

# Statistical Machine Learning Project

Vencel David Koczka, Aiman Baig, Giorgos Tsouderos, Tomas Grahn

March 16, 2024

## Abstract

The goal of this project was to find a model to predict increased bike demand. Feature engineering was used to highlight features that showed high proportions of increased demand and remove features which showed low information about high demand times. From these features four models were trained and optimised with respect to their F-Beta score ( $\beta=1.5$ ). These models included logistic regression, quadratic discriminant analysis, K-nearest neighbors and random forest. Of these models quadratic discriminant analysis proved the best with a F-Beta score ( $\beta=1.5$ ) of 0.5767 on a hold-out data set.

## 1 Introduction

Capital Bikeshare is a 24-hour public bicycle-sharing. The problem arises is that there are certain occasions when, due to various circumstances, there are not as many bikes available as there are demands. The District Department of Transportation in the city wants to know if at certain hours an increase in the number of bikes available will be necessary.

The goal of the project is to predict whether an increase in the number of bikes is necessary or not based on various time and weather data. To accomplish this, first the data was explored. Then feature engineering applied to best expose times of increased demand. Then four models were explored to determine which is best able to predict an increase in demand given a set test of data. Finally the trained models were tested on a set of held out test data to determine which model is best suited for the problem.

## 2 Data Analysis

### 2.1 Numerical and categorical

Reviewing the data in each column revealed that the following categories were categorical: Hour of day, Day of week, Month, Holiday, Weekday, Summertime, Snow, Increase stock. These were numerical: Temperature, Dew, Humidity, Precipitation, Snow depth, Wind speed, Cloud cover, Visibility.

## 2.2 Trends in time

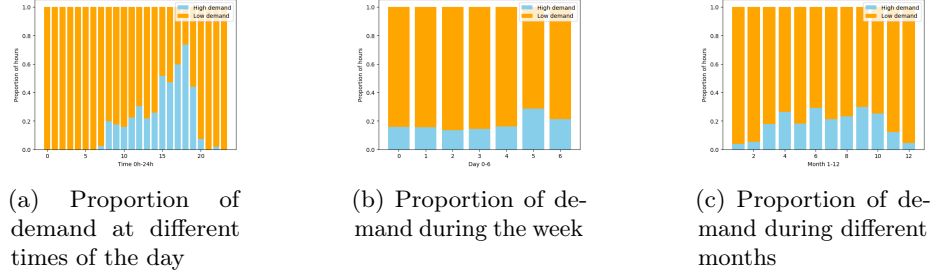


Figure 1: Trends in time,days of the week and months

Figure 1a shows 2 separate areas of high demand. A lower one from 8am to 2pm and a larger one from 3pm to 7pm. Figure 1b shows a slight peak in high demand on Saturday and Sunday, with all other days being fairly equal. Figure 1c see that April, June, September have the highest demand and January, February and December have the least.

### 2.2.1 Weekdays & Holidays

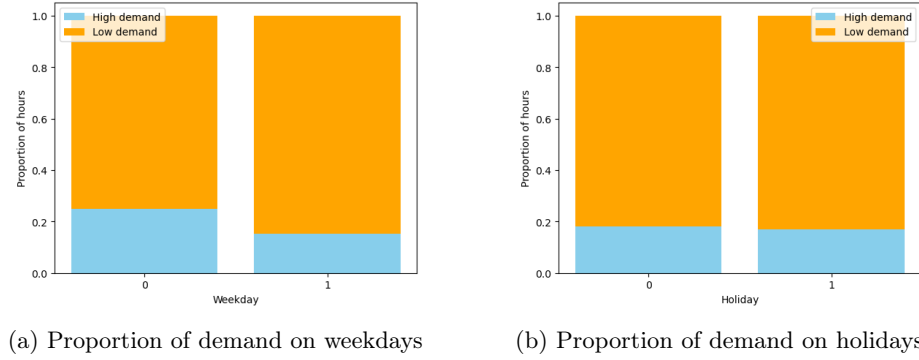


Figure 2: Trends in weekdays and holidays

Figure 2a shows that there is slightly more demand on weekends. Figure 2b shows that there is no significant difference between holiday and non-holiday on demand.

