
Do we need more bikes? Project in Statistical Machine Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 In this project, we aim to create a classification model that can predict whether
2 an increase in the number of bikes is needed in Washington D.C. on a specific
3 temporal and meteorological conditions. We will compare various classification
4 models which consists of Logistic Regression, Linear Discriminant Analysis (LDA),
5 K-Nearest Neighbors (KNN), Random Forest, and Gradient Boosting. The models
6 are evaluated using Accuracy, Precision, Recall, and F1-Score metrics with 10-Fold
7 Cross Validation. The results indicate that the Random Forest model outperforms
8 the other models in all evaluation metrics, with an accuracy of 94.44%, recall of
9 94.99%, precision of 94.00%, and F1-score of 94.47%. This suggests that Random
10 Forest is the most effective model for predicting bike demand in this context.
11 Number of group member: 4

12 **1 Problem Description**

13 Capital Bikeshare is a 24-hour public bicycle-sharing system that serves Washington, D.C., and offers
14 transportation for thousands of people throughout the city. The problem that arises is that there are
15 certain occasions when, due to various circumstances, there are not as many bikes available as there
16 are demands. In the long term, this situation will result in more people taking the car instead of
17 the bicycle, increasing CO2 emissions in the city. To tackle this situation, the District Department
18 of Transportation in the city wants to know if at certain hours an increase in the number of bikes
19 available will be necessary.

20 In this Project, we aim to analyze whether the increase in the number of bikes is necessary or not
21 based on the various temporal and meteorological data provided in the dataset.

22 **2 Data Analysis**

23 The Training Dataset training.csv consists of 1600 randomly selected observation over the period
24 of three years in the city of Washington D.C. The dataset contains 16 features and 1 target variable.
25 The features are: *hour_of_day*, *day_of_week*, *month*, *holiday*, *weekday*, *summertime*, *temp*, *dew*,
26 *humidity*, *precip*, *snow*, *snow_depth*, *windspeed*, *cloudcover*, and *visibility*. And the target variable is
27 *increase_stock*.

28 **2.1 Variable Types and Processing**

29 The target variable *increase_stock* indicates whether an increase in the number of bikes is needed
30 at a particular hour. the value '*low_bike_demand*' indicates that no increase is needed, while
31 '*high_bike_demand*' indicates that an increase is necessary. For the analysis, we will convert these
32 categorical values into binary numerical values, where '*low_bike_demand*' is represented as 0 and

33 'high_bike_demand' as 1. Since the target variable is binary, this problem can be treated as a binary
34 classification task.

35 For Binary features such as *holiday*, *weekday*, and *summertime*, they will be counted as categorical
36 variables with values 0 and 1.

37 The Features *temp*, *dew*, *humidity*, *precip*, *snow*, *snow_depth*, *windspeed*, *cloudcover*, and *visibility*
38 will be treated as numerical variables as they represent continuous measurements.

39 Regarding Ordinal features such as *hour_of_day*, *day_of_week*, and *month*, special attention is needed.

40 For these features, we can't simply treat them as numerical values due to their cyclical nature. For
41 example, after hour 23 comes hour 0 again. However, treating them as categorical variables may lead
42 to loss of information regarding their order and cyclical patterns. Therefore, there are several possible
43 approaches to handle them:

44 • One-Hot Encoding: Convert each of these features into multiple binary features, each
45 representing a specific category. For example, *hour_of_day* would be converted into 24
46 binary features.

47 • Cyclical Transformation: Transform these features using sine and cosine functions to capture
48 their cyclical nature. For example, for *hour_of_day*, we can create two new features:

$$\text{hour_sin} = \sin\left(2\pi \cdot \frac{\text{hour_of_day}}{24}\right) [4]$$

49

$$\text{hour_cos} = \cos\left(2\pi \cdot \frac{\text{hour_of_day}}{24}\right) [4]$$

50 Since One-Hot Encoding wouldn't effectively capture the cyclical nature of these features and may
51 lead to high dimensionality, for this analysis, we will use the Cyclical Transformation approach to
52 handle these ordinal features. [4]

53 After processing, the dataset will consist of 18 features and 1 target variable. Those features are
54 shown in Table 1.

Table 1: Processed Features in the Dataset

Feature	Type	Description
hour_sin	Numerical	Sine transformation of hour of the day
hour_cos	Numerical	Cosine transformation of hour of the day
day_sin	Numerical	Sine transformation of day of the week
day_cos	Numerical	Cosine transformation of day of the week
month_sin	Numerical	Sine transformation of month of the year
month_cos	Numerical	Cosine transformation of month of the year
holiday	Binary / Categorical	Whether the day is a holiday or not (0 or 1)
weekday	Binary / Categorical	Whether the day is a weekday or not (0 or 1)
summertime	Binary / Categorical	Whether the day is in the summer time period or not (0 or 1)
temp	Numerical	Temperature in Celsius
dew	Numerical	Dew point temperature in Celsius
humidity	Numerical	Relative Humidity in percentage
precip	Numerical	Precipitation in mm
snow	Numerical	Amount of snow in the last hour in cm
snow_depth	Numerical	Accumulated snow depth in cm
windspeed	Numerical	Wind speed in km/h
cloudcover	Numerical	Percentage of cloud cover
visibility	Numerical	Distance in km at which objects or landmarks can be clearly seen and identified
increase_stock (Target)	Binary / Categorical	Whether an increase in bike stock is needed (0 or 1)

