

**A Project Report submitted  
For Simulation & Modelling (UCS751)**

by

**Anirudh Anand**

**101903004**

**Yuvraj Gupta**

**101903549**

**Submitted to  
Mr. Sanjeev Rao**



**THAPAR INSTITUTE**  
OF ENGINEERING & TECHNOLOGY  
(Deemed to be University)

**DEPARTMENT OF COMPUTER SCIENCE AND  
ENGINEERING**  
**THAPAR INSTITUTE OF ENGINEERING AND TECHNOLOGY,**  
**(A DEEMED TO BE UNIVERSITY),**  
**PATIALA, PUNJAB**  
**INDIA**  
**Nov 2022**

# **Traffic Flow Prediction**

## **Objective**

Traffic flow prediction is an essential part of the intelligent transport system. This is the accurate estimation of traffic flow in a given region at a particular interval of time in the future. The study of traffic forecasting is useful in mitigating congestion and make safer and cost-efficient travel.

While traditional models use shallow networks, there has been an exponential growth in the number of vehicles in recent times and these traditional machine learning models fail to work in current scenarios. In our paper, we review some of the latest works in deep learning for traffic flow prediction.

Many deep learning architectures include Convolutional Neural Network (CNN), Recurrent Neural Network (RNN), Long Short-Term Memory (LSTM), Restricted Boltzmann Machines (RBM), and Stacked Auto Encoder (SAE). These deep learning models use multiple layers to extract higher level of features from raw input progressively. The latest deep learning models developed to tackle this very problem are reviewed and due to the complexity of transport networks, this review gives the reader information about how various factors influence these models and what models work best in different scenarios.

Machine Learning technology has developed so much that we can now help the community by detecting breast cancer on the basis of various features such as radius, texture, compactness etc. We have used the K-nearest neighbor (KNN) algorithm to develop a data model which predicts whether the tumor is benign or malignant.

The Traffic flow dataset by Coplin has been used for the training and testing of the data model created by us.

## Dataset

The dataset used in this project is the Coplin Traffic flow Dataset. The dataset has 30 parameters defined on the basis of which we determine whether the result is malignant or benign.

The different parameters are states as follows:

radius\_mean  
texture\_mean  
perimeter\_mean  
area\_mean  
smoothness\_mean  
compactness\_mean  
concavity\_mean  
concave points\_mean  
symmetry\_mean  
fractal\_dimension\_mean  
radius\_se  
texture\_se  
perimeter\_se  
area\_se  
smoothness\_se  
compactness\_se  
concavity\_se  
concave points\_se  
symmetry\_se  
fractal\_dimension\_se  
radius\_worst  
texture\_worst  
perimeter\_worst  
area\_worst  
smoothness\_worst  
compactness\_worst  
concavity\_worst  
concave points\_worst  
symmetry\_worst  
fractal\_dimension\_worst

## Methodology

The Data Model in this project is made using the K-nearest neighbor algorithm. Among the supervised machine learning algorithms, K-nearest neighbors (KNN) is one of the most effective techniques. It performs classification on certain data points. The KNN algorithm is a type of supervised ML algorithm that can be used for both classifications as well as regression predictive problems. It uses 'attribute similarity' to predict the values of new data-points and then the new data point will be assigned a value based on how closely it matches the points in the training set.

### Algorithm

**Step-1: Select the number K of the neighbors**

**Step-2: Calculate the Euclidean distance of K number of neighbors**

**Step-3: Take the K nearest neighbors as per the calculated Euclidean distance.**

**Step-4: Among these k neighbors, count the number of the data points in each category.**

**Step-5: Assign the new data points to that category for which the number of the neighbor is maximum.**

**Step-6: Our model is ready.**

The K value for our dataset is 23, because the optimal K values is the square root of the number of entries in the dataset, i.e., 569, which comes out to be 23.85372. Rounding it down to 23 is better because we want an odd K value as an odd number so that we can calculate a clear majority in the case where only two groups are possible.