## 2.3 Trends for weather

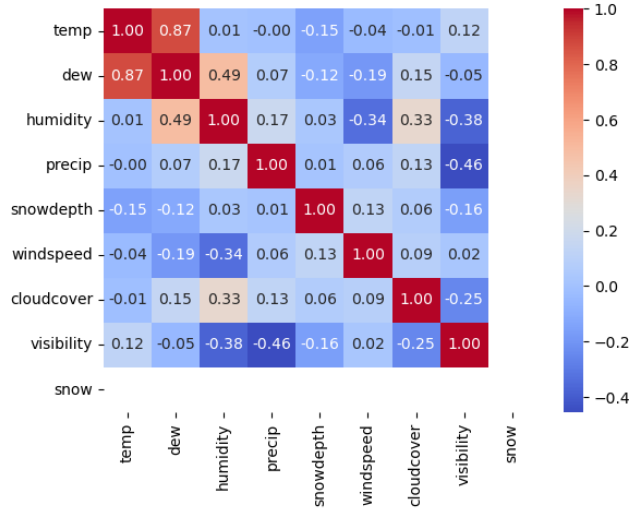


Figure 3: Correlation between features

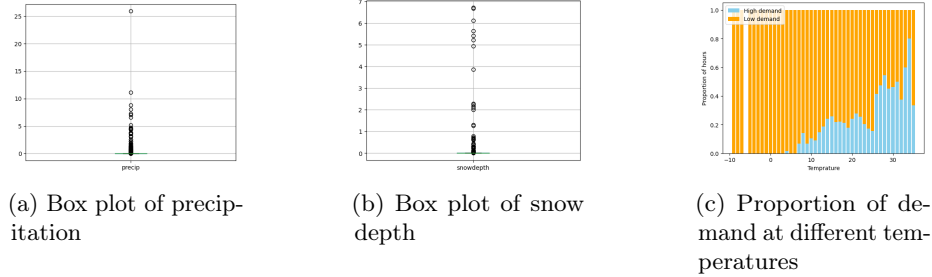


Figure 4: Trends in precipitation, snow depth and temperature

Figure 3 shows us that Snow has no correlation to anything, This is because it is only 0 or NaN in the dataset. It also shows us that dew is highly correlated to temperature. Figure 4a shows that a majority of the values are 0. Furthermore looking at the values that are not 0 shows they are all low demand hours. Figure 4b shows that a majority of the values are 0. Furthermore looking at the values that are not 0 shows they are all low demand hours. Figure 4c there is a spike in demand past 25 degrees. There is also increased demand past 12 degrees with central spikes at 16 and 21.

## 2.4 Basic analysis and naive classifier

By looking at the proportion of high and low demand samples a naive classifier can be constructed. The data is 82% low demand and 18% high demand. Hence a classifier that always classifies as low demand would have an accuracy of 82%. This classifier however has 0 F-Beta score ( $\beta=1.5$ ) which is our key metric for finding an appropriate classifier.

## 2.5 Feature engineering

Feature engineering was used to enhance as much information relating to high demand as possible. Snow was removed as it was only 0 or Nan. Holiday was removed as there was no greater proportion of high demand whether there was a holiday or not. Snow depth and precipitation were removed as they are almost always 0 and when it they are non zero they are low demand. Dew was removed as it was highly correlated with temperature and it so did not add extra information regarding the high demand times. Finally the features were normalised so that k-NN could use.

Due to the high demand between 3pm and 7pm an additional feature was constructed called rush hour to prioritise this information. Finally another feature was added for good temperature to reflect the high demand when the temperature was greater than 25 degrees.

# 3 Implementation of Methods

Based on the direction given by Capital Bikeshare to look for hours when an increased bike demand would occur, it was decided that the key metric for optimisation of models would be a F-beta score ( $\beta=1.5$ ) to give more weight to recall. Recall was seen as more important as it was assumed missing high demand days was the biggest loss for the company. Accuracy, precision and recall are also reported.

## 3.1 Logistic Regression

Logistic Regression is a binary classification algorithm that models the probability of an instance belonging to a certain class. The Sigmoid function is used in logistic regression. The sigmoid function is defined as:  $\sigma(z) = \frac{1}{1+e^{-z}}$  where  $\sigma(z)$  is the output between 0 and 1, and  $z$  is linear combination of the input.