55 **2.2 Exploratory Data Analysis**

- 56 For the initial stage, we will perform Exploratory Data Analysis (EDA) to understand the distribution
57 and trends that arises in the dataset. Including which features are more correlated with the target
58 variable *increase_stock*.
- 59 The feature *snow* only contains zero values in all observations, therefore it will be removed from the
60 dataset as it doesn't provide any useful information for the analysis. Upon analyzing the dataset, we
61 found that there are no missing values in any of the features or the target variable. Therefore, no
62 handling is required.
- 63 Using Pearson correlation coefficient, we found correlation values between each feature and the target
64 variable as shown in Table 2.

Table 2: Ordered Correlation between Features and Target Variable

Feature	Correlation Coefficient
hour_of_day_cos	-0.339960
temp	0.336981
humidity	-0.308726
hour_of_day_sin	-0.308121
summertime	0.216052
month_cos	-0.169059
dew	0.132663
weekday	-0.116446
visibility	0.113443
windspeed	0.096011
month_sin	-0.092078
day_of_week_sin	-0.088152
precip	-0.059304
snowdepth	-0.047526
cloudcover	-0.045534
day_of_week_cos	-0.031473
holiday	-0.004909

- 65 As the table suggests, the feature *hour_of_day_cos* has the highest positive correlation with the target
66 variable *increase_stock*, indicating that the time of day plays a significant role in determining whether
67 an increase in bike stock is needed. On the other hand, the feature *holiday* has the lowest correlation
68 with the target variable, suggesting that whether a day is a holiday or not has minimal impact on bike
69 demand.

70 **2.3 Imbalance in the Dataset**

- 71 Upon analyzing the target variable *increase_stock*, we found that there is 1312 instances of class 0
72 (low bike demand) and 288 instances of class 1 (high bike demand). This indicates a significant class
73 imbalance in the dataset [2], with class 0 being the majority class.
- 74 To address this class imbalance, we will employ the use of Synthetic Minority Over-sampling
75 Technique (SMOTE). SMOTE works by generating synthetic samples for the minority class (class 1
76 in this case) by interpolating between existing minority class instances. This helps to balance the
77 class distribution and provides the model with more representative samples of the minority class
78 during training. [1]
- 79 The interpolation is done by selecting a minority class instance and finding its k-nearest neighbors.
80 A synthetic sample is then created by randomly selecting one of the neighbors and interpolating
81 between the two instances. This process is repeated until the desired balance between the classes is
82 achieved.

83 **3 Models and Methods**

84 In this experiment, we will compare various classification models to determine which one performs
85 best for predicting whether an increase in bike stock is needed using the provided dataset. The models
86 we will consider includes Logistic Regression, Linear Discriminant Analysis (LDA), K-Nearest
87 Neighbors (KNN), Random Forest, and Gradient Boosting.

88 For each model, we will perform hyperparameter tuning using techniques such as Grid Search or
89 Random Search combined with cross-validation to find the optimal set of hyperparameters that yield
90 the best performance on the validation set.

91 **3.1 Benchmark Model**

92 As the benchmark model, we will use a naive model that predict each instance using stratified random
93 sampling based on the training set's class distribution. This means that for each instance, the model
94 will randomly assign a class label (0 or 1) based on the proportion of each class in the training data
95 [3]. This will provide a baseline accuracy to compare the performance of more sophisticated models.

96 **3.2 Evaluation Metrics**

97 To evaluate the performance of each classification model, we will use several metrics including
98 Accuracy, Precision, Recall, and F1-Score. These metrics will provide a comprehensive understanding
99 of how well each model performs in predicting the target variable. [2]

100 We will also use K-Fold Cross Validation to ensure that our evaluation metrics are robust and not
101 overly dependent on a particular train-test split.

102 K-Fold Cross Validation involves dividing the dataset into K subsets, using one of the subsets as the
103 test set and the remaining K-1 subsets as the training set. This process is repeated K times, with each
104 subset used as the test set once. The final evaluation metrics are then averaged over all K iterations to
105 provide a more reliable estimate of model performance. For this experiment, we will use K=10. [2]

106 **3.3 Logistic Regression**

107 Logistic Regression is a linear model used for binary classification tasks. It models the probability of
108 the target variable being in a particular class using the logistic function. The model estimates the
109 coefficients for each feature, which represent the impact of each feature on the log-odds of the target
110 variable [2].

