

Configuration Manual

Demand prediction in a bike-sharing system using machine learning techniques Masters in Data Analytics (B)

> Ashish Rawat X18185801

School of Computing National College of Ireland

Supervisor: Paul Stynes

National College of Ireland Project Submission Sheet School of Computing



Student Name:	Ashish Rawat
Student ID:	X18185801
Programme:	Masters in Data Analytics (B)
Year:	2020
Module:	Demand prediction in a bike-sharing system using machine
	learning techniques
Supervisor:	Paul Stynes
Submission Due Date:	17/08/2020
Project Title:	Configuration Manual
Word Count:	901
Page Count:	11

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

<u>ALL</u> internet material must be referenced in the bibliography section. Students are required to use the Referencing Standard specified in the report template. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

Signature:	
Date:	17th August 2020

PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST:

Attach a completed copy of this sheet to each project (including multiple copies).	
Attach a Moodle submission receipt of the online project submission, to	
each project (including multiple copies).	
You must ensure that you retain a HARD COPY of the project, both for	
your own reference and in case a project is lost or mislaid. It is not sufficient to keep	
a copy on computer.	

Assignments that are submitted to the Programme Coordinator office must be placed into the assignment box located outside the office.

Office Use Only	
Signature:	
Date:	
Penalty Applied (if applicable):	

Configuration Manual

Ashish Rawat X18185801

1 Introduction

The following configuration manual discusses the hardware and software requirements for this research. It also describes the programming stages of the implementation in detail.

2 System Setup

2.1 Hardware specifications

- Processor: Intel(R) Core(TM) i7-8565U CPU @ 1.80GHz 1.99GHz
- RAM: 8 GB
- System Type: 64-bit Operating System
- Storage: 512 SSD with 1 TB HDD

2.2 Software required

- Microsoft Excel: It is a spreadsheet program, which is used to store data in a grid of rows and columns. And in context to this study, has been used to access CSV (comma separated files) format files.
- Google Colaboratory: It is a free cloud service that has allowed continuous accessibility of 12 GB NVIDIA Tesla K80 GPU for 12 hours (Bar-El; 2018). This study uses the environment of colab to examine metro bike data and create models for machine learning. GPU setting is activated for some models to execute the code in less time.
- Tableau: Some visualizations in the final report were made through Tableau. It is a free visualization software which can be downloaded from Tableau website¹.

3 Project Development

The research work was conducted using version 3.6.7 of Python. And it was implemented in four phases : Data pre-processing and feature engineering, data preparation, modelling machine learning techniques and finally, evaluation.

¹https://www.tableau.com/

3.1 Data Collection

The data for the study is obtained from the Metro Bike website². It consists of 4 csv files containing quarterly data of bike trips made in 2019 and 1 csv file with station details. The data files are downloaded and stored in Google Drive for it to be accessed by the Google Colab.

3.2 Data Processing and Feature Engineering

First of all, Google Drive is connected to the Colab and the necessary libraries are imported as shown in figure 1.

Link the Google Drive with the Colab, to access data files.

```
[ ] from google.colab import drive
    drive.mount('/content/drive')
```

Installing the libraries required.

```
[ ] import os
from os import listdir, makedirs, remove
import os.path as osp
import pandas as pd
from datetime import datetime, timedelta
import numpy as np
from numpy.random import Generator, PCG64
```

Figure 1: Data Setup

To facilitate the processing and achieve the desired data set following functions are developed:

1. "time_bins" function as shown in figure 2 creates a dataframe of start date time and end date time with a fixed interval value and returns it.



Figure 2: time_bin function

²https://bikeshare.metro.net/about/data/

2. "format_data" function loads the bike trip data, one file at a time and performs pre-preprocessing on it. And at end combines it with the time interval dataset obtained from "time_bins" function to develop transformed dataset and weights for graph (figure 3).