The Logistic Regression model is trained by adjusting the coefficients to minimize the difference between the predicted probabilities and the actual class labels in the training data. The model parameters are adjusted iteratively to minimize the overall loss across all training instances. Grid search with 5 folds was used to search over the possible hyper-parameters. The search iterated over different parameters for regularization amount, regularization method and the maximum iterations. The search was optimised for F-Beta ( $\beta=1.5$ ) and the following score was achieved.

Table 1: Logistic Regression Performance

Accuracy	Precision	Recall	$F_\beta$
0.8469	0.5128	0.4000	0.4290

The method performed quite well with respect to accuracy but was lacking in the ability to separate the positive samples. This was likely due to the model’s linear nature which limits the complexity of the boundaries it can find.

### 3.2 Discriminant Analysis: QDA

Discriminant analysis is a method of classification in which density functions are constructed based on the training data and are used to classify new points. The density functions are assumed to be normally distributed. The maximum likelihood distribution is found during training by optimising the following

$$f(x; \theta) = \arg \max_{\theta} \sum_{i=1}^n \ln(\mathbb{P}(x_i|y_i, \theta)) + \ln(\mathbb{P}(y_i|\theta))$$

If you assume that the functions will share the same covariance matrix then you will be performing linear discriminant analysis. If you allow them to have differing covariance matrices then you arrive at quadratic discriminant analysis. If your data is not linearly separable it is favorable to use quadratic discriminant analysis as it can create non linear boundaries. In the case of our problem quadratic discriminate analysis was seen as favorable as there are some features such as temperature that could yield additional information if a more complex boundary is used.

Sklearn offers the ability to regularise a QDA model with shrinkage. This can assist with data that has a high number of features compared to the amount of data as it helps to compensate for the fact the empirical covariance matrix is a poor estimate for the true covariance matrix. Table 2 shows that shrinkage only made the model perform worse with respect to F Beta ( $\beta=1.5$ ) so for the final QDA no shrinkage was used.

Table 2: QDA Performance

Version	Accuracy	Precision	Recall	$F_\beta$
QDA	0.8203	0.5443	0.8113	0.7049
QDA with shrinkage	0.8320	0.7083	0.3208	0.3857

The QDA method performed well particularly with respect to recall which also gave it a high F-beta score. QDA seems to perform well due to it’s ability to fit non-linear boundaries between the data. This may be due to how the model can better reflect the relationship between high demand and features such as

temperature and hour of the day. These features don't appear linear in when there is or isn't demand and the QDA can likely model the relationship better.

### 3.3 K-nearest neighbor

The k-nearest neighbors algorithm is a supervised machine learning algorithm used for classification and regression problems. A data point is a vector in a multidimensional feature space. When classifying a new data point, its assigned a class label is based on the classes of its k closest neighbors. There are several ways to calculate the distance between each data point (Euclidean, Manhattan, Hamming) but in this case the Minkowski distance was used.

$$MinkowskiDistance = \sum_{i=1}^n (|x_i - y_i|)^{\frac{1}{p}}$$

This distance metric is a generalised form of the Euclidean and Manhattan distances, where the change of the p parameter enables the creation of other distance metrics. The k parameter allows the selection of the amount of neighbors checked when deciding to assign a label to a new data point. To find the optimal k parameter, the elbow method was used where the k-NN algorithm is executed for a range of k values and the k that minimizes the loss of the optimizing metric we chose, is the that is selected. In our case, to find the optimal number of neighbors we averaged the F-Beta error rate over 5 random splits which resulted in k=2. When a small value of k is used it usually indicates overfitting, however since KNN is a non parametric algorithm, it can be difficult to avoid these cases. The results after running the KNN algorithm for k=2 are highlighted in the table below.

Table 3: KNN Performance

Accuracy	Precision	Recall	$F_{\beta}$
0.7406	0.3225	0.6000	0.4745

### 3.4 Tree-based methods

We experimented with three different kind of Tree-based methods, namely Decision Trees, Bagging, and Random Forest. Our best results were achieved with Random Forest, out of these methods. Random forest trains multiple decision trees and averages the results which leads to better performance.