111 In Logistic Regression, we use the sigmoid function to map the linear combination of features to a
112 probability value between 0 and 1. With the formulation:

$$P(Y = 1|X) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n)}} [3]$$

113 Where $P(Y = 1|X)$ is the probability of the target variable being in class 1 given the features X , β_0
114 is the intercept, and $\beta_1, \beta_2, \dots, \beta_n$ are the coefficients for each feature X_1, X_2, \dots, X_n .

115 The parameters of the model are estimated using Maximum Likelihood Estimation (MLE), which
116 finds the set of coefficients that maximize the likelihood of the observed data given the model.

117 **3.4 Linear Discriminant Analysis (LDA)**

118 Linear Discriminant Analysis (LDA) is a classification method that finds a linear combination of
119 features that best separates the classes. [2].

120 LDA is derived from the probabilistic model which models the class-conditional distributions of the
121 data $P(X|y=k)$ for each class k [3]. Predictions are made by applying Bayes theorem for each training
122 sample $x \in \mathbb{R}^d$:

$$P(y = k|x) = \frac{P(x|y = k)P(y = k)}{P(x)} [3]$$

123 Then we select the class with the highest posterior probability. For LDA, $P(X|y=k)$ is modeled as a
 124 multivariate Gaussian distribution with density function as follows:

$$P(x|y = k) = \frac{1}{(2\pi)^{d/2}|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu_k)^T \Sigma^{-1}(x - \mu_k)\right) [3]$$

125 where μ_k is the mean vector of class k, and Σ is the shared covariance matrix across all classes, and d
 126 is the number of features.

127 In LDA, we assume that the covariance matrices of all classes are equal, i.e., $\Sigma_k = \Sigma$ for all k. This
 128 reduces the log posterior to:

$$\log P(y = k|x) = -\frac{1}{2}(x - \mu_k)^T \Sigma^{-1}(x - \mu_k) + \log P(y = k) + Constant [3]$$

129 which can be simplified to:

$$\log P(y = k|x) = w_k^T x + w_{k0} + Constant [3]$$

130 where $w_k = \Sigma^{-1}\mu_k$ and $w_{k0} = -\frac{1}{2}\mu_k^T \Sigma^{-1}\mu_k + \log P(y = k)$. [3]

131 3.5 K-Nearest Neighbors (KNN)

132 K-Nearest Neighbors (KNN) is a non-parametric classification method that classifies new instances
 133 based on the majority class of their K nearest neighbors in the feature space. The distance metric used
 134 to determine the nearest neighbors can be Euclidean distance, Manhattan distance, or other distance
 135 measures. [3]

136 KNN algorithm are usually used for unsupervised learning tasks such as clustering. However, it can
 137 also be adapted for supervised learning tasks by using the labels of the nearest neighbors to make
 138 predictions.

139 For Supervised learning tasks, KNN works using these following steps:

- 140 • Choose the number of neighbors K.
- 141 • For each new instance to be classified, calculate the distance between the new instance and
 142 all instances in the training dataset.
- 143 • Identify the K nearest neighbors based on the calculated distances.
- 144 • Determine the majority class among the K nearest neighbors.
- 145 • Assign the majority class as the predicted class for the new instance.

146 For the distance metric, we will use Manhattan distance, which is defined as:

$$d(p, q) = \sum_{i=1}^n |p_i - q_i| [3]$$

147 where p and q are two instances in the feature space, and n is the number of features.

148 3.6 Random Forest

149 Random Forest is an ensemble learning method that consists of multiple decision trees. Each tree is
 150 trained on a random subset of the training data and a random subset of features. The final prediction
 151 is made by aggregating the predictions from all trees. The method to aggregate the results in this
 152 experiment will use majority voting. [3]

153 Decision trees in principle work by recursively splitting the data based on feature values to create
 154 branches that lead to leaf nodes representing class labels. The splits are chosen based on the feature
 155 that maximizes the information gain or minimizes the impurity at each node. For the metric to
 156 measure impurity, we will use Entropy, which is defined as:

$$H(X) = - \sum_{i=1}^c p_i \log_2(p_i) [3]$$

157 where p_i is the proportion of instances belonging to class i in the node, and c is the number of classes.

158 Decision trees can be prone to overfitting, especially when they are deep and complex. To mitigate
159 this, we can use Random Forest.

160 Random Forest Algorithm, works using these following steps:

161 • For each tree in the forest:

- Randomly sample the training data with replacement (bootstrap sampling).
- Randomly select a subset of features to consider for splitting at each node.
- Train a decision tree on the sampled data using the selected features.

165 • For making predictions:

- For each new instance, pass it through each tree in the forest to obtain the predicted class.
- Aggregate the predictions from all trees using majority voting to determine the final predicted class.

170 Due to the randomness introduced in the training process, Random Forests are less prone to overfitting
171 compared to individual decision trees and often achieve better generalization performance. [3]

172 3.7 Gradient Boosting

173 Gradient Boosting is an ensemble learning method that builds a series of weak learners in a sequential
174 manner. Each weak learner is trained to correct the errors made by the previous learners. The final
175 prediction is made by combining the predictions from all weak learners. [3]

176 Gradient Boosting Algorithm works using these following steps:

177 • Initialize the model with a constant value, typically the mean of the target variable.

