

---

# Do we need more bikes?

## Project in Statistical Machine Learning

---

Erik Björk

Simona Stoyanoska

Olov Rahm

Jakob Hanson

### Abstract

In this project we develop, and study different statistical machine learning models for predicting whether the number of available bikes at a given hour should be increased, a project by the District Department of Transportation in Washington D.C. The training data set consists of 1600 instances of hourly bike rentals, and a test set of 400 instances. The models for prediction we have used are: *Logistic regression*, *Discriminant methods: LDA, QDA*, *k- Nearest Neighbour*, and *Tree Based Methods*. We have found that *k- Nearest Neighbour* gives best prediction, with accuracy 92%.  
The group consists of 4 students.

# 1 Introduction

Statistical machine learning is a subject that aims to build and train algorithms, that analyse large amount of data, and make predictions for the future, which are computed by using established statistical models, and tools from functional analysis. This is a project in supervised, statistical machine learning, where several models were created, and trained, in order to analyse which one of them gives best prediction for the project "Do we need more bikes", where we want to understand, and predict if there is a high, or low demand of city bikes in the public transportation of Washington, a project by the District Department of Transportation in Washington D.C..

The data set used for training our models, consist of 15 variables, containing quantitative/qualitative data. We developed several models, and evaluated them with cross-validation, in order to understand which algorithm gives the best prediction.

## 2 Theoretical Background

Here we review the theoretical background of the models. We follow mostly [1].

Let  $\{\mathbf{x}_i, y_i\}_{i \in \mathbb{N}}$  be the training data set, where  $\mathbf{x}_i$  is a matrix representing the input, and  $y_i$  is the output.

### 2.1 Mathematical Overview of the Models

#### 2.1.1 Logistic Regression

The backbone of logistic regression is linear regression, i.e. finding the least-squares solution to an equation system

$$X\theta = y$$

given by the normal equations

$$X^T X \theta = X^T y$$

where  $X$  is the training data matrix,  $\theta$  is the coefficient vector and  $y$  is the training output. The parameter vector is then used in the sigmoid function  $\sigma(z) : \mathbb{R} \rightarrow [0, 1]$

$$\sigma(z) = \frac{e^z}{1 + e^z} :$$

where  $z = x^T \theta$ , and  $x$  is the testing input. This gives a statistical interpretation of the input vector. In the case of a binary True/False classification, the value of the sigmoid function then determines the class. The model was tuned using the `lbfgs` algorithm, which is standard in *SKLearn*. It optimizes the weights  $\theta_i$  to minimize the log-loss function:

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y_i \log(h_\theta(x_i)) + (1 - y_i) \log(1 - h_\theta(x_i))] + \frac{\lambda}{2} \|\theta\|^2, \quad (1)$$

where  $\lambda$  is the regularization strength (default for *SKLearn*:  $\lambda = 1.0$ ),  $h_\theta$  is the sigmoid function for each of the  $m$  input vectors  $x_i$ , and  $\theta$  is the parameter vector being optimized. It is being optimized through a gradient method.

#### 2.1.2 Random forest

The random forest method is a based upon decision trees, i.e. dividing the data point into binary groups based on Gini-impurity, entropy or classification error, Gini being the most common. These divisions are then used to create a binary tree shown in figure ??Tree) and where the leaf-nodes are used to classify the target variables based on the input. As of itself the decision tree tends to have unsatisfying results which leads to methods like random forest and sandbagging that boost its accuracy. Sandbagging is a way to sample the data in order to get multiple datasets from the same data. One then creates a decision-tree for every subset data to then combine them into one model. This lessens the variance of the model but increases bias. This means that sandbagging can increase false negatives which in this application makes it nonviable. Random forest on the otherhand is viable, it creates multiple trees while discarding random input variable this randomness decreases overfitting creating a more robust model.

### 2.1.3 Non-parametric method: k-Nearest Neighbour

*k-Nearest Neighbour* ( $k$ -NN) is a distance based method that takes a  $k$  amount of points from the training data set, called *neighbours*, computes the distance between them, then assumes that the predicted value  $\hat{y}(x_*)$  follows the trend of the  $k$ -nearest neighbours. Since  $k$ -NN uses the training data explicitly it is also called a *nonparametric* method.

