
Do we need more bikes?

Project in Statistical Machine Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

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

10 **1 Plan**

11 **1.1 From Intro**

- 12 (i) Explore and preprocess data
- 13 (ii) try some or all classification methods, which are these?
 - 14 • Logistic Regression
 - 15 • Discriminant analysis: LDA, QDA
 - 16 • K-nearest neighbor
 - 17 • Tree-based methods: classification trees, random forests, bagging
 - 18 • Boosting
- 19 (iii) Which of these are to be "put in production"?

20 **1.2 From Data analysis task**

- 21 • Can any trend be seen comparing different hours, weeks, months?
- 22 • Is there any difference between weekdays and holidays?
- 23 • Is there any trend depending on the weather?

24 **1.3 From Implementation of methods**

- 25 Each group member should implement one family each, who did what shall be clear!
- 26 DNNs are encouraged to be implemented, do this if there is time. (DNN is not a thing a group member can claim as their family.)
- 27
- 28 Implement a naive version, let's do: *Always low_bike_demand*

29 **1.3.1 What to do with each method**

- 30 1. Implement the method (each person individually)
- 31 2. Tune hyper-parameters, discuss how this is done (each person individually)
- 32 3. Evaluate with for example cross-validation. Don't use E_{k-fold} (what is that?) (need to do
- 33 together)
- 34 4. (optional) Think about input features, are all relevant? (together)
- 35 Before training, unify pre-processing FOR ALL METHODS and choose ONE OR MULTIPLE
- 36 metrics to evaluate the model. (is it necessary to have the same for all?, is it beneficial?) Examples:
 - 37 • accuracy
 - 38 • f1-score
 - 39 • recall
 - 40 • precision

41 Use same test-train split for ALL MODELS

42 **2 Introduction**

- 43 Statistical machine learning is a subject that aims to build and train algorithms, that analyse large
- 44 amount of data, and make predictions for the future, which are computed by using established
- 45 statistical models, and tools from functional analysis. This is a project in supervised, statistical
- 46 machine learning, where several models were created, and trained, in order to analyse which one of
- 47 them gives best prediction for the project "Do we need more bikes", where we want to understand,
- 48 and predict if there is a high, or low demand of city bikes in the public transportation of Washington,
- 49 a project by the District Department of Transportation in Washington D.C..
- 50 The data set used for training our models, consist of 15 variables, containing quantitative/qualitative
- 51 data. We developed several models, and evaluated them with cross-validation, in order to understand
- 52 which algorithm gives the best prediction.

3 Theoretical Background

3.1 Mathematical Overview of the Models

3.1.1 Logistic Regression

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

$$X\theta = y \quad (1)$$

given by the normal equations

$$X^T X \theta = X^T y \quad (2)$$

where X is the training data matrix, θ is the coefficient vector and b is the training output. The parameter vector is then used in the sigmoid function:

$$\sigma(z) = \frac{e^z}{1 + e^z} : \mathbb{R} \rightarrow [0, 1], \quad (3)$$

$$z = x^T \theta, \quad (4)$$

where 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.

3.1.2 Random forest

The random forest method is 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 other hand is viable, it creates multiple trees while discarding random input variable this randomness decreases overfitting creating a more robust model.

3.1.3 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 alia *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 belongs to the high demand-class. Denote c = high demand class. For simplicity, assume Euclidean distance. 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, χ 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%.