178 • For each iteration $m = 1$ to M :

- Compute the pseudo-residuals, which are the negative gradients of the loss function with respect to the current model's predictions.
- Train a weak learner (e.g., decision tree, logistic regression) on the pseudo-residuals.
- Compute the optimal step size (learning rate) for the weak learner.
- Update the model by adding the weighted predictions of the weak learner to the current model.

185 • For making predictions:

- For each new instance, pass it through all weak learners and sum their weighted predictions to obtain the final predicted value.

188 Since the problem we're working on is a binary classification task, we will use logistic loss as the
189 loss function for Gradient Boosting. The logistic loss is defined as:

$$L(y, \hat{y}) = -\frac{1}{N} \sum_{i=1}^N [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)] [3]$$

190 where y_i is the true label, \hat{y}_i is the predicted probability, and N is the number of instances. For the
191 weak learners, we will use decision trees with a maximum depth of 5.

192 3.8 Hyperparameter Tuning

193 For each classification model, we will perform hyperparameter tuning using Grid Search com-
194 bined with 10-Fold Cross Validation to find the optimal set of hyperparameters that yield the best
195 performance on the validation set. The hyperparameters to be tuned for each model are as follows:

196 Grid Search works by exhaustively searching through a specified subset of hyperparameters for
197 each model. For each combination of hyperparameters, the model is trained and evaluated using
198 cross-validation. The combination that yields the best average performance across the folds is selected
199 as the optimal set of hyperparameters. [2]

200 **4 Experiment and Results**

201 After performing the experiments using the described models and methods. We obtained the results
202 as shown in Table 3.

Table 3: Results of Classification Models (rounded to 4 decimal places)

Model	Accuracy	Precision	Recall	F1-Score
Random Forest	0.9037	0.9065	0.9037	0.9043
Gradient Boosting	0.8994	0.9015	0.8994	0.8999
Logistic Regression	0.8100	0.8696	0.8100	0.8271
LDA	0.7956	0.8676	0.7956	0.8155
K-Nearest Neighbors	0.7675	0.8233	0.7675	0.7863
Benchmark Model	0.4781	0.6889	0.4781	0.5375

203 Using Grid Search with 10-Fold Cross Validation, we found the optimal hyperparameters for each
204 model is as follows:

- **Logistic Regression:** Regularization strength $C = 1.0$, Solver type using Liblinear, and Penalty type using Ridge Regularization (L2).
- **LDA:** Solver Type using SVD (Singular Value Decomposition)
- **K-Nearest Neighbors:** Number of neighbors $K = 3$, Distance metric using Manhattan distance, and Weighting type using the inverse of their distance.
- **Random Forest:** Number of trees $n_estimators = 250$, Maximum depth of each tree $max_depth = 15$, Minimum samples per leaf $min_samples_leaf = 1$, Minimum samples per split $min_samples_split = 2$, and Criterion using Entropy.
- **Gradient Boosting:** Number of estimators $n_estimators = 100$, Learning rate = 0.1, Maximum Depth for each estimator = 10, and Subsample = 0.6

215 Based on the results, we can see that the more complex models such as Random Forest and Gradient
216 Boosting outperformed the simpler models like Logistic Regression, LDA, and KNN.

217 Random Forest achieved an accuracy of 90.37%, indicating that it correctly classified a high percent-
218 age of instances in the dataset. The F1-Score of 90.43% also suggests that the model has a good
219 balance between precision and recall, suggesting that it is effective in identifying both positive and
220 negative instances. Gradient Boosting also performed comparably well, achieving an accuracy and
221 F1-Score of 89.94% and 89.99% respectively, which is slightly lower than Random Forest but still
222 significantly better than the other models.

223 If we compare the simpler models such as Logistic Regression, LDA, and KNN, we can see that
224 they achieved lower accuracy and F1-scores. Logistic Regression achieved an accuracy of 81.00%,
225 LDA achieved 79.56%, and KNN achieved 76.75%. This indicates that these models may not be as
226 effective in capturing the complex relationships in the dataset compared to the ensemble methods.

227 Due to its superior performance, for production deployment, we recommend using either the Random
228 Forest or Gradient Boosting model for predicting whether an increase in bike stock is needed.

229 **5 Conclusion**

230 Based on the analysis and experiments conducted in this project, we can conclude that the Random
231 Forest is the most effective model for predicting whether an increase in bike stock is needed in
232 Washington D.C. The model achieved the highest accuracy, recall, precision, and F1-score among all
233 the models evaluated.