The training process for decision trees splits the feature space into regions. Step by step we always choose the best available feature and value to split on, and this way the tree becomes deeper and deeper. By setting parameters for the trees we can restrict how deep a tree can get, and also manipulate the splitting logic itself, by for example setting the minimum samples a leaf should contain. When predicting new datapoints we use  $\hat{y}_l = MajorityVote\{y_i : \mathbf{x}_i \in R_l\}$ , Which

means that each final region (leaf) has a class determined by the dominant data points contained in that region. The new point will be classified based on that.

For finding the best hyperparameters we used RandomGridSearch method provided by sklearn, as it is faster than simple GridSearch. This method tries out different combination of hyperparameters, and compares the results with cross-validation. For random forest classifiers we found that we should set the number of estimators to 334, the maximum depth of the trees to 50, and minimum samples in each leaf to 1, minimum samples in each split to 5, with Bootstrap set to False, and max\_features parameter set to 'sqrt'.

When using decision trees, we can easily overfit on our data if the maximum tree depth is too large. In the case of random forests this is not a bad behaviour, because we average the results out on numerous trees, consequently regularizing the model. For this method our feature engineering does not affect the final performance. This is because the splitting happens on the same features. By training a simple decision tree, we can see that these parameters are hour\_of\_day, and temp.

Table 4: Random Forest Performance

Accuracy	Precision	Recall	$F_\beta$
0.8531	0.5294	0.5400	0.5366

## 4 Conclusion

Table 5: Results of evaluation

Model	Accuracy	Precision	Recall	$F_\beta$
Logistic Regression	0.8469	0.5128	0.4000	0.4290
QDA	0.7750	0.3854	0.7400	0.5767
K-nn	0.7406	0.3225	0.6000	0.4745
Random Forest	0.8531	0.5294	0.5400	0.5366

Figure 5 shows performance of the tuned methods with respect to the held out test set. In terms of our key metric F-beta ( $\beta=1.5$ ) the QDA method performed the best and will be chosen as the model for production. Random forest had a better performance in accuracy and a strong performance in precision. Should the users require a model with higher precision and accuracy it is recommended to switch to the random forest model. The F-beta score could be further improved by exploring the use of neural networks, more data exploration and tailored feature engineering for the methods.

