
lightGBMBoosting
default model fitting...
Now tuning LightGBMRegressor('boosting_type':'gbdt', 'class_weight':None, 'colsample_bytree':1.0, 'importance_type':'split', 'learning_rate':0.001, 'max_depth':-1, 'min_child_samples':20, 'min_child_weight':0.001, 'min_split_gain':0.0, 'n_estimators':100, 'n_jobs':-1, 'num_leaves':31, 'objective':None, 'random_state':None, 'reg_alpha':0.0, 'reg_lambda':0.0, 'silent':True, 'subsample':1.0, 'subsample_for_bin':200000, 'subsample_freq':0).
Cross CV validation...
=====
CatBoostRegressor
default model fitting...
Now tuning catboost.core.CatBoostRegressor object at 0x7f9278bf8940>.
Cross CV validation...
=====
NeuralNet_MLP
default model fitting...
Now tuning MuRegression('activation':'relu', 'alpha':0.0001, 'batch_size':'auto', 'beta_1':0.9, 'beta_2':0.999, 'early_stopping':False, 'epsilon':1e-09, 'hidden_layer_sizes':(100,), 'learning_rate':'constant', 'loss':'logistic', 'max_iter':100, 'max_fun':100000, 'momentum':0.9, 'n_iter_no_change':10, 'nesterovs_momentum':True, 'power_t':0.5, 'random_state':None, 'shuffle':True, 'solver':'adam', 'tol':0.0001, 'validation_fraction':0.1, 'verbose':False, 'warm_start':False).
Cross CV validation...
```



```
LightGBMRegressor
default model fitting...
Now tuning LightGBMRegressor
{'objective_type': 'lgbt', 'class_weight': None, 'colsample_bytree': 1.0,
 'learning_rate': 0.01, 'max_depth': 1,
 'min_child_samples': 10, 'min_child_weight': 0.001, 'min_split_gain': 0.0,
 'n_estimators': 100, 'n_jobs': -1, 'num_leaves': 31, 'objective': 'None',
 'random_state': None, 'reg_alpha': 0.0, 'reg_lambda': 0.0, 'silent': True,
 'subsample': 1.0, 'subsample_for_bin': 200000, 'subsample_freq': 0}.
Cross CV validation...
```

```
CategoricalGradientBoosting
default model fitting...
Now tuning CatBoost.core.CatBoostRegressor object at 0x7f92788f8940>.
Cross CV validation...
```