234 **References**

- 235 [1] G. Lemaître, F. Nogueira, and C. K. Aridas. Imbalanced-learn: A python toolbox to tackle the
236 curse of imbalanced datasets in machine learning. *Journal of Machine Learning Research*, 18
237 (17):1–5, 2017. URL <http://jmlr.org/papers/v18/16-365.html>.
- 238 [2] A. Lindholm, N. Wahlström, F. Lindsten, and T. B. Schön. *Machine Learning - A First Course
239 for Engineers and Scientists*. Cambridge University Press, 2022. URL <https://smlbook.org>.
- 240 [3] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Pret-
241 tenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot,
242 and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning
Research*, 12:2825–2830, 2011.
- 243 [4] D. Radečić. How to handle cyclical data in machine
244 learning, 2020. URL [https://towardsdatascience.com/
245 how-to-handle-cyclical-data-in-machine-learning-3e0336f7f97c](https://towardsdatascience.com/how-to-handle-cyclical-data-in-machine-learning-3e0336f7f97c).

247 A Appendix

```
248 1 # %% [markdown]
249 2 # # Appendix A
250 3
251 4 # %%
252 5 import pandas as pd
253 6 import pycaret
254 7 import numpy as np
255 8 import copy
256 9
257 10 # %% [markdown]
258 11 # ## Read Raw Data
259 12
260 13 # %%
261 14 df_raw = pd.read_csv('training_data_ht2025.csv')
262 15
263 16 # %%
264 17 df_raw
265 18
266 19 # %% [markdown]
267 20 # ### Check Data Contents
268 21
269 22 # %%
270 23 df_raw['increase_stock'].value_counts()
271 24
272 25 # %% [markdown]
273 26 # Based on above code that shows the values on the target variable *increase_stock*, there's a problem of Imbalance Dataset, where the records with low bike demand are much more compared to records with high bike demands.
274 27
275 28 # %%
276 29 df_raw['holiday'].value_counts()
277 30
278 31 # %%
279 32 df_raw['weekday'].value_counts()
280 33
281 34 # %%
282 35 df_raw['summertime'].value_counts()
283 36
284 37 # %%
285 38 df_raw['increase_stock'].value_counts()
286 39
287 40 # %%
288 41 df_raw['hour_of_day'].value_counts()
289 42
290 43 # %%
291 44 df_raw['day_of_week'].value_counts()
292 45
293 46 # %%
294 47 df_raw['month'].value_counts()
295 48
296 49 # %%
297 50 df_raw['snow'].value_counts()
298 51
299 52 # %% [markdown]
300 53 # Note: The Snow Features only have one value: 0, and therefore isn't
301 54     really useful for this dataset. And In Feature selection we can
302 55     drop them.
303 56
304 57 # %%
305 58 df_raw.isnull().sum()
306 59
307 60 # %% [markdown]
```

```

31159 # Based on the checking results, there're no NaN values, so we can
312     proceed with feature processing.
313
3140
3141 # %% [markdown]
3152 # ##### Cyclical Encoding
3163
3174 # %% [markdown]
3185 # Function to encode Ordinal Variables into Cyclical encoding using
319     Sine and Cosine
320
3217 # %%
3228 def cyclical_encode(x, max_val, start_val=0):
3239     x_arr = np.asarray(x, dtype=float)
3240
3251     # If values are 1..period (e.g. months 1..12), shift to 0..period
326     -1
3272     if not start_val:
3283         x_arr = x_arr - 1
3294
3305     angle = 2 * np.pi * x_arr / max_val
3316     sin_x = np.sin(angle)
3327     cos_x = np.cos(angle)
3338
3349     if np.isscalar(x):
3350         return float(sin_x), float(cos_x)
3361
3372     return sin_x, cos_x
3383
3394 # %%
3405
3416
3427 # %%
3438 # Transform Cyclical Values for hour_of_day
3449
3450 # Transform Cyclical Values for hour_of_day
3461 # Hour of day ranges from 0 to 23
3472 df_hour_feat = pd.DataFrame(df_raw['hour_of_day'].apply(lambda x:
348     cyclical_encode(x, max(df_raw['hour_of_day']), min(df_raw['
349         hour_of_day']))).to_list())
3503 df_hour_feat.columns = ['hour_of_day_sin', 'hour_of_day_cos']
3514 df_hour_feat
3525
3536
3547 # Transform Cyclical Values for day_of_week
3558 # Day of week ranges from 0 (Monday) to 6 (Sunday)
3569 df_day_feat = pd.DataFrame(df_raw['day_of_week'].apply(lambda x:
357     cyclical_encode(x, 6, 0)).to_list())
3580 df_day_feat.columns = ['day_of_week_sin', 'day_of_week_cos']
3591 df_day_feat
3602
36103 # Transform Cyclical Values for month
36214 # Month ranges from 1 to 12
36315 df_month_feat = pd.DataFrame(df_raw['month'].apply(lambda x:
364     cyclical_encode(x, 12, 1)).to_list())
36516 df_month_feat.columns = ['month_sin', 'month_cos']
36617 df_month_feat
36718
36819 # Concat Cyclical Features
36910 df_cyclical = pd.concat([df_hour_feat, df_day_feat, df_month_feat],
370     axis=1)
37111 df_cyclical
37212
37313 # %%
37414 # Creating Final Feature and Target DataFrames
37515 df_features = copy.deepcopy(df_raw)