## A Appendix

```
1 # -*- coding: utf-8 -*-
2 """SML_Project_Final_corrected.ipynb
3
4 Automatically generated by Colaboratory.
5
6 Original file is located at
7     https://colab.research.google.com/drive/1
8     QTEIux5jWMkKk9PVBq7J5F8Bs9N638if
9
10 # Data Analysis
11 i) Which are the numerical and which are the categorical?
12 ii) Is there a greater trend to need an increase in the
13     availability of bicycles? Study this question
14     from various perspectives:
15
16     Can any trend be seen comparing different hours, weeks, and
17     months?
18
19     Is there any difference between weekdays and holidays?
20
21     Is there any trend depending on the weather? Rainy days, snowy
22     days, etc.
23
24 Write concise answers to each question and support your findings
25 with evidence (statistics, plots,
26 etc.). Discuss the results. Additionally, you can explore the
27 correlation of features, outliers, range of
28 values, and many more aspects.
29 """
30
31 import pandas as pd
32 import numpy as np
33 import matplotlib.pyplot as plt
34 data = pd.read_csv('data/training_data.csv')
35 data
36
37 data.info()
38
39 categorical = data[['hour_of_day', 'day_of_week', 'month', 'holiday', '
40     weekday', 'summertime', 'snow', 'increase_stock']]
41 # Snow is a very bad variable it's only either 0 or nan
42 numerical = data[['temp', 'dew', 'humidity', 'precip', 'snowdepth', '
43     windspeed', 'cloudcover', 'visibility']]
44
45 # Naive classifier and proportion of data
46 total = len(data)
47 increase = len(data[data["increase_stock"]=='high_bike_demand'])
48 decrease = len(data[data["increase_stock"]=='low_bike_demand'])
49 proportion_increase = increase/total #18%
50 proportion_decrease = decrease/total #82%
51 # Hence a naive classifier that always predicts decrease achieve
52     82% missclassification
53
54 # Different hours
```



```

47 hours = sorted(data['hour_of_day'].unique())
48 hours_high_demand = []
49 hours_low_demand = []
50 for hour in hours:
51     hours_high_demand.append(len(data.query("hour_of_day == @hour and
52         increase_stock == 'high_bike_demand'")))
53     hours_low_demand.append(len(data.query("hour_of_day == @hour and
54         increase_stock == 'low_bike_demand'")))
55 plt.figure()
56 plt.bar(hours, hours_high_demand)
57 plt.xlabel('Time 0h-24h')
58 plt.ylabel('Number of hours with high demand')
59 plt.figure()
60 plt.bar(hours, hours_low_demand)
61 plt.xlabel('Time 0h-24h')
62 plt.ylabel('Number of hours with low demand')
63 plt.figure()
64 plt.bar(hours, hours_high_demand, color='skyblue', label='High demand')
65 plt.bar(hours, hours_low_demand, bottom=hours_high_demand, color='
66     orange', label='Low demand')
67 plt.xlabel('Time 0h-24h')
68 plt.ylabel('Hours')
69 plt.legend()
70 plt.figure()
71 total_hours = np.add(hours_high_demand, hours_low_demand)
72 plt.bar(hours, hours_high_demand/total_hours, color='skyblue', label='
73     High demand')
74 plt.bar(hours, hours_low_demand/total_hours, bottom=hours_high_demand
75     /total_hours, color='orange', label='Low demand')
76 plt.xlabel('Time 0h-24h')
77 plt.ylabel('Proportion of hours')
78 plt.legend()
79
80 # Different weeks
81 dws = sorted(data['day_of_week'].unique())
82 dw_high_demand = []
83 dw_low_demand = []
84 for dw in dws:
85     dw_high_demand.append(len(data.query("day_of_week == @dw and
86         increase_stock == 'high_bike_demand'")))
87     dw_low_demand.append(len(data.query("day_of_week == @dw and
88         increase_stock == 'low_bike_demand'")))
89
90 plt.figure()
91 plt.bar(dws, dw_high_demand)
92 plt.xlabel('Day 0-6')
93 plt.ylabel('Hours in high demand')
94 plt.figure()
95 plt.bar(dws, dw_low_demand)
96 plt.xlabel('Day 0-6')
97 plt.ylabel('Hours in low demand')
98 plt.figure()
99 plt.bar(dws, dw_high_demand, color='skyblue', label='High demand')

```

```

196 plt.bar(dws,dw_low_demand,bottom=dw_high_demand,color='orange',
197         label='Low demand')
198 plt.xlabel('Day 0-6')
199 plt.ylabel('Hours')
200 plt.legend()
201 plt.figure()
202 total_hours = np.add(dw_high_demand,dw_low_demand)
203 plt.bar(dws,dw_high_demand/total_hours,color='skyblue',label='High
204         demand')
205 plt.bar(dws,dw_low_demand/total_hours,bottom=dw_high_demand/
206         total_hours,color='orange',label='Low demand')
207 plt.xlabel('Day 0-6')
208 plt.ylabel('Proportion of hours')
209 plt.legend()
210
211 # Different months
212 months = sorted(data['month'].unique())
213 month_high_demand = []
214 month_low_demand = []
215 for month in months:
216     month_high_demand.append(len(data.query("month == @month and
217         increase_stock == 'high_bike_demand'")))
218     month_low_demand.append(len(data.query("month == @month and
219         increase_stock == 'low_bike_demand'")))
220
221 plt.figure()
222 plt.bar(months,month_high_demand)
223 plt.xlabel('Month 1-12')
224 plt.ylabel('Hours with high demand')
225 plt.figure()
226 plt.bar(months,month_low_demand)
227 plt.xlabel('Month 1-12')
228 plt.ylabel('Hours with low demand')
229 plt.figure()
230 plt.bar(months,month_high_demand,color='skyblue',label='High demand
231         ')
232 plt.bar(months,month_low_demand,bottom=month_high_demand,color='
233         orange',label='Low demand')
234 plt.xlabel('Month 1-12')
235 plt.ylabel('Hours')
236 plt.legend()
237 plt.figure()
238 total_hours = np.add(month_high_demand,month_low_demand)
239 plt.bar(months,month_high_demand/total_hours,color='skyblue',label=
240         'High demand')
241 plt.bar(months,month_low_demand/total_hours,bottom=
242         month_high_demand/total_hours,color='orange',label='Low demand'
243         )
244 plt.xlabel('Month 1-12')
245 plt.ylabel('Proportion of hours')
246 plt.legend()
247
248 """In the next cell we test for the diferent correlations.
249
250 We remove three columns [dew, snow, weekday] because they are
251 highly correlated with other columns.