O

def	<pre>format_data(file_path, tinterval, interval_size, duration_upper, station_exclude, station_const_ids):</pre>
	<pre>def round_minutes(dt, direction, resolution): new_minute = (dt.minute // resolution + (1 if direction == 'up' else 0)) * resolution return dt + timedelta(minutes-new_minute - dt.minute)</pre>
	<pre># Read bike trip data and load it into a pandas dataframe. df_raw = pd.read_csv(file_path, low_memory= False) df_raw[=Date] = pd.to_datetime(df_raw['idm_time'], dayfirst=True) df_raw[=Traw.=Traw_memame(columns=('end_station : 'EndStation Id', 'start_station': 'StartStation Id', 'duration': 'Duration']) df_raw = df_raw_memame(columns=('end_station Id'), isnd()] df_raw = df_raw_memame(tradiction Id'), isna()] df_raw['EndStation Id'] = df_raw['EndStation Id'].astype(int) df_raw = df_raw.loc[df_raw['Duration'] < duration_upper]</pre>
	<pre># Check to confirm no unknown station value is present in the data which is not recorded in bike station details dataset. s_diff = set(df_raw['startStation Id'].to_list()) - station_const_ids e_diff = set(df_raw['endStation Id'].to_list()) - station_const_ids if len(s_diff) > 0:</pre>
	<pre>raise RuntimeError('In file {} unknown start station IDs: {}'.format(file_path, s_diff)) if len(e_diff) > 0: raise RuntimeError('In file {} unknown end station IDs: {}'.format(file_path, e_diff))</pre>
	<pre># Remove unwanted stations like virtual stations df_raw = df_raw.loc[~df_raw['StartStation Id'].isin(station_exclude)] df_raw = df_raw.loc[~df_raw['EndStation Id'].isin(station_exclude)]</pre>
	<pre># Count the number of bike trips made for station-to-station. df_graph_weight = df_raw.groupby(['StartStation Id', 'EndStation Id']).size()</pre>
	<pre># Maps time of bike trip to a starting time of the interval, as defined by the dataset created by time_bins(). df_raw['End Date Lower Bound'] = df_raw['End Date'].apply(round_minutes,</pre>
	<pre>#Werge df_raw dataset with time interval dataset obtained from the time_bins(). df1 = pd.merge(df_raw, tinterval, left_on='End Date Lower Bound', right_on='inner_limit').drop(columns=['innerval_id': 'End Date ID'), axis=1) df1 = df1.rename({'interval_id': 'End Date ID'), axis=1) df1 = df1.rename({'interval_id': 'Start Date Lower Bound', right_on='inner_limit').drop(columns=['interval_id': 'Start Date ID'), axis=1) df1 = df1.rename({'interval_id': 'Start Date ID', 'rename('interval_id', 'start_lat', 'rend_id', 'start Date Lower Bound', 'Start Date Lower Bound']) </pre>
	<pre># raises an error if len(df) == 0: raise RuntimeError('Incorrect time intervals for file {}'.format(file_path))</pre>
	<pre>g_arrival = df1.groupby(by=['EndStation Id', 'End Date ID']) df_arrival = g_arrival.size() g_departure = df1.groupby(by=['StartStation Id', 'Start Date ID']) df_departure = g_departure.size()</pre>
	<pre># Zero is substituted in place of records with missing value of time ids, s1 = df_arrival.index.get_level_values('EndStation Id').unique() s2 = df_departure.index.get_level_values('StartStation Id').unique() s_index = s1.union(s2) t1 = df_arrival.index.get_level_values('End Date ID').min() t2 = df_departure.index.get_level_values('Start Date ID').max() t_index = pd.Int64Index(list(range(t1, t2 + 1)), dtype='int64') new_index = pd.MultiIndex.from_product([s_index, t_index],</pre>
	<pre>df_all = df_arrival.join(df_departure, how='outer')</pre>
	return df_all, df_graph_weight

Figure 3: format_data function

- 3. "merge_data" function combines individual formatted files into a single file and returns the complete station Id list (figure 4).
- 4. "make_graph_weights" function uses the weights generated from "format_data" function to create a new features "average percent" which is used to assign weights to edges (figure 5).



Figure 4: merge_data function

```
def make_graph_weights(graphs, stations_include, norm='median'):
    # Create index that contains all station pairs that ever appears in the raw data files
    # and apply this index to all graphs.
    mega_index = graphs[0].index
    for graph in graphs[1:]:
        mega_index = mega_index.union(graph.index)
    g_expand = [g.reindex(mega_index, fill_value=0) for g in graphs]
    cat_graph = pd.concat(g_expand).reset_index()
    cat graph = cat graph.loc[cat graph['StartStation Id'].isin(stations include)]
    cat graph = cat graph.loc[cat graph['EndStation Id'].isin(stations include)]
    cat_graph_all = cat_graph.groupby(['StartStation Id', 'EndStation Id']).sum()
    cat_graph_all = cat_graph_all.rename(columns={0: 'total'})
    all_df = []
    #Create another feature in the form of percentage of usage.
    for s_ind, df_slice in cat_graph_all.groupby('StartStation Id'):
        mean slice = df slice['total'].sum()
        df slice.loc[:, 'average percent'] = 100.0 * df slice['total'] / mean slice
        all df.append(df slice)
    cat_graph_all = pd.concat(all_df)
    return cat graph all
```



5. "transform_data" function as shown in figure 6 is the parent function which calls all the above functions and saves the datasets formed into a directory in google drive.