The  $k$ -NN method can be divided into several subcategories, inter alias *classification*  $k$ -NN method, *regression*  $k$ -NN method. In this project, we are using the classification method, since we are trying to predict in which of the two classes low, or high demand, the given, and predicted data points belong.

The classification  $k$ -NN algorithm evaluates  $\hat{y}(x_*)$  by computing the most frequently occurring class among the  $k$  nearest neighbours. Here, we try to identify whether a data point belong to the high demand-class. Denote  $c$  = high demand class. For simplicity, assume Euclidean ambience. Then

$$\hat{y}(x_*) = \arg \max_c \sum_{n \in \mathbb{N}} \chi_{(y_n=c)},$$

where  $y_i$  is the class of the nearest neighbour,  $\chi$  is the characteristic function

$$\chi_{(y_i=c)} = \begin{cases} 1 & \text{if } y_n = c, \\ 0 & \text{otherwise.} \end{cases}$$

It is very common to use a weighted sum to predict the next value, i.e.

$$\hat{y}(x_*) = \arg \max_c \sum_{n \in \mathbb{N}} \frac{\chi_{(y_n=c)}}{d(x, x_n)},$$

where  $d$  is the standard Euclidean metric, computing the distance between an input  $x$ , and a neighbour  $x_n$ .

When using this model it is important to choose an optimal  $k$ -value. There are several tests for this, here we implement *uniform weighting*, and *distance weighting*. The first algorithm creates a  $k$ -NN model for each new  $k \in [1, 500]$ , and trains the model with uniform weights, i.e. the contribution of all neighbours is equal. Similarly, the latter trains a  $k$ -NN classifier for each  $k \in [1, 500]$ , with the difference that it uses distance based weighting, i.e. closer neighbours have greater influence. After testing different upper boundaries for  $k$ , the two models gave good results in the interval  $[1, 500]$ , see Figure 1. From the figures, we can see that the second test gives a better value for  $k$ , since the plot follows smoother trend, in comparison to the uniform weighting test, which makes it easier to identify an optimal  $k$  value ( $k = 120$ ). Moreover, the distance weighting algorithm is providing results for larger values of  $k$ , that is for  $k \in [1, 400)$  before the curve converges, while the uniform weighting algorithm converges earlier, when  $k = 120$ . This means that for large  $k$ , both test algorithms make prediction based on the most common class in the data set, instead of making prediction based on the behaviour of the neighbours. Thus for sufficiently large  $k$ , for any given data point, the model will consider unnecessarily large amount of neighbours, and the prediction will be evaluated to belong to the most frequent class. Since the distance weighting has a larger range of  $k$ -value, it should be more trustworthy.

When  $k = 120$ , the accuracy of the model is 92%.

### 2.1.4 Discriminant analysis: LDA and QDA

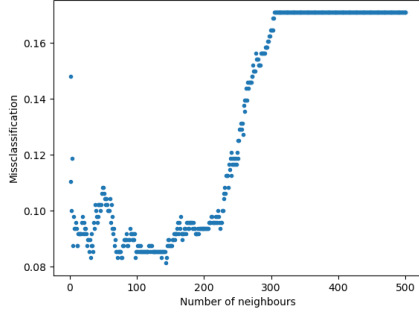
Linear Discriminant Analysis is a generative model, which means it is a model that's creating and using a probability distribution  $P(\mathbf{x}, y)$  to create an estimation for the probability  $P(y = m|\mathbf{x})$  using bayes theorem.

Bayes theorem is:

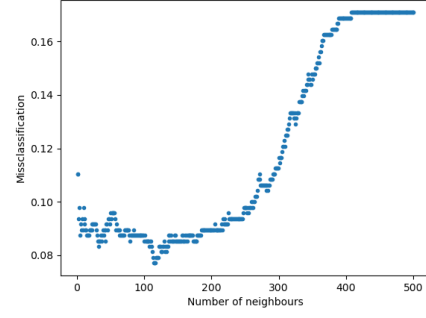
$$p(y|\mathbf{x}) = \frac{p(y, \mathbf{x})}{p(\mathbf{x})} = \frac{p(y)p(\mathbf{x}|y)}{\int_y p(y, \mathbf{x})}$$

For the discrete version it is obtained:

$$p(y = m|\mathbf{x}) = \frac{p(y = m)p(\mathbf{x}|y = m)}{\sum_{m=1}^M p(y = m)p(\mathbf{x}|y = m)}$$



(a) Uniform distance test for  $k$ .