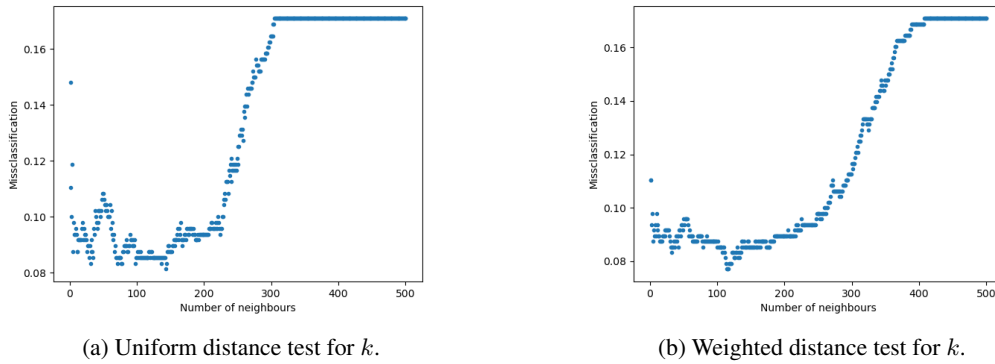


Figure 1: Test for choosing an optimal k -value.

3.1.4 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:

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

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:

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

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

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

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

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

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

125
 126 Quadratic discriminant analysis (QDA) is heavily based on LDA with the sole difference
 127 being how the covariance matrix Σ is created. In LDA, the covariance matrix is assumed to be the
 128 same for data in each and every class. In QDA however, the covariance matrix is calculated for each
 129 class as follows:

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

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

136 3.2 Input Data Modification

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

- 139 • *Kept as-is:* weekday, windspeed, visibility, temp
- 140 • *Modified:*
 - 141 – month - split into two inputs, one cosine and one sine part. This makes the new inputs
 - 142 linear and can follow the fluctuations of the year. The original input was discarded.
 - 143 – hour_of_day - split into three boolean variables: demand_day, demand_evening,
 - 144 and demand_night, reflecting if the time was between 08-14, 15-19 or 20-07 respec-
 - 145 tively. This was done because plotting the data showed three different plateaus of
 - 146 demand for the different time intervals. The original input was discarded.
 - 147 – snowdepth, precip were transformed into booleans, reflecting if it was raining or
 - 148 if there was snow on the ground or not. This was done as there were no times where
 - 149 demand was high when it was raining or when there was snow on the ground.
- 150 • *Removed:* cloudcover, day_of_week, snow, dew, holiday, summertime. These were
- 151 removed due to being redundant (e.g. summertime), not showing a clear trend (e.g.
- 152 cloudcover), giving a worse score when used, or all three (e.g. day_of_week).

153 4 Data Analysis

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

- 155 • *Numerical:* temp, dew, humidity, precip, snow, snowdepth, windspeed, cloudcover
- 156 and visibility.
- 157 • *Categorical:* hour_of_day, day_of_week, month, holiday, weekday, summertime, and
- 158 increase_stock