```

```

376|6
377|7 # Creating Target Variable and Mapping text to binary
378|8 df_target = df_features['increase_stock'].map({'low_bike_demand': 0, ,
379|9     'high_bike_demand': 1})
380|9
381|0 # Dropping original cyclical columns and target from features
382|1 df_features = df_features.drop(['hour_of_day', 'day_of_week', 'month',
383|2     'increase_stock'], axis=1)
384|2
385|3 # Concatenating Cyclical Features to Features DataFrame
386|4 df_features = pd.concat([df_cyclical, df_features], axis=1)
387|5
388|6 df_features
389|7
390|8 # %%
391|9 df_target
392|0
393|1 # %%
394|2 # Correlation Analysis between Features and Target
395|3 df_corr = pd.concat([df_features, df_target], axis=1).corr()
396|4 df_corr['increase_stock'].sort_values(key=abs, ascending=False)
397|5
398|6 # %% [markdown]
399|7 # Based on above result, The variables hour_of_day_cos, temp, humidity
400|8     , hour_of_day_sin, and summertime held the biggest correlation to
401|9     the target variable. Because of this, we could check the trends
402|0     more based on those attributes
403|1
404|2 # %%
405|3 from imblearn.over_sampling import SMOTE
406|4
407|5 # Applying SMOTE to Balance the Dataset
408|6 smote = SMOTE(sampling_strategy='minority')
409|7
410|8 print("Before SMOTE:\n", df_target.value_counts())
411|9
412|0 print()
413|1
414|2 X_resampled, y_resampled = smote.fit_resample(df_features, df_target)
415|3
416|4 print("After SMOTE:\n", y_resampled.value_counts())
417|5
418|6 # Note: Don't use SMOTE outside training data to avoid data leakage,
419|7     this block is just for experimentation purposes. In practice,
420|8     SMOTE should only be applied to the training set within each cross
421|9     -validation folds.
422|0
423|1
424|2 # %% [markdown]
425|3 # ## Model Experiment
426|4
427|5 # %%
428|6 from sklearn.metrics import accuracy_score, precision_score,
429|7     recall_score, f1_score, roc_auc_score, classification_report
430|8 from sklearn.model_selection import KFold, train_test_split
431|9 from sklearn.model_selection import GridSearchCV
432|0
433|1
434|2 # %% [markdown]
435|3 # #### K-Fold Classification Report
436|4
437|5 # %%
438|6 ## Function to do K-Fold Cross Validation and return the
439|7     classification report
440|8 def kfold_classification_reports(clf, X, y, n_splits=10, shuffle=True,
441|9     random_state=42):

```

```

44170 kf = KFold(n_splits=n_splits, shuffle=shuffle, random_state=
442 random_state)
44371 reports = []
44472
44573 for train_index, test_index in kf.split(X):
44674     # Applying SMOTE to Balance the Dataset
44775     smote = SMOTE(sampling_strategy='minority')
44876
44977
45078     X_train, X_test = X.iloc[train_index], X.iloc[test_index]
45179     y_train, y_test = y.iloc[train_index], y.iloc[test_index]
45280
45381     X_resampled, y_resampled = smote.fit_resample(X_train, y_train)
45482
45583     clf.fit(X_resampled, y_resampled)
45684     y_pred = clf.predict(X_test)
45785
45886     report = classification_report(y_test, y_pred, output_dict=
45987     True)
46088     reports += [report]
46289
46389 detailed = []
46490 for i, rep in enumerate(reports):
46591     fold_result = {
46692         "fold": i + 1,
46793         "accuracy": rep["accuracy"],
46894         "precision": rep["weighted avg"]["precision"],
46995         "recall": rep["weighted avg"]["recall"],
47096         "f1": rep["weighted avg"]["f1-score"],
47197     }
47298     detailed.append(fold_result)
47399
47400 # compute averaged (generalized) metrics
47591 accuracies = [d["accuracy"] for d in detailed]
47692 precisions = [d["precision"] for d in detailed]
47793 recalls = [d["recall"] for d in detailed]
47894 f1s = [d["f1"] for d in detailed]
47995
48096 report_dict = {
48197     "general": {
48298         "accuracy_mean": float(np.mean(accuracies)),
48399         "accuracy_std": float(np.std(accuracies)),
48410         "precision_mean": float(np.mean(precisions)),
48511         "precision_std": float(np.std(precisions)),
48612         "recall_mean": float(np.mean(recalls)),
48713         "recall_std": float(np.std(recalls)),
48814         "f1_mean": float(np.mean(f1s)),
48915         "f1_std": float(np.std(f1s)),
49016     },
49117     "detailed": detailed
49218 }
49399
49420     return report_dict
49599
49622 # %% [markdown]
49723 # ## Benchmark Model
49899
49925 # %% [markdown]
50026 # For the benchmark model, we use a naive model that predict each
501 instance as the majority class in the training dataset. This will
502 provide a baseline accuracy to compare the performance of more
503 sophisticated models.
50427
50528 # %%