```
NeuralNet_MLP
default model fitting...
Now tuning MLPRegressor(activation='relu', alpha=0.0001, batch_size='auto', beta_1=0.9,
beta_2=0.999, early_stopping=False, epsilon=1e-08,
hidden_layer_sizes=(100,), learning_rate='constant',
learning_rate_init=0.001, max_fun=15000, max_iter=200,
momentum=0.9, n_iter_no_change=10, nesterovs_momentum=True,
power_t=0.5, random_state=None, shuffle=True, solver='adam',
tol=0.0001, validation_fraction=0.1, verbose=False,
warm_start=False).
Cross CV validation...
```

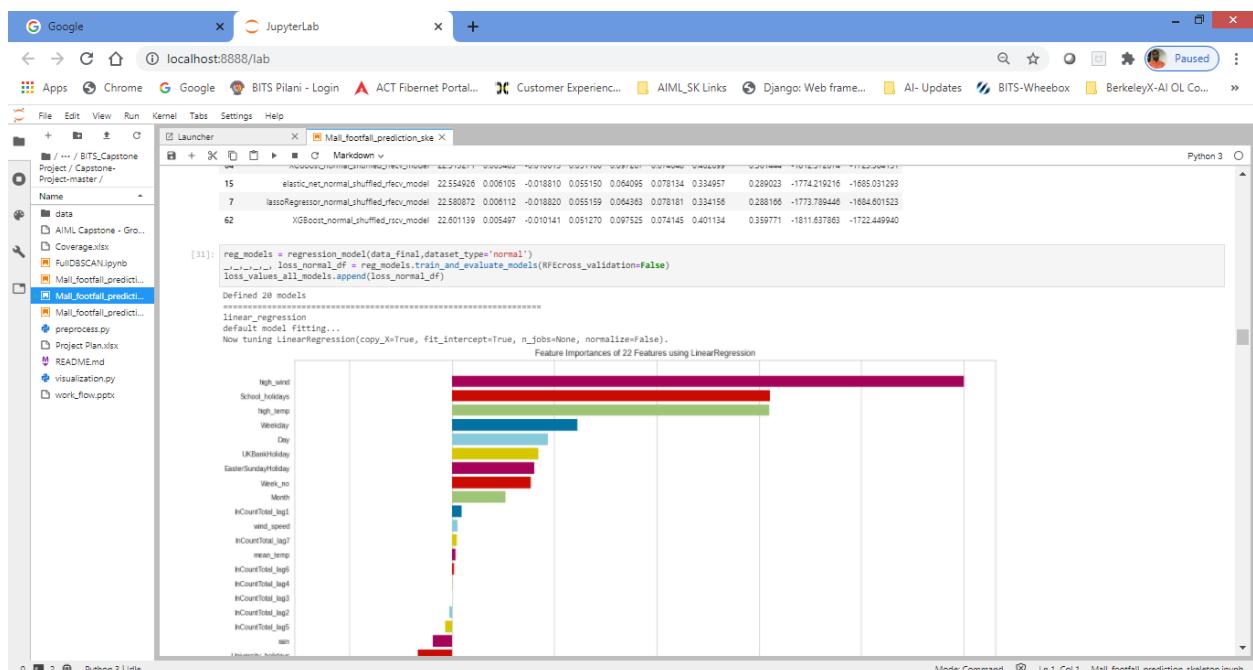
Google JupyterLab localhost:8888/lab

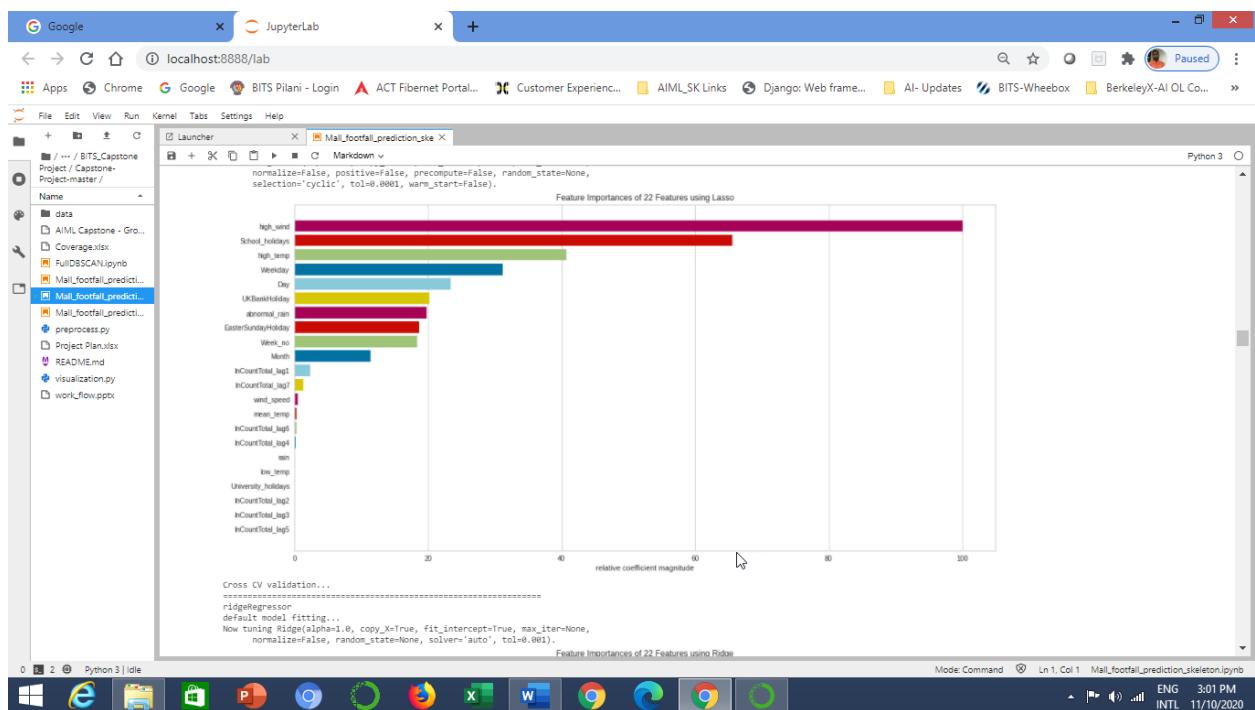
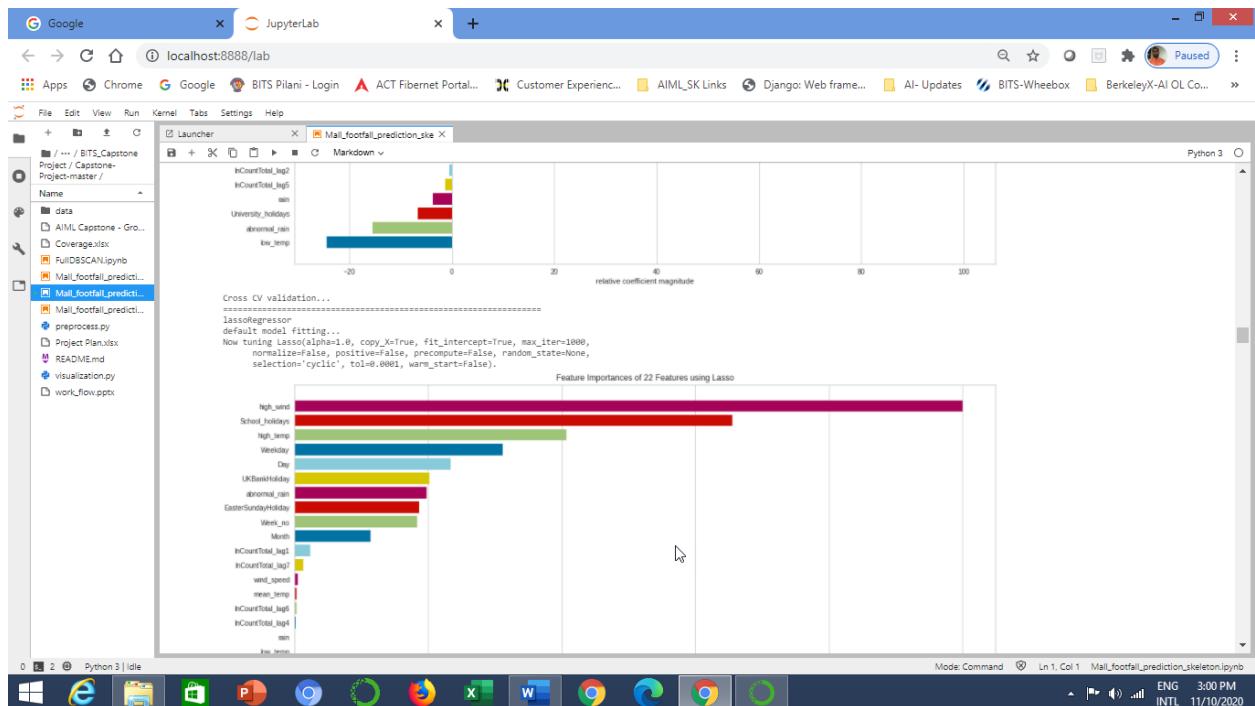
File Edit View Run Kernel Tabs Settings Help

Launcher Mail_footfall_prediction_ske