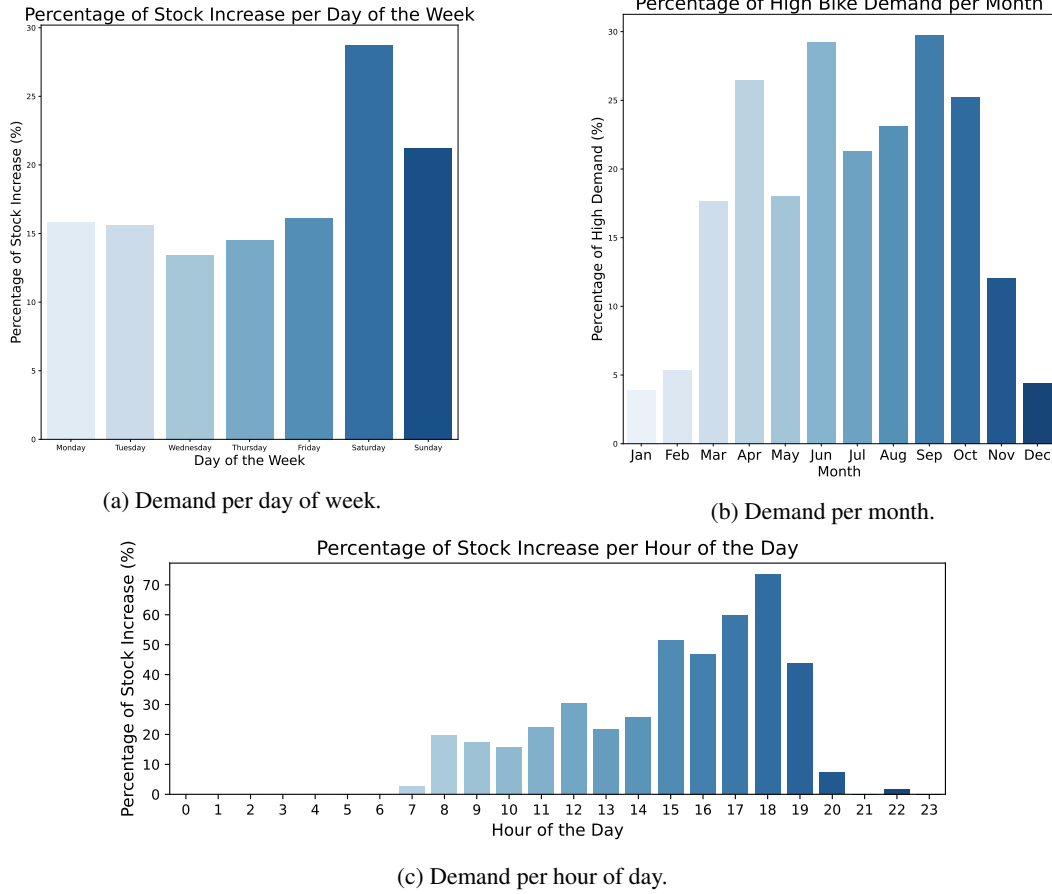


Figure 2: Bike demand vs. day of week and month.

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 CO₂ emissions.

5 Result

The method used to evaluate the different models were simply chosen to be the accuracy defined by:

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

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:

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

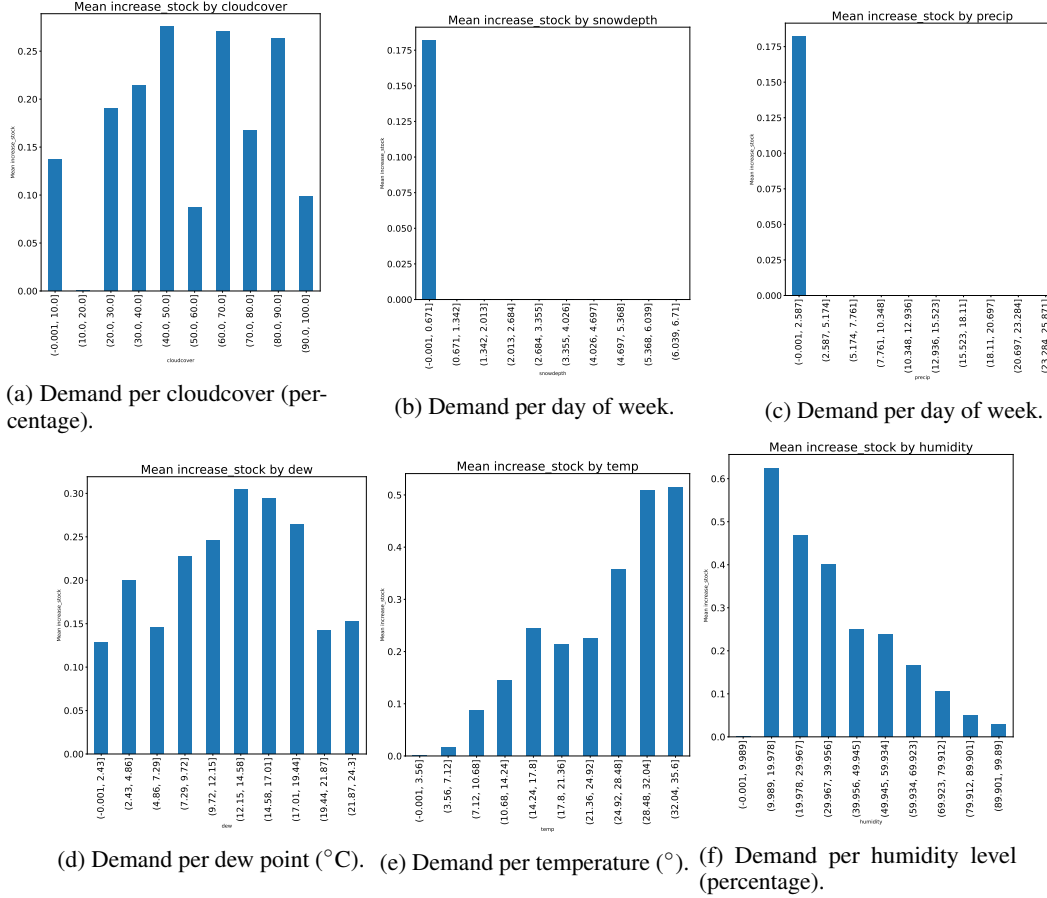


Figure 3: Bike demand vs. various weather parameters.

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%
Naive	83%	0%	0%

177 and kNN the group would proceed with the kNN method, its higher accuracy and precision score out
 178 waying the slightly better recall score of random forest. This will mean a slight loss in income caused
 179 by increasing false negatives but is thought to be covered by fewer false positives.

180 6 Conclusion

181 From the evaluation the two models k-nn performed the best with the highest accuracy and precision.
 182 As for the recall this is not the highest but is considered adequate and hence this method is chosen.

183

184 One reason for the discriminant analysis falling short of the other models is likely due to
 185 these models being designed with the assumption of variables being normally distributed. This is not
 186 the case for this particular data set.

187 A Appendix

```

188 1 import pandas as pd
189 2 import numpy as np
190 3 from sklearn.model_selection import train_test_split
191 4 from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
192 5 from sklearn.linear_model import LogisticRegression
193 6 from sklearn.metrics import accuracy_score
194 7 from sklearn.metrics import classification_report
195 8
196 9 df = pd.read_csv('training_data_vt2025.csv')
197 10
198 11 # modify the month to represent the periodicity that is observed in
199 12 data.
200 13 df['month_cos'] = np.cos(df['month']*2*np.pi/12)
201 14 df['month_sin'] = np.sin(df['month']*2*np.pi/12)
202 15
203 16 # time of day, replaced with 3 bool values: is_night, is_day and
204 17 is_evening,
205 18 # adding the new categories back in the end.
206 19 def categorize_demand(hour):
207 20     if 20 <= hour or 7 >= hour:
208 21         return 'night'
209 22     elif 8 <= hour <= 14:
210 23         return 'day'
211 24     elif 15 <= hour <= 19:
212 25         return 'evening'
213 26
214 27 df['time_of_day'] = df['hour_of_day'].apply(categorize_demand)
215 28 df_dummies = pd.get_dummies(df['time_of_day'], prefix='is', drop_first
216 29 =False)
217 30 df = pd.concat([df, df_dummies], axis=1)
218 31
219 32 # Create bool of snowdepth and percipitation
220 33 df['snowdepth_bool'] = df['snowdepth'].replace(0, False).astype(bool)
221 34 df['precip_bool'] = df['precip'].replace(0, False).astype(bool)
222 35
223 36 # Seperate training data from target
224 37 X=df[['#holiday',
225 38     'weekday',
226 39     '#summertime',
227 40     'temp',
228 41     '#dew',
229 42     '#humidity',
230 43     '#visibility',
231 44     '#windspeed',
232 45     '#month',
233 46     'month_cos',
234 47     'month_sin',
235 48     '#hour_of_day',
236 49     'is_day',
237 50     'is_evening',
238 51     'is_night',
239 52     '#hour_cos',
240 53     '#hour_sin',
241 54     'snowdepth_bool',
242 55     'precip_bool'
243 56 ]]
244 57
245 58 y=df['increase_stock']
246 59
247 60 # Split dataset into training and test sets
248 61 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size
249 62 =0.2, random_state=42)
250 63