```

```

50629 from sklearn.dummy import DummyClassifier
50730
50831 clf_dummy = DummyClassifier(strategy="stratified")
50932
51033 # %%
51134 reports = kfold_classification_reports(clf_dummy, df_features,
512     df_target)
51335 reports['general']
51436
51537 # %% [markdown]
51638 # ## Logistic Regression
51739
51840 # %% [markdown]
51941 # For Logistic Regression, we will do hyperparameter tuning on
520     Regularization Strength (*C*), Penalty type (*penalty*), and
521     Solver type (*solver*)
52242
52343 # %%
52444
52545
52646 # %%
52747
52848 from sklearn.linear_model import LogisticRegression
52949
53050 clf = LogisticRegression()
53151
53252 # Grid Search to find the best hyperparameters for Logistic Regression
53353 grid_search = GridSearchCV(
53454     estimator=clf,
5355     param_grid={
53656         'C': [0.01, 0.1, 1, 10, 100],
53757         'penalty': ['l1', 'l2'],
53858         'solver': ['liblinear', 'lbfgs', 'saga', 'newton-cg']
53959     },
54060     scoring='f1',
54161     cv=10,
54262     n_jobs=-1,
54363     verbose=1
54464 )
54565 clf_grid_lg = grid_search.fit(df_features, df_target)
54666
54767 # %%
54868 clf_grid_lg.best_params_
54969
55070 # %%
55171 clf_lg = LogisticRegression(**clf_grid_lg.best_params_)
55272
55373 # %%
55474 reports = kfold_classification_reports(clf_lg, df_features, df_target)
55575 reports['general']
55676
55777 # %% [markdown]
55878 # ## Linear Discriminant Analysis (LDA)
55979
56080 # %% [markdown]
56181 # For LDA, we us Hyperparameter tuning for Solver Type (*solver*)
56282
56383 # %%
56484 from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
56585
56686 clf = LinearDiscriminantAnalysis()
56787
56888 # Grid Search to find the best hyperparameters for LDA
56989 grid_search = GridSearchCV(
57089     estimator=clf,

```

```

57191     param_grid={
57192         'solver': ['svd', 'lsqr', 'eigen']
57193     },
57194     scoring='f1',
57195     cv=10,
57196     n_jobs=-1,
57197     verbose=1
57198 )
57199 clf_grid_lda = grid_search.fit(df_features, df_target)
5800
5810
5811 # %%
5812 clf_grid_lda.best_params_
5813
5814 # %%
5815 clf_lda = LinearDiscriminantAnalysis(**clf_grid_lda.best_params_)
5816 clf_lda
5817
5818 # %%
5819 reports = kfold_classification_reports(clf_lda, df_features, df_target
5820     )
5821 reports['general']
5822
5823 # %% [markdown]
5824 # ## K Nearest Neighbor (KNN)
5825
5826 # %% [markdown]
5827 # For K Nearest Neighbor, we use Hyperparameter tuning for selecting
5828 # the numer of neighbors (*n_neighbors*), weight function used in
5829 # prediction (*weights*), and the metric used for distance
5830 # computation (*metric*)
5831
5832 # %%
5833 from sklearn.neighbors import KNeighborsClassifier
5834
5835 clf = KNeighborsClassifier()
5836
5837 #Grid Search to find the best hyperparameters for KNN
5838 grid_search = GridSearchCV(
5839     estimator=clf,
5840     param_grid={
5841         'n_neighbors': range(3, 21, 2),
5842         'weights': ['uniform', 'distance'],
5843         'metric': ['euclidean', 'manhattan', 'minkowski']
5844     },
5845     scoring='f1',
5846     cv=10,
5847     n_jobs=-1,
5848     verbose=1
5849 )
5850 clf_grid_knn = grid_search.fit(df_features, df_target)
5851
5852 # %%
5853 clf_grid_knn.best_params_
5854
5855 # %%
5856 clf_knn = KNeighborsClassifier(**clf_grid_knn.best_params_)
5857 clf_knn
5858
5859 # %%
5860 reports = kfold_classification_reports(clf_knn, df_features, df_target
5861     )
5862 reports['general']
5863
5864 # %% [markdown]