(b) Weighted distance test for  $k$ .

Figure 1: Test for choosing an optimal  $k$ -value.

For this form of the equation to be useful, it is necessary to obtain an accurate estimation of  $p(y = m)$  and  $p(\mathbf{x}|y = m)$  for all classes  $m$ .

In LDA,  $p(y = m)$  is estimated by counting the percentage of data points (in the training data) being in each of the classes and using that percentage as the probability of a data point being in that class. In mathematical terms:

$$p(y = m) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}\{y_i = m\} = \frac{n_m}{n}$$

To estimate the probability distribution  $p(\mathbf{x}|y = m)$ , a multi-dimensional gaussian distribution is used:

$$\mathcal{N}(\mathbf{x}|\mu, \Sigma) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1}(\mathbf{x} - \mu)\right)$$

Where  $\mathbf{x}$  is the  $d$ -dimensional data point,  $\mu$  is the ( $d$ -dimensional) mean of the random variable.  $\Sigma$  is the symmetric, positive definite covariance matrix defined by:

$$\Sigma = \frac{1}{n - M} \sum_{m=1}^M \sum_{i:y_i=m} (\mathbf{x}_i - \mu_m)(\mathbf{x}_i - \mu_m)^T$$

Using these estimations results in an expression for the quantity  $p(y = m|\mathbf{x}) \forall m$ . LDA then uses maximum likelihood to categorize an input  $\mathbf{x}$  into a class  $m$ .

Quadratic discriminant analysis (QDA) is heavily based of LDA with the sole difference being how the covariance matrix  $\Sigma$  is created. In LDA, the covariance matrix is assumed to be the same for data in each and every class. In QDA however, the covariance matrix is calculated for each class as follows:

$$\Sigma_m = \frac{1}{n_m - 1} \sum_{i:y_i=m} (\mathbf{x}_i - \mu_m)(\mathbf{x}_i - \mu_m)^T$$

To optimize LDA and QDA, it's all about estimating  $\Sigma_m$  and  $\mu_m$  for all classes  $m$ . The mean  $\mu$  is estimated as the mean of all the input data  $\mu = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$ . The covariance matrix  $\Sigma_m$  is defined previously in the section.

One thing to note about LDA and QDA is that the use of a multi-variable gaussian distribution benefits normally distributed variables. In this project however, there is a dependance on positive definite values which are not normally distributed by nature. This is an issue when using QDA since in the class of *high\_bike\_demand*, all data points have a snow depth of 0 and has hence no variance. This results in this class having a undefined inverse for the covariance matrix. The solution used was to exclude this variable from this model.

## 2.2 Input Data Modification

By plotting the data and analyzing the .csv file, some observations were made. The different inputs were then changed accordingly:

- *Kept as-is*: weekday, windspeed, visibility, temp
- *Modified*:
  - month - split into two inputs, one cosine and one sine part. This make the new inputs linear and can follow the fluctuations of the year. The original input was discarded.
  - hour\_of\_day - split into three boolean variables: demand\_day, demand\_evening, and demand\_night, reflecting if the time was between 08-14, 15-19 or 20-07 respectively. This was done because plotting the data showed three different plateaus of demand for the different time intervals. The original input was discarded.
  - snowdepth, precip were transformed into booleans, reflecting if it was raining or if there was snow on the ground or not. This was done as there was no times where demand was high when it was raining or when there was snow on the ground.
- *Removed*: cloudcover, day\_of\_week, snow, dew, holiday, summertime. These were removed due to being redundant (e.g. summertime), not showing a clear trend (e.g. cloudcover), giving a worse score when used, or all three (e.g. day\_of\_week).

## 3 Data Analysis

In the given data, there are some numerical and categorical features:

- *Numerical*: temp, dew, humidity, precip, snow, snowdepth, windspeed, cloudcover and visibility.
- *Categorical*: hour\_of\_day, day\_of\_week, month, holiday, weekday, summertime, and increase\_stock

There are some trends seen in the data when it comes to time and weather. From figure 2, one can see a periodic relationship for the months, where there is a higher demand during the warmer months, loosely following a trigonometric curve. Over the week, the demand is rather stable, with a peak on the weekend, especially Saturdays.