```



```

25160 # Apply Linear Discriminant Analysis (LDA)
25261 lda = LinearDiscriminantAnalysis(n_components=1)
25362 X_train_lda = lda.fit_transform(X_train, y_train)
25463 X_test_lda = lda.transform(X_test)
25564
25665 # Train a classifier (Logistic Regression)
25766 clf = LogisticRegression()
25867 clf.fit(X_train_lda, y_train)
25968
26069 # Make predictions
26170 y_pred = clf.predict(X_test_lda)
26271
26372 # Evaluate accuracy
26473 accuracy = accuracy_score(y_test, y_pred)
26574 print(f"Model Accuracy: {accuracy:.2f}")
26675
26776 print(classification_report(y_test, y_pred))

```

Listing 1: Code for LDA

```

268 1 import pandas as pd
269 2 import numpy as np
270 3 from sklearn.model_selection import train_test_split
271 4 from sklearn.discriminant_analysis import
272     QuadraticDiscriminantAnalysis
273 5 from sklearn.metrics import accuracy_score
274 6 from sklearn.metrics import classification_report
275 7
276 8 df = pd.read_csv('training_data_vt2025.csv')
277 9
27810 # modify the month to represent the periodicity that is observed in
279     data.
28011 df['month_cos'] = np.cos(df['month']*2*np.pi/12)
28112 df['month_sin'] = np.sin(df['month']*2*np.pi/12)
28213
28314 # time of day, replaced with 3 bool values: is_night, is_day and
284     is_evening,
28515 # adding the new categories back in the end.
28616 def categorize_demand(hour):
28717     if 20 <= hour or 7 >= hour:
28818         return 'night'
28919     elif 8 <= hour <= 14:
29020         return 'day'
29121     elif 15 <= hour <= 19:
29222         return 'evening'
29323
29424 df['time_of_day'] = df['hour_of_day'].apply(categorize_demand)
29525 df_dummies = pd.get_dummies(df['time_of_day'], prefix='is', drop_first
296     =False)
29726 df = pd.concat([df, df_dummies], axis=1)
29827
29928 # Create bool of snowdepth and percipitation
30029 df['snowdepth_bool'] = df['snowdepth'].where(df['snowdepth'] == 0, 1)
30130 df['precip_bool'] = df['precip'].where(df['precip'] == 0, 1)
30231
30332 # Seperate training data from target
30433 X=df[['#holiday',
30534     'weekday',
30635     '#summertime',
30736     'temp',
30837     '#dew',
30938     '#humidity',
31039     '#visibility',
31140     '#windspeed',
31241     '#month',

```

```

31342         'month_cos',
31443         'month_sin',
31544         #'hour_of_day',
31645         'is_day',
31746         'is_evening',
31847         'is_night',
31948         #'snowdepth_bool',
32049         'precip_bool'
32150     ]
32251
32352 y=df['increase_stock']
32453
32554 # Split dataset into training and test sets
32655 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size
327         =0.2, random_state=42)
32856
32957 # Apply Quadratic Discriminant Analysis (QDA)
33058 qda = QuadraticDiscriminantAnalysis()
33159 X_train_lda = qda.fit(X_train, y_train)
33260
33361 # Make predictions
33462 y_pred = qda.predict(X_test)
33563
33664 # Evaluate accuracy
33765 accuracy = accuracy_score(y_test, y_pred)
33866 print(f"Model Accuracy: {accuracy:.2f}")
33967
34068 print(classification_report(y_test, y_pred))

```

Listing 2: Code for QDA