```
[30]: loss_normal_df.sort_values('mape').head(25)
```

	Model_type	mape	mse	me	mae	mpe	rmse	R2 score	Adj R2 score	AIC	BIC
69	CategoricalGradientBoosting_normal_shuffled_rfcv	0.005194	-0.015631	0.049555	0.038009	0.072067	0.434226	0.395149	-1831.931066	-1742.743143	
57	gradient_boosting_regressor_normal_shuffled_rf	0.005815	-0.019362	0.052390	0.027793	0.074909	0.388287	0.346036	-1804.060158	-1714.872235	
68	CategoricalGradientBoosting_normal_shuffled_rf	0.004999	-0.014414	0.048720	0.035079	0.070706	0.455398	0.417783	-1845.546436	-1756.358513	
49	random_forest_normal_shuffled_default_model	0.004767	-0.012446	0.047276	0.064320	0.069042	0.480725	0.444489	-1882.547285	-1773.359362	
46	bagging_normal_shuffled_default_model	0.004732	-0.012377	0.048501	0.064547	0.068793	0.484467	0.448860	-1865.129562	-1775.941659	
61	XGBoost_normal_shuffled_default_model	0.005281	0.005520	-0.017470	0.052288	0.044728	0.074925	0.357115	-1810.193559	-1721.005456	
60	gradient_boosting_regressor_normal_shuffled_rf	0.005433	-0.017211	0.051460	0.045572	0.073708	0.080174	0.367297	-1815.859485	-1726.671582	
53	extra_trees_regressor_normal_shuffled_rf	0.005094	-0.014785	0.049404	0.058715	0.071512	0.045091	0.406764	-1838.859078	-1749.665155	
45	adaBoost_normal_shuffled_rf	0.005237	-0.005556	0.050943	0.038739	0.072364	0.049550	0.390149	-1628.992344	-1739.804421	
65	LightGradientBoosting_normal_shuffled_rf	0.005345	-0.015181	0.050980	0.055039	0.073246	0.415568	0.375202	-1820.347939	-1731.160007	
66	LightGradientBoosting_normal_shuffled_rf	0.005155	-0.015762	0.049409	0.059171	0.071660	0.440596	0.401959	-1835.973252	-1746.785329	
43	adaBoost_normal_shuffled_rf	0.005345	-0.015170	0.058657	0.072488	0.427593	0.380800	-1827.799953	-1738.582030		
58	gradient_boosting_regressor_normal_shuffled_rf	0.005532	-0.018023	0.052252	0.055551	0.074380	0.397331	0.355770	-1809.377964	-1720.190041	
54	extra_trees_regressor_normal_shuffled_rf	0.005201	-0.015120	0.048837	0.065786	0.072115	0.433409	0.394240	-1831.453730	-1742.265807	
52	random_forest_normal_shuffled_rf	0.005072	-0.019348	0.048512	0.074558	0.071221	0.447441	0.409276	-1840.386442	-1751.180519	
50	random_forest_normal_shuffled_rf	0.005154	-0.014407	0.048533	0.075950	0.071792	0.438538	0.399759	-1834.662278	-1745.474355	
3	linear_regression_normal_shuffled_rf	0.005224	0.006029	-0.019224	0.054684	0.054663	0.077445	0.432304	-1778.706762	-1699.510839	
56	extra_trees_regressor_normal_shuffled_rf	0.005250	-0.015162	0.050158	0.073199	0.072453	0.428148	0.388648	-1828.114691	-1738.926768	
47	bagging_normal_shuffled_rf	0.005581	-0.015118	0.057781	0.074570	0.094243	0.352404	0.362404	-1807.535365	-1718.365642	
11	ridgeRegressor_normal_shuffled_rf	0.006158	-0.020317	0.035581	0.052538	0.078474	0.329157	0.283233	-1771.119537	-1681.931014	
72	NeuralNet_MLP_normal_shuffled_rf	0.005971	-0.016289	0.055082	0.064159	0.077274	0.340914	0.304585	-1782.120210	-1692.932287	
64	XGBoost_normal_shuffled_rf	0.005483	-0.010013	0.051166	0.097267	0.074048	0.402699	0.361444	-1812.872074	-1723.384151	





Google JupyterLab

localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

Launcher Mail_footfall_prediction_ske