Looking at the weather (figure 3); if there is rain or if there is snow on the ground, there is close to always low demand. Cloudcover did not make a big impact, which is also intuitive, as a cloudy day does not make biking more difficult. Dew point also does not have a clear trend, while humidity however has a clear trend downwards as the humidity increases. Temperature had a more clear impact, where more people wanted to bike the warmer it got.

The overall trend is that about one eighth of observations correspond to a high bike demand. During the night, or in bad weather, the demand is (intuitively) low. But during rush hour (figure 2c), the demand is very high, and should probably be increased in order to minimize excessive CO<sub>2</sub> emissions.

## 4 Result

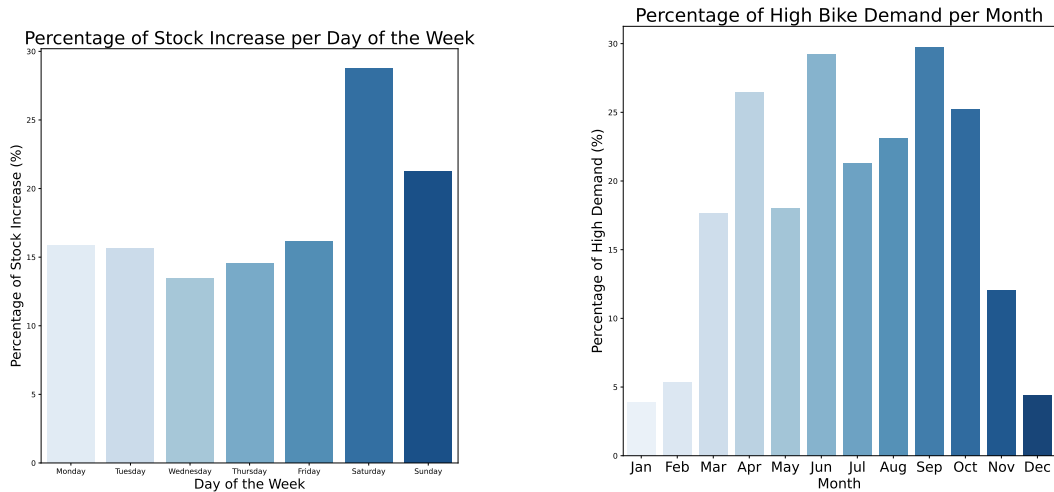
The method used to evaluate the different models where chosen to be the accuracy as well as the precision and recall of the class "high bike demand". The accuracy is defined simply as:

$$\text{Accuracy} = \frac{n_{\text{correct}}}{n_{\text{tot}}}$$

And the precision and recall is defined as:

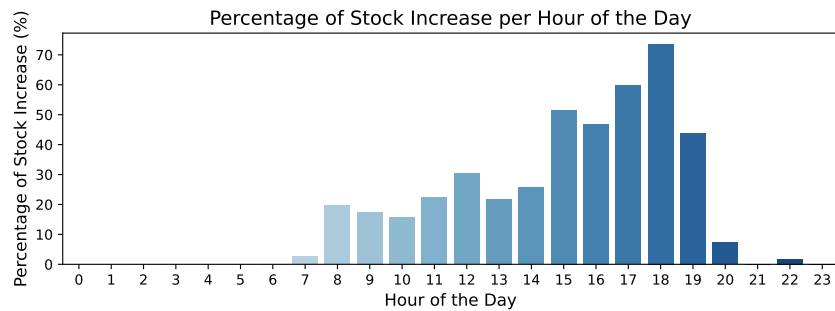
$$\text{Precision} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}} \quad \text{Recall} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$

Furthermore, a naive model that only guessed there is a low demand was compared to the rest of the models. The different models were tested and the accuracy where:



(a) Demand per day of week.

(b) Demand per month.



(c) Demand per hour of day.

Figure 2: Bike demand vs. day of week and month.