```

```

142 We keep all the others for now, we may remove ones later.
143 """
144
145 import seaborn as sns
146 # Weekdays and holidays
147 # Weekdays correlated to days as weekdays are just 0-5, weekends
   are 0-2
148 correlation_matrix = data.corr()
149 # Plotting the heatmap
150 sns.heatmap(correlation_matrix, annot=True, cmap="coolwarm", fmt="
   .2f")
151 plt.show()
152
153 # Do some bar plots of numerical values
154 numerical.boxplot()
155
156 # Temperature
157 data['itemp'] = data['temp'].astype(int)
158 temps = data['itemp']
159 temps_high_demand = []
160 temps_low_demand = []
161 for temp in temps:
162     temps_high_demand.append(len(data.query("itemp == @temp and
   increase_stock == 'high_bike_demand'")))
163     temps_low_demand.append(len(data.query("itemp == @temp and
   increase_stock == 'low_bike_demand'")))
164 plt.figure()
165 total_hours = np.add(temps_high_demand, temps_low_demand)
166 plt.bar(temps, temps_high_demand/total_hours, color='skyblue', label='
   High demand')
167 plt.bar(temps, temps_low_demand/total_hours, bottom=temps_high_demand
   /total_hours, color='orange', label='Low demand')
168 plt.xlabel('Temperature')
169 plt.ylabel('Proportion of hours')
170 plt.legend()
171
172 # Weekdays
173 weekdays = data['weekday']
174 weekday_high_demand = []
175 weekday_low_demand = []
176 for weekday in weekdays:
177     weekday_high_demand.append(len(data.query("weekday == @weekday
   and increase_stock == 'high_bike_demand'")))
178     weekday_low_demand.append(len(data.query("weekday == @weekday and
   increase_stock == 'low_bike_demand'")))
179 plt.figure()
180 total_hours = np.add(weekday_high_demand, weekday_low_demand)
181 plt.bar weekdays, weekday_high_demand/total_hours, color='skyblue',
   label='High demand')
182 plt.bar weekdays, weekday_low_demand/total_hours, bottom=
   weekday_high_demand/total_hours, color='orange', label='Low
   demand')
183 plt.xlabel('Weekday')
184 plt.ylabel('Proportion of hours')
185 plt.xticks([0,1])
186 plt.legend()
187

```

```

188 #Holiday
189 holidays = data['holiday']
190 holiday_high_demand = []
191 holiday_low_demand = []
192 for holiday in holidays:
193     holiday_high_demand.append(len(data.query("holiday == @holiday
194         and increase_stock == 'high_bike_demand'")))
195     holiday_low_demand.append(len(data.query("holiday == @holiday and
196         increase_stock == 'low_bike_demand'")))
197 plt.figure()
198 total_hours = np.add(holiday_high_demand, holiday_low_demand)
199 plt.bar(holidays, holiday_high_demand/total_hours, color='skyblue',
200     label='High demand')
201 plt.bar(holidays, holiday_low_demand/total_hours, bottom=
202     holiday_high_demand/total_hours, color='orange', label='Low
203     demand')
204 plt.xlabel('Holiday')
205 plt.ylabel('Proportion of hours')
206 plt.xticks([0,1])
207 plt.legend()
208
209 """# Importing packages"""
210
211 import numpy as np
212 import pandas as pd
213
214 import sklearn.preprocessing as skl_pre
215 import sklearn.linear_model as skl_lm
216 import sklearn.discriminant_analysis as skl_da
217 import sklearn.neighbors as skl_nb
218
219 from sklearn import tree
220 from sklearn.ensemble import BaggingClassifier,
221     RandomForestClassifier
222
223
224 import matplotlib.pyplot as plt
225
226 from sklearn.model_selection import train_test_split
227
228 from sklearn.linear_model import LogisticRegression
229 from sklearn.model_selection import GridSearchCV
230 from sklearn.datasets import make_classification
231 from sklearn.metrics import fbeta_score, make_scorer, recall_score,
232     precision_score
233 from sklearn.preprocessing import normalize
234
235 """# Train - valid - test split"""
236
237 #just importing the feature manipulation here
238 data = pd.read_csv('data/training_data.csv')
239 filtered = data.drop(['precip', 'holiday', 'snowdepth', 'snow', '
240     dew'], axis=1)
241 data_set = filtered
242 data_set['rushhour'] = filtered['hour_of_day'].apply(lambda x: 1 if
243     x > 2 and x < 8 else 0)