Figure 6: transform_data function

Finally, the values of variables are passed to the "transform_data" function to start the execution process as shown in the figure 7.



Figure 7: Variables used to generate transformed datasets

3.3 Data Preparation

1. For ARIMA and LSTM Model the following data preparation followed in the research. The datasets obtained from processing are utilised to generate individual bike station dataset for both bike demands and docks demands as shown in figure 8.



Figure 8: Data preparation for ARIMA and LSTM models

2. However, for STGCN and TAGCN model graph data is generated using Pytorch Geometric library. And libraries shown in figure 9 are mandatory for the process. And To ease the implementation following functions are developed.

```
import torch
import torch.nn.functional as F
from torch_geometric.data import Data
import os
from os import listdir, makedirs, remove
import os.path as osp
import pandas as pd
from datetime import datetime, timedelta
import numpy as np
from numpy.random import Generator, PCG64
from torch_geometric.data import Dataset, DataLoader
from torch_geometric.data import Data
import torch.nn.functional as funct
```

Figure 9: Libraries necessary for graph data preparation

• "make_edges" function converts the graph weight dataset into edge data compatible with Pytorch geometric (figure 10).



Figure 10: make_edges function

• "time_windows" functions slice the dependent features from the data and for it to be converted into tensors as shown in figure 11.



Figure 11: time_windows function

• Finally, "create_torch_data" function creates a graph structure of nodes and edges from the transformed bikes data (figure 12).

3.4 Modelling

- To train the model 80% of data has been utilised in all the models and rest 20% is used for training. In the case of LSTM model 20% of data from train set id used for validation.
- The model has been implemented as following:



Figure 12: create_torch_data function

1. ARIMA Model

The ARIMA model is implemented using Statsmodel library. To utilise the data for ARIMA stationary test has also been performed and model is implemented as shown in figure 13.



Figure 13: ARIMA model train and test implementation

2. LSTM Model

A five LSTM model has been implemented using Keras library and the data is scaled using MinMaxScaler before training. The code can be referred from figure 14.

0	<pre>model = Sequential() #Adding the first LSTM layer and some Dropout regularisation model.add(LSTM(units = 50, return_sequences = True, input_shape = (X_train.shape[1], 1))) #Adding a second LSTM layer and some Dropout regularisation model.add(Dropout(0.2)) # Adding a third LSTM layer and some Dropout regularisation model.add(Dropout(0.2)) # Adding a third LSTM layer and some Dropout regularisation model.add(Dropout(0.2)) # Adding a third LSTM layer and some Dropout regularisation model.add(Dropout(0.2)) # Adding a fourth LSTM layer and some Dropout regularisation model.add(Dropout(0.2)) # Adding a fourth LSTM layer and some Dropout regularisation model.add(Dropout(0.2)) # Adding the output layer model.add(Dense(units = 1)) model.summary()</pre>
[]	<pre>model.compile(optimizer = 'adam', loss = 'mean_squared_error') # Fitting the model to the Training set history = model.fit(X_train, y_train, epochs = 30, batch_size = 32, validation_split= 0.2, callbacks = EarlyStopping(monitor='val_loss', patience = 5))</pre>

Figure 14: LSTM model implementation

3. STGCN Model

The implementation of STGCN model is a intricate process and was referred from (Ohrn; 2020). The model consists of 4 classes which are referenced in the main "STGCN" class. The figure 15 shows the 3 individual classes and figure 16, the main STGCN class.