## Code

```
traffice_flow_prediction.py •
C: > Users > pnb > Downloads > traffice_flow_prediction.py
1  # -*- coding: utf-8 -*-
2  """traffice_flow_prediction.ipynb
3
4  Automatically generated by Colaboratory.
5
6  Original file is located at
7  | https://colab.research.google.com/drive/1pKuw7QWql39bdV88hWsu29zExOMBeI72
8  | """
9
10 !unzip archive.zip
11 |
12
13 import pandas as pd
14 import glob
15
16 path = '/content/DataSet'
17 all_files = glob.glob(path + "/*.csv")
18
19 li = []
20
21 for filename in all_files:
22     df = pd.read_csv(filename, index_col=None, header=0)
23     li.append(df)
24
25 df = pd.concat(li, axis=0, ignore_index=True)
26 print(df.head(3))
27
28 """### Descriptive Analysis"""
29
30 df.describe()
31
```

```


30 df.describe()
31
32 """### Getting new features"""
33
34 df['timestamp'] = df['timestamp'].astype('datetime64[ns]')
35 df['weekday'] = df['timestamp'].dt.weekday
36
37 # Feature engineering with the date
38 df['year'] = df['timestamp'].dt.year
39 df['month'] = df['timestamp'].dt.month
40 df['day'] = df['timestamp'].dt.day
41
42 df.head(3)
43
44 original_df = df.copy()
45
46 import matplotlib.pyplot as plt
47 plt.style.use('ggplot')
48
49 plt.figure(figsize=(50, 8))
50 mean_group = df[['timestamp', 'hourly_traffic_count']].groupby(['timestamp'])['hourly_traffic_count'].mean()
51 plt.plot(mean_group)
52 plt.title('Time Series - Average')
53 plt.show()
54
55 plt.figure(figsize=(50, 8))
56 median_group = df[['timestamp', 'hourly_traffic_count']].groupby(['timestamp'])['hourly_traffic_count'].median()
57 plt.plot(median_group, color = 'b')
58 plt.title('Time Series - median')
59 plt.show()
60
61 df['weekday_num'] = df['weekday']
62 df['weekday'].replace(0, '01 - Monday', inplace=True)
63 df['weekday'].replace(1, '02 - Tuesday', inplace=True)
64 df['weekday'].replace(2, '03 - Wednesday', inplace=True)
65 df['weekday'].replace(3, '04 - Thursday', inplace=True)
66 df['weekday'].replace(4, '05 - Friday', inplace=True)
67 df['weekday'].replace(5, '06 - Saturday', inplace=True)
68 df['weekday'].replace(6, '07 - Sunday', inplace=True)
69
70 train_group = df.groupby(["month", "weekday"])['hourly_traffic_count'].mean().reset_index()
71 train_group = train_group.pivot('weekday', 'month', 'hourly_traffic_count')
72 train_group.sort_index(inplace=True)
73
74 import seaborn as sns
75 sns.set(font_scale=1.2)
76
77 f, ax = plt.subplots(figsize=(8, 8))
78 sns.heatmap(train_group, annot=False, ax=ax, fmt="d", linewidths=2)
79 plt.title('Web Traffic Months cross Weekdays')
80 plt.show()
81
82 """##### We can clearly see that there is no activity during the December and people during October contributed the highest activity """
83
84 times_series_means = pd.DataFrame(mean_group).reset_index(drop=False)
85 times_series_means['weekday'] = times_series_means['timestamp'].apply(lambda x: x.weekday())
86 times_series_means['Date_str'] = times_series_means['timestamp'].apply(lambda x: str(x))
87 times_series_means[['year', 'month', 'day']] = pd.DataFrame(times_series_means['Date_str'].str.split('-', 2).tolist(), columns = ['year', 'month', 'day'])
88 date_staging = pd.DataFrame(times_series_means['day'].str.split(' ', 2).tolist(), columns = ['day', 'other'])
89 times_series_means['day'] = date_staging['day']*1

```

```

88 date_staging = pd.DataFrame(times_series_means['day'].str.split(' ',2).tolist(), columns = ['day','other'])
89 times_series_means['day'] = date_staging['day']*1
90 times_series_means.drop('Date_str',axis = 1, inplace =True)
91 del times_series_means['timestamp']
92 times_series_means.head()
93
94 """### Train/Test Preparation"""
95
96 from sklearn.model_selection import train_test_split
97
98 X, y = times_series_means.drop(['hourly_traffic_count','year'],axis=1), times_series_means['hourly_traffic_count']
99 trainx, testx, trainy, testy = train_test_split(X, y, test_size=0.2)
100
101
102 # Linear Model
103 from sklearn.ensemble import GradientBoostingRegressor, AdaBoostRegressor
104 from sklearn.metrics import mean_absolute_error, r2_score
105
106 def modelisation(x_tr, y_tr, x_ts, y_ts, model):
107     # Modelisation with all product
108     model.fit(x_tr, y_tr)
109
110     prediction = model.predict(x_ts)
111     r2 = r2_score(y_ts.to_numpy(), model.predict(x_ts))
112     mae = mean_absolute_error(y_ts.to_numpy(), model.predict(x_ts))
113     print ("-----")
114     print ("mae with 80% of the data to train:", mae)
115     print ("-----")
116
117     return prediction, model