Accuracy of the models

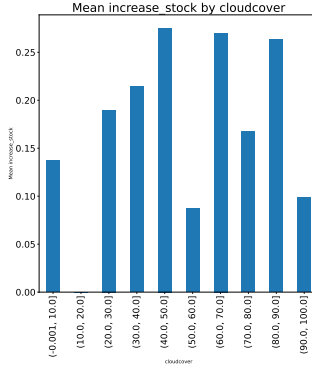
Model	Accuracy	Precision	Recall
LDA	85%	53%	50%
QDA	87%	67%	36%
k-nearest neighbour	92%	81%	70%
Random Forest	91%	77%	71%
Logistic Regression	90%	73%	63%
Naive	83%	0%	0%

Here you can clearly see random forest and k-nearest neighbour are the best classifiers both outperforming linear and quadratic regression on accuracy, precision and recall. Out of random forest and kNN the group would proceed with the kNN method, its higher accuracy and precision score outweighing the slightly better recall score of random forest. This will mean a slight loss in income caused by increasing false negatives but is thought to be covered by fewer false positives.

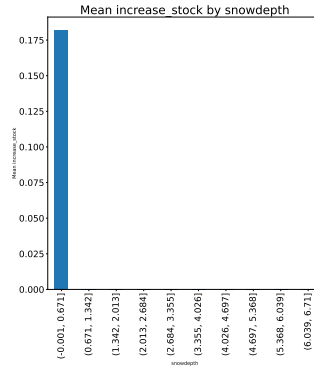
## 5 Conclusion

From the evaluation the models, k-nn performed the best with the highest accuracy and precision. As for the recall, k-nn did not perform the best but is considered adequate and hence this method is chosen as the best one.

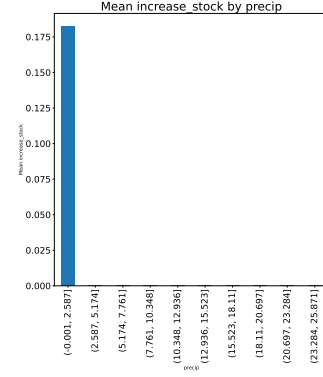
One reason for the discriminant analysis falling short of the other models is likely due to



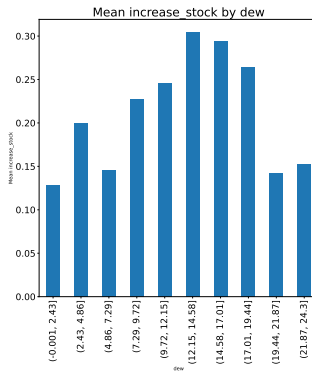
(a) Demand per cloudcover (percentage).



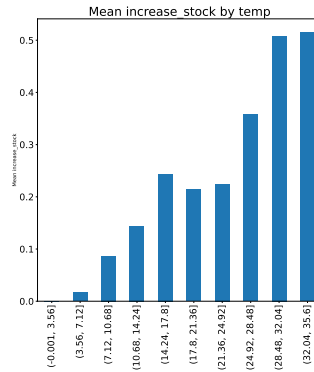
(b) Demand per snowdepth (mm).



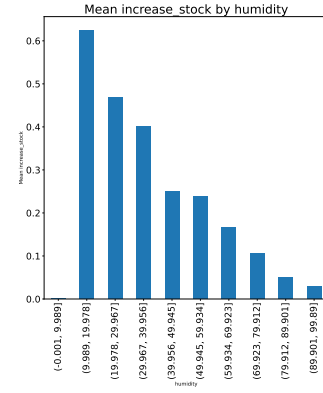
(c) Demand per precipitation (mm/hour).



(d) Demand per dew point ( $^{\circ}\text{C}$ ).



(e) Demand per temperature ( $^{\circ}$ ).



(f) Demand per humidity level (percentage).

Figure 3: Bike demand vs. various weather parameters.

these models being designed with the assumption of variables being normally distributed. This is not the case for this particular data set.

## References

- [1] A. Lindholm, N. Wahlström, F. Lindsten, and T. B. Schön. *Machine learning: a first course for engineers and scientists*. Cambridge University Press, 2022.