```

```

235 data_set['goodtemp'] = filtered['temp'].apply(lambda x: 1 if x>25
        else 0)
236
237 X = normalize(data_set.drop('increase_stock',axis=1))
238 y = data_set['increase_stock']
239
240 from sklearn.model_selection import train_test_split
241 X_in, X_test, y_in, y_test = train_test_split(X, y, test_size=0.20,
        random_state=42)
242
243 def print_metrics(model,X_test,y_test):
244     y_preds = model.predict(X_test)
245     y_test = np.array(y_test).reshape(-1)
246     #Confusion Matrix
247     cross_vals=pd.crosstab(y_preds, y_test)
248     print(cross_vals)
249     missclassification_rate = np.mean([y_preds != y_test])
250     print(f"Missclassification {missclassification_rate}")
251     accuracy = 1-missclassification_rate
252     print(f"Accuracy {accuracy}")
253     recall = recall_score(y_test, y_preds,pos_label='high_bike_demand'
        ')
254     print(f"Recall {recall}")
255     precision = precision_score(y_test, y_preds,pos_label='
        high_bike_demand')
256     print(f"Precision {precision}")
257     beta = 1.5
258     F_beta = ((1+beta**2) *precision * recall) / (beta**2 * precision
        + recall)
259     print(f"F_beta {F_beta}")
260
261 """# Logistic Regression"""
262
263 model = LogisticRegression(solver='liblinear')
264 param_grid = {
265     'C': [0.001, 0.01, 0.1, 1, 10, 100, 1000],
266     'penalty': [None,'l1','l2'],
267     'max_iter': [10, 50 ,100, 500, 1000,5000],
268     'solver': ['liblinear']
269 }
270 fb_score = make_scorer(fbeta_score, beta=1.5, pos_label='
        high_bike_demand')
271 grid_search = GridSearchCV(model, param_grid, cv=5, scoring=
        fb_score)
272 grid_search.fit(X_in, y_in)
273 best_model = grid_search.best_estimator_
274 #print_metrics(best_model, X_test, y_test)
275
276 grid_search.best_params_
277
278 model = LogisticRegression(solver='liblinear',C=1000,max_iter=50,
        penalty='l1')
279 model.fit(X_in, y_in)
280 print_metrics(model, X_test, y_test)
281
282 """# Discriminate Analysis"""
283

```

```

284 X_train, X_valid, y_train, y_valid = train_test_split(X_in, y_in,
    test_size=0.20, random_state=1337)
285
286 model1 = skl_da.QuadraticDiscriminantAnalysis(reg_param=0.005)
287 model2 = skl_da.QuadraticDiscriminantAnalysis(reg_param=0)
288 model1.fit(X_train, y_train)
289 model2.fit(X_train, y_train)
290 print_metrics(model1, X_valid, y_valid)
291 print_metrics(model2, X_valid, y_valid)
292
293 model = skl_da.QuadraticDiscriminantAnalysis(reg_param=0)
294 model.fit(X_in, y_in)
295 print_metrics(model, X_test, y_test)
296
297 """# K nearest neighbours"""
298
299 #elbow method for figuring out optimal number of k neighbours
300 from sklearn.metrics import precision_score
301 from sklearn.model_selection import KFold
302
303 def metric_f_beta(y_test, y_preds):
304     y_test = np.array(y_test).reshape(-1)
305     recall = recall_score(y_test, y_preds, pos_label='high_bike_demand')
306     precision = precision_score(y_test, y_preds, pos_label='high_bike_demand')
307     beta = 1.5
308     F_beta = ((1+beta*2) * precision * recall) / (beta*2 * precision + recall)
309     return F_beta
310 n = [x for x in range(1,50)]
311 error_rates = []
312
313 kf = KFold(n_splits=5, random_state=42, shuffle=True)
314 y_in = np.array(y_in)
315 for j, (train_index, test_index) in enumerate(kf.split(X_in)):
316     error_rate_i = []
317     for i in n:
318         model = skl_nb.KNeighborsClassifier(n_neighbors=i)
319         model.fit(X_in[train_index], y_in[train_index])
320         y_pred = model.predict(X_in[test_index])
321         error_rate = 1-metric_f_beta(y_in[test_index], y_pred)
322         error_rate_i.append(error_rate)
323     error_rates.append(error_rate_i)
324
325 avg_errors = np.mean(np.array(error_rates), axis=0)
326
327 plt.plot(n, avg_errors)
328 plt.xlabel('Number of neighbors (k)')
329 plt.ylabel('Error f_beta rate')
330 plt.show()
331 print(f'minimum is: {np.min(avg_errors)} at k={np.argmin(avg_errors)+1}')
332
333 n = 2
334 model = skl_nb.KNeighborsClassifier(n_neighbors=n)
335 model.fit(X_in, y_in)