Figure 15: Three convolution layer classes of the STGCN model

```
class STGCN(torch.nn.Module):
       graph_conv_kwargs={}):
          super(STGCN, self).__init__()
          self.n_temporal = n_temporal_dim
          self.n_spatial = n_spatial_dim
          self.n_input_channels = n_input_channels
          self.co_temporal = co_temporal
          self.co_spatial = co_spatial
          self.permute_norm = (IDX_CHANNEL, IDX_SPATIAL, IDX_TEMPORAL)
assert n_temporal_dim - 4 * time_conv_length + 4 == 1
          self.model_t_1a = Time1dConvGLU(n_spatial_dim, n_temporal_dim,
                                                  channel_inputs=n_input_channels,
channel_outputs=co_temporal,
                                                  time_convolution_length=time_conv_length)
          self.model_s_1 = SpatialGraphConv(n_spatial_dim, n_temporal_dim - time_conv_length + 1,
                                                    channel_inputs=co_temporal,
                                                    channel_outputs=co_spatial,
                                                    graph_conv_kwargs=graph_conv_kwargs)
          channel_outputs=co_temporal,
time_convolution_length=time_conv_length)
          self.layer_norm_1 = torch.nn.LayerNorm([n_spatial_dim, n_temporal_dim - 2 * time_conv_length + 2])
          self.model_t_2a = Time1dConvGLU(n_spatial_dim, n_temporal_dim - 2 * time_conv_length + 2,
                                                  channel inputs=co temporal,
                                                  channel_outputs=co_temporal,
          time_convolution_length=time_conv_length)
self.model_s_2 = SpatialGraphConv(n_spatial_dim, n_temporal_dim - 3 * time
                                                                                                time conv length + 3,
                                                    channel inputs=co temporal,
                                                    channel_outputs=co_spatial,
          graph_conv_kwargs=graph_conv_kwargs)
self.model_t_2b = Time1dConvGLU(n_spatial_dim, n_temporal_dim - 3 * time_conv_length + 3,
                                                  channel_inputs=co_spatial,
                                                  channel_outputs=co_temporal,
                                                  time_convolution_length=time_conv_length)
          self.layer_norm_2 = torch.nn.LayerNorm([n_spatial_dim, n_temporal_dim - 4 * time_conv_length + 4])
          self.model output = torch.nn.Sequential(torch.nn.Linear(in features=n spatial dim * co temporal,
                                                                               out_features=n_spatial_dim * n_input_channels),
                                                            torch.nn.ReLU().
                                                            def forward(self, data_graph):
         ctronomodectri, dote_graph.
data_step_0 = data_graph
edge_index = data_graph.edge_index
edge_attr = data_graph.edge_attr
self.n_graphs_in_batch = self._infer_batch_size(data_step_0.x)
          data_step_0_x = data_step_0.x
         data_step_1_x = self.model_t_1a(data_step_0_x)
data_step_1 = Data(data_step_1_x, edge_index, edge_attr)
         data_step_1x = self.model_s_1(data_step_1)
data_step_3x = self.model_t_1b(data_step_2x)
data_step_3_x = self.compute_layer_norm(data_step_3_x, self.layer_norm_1)
         data_step_4_x = self.model_t_2a(data_step_3_x)
data_step_4 = Data(data_step_4_x, edge_index, edge_attr)
data_step_5_x = self.model_s_2(data_step_4)
data_step_6_x = self.model_t_2b(data_step_5_x)
data_step_6_x = self.compute_layer_norm(data_step_6_x, self.layer_norm_2)
          data_out_reshape = data_step_6_x.reshape(self.n_graphs_in_batch, self.n_spatial * self.co_temporal)
          data_output_x = self.model_output(data_out_reshape)
data_output_x = data_output_x.reshape(self.n_graphs_in_batch * self.n_spatial, self.n_input_channels)
          return Data(y=data_output_x, edge_index=edge_index, edge_attr=edge_attr)
       def __infer_batch_size(self, data_x):
    if data_x.shape[IDX_SPATIAL] % self.n_spatial == 0:
            n_batches = data_x.shape[IDX_SPATIAL] // self.n_spatial
         else:
            raise RuntimeError('Incorrect batch size')
         return n_batches
        def _compute_layer_norm(self, data_x, norm_func):
          ret = []
          for k_batch in range(self.n_graphs_in_batch):
           data_x_pergraph = data_x.narrow(IDX_SPATIAL, k_batch * self.n_spatial, self.n_spatial)
data_x_pergraph_norm = norm_func(data_x_pergraph.permute(self.permute_norm)).permute(self.permute_norm)
ret.append(data_x_pergraph_norm)
          return torch.cat(ret)
```



4. TAGCN Model

Implementation of TAGCN model is similar to STCGN model, with the only difference in the "SpatialGraphConv" class in the main TAGCN class. Figure 17 represents the SpatialGraphConv class in TAGCN.



Figure 17: SpatialGraphConv class in TAGCN model

3.5 Evaluation

To evaluate the models, three metrics; RMSE, MSE and MAE are utilised which are implemented using sklearn library. The code can be referenced through figure 18



Figure 18: Evaluation metrics used in the research.

References

Bar-El, O. (2018). Getting the most out of your google colab. URL: https://medium.com/@oribarel/getting-the-most-out-of-your-google-colab-2b0585f82403

Ohrn, A. (2020). London bike ride forecasting with graph convolutional networks. URL: https://towardsdatascience.com/london-bike-ride-forecasting-with-graphconvolutional-networks-aee044e48131