## A Appendix

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import sklearn.linear_model as skl_lm
5 import sklearn.discriminant_analysis as skl_da
6 import sklearn.preprocessing as pp
7 import sklearn.metrics as skl_m
8
9 import sklearn.neighbors as skl_nb
10
11
12 # ----- DATA MODIFICATION AND SETUP -----
13
14 df = pd.read_csv('training_data_vt2025.csv')
15
16 # Modify the dataset, emphasizing different variables
17
18 df['month_cos'] = np.cos(df.month*np.pi/12)
19 df['month_sin'] = np.sin(df.month*np.pi/12)
20
21 # time of day, replaed with low,medium and high demand,
22 # adding the new categories back in the end.
23 def categorize_demand(hour):
24     if 20 <= hour or 7 >= hour:
25         return 'night'
26     elif 8 <= hour <= 14:
27         return 'day'
28     elif 15 <= hour <= 19:
29         return 'evening'
30
31 df['demand_category'] = df['hour_of_day'].apply(categorize_demand)
32 df_dummies = pd.get_dummies(df['demand_category'], prefix='demand',
33                             drop_first=False)
34 df = pd.concat([df, df_dummies], axis=1)
35
36 # converting to bools
37 def if_zero(data):
38     if data == 0:
39         return True
40     else:
41         return False
42
43 # temperature
44
45 df['snowdepth_bool'] = df['snowdepth'].replace(0, False).astype(bool)
46 df['precip_bool'] = df['precip'].replace(0, False).astype(bool)
47
48 # Split into train and test:
49 np.random.seed(0)
50 df_modified=df[['holiday',
51                 'weekday',
52                 'summertime',
53                 'temp',
54                 'dew',
55                 'humidity',
56                 'visibility',
57                 'windspeed',
58                 'month_cos',
59                 'month_sin',
60                 'demand_day',
61                 'demand_evening',
62                 'demand_night',
```



```

63         'snowdepth_bool',
64         'precip_bool',
65         'increase_stock']]
66
67 N = df_modified.shape[0]
68 n = round(0.7*N)
69 trainI = np.random.choice(N,size=n,replace=False)
70 trainIndex = df_modified.index.isin(trainI)
71 train = df_modified.iloc[trainIndex]
72 test = df_modified.iloc[~trainIndex]
73
74 # Set up X,Y
75
76 # Train data
77 X = train.iloc[:,0:-2]
78 Y = train['increase_stock']
79
80 # Test data
81 X_test = test.iloc[:,0:-2]
82 Y_test = test['increase_stock']
83
84
85 # ----- K-NN METHOD -----
86
87 """
88 # Tests for k-value
89 # TEST 1 - uniform distance
90 missclassification = []
91 for k in range(500): # Try n_neighbours = 1, 2, ...,
92
93     #kNN method
94     scaler = pp.StandardScaler().fit(X)
95     model = skl_nb.KNeighborsClassifier(n_neighbors = k+1, weights = '
uniform')
96     model.fit(scaler.transform(X),Y)
97
98     # Prediction
99     y_hat = model.predict(scaler.transform(X_test))
100     missclassification.append(np.mean(y_hat != Y_test))
101
102 K = np.linspace(1, 500, 500)
103 plt.plot(K, missclassification, '.')
104 plt.ylabel('Missclassification')
105 plt.xlabel('Number of neighbours')
106 plt.show()
107
108 #TEST 2
109 missclassification = []
110 for k in range(500): # Try n_neighbours = 1, 2, ...,
111
112     #kNN method
113     scaler = pp.StandardScaler().fit(X)
114     model = skl_nb.KNeighborsClassifier(n_neighbors = k+1, weights = '
distance')
115     model.fit(scaler.transform(X),Y)
116
117     # Prediction
118     y_hat = model.predict(scaler.transform(X_test))
119     missclassification.append(np.mean(y_hat != Y_test))
120
121 K = np.linspace(1, 500, 500)
122 plt.plot(K, missclassification, '.')
123 plt.ylabel('Missclassification')
124 plt.xlabel('Number of neighbours')
125 plt.show()

```