```

```

336 print_metrics(model, X_test, y_test)
337
338 """# Tree based methods"""
339
340 from sklearn.model_selection import RandomizedSearchCV
341 from sklearn.ensemble import RandomForestClassifier
342 model = RandomForestClassifier()
343 # Number of trees in random forest
344 n_estimators = [int(x) for x in np.linspace(start = 1, stop = 1000,
345 num = 10)]
346 # Number of features to consider at every split
347 max_features = ['auto', 'sqrt']
348 # Maximum number of levels in tree
349 max_depth = [int(x) for x in np.linspace(1, 100, num = 5)]
350 max_depth.append(None)
351 # Minimum number of samples required to split a node
352 min_samples_split = [2, 5, 8, 10]
353 # Minimum number of samples required at each leaf node
354 min_samples_leaf = [1, 2, 4, 8, 15]
355 # Method of selecting samples for training each tree
356 bootstrap = [True, False]
357 # Create the random grid
358 random_grid = {'n_estimators': n_estimators,
359 'max_features': max_features,
360 'max_depth': max_depth,
361 'min_samples_split': min_samples_split,
362 'min_samples_leaf': min_samples_leaf,
363 'bootstrap': bootstrap}
364 fb_score = make_scoring(fbeta_score, beta=1.5, pos_label='
365 high_bike_demand')
366 rf_random = RandomizedSearchCV(estimator=model, param_distributions
367 =random_grid, n_iter=100, cv=3, verbose=2, random_state=42,
368 n_jobs=-1, scoring=fb_score)
369 rf_random.fit(X_in, y_in)
370
371 rf_random.best_params_
372
373 model = RandomForestClassifier(n_estimators=334, max_depth=50,
374 min_samples_leaf=1, min_samples_split=5, bootstrap=False,
375 max_features='sqrt')
376 model.fit(X_in, y_in)
377 print_metrics(model, X_test, y_test)
378
379 """# Final test"""
380
381 final_test = pd.read_csv('test_data.csv')
382 final_test = final_test.drop(['precip', 'holiday', 'snowdepth', '
383 snow', 'dew'], axis=1)
384 final_test['rushhour'] = final_test['hour_of_day'].apply(lambda x:
385 1 if x > 2 and x < 8 else 0)
386 final_test['goodtemp'] = final_test['temp'].apply(lambda x: 1 if x
387 >25 else 0)
388 final_test = normalize(final_test)
389
390 # Run final test using random forest as it had highest accuracy
391 # and we assume this will be the metric used to measure performance

```

```

383 model = RandomForestClassifier(n_estimators=334, max_depth=50,
    min_samples_leaf=1, min_samples_split=5, bootstrap=False,
    max_features='sqrt')
384 model.fit(X, y)
385
386 preds = model.predict(final_test)
387 np.unique(preds, return_counts=True)
388
389 import csv
390 out_preds = [0 if pred == 'low_bike_demand' else 1 for pred in
    preds]
391 with open(r'out.csv', 'w') as fp:
392     for pred in out_preds:
393         # write each item on a new line
394         fp.write("%s," % pred)

```