```

C:\> Users > pnb > Downloads >  traffic\_flow\_prediction.py

```

117     return prediction, model
118
119 model = AdaBoostRegressor(n_estimators = 5000, random_state = 42, learning_rate=0.01)
120
121 prediction, clr = modelisation(trainx, trainy, testx, testy, model)
122
123 import numpy as np
124 plt.figure(figsize=(16, 4))
125 line_up, = plt.plot(prediction,label='Prediction', color="red")
126 plt.ylabel('Series')
127 plt.legend(handles=[line_up])
128 plt.title('Performance of predictions - Benchmark Predictions')
129 plt.show()
130
131 plt.figure(figsize=(16, 6))
132 line_down, = plt.plot(np.array(testy),label='Reality')
133 plt.ylabel('Series')
134 plt.legend(handles=[line_down])
135 plt.title('Performance of predictions - Benchmark Reality')
136 plt.show()
137
138 trainx.shape
139
140 """### Keras LSTM"""
141
142 trainx.head()
143
144 from keras.models import Sequential
145 from keras.layers import Dense, LSTM
146 from pandas import DataFrame, concat
147 from sklearn.preprocessing import MinMaxScaler

```



C: > Users > pnb > Downloads > traffic\_flow\_prediction.py

```
117     return prediction, model
118
119     model = AdaBoostRegressor(n_estimators = 5000, random_state = 42, learning_rate=0.01)
120
121     prediction, clr = modelisation(trainx, trainy, testx, testy, model)
122
123     import numpy as np
124     plt.figure(figsize=(16, 4))
125     line_up, = plt.plot(prediction, label='Prediction', color="red")
126     plt.ylabel('Series')
127     plt.legend(handles=[line_up])
128     plt.title('Performance of predictions - Benchmark Predictions')
129     plt.show()
130
131     plt.figure(figsize=(16, 6))
132     line_down, = plt.plot(np.array(testy), label='Reality')
133     plt.ylabel('Series')
134     plt.legend(handles=[line_down])
135     plt.title('Performance of predictions - Benchmark Reality')
136     plt.show()
137
138     trainx.shape
139
140     """### Keras LSTM"""
141
142     trainx.head()
143
144     from keras.models import Sequential
145     from keras.layers import Dense, LSTM
146     from pandas import DataFrame, concat
147     from sklearn.preprocessing import MinMaxScaler
148
149     from keras.models import Sequential
150     from keras.layers import Dense, LSTM
151     from pandas import DataFrame, concat
152     from sklearn.preprocessing import MinMaxScaler
153
154     X, y = times_series_means.drop(['hourly_traffic_count', 'year'], axis=1), times_series_means['hourly_traffic_count']
155     scaler = MinMaxScaler(feature_range=(0, 1))
156     scaled = scaler.fit_transform(X)
157
158     trainx, testx, trainy, testy = train_test_split(scaled, y, test_size=0.2)
159
160     trainx = trainx.reshape((trainx.shape[0], 1, trainx.shape[1]))
161     testx = testx.reshape((testx.shape[0], 1, testx.shape[1]))
162
163     # design network
164     model = Sequential()
165     model.add(LSTM(50, input_shape=(trainx.shape[1], trainx.shape[2])))
166     model.add(Dense(1))
167     model.compile(loss='mae', optimizer='adam')
168     # fit network
169     history = model.fit(trainx, trainy, epochs=50, batch_size=8, validation_data=(testx, testy), verbose=2, shuffle=False)
170
171     """### Plot history and Evaluation"""
172
173     from sklearn.metrics import mean_squared_error
174
175     plt.plot(history.history['loss'], label='train')
176     plt.plot(history.history['val_loss'], label='test')
```



```
171
172 plt.plot(history.history['loss'], label='train')
173 plt.plot(history.history['val_loss'], label='test')
174 plt.legend()
175 plt.show()
176
177 # make a prediction
178 yhat = model.predict(testx)
179 # calculate RMSE
180 rmse = np.sqrt(mean_squared_error(testy, yhat))
181 print('Test RMSE: %.3f' % rmse)
```

## Output

```
PROBLEMS    OUTPUT    DEBUG CONSOLE    TERMINAL

Accuracy : 91.23711340206185%
PS C:\Users\A\Desktop\ML projects> & C:/Users/A/Desktop/ML projects/output.py
Accuracy : 94.53551912568307%
PS C:\Users\A\Desktop\ML projects> & C:/Users/A/Desktop/ML projects/output.py
Accuracy : 93.92265193370166%
PS C:\Users\A\Desktop\ML projects>
```

## Conclusion

This project helped us understand the working of data models using KNN algorithm. Machine learning approaches have been increasing rapidly in the practical life due to their monumental performance in predicting and classifying problems.

The average accuracy of the data model is 93.23% with the maximum accuracy reached 94.5%.

## References

The dataset has been taken from the following link:

<https://www.kaggle.com/datasets/coplin/traffic>

The data Model has been made with the help of the following videos:

<https://www.youtube.com/watch?v=4HKqjENq9OU>

<https://www.youtube.com/watch?v=2btpg2xmq80>