```

126 """
127
128 # creating the model
129 model = skl_nb.KNeighborsClassifier(n_neighbors = 120, weights = '
      distance')
130
131 # Scaling the data, otherwise
132 scaler = pp.StandardScaler().fit(X)
133 model.fit(scaler.transform(X),Y)
134 y_hat = model.predict(scaler.transform(X_test))
135
136 # Get confusion matrix
137 diff = pd.crosstab(y_hat, Y_test)
138 print(f'Confusion matrix for K-NN: \n {diff}')
139
140 # Get metrics:
141 print(f'Metrics for K-NN: \n{skl_m.classification_report(Y_test, y_hat
      )}')
142
143
144
145 # ----- LDA & QDA -----
146
147 # LDA:
148
149 # Apply Linear Discriminant Analysis (LDA)
150 lda = skl_da.LinearDiscriminantAnalysis(n_components=1)
151 X_train_lda = lda.fit_transform(X, Y)
152 X_test_lda = lda.transform(X_test)
153
154 # Train a classifier (Logistic Regression)
155 clf = skl_lm.LogisticRegression()
156 clf.fit(X_train_lda, Y)
157
158 # Make predictions
159 y_pred = clf.predict(X_test_lda)
160
161 # Evaluate accuracy
162 accuracy = skl_m.accuracy_score(Y_test, y_pred)
163 print(f"Model Accuracy for LDA: \n {accuracy:.2f}")
164
165 # Get confusion matrix
166 diff = pd.crosstab(y_pred, Y_test)
167 print(f'Confusion matrix for LDA: \n {diff}')
168
169 # QDA:
170
171 # Has a variance of 0
172 X = X.drop(columns='snowdepth_bool')
173 X_test = X_test.drop(columns='snowdepth_bool')
174
175 # Apply Quadratic Discriminant Analysis (QDA)
176 qda = skl_da.QuadraticDiscriminantAnalysis()
177 X_train_lda = qda.fit(X, Y)
178
179 # Make predictions
180 y_pred = qda.predict(X_test)
181
182 # Evaluate accuracy
183 accuracy = skl_m.accuracy_score(Y_test, y_pred)
184 print(f"Model Accuracy: {accuracy:.2f}")
185
186 print(f'Metrics for QDA: \n{skl_m.classification_report(Y_test, y_pred
      )}')
187

```

```

188 # ----- LOGISTIC REGRESSION -----
189
190 # Scaling the data, otherwise
191 scaler = pp.StandardScaler().fit(X)
192 model.fit(scaler.transform(X),Y)
193 y_hat = model.predict(scaler.transform(X_test))
194
195 # Get confusion matrix
196 diff = pd.crosstab(y_hat, Y_test)
197 print(f'Confusion matrix: \n {diff}')
198
199 # Get metrics:
200 print(f'metrics for Logistic regression: \n {skl_m.
      classification_report(Y_test, y_hat)}')
201
202
203 # ----- TREE-BASED METHODS -----
204
205 from sklearn.ensemble import BaggingClassifier, RandomForestClassifier
206 from sklearn.model_selection import GridSearchCV, RandomizedSearchCV
207
208 model = RandomForestClassifier(random_state=42)
209 param_grid = {
210     'n_estimators': [100, 200, 300],
211     'max_depth': [10, 20, None],
212     'min_samples_split': [2, 5, 10],
213     'min_samples_leaf': [1, 2, 4]
214 }
215
216 # Set up Grid Search
217 random_search = RandomizedSearchCV(model, param_grid, cv=5, scoring='
      accuracy', n_jobs=-1, verbose=2)
218
219 # Fit on training data
220 random_search.fit(X, Y)
221
222 # Get the best hyperparameters
223 print("Best Parameters: ", random_search.best_params_)
224 print("Best Accuracy: %.2f" % random_search.best_score_)
225
226 # Update the model with the best parameters
227 best_model = random_search.best_estimator_
228
229 # Fit the best model on the training data
230 best_model.fit(X, Y)
231
232 # Make predictions using the optimized model
233 y_predict = best_model.predict(X_test)
234
235 print(f'Metrics for tree method: \n{skl_m.classification_report(Y_test
      , y_predict)}')
236
237 # ----- NAIVE MODEL -----
238
239 # We don't train any model, we just guess there is always a low demand
240 y_test = test['increase_stock']
241 y_predict = np.array(["low_bike_demand"]*len(y_test))
242
243 print(f'Metrics for naive model: \n{skl_m.classification_report(y_test
      , y_predict)}')

```

Listing 1: Code for all methods.