```

ridgeRegressors
default model fitting...
Now tuning Ridge(alpha=1.0, copy_X=True, fit_intercept=True, max_iter=None,
normalize=False, random_state=None, solver='auto', tol=0.001)

Feature Importances of 22 Features using Ridge

high_wind
School_holiday
high_temp
Weekday
Day
UKBankHoliday
EasterSundayHoliday
Week_no
Month
InCountTotal_lag1
InCountTotal_lag2
InCountTotal_lag3
InCountTotal_lag4
InCountTotal_lag5
University_holiday
InCountTotal_lag6
rain
low_temp
abnormal_rain

```

Cross CV validation...

```

elasticNet
default model fitting...
Now tuning ElasticNet(alpha=1.0, copy_X=True, fit_intercept=True, l1_ratio=0.5,
max_iter=1000, normalize=False, positive=False, precompute=False,
random_state=None, selection='cyclic', tol=0.0001, warm_start=False).

```

Mode Command Ln 1, Col 1 Mail_footfall_prediction_skeleton.ipynb

Python 3

Google JupyterLab

localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

Launcher Mail_footfall_prediction_ske

```

elasticNet
default model fitting...
Now tuning ElasticNet(alpha=1.0, copy_X=True, fit_intercept=True, l1_ratio=0.5,
max_iter=1000, normalize=False, positive=False, precompute=False,
random_state=None, selection='cyclic', tol=0.0001, warm_start=False).

Feature Importances of 22 Features using ElasticNet

high_wind
School_holiday
high_temp
Weekday
Day
UKBankHoliday
abnormal_rain
EasterSundayHoliday
Week_no
Month
InCountTotal_lag1
InCountTotal_lag2
InCountTotal_lag3
InCountTotal_lag4
InCountTotal_lag5
rain
low_temp
University_holiday
InCountTotal_lag6
InCountTotal_lag7

```

Cross CV validation...

```

huberRegressor

```

Mode Command Ln 1, Col 1 Mail_footfall_prediction_skeleton.ipynb

Python 3

Google JupyterLab

localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

Launcher Mail_footfall_prediction_ske

```
huberregressor
default model fitting...
Now tuning HuberRegressor(alpha=0.0001, epsilon=1.35, fit_intercept=True, max_iter=100,
                           tol=1e-05, warm_start=False).

Feature Importances of 22 Features using HuberRegressor
```

```
Cross CV validation...
passive_aggressor
default model fitting...
Now tuning PassiveAggressiveRegressor(C=1.0, average=False, early_stopping=False,
```

Mode Command Ln 1, Col 1 Mail_footfall_prediction_skeleton.ipynb

Python 3

Google JupyterLab

localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

Launcher Mail_footfall_prediction_ske

```
passive_aggressor
default model fitting...
Now tuning PassiveAggressiveRegressor(C=1.0, average=False, early_stopping=False,
                                         epsilon=0.1, fit_intercept=True,
                                         loss='epsilon_insensitive', max_iter=1000,
                                         n_iter_no_change=5, random_state=None, shuffle=True,
                                         tol=0.001, validation_fraction=0.1, verbose=0,
                                         warm_start=False).
```

```
Cross CV validation...
```

Mode Command Ln 1, Col 1 Mail_footfall_prediction_skeleton.ipynb

Python 3

The screenshot shows a Jupyter Notebook interface with a Python 3 kernel. The code cell displays the following SGDRegressor configuration:

```
sgd
default model fitting...
Now tuning SGDRegressor(alpha=0.0001, average=False, early_stopping=False, epsilon=0.1,
eta0=0.01, fit_intercept=True, l1_ratio=0.15,
learning_rate='invscaling', loss='squared_loss', max_iter=1000,
n_iter_no_change=5, penalty='l2', power_t=0.25, random_state=None,
shuffle=True, tol=0.001, validation_fraction=0.1, verbose=0,
warm_start=False).
```

The notebook also contains a bar chart titled "Feature Importances of 22 Features using SGDRegressor". The x-axis represents the "relative coefficient magnitude" from 0 to 100. The y-axis lists 22 features. The bars are color-coded by category: red for InCountTotal_1ag7, green for Weekday, Week_no, Day, InCountTotal_1ag1, Month, UKBankHoliday, EasterSundayHoliday, abnormal_rain, high_temp, high_wind, School_holidays, University_holidays, low_temp, InCountTotal_1ag8, InCountTotal_1ag9, InCountTotal_1ag10, InCountTotal_1ag11, InCountTotal_1ag12, and InCountTotal_1ag13; blue for InCountTotal_1ag6; and yellow for InCountTotal_1ag5.

Feature	Relative Coefficient Magnitude (approx.)
InCountTotal_1ag7	95
Weekday	90
Week_no	85
Day	75
InCountTotal_1ag1	70
Month	65
UKBankHoliday	60
EasterSundayHoliday	55
abnormal_rain	50
high_temp	45
high_wind	45
School_holidays	45
University_holidays	45
low_temp	40
InCountTotal_1ag8	35
InCountTotal_1ag9	30
InCountTotal_1ag10	30
InCountTotal_1ag11	30
InCountTotal_1ag12	25
InCountTotal_1ag13	20
InCountTotal_1ag6	10
InCountTotal_1ag5	10

Google JupyterLab

localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

+ - +

Launcher Mail_footfall_prediction_ske

Markdown

Cross CV validation.....

```
kon_regressor
default model fitting...
Now tuning NeighborsRegressor(algorithm='auto', leaf_size=30, metric='minkowski',
                               metric_params=None, n_jobs=None, n_neighbors=7, p=2,
                               weights='uniform').
```

Cross CV validation.....

```
decisiontree_regressor
default model fitting...
Now tuning DecisionTreeRegressor(ccp_alpha=0.0, criterion='mse', max_depth=None,
                                   max_features=None, max_leaf_nodes=None,
                                   min_impurity_decrease=0.0, min_impurity_split=None,
                                   min_samples_leaf=1, min_samples_split=2,
                                   min_weight_fraction_leaf=0.0, presort='deprecated',
                                   random_state=None, splitter='best').
```

Feature Importances of 22 Features using DecisionTreeRegressor

relative coefficient magnitude

Python 3

Google JupyterLab

localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

Launcher Mail_footfall_prediction_skeleton.py

```
decisiontree_regressor
default model fitting...
Now tuning DecisionTreeRegressor(ccp_alpha=0.0, criterion='mse', max_depth=None,
                                    max_features=None, max_leaf_nodes=None,
                                    min_impurity_decrease=0.0, min_impurity_split=None,
                                    min_samples_leaf=1, min_samples_split=2,
                                    min_weight_fraction_leaf=0.0, presort='deprecated',
                                    random_state=None, splitter='best').
```

Feature Importances of 22 Features using DecisionTreeRegressor

Cross CV validation...

Python 3 | idle

Mode Command Ln 1, Col 1 Mail_footfall_prediction_skeleton.py

ENGLISH 3:06 PM INTL 11/10/2020

Google JupyterLab

localhost:8888/lab

File Edit View Run Kernel Tabs Settings Help

Launcher Mail_footfall_prediction_skeleton.py

```
Cross CV validation...
=====
extra_tree_regressor
default model fitting...
Now tuning ExtraTreeRegressor(ccp_alpha=0.0, criterion='mse', max_depth=None,
                                max_features='auto', max_leaf_nodes=None,
                                min_impurity_decrease=0.0, min_impurity_split=None,
                                min_samples_leaf=1, min_samples_split=2,
                                min_weight_fraction_leaf=0.0, random_state=None,
                                splitter='random').
```

Feature Importances of 22 Features using ExtraTreeRegressor

Cross CV validation...

Python 3 | idle

Mode Command Ln 1, Col 1 Mail_footfall_prediction_skeleton.py

ENGLISH 3:07 PM INTL 11/10/2020