```

```

636:1 # ## Random Forest Classifier
637:2
638:3 # %%
639:4 # For Random Forest Classifier, we use Hyperparameter tuning for
640:5     selecting number of trees (*n_estimators*), maximum depth of each
641:6     trees (*max_depth*), minimum number of samples required to split
642:7     (*min_samples_split*), and minimum samples required to be a leaf
643:8     node (*min_samples_leaf*)
644:9
645:10 # %%
646:11 from sklearn.ensemble import RandomForestClassifier
647:12 clf = RandomForestClassifier()
648:13 grid_search = GridSearchCV(
649:14     estimator=clf,
650:15     param_grid={
651:16         'n_estimators': [20, 50, 100, 200, 250],
652:17         'max_depth': [3, 5, 8, 10, 15, 20],
653:18         'min_samples_split': [2, 5, 10],
654:19         'min_samples_leaf': [1, 2, 4],
655:20         'criterion': ['gini', 'entropy']
656:21     },
657:22     scoring='f1',
658:23     cv=10,
659:24     n_jobs=-1,
660:25     verbose=1
661:26 )
662:27 clf_grid_rf = grid_search.fit(df_features, df_target)
663:28
664:29 # %%
665:30 clf_grid_rf.best_params_
666:31
667:32 # %%
668:33 clf_rf = RandomForestClassifier(**clf_grid_rf.best_params_)
669:34 clf_rf
670:35
671:36 # %%
672:37 reports = kfold_classification_reports(clf_rf, df_features, df_target)
673:38 reports['general']
674:39
675:40 # %% [markdown]
676:41 # ## Gradient Boosting
677:42
678:43 # %% [markdown]
679:44 # For Gradient Boosting, we use Hyperparameter tuning to determine,
680:45     the number of sequential estimator (*n_estimators*), the learning
681:46     rate (*learning_rate*), Maximum depth of the individual regression
682:47     estimators (*max_depth*), and the fraction of samples to be used
683:48     for fitting the individual estimators (*subsample*).
684:49
685:50 # %%
686:51 from sklearn.ensemble import GradientBoostingClassifier
687:52 clf = GradientBoostingClassifier()
688:53 grid_search = GridSearchCV(
689:54     estimator=clf,
690:55     param_grid={
691:56         'n_estimators': [50, 100, 200],
692:57         'learning_rate': [0.01, 0.05, 0.1, 0.2, 0.5],
693:58         'max_depth': [3, 5, 8, 10],
694:59         'subsample': [0.6, 0.8, 1.0]
695:60     },
696:61     scoring='f1',
697:62     cv=10,
698:63     n_jobs=-1,
699:64     verbose=1
700:65 )

```

```

70108 clf_grid_gb = grid_search.fit(df_features, df_target)
70109 # %%
70110 clf_grid_gb.best_params_
70111
70112 # %%
70113 clf_gb = GradientBoostingClassifier(**clf_grid_gb.best_params_)
70114 clf_gb
70115
70116 # %%
70117 reports = kfold_classification_reports(clf_gb, df_features, df_target)
70118 reports['general']
70119
70120 # %% [markdown]
70121 # ### Saving Model
70122
70123 # %% [markdown]
70124 # Saving all model so we can use it later
70125
70126 # %%
70127 # Applying SMOTE to Balance the Full Training Dataset
70128 smote = SMOTE(sampling_strategy='minority')
70129
70130 print("Before SMOTE:\n", df_target.value_counts())
70131
70132 print()
70133
70134 X_resampled, y_resampled = smote.fit_resample(df_features, df_target)
70135
70136 print("After SMOTE:\n", y_resampled.value_counts())
70137
70138 # %%
70139 # import pickle
70140
70141 # Save the Dummy Classifier model used as benchmark
70142 with open('clf_dummy.pkl', 'wb') as f:
70143     pickle.dump(clf_dummy, f)
70144
70145 # %%
70146 # Train the Logistic Regression model on the full training dataset
70147 clf_lg = LogisticRegression(**clf_grid_lg.best_params_)
70148 clf_lg.fit(X_resampled, y_resampled)
70149
70150 # Save the trained Logistic Regression model
70151 with open('clf_lg.pkl', 'wb') as f:
70152     pickle.dump(clf_lg, f)
70153
70154 # %%
70155 # Train the LDA model on the full training dataset
70156 clf_lda = LinearDiscriminantAnalysis(**clf_grid_lda.best_params_)
70157 clf_lda.fit(X_resampled, y_resampled)
70158
70159 # Save the trained LDA model
70160 with open('clf_lda.pkl', 'wb') as f:
70161     pickle.dump(clf_lda, f)
70162
70163 # %%
70164 # Train the KNN model on the full training dataset
70165 clf_knn = KNeighborsClassifier(**clf_grid_knn.best_params_)
70166 clf_knn.fit(X_resampled, y_resampled)
70167
70168 # Save the trained KNN model
70169 with open('clf_knn.pkl', 'wb') as f:

```

```
766:3     pickle.dump(clf_knn, f)
767:4 # %%
768:5 # Train the Random Forest Classifier model on the full training
769:6     dataset
770:7 clf_rf = RandomForestClassifier(**clf_grid_rf.best_params_)
771:8 clf_rf.fit(X_resampled, y_resampled)
772:9
773:10 # Save the trained Random Forest Classifier model
774:11 with open('clf_rf.pkl', 'wb') as f:
775:12     pickle.dump(clf_rf, f)
776:13
777:14 # %%
778:15 # Train the Gradient Boosting Classifier model on the full training
779:16     dataset
780:17 clf_gb = GradientBoostingClassifier(**clf_grid_gb.best_params_)
781:18 clf_gb.fit(X_resampled, y_resampled)
782:19
783:20 # Save the trained Gradient Boosting Classifier model
784:21 with open('clf_gb.pkl', 'wb') as f:
785:22     pickle.dump(clf_gb, f)
786:23
787:24 # %%
```