
Do we need more bikes?

Project in Statistical Machine Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

In this project, four statistical machine learning classification methods were trained and evaluated using a dataset containing hourly bike rental data from the District Department of Transportation in Washington, D.C. The dataset includes various temporal and meteorological features, such as time of day, weather conditions and temperature, to determine factors influencing bike demand. The classification models implemented included Logistic regression, Linear discriminant analysis (LDA), K-nearest neighbours (k-NN) and Random forest. The objective was to identify the model that most accurately predicts whether an increase in the number of bikes is necessary based on the given features. The models were assessed based on accuracy, precision, recall and F1-score. Hyperparameter tuning and cross-validation were performed to optimize performance. Among the tested models, Random forest demonstrated the best predictive capability by performing best on key evaluation metrics, balancing precision and recall effectively. The findings highlight the potential of machine learning models to support demand forecasting in bike-sharing systems. Number of group member: 4

1 Introduction

Classification is a machine learning technique widely used in areas such as anomaly detection. It automatically categorises data into predefined classes based on various input features and can be applied to demand forecasting, such as predicting the need for additional bikes in a bike-sharing system. Ensuring sufficient bike availability is crucial as shortages could discourage users and lead to increased car usage and carbon emissions. The District Department of Transportation in Washington, D.C. wants to address this with a predictive model to determine when more bikes are needed. A dataset of 1600 observations will be preprocessed and used to train four classification models: Logistic regression, Linear discriminant analysis (LDA), K-nearest neighbours (k-NN) and Random forest. The best model will be selected based on precision, recall, F1-score and accuracy.

2 Data analysis

Prior to model development, an exploratory data analysis (EDA) was performed to gain insights into the data. The EDA was meant to answer questions on which features are numerical and categorical as well as examine greater trends between bike availability and time and weather features. See A for the data analysis code.

The categorical features are hour_of_day, day_of_week, month, holiday, weekday, summertime, and increase_stock. The rest are numerical.

The trends when comparing different hours, weeks, and months can be seen in Figure 1.

34 For the hours of day, the frequency to increase availability is highest between 15-19, while it is lowest
 35 between 23-06.

36 For the day of week, the need to increase availability is slightly higher during the weekend, especially
 37 on Saturdays.

38 For the months, the need to increase availability is lowest during winter months and peaks in April,
 39 June, September.

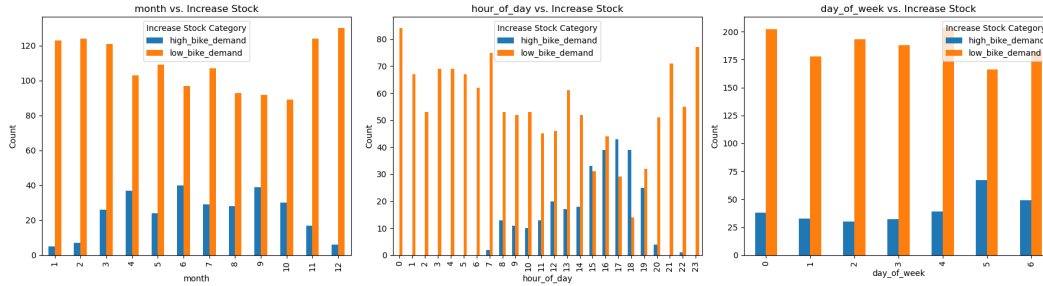


Figure 1: Time features plotted against `increase_stock`.

40 The relationships between the binary holidays and weekdays and the corresponding need to increase
 41 availability are displayed in Table 1 and Table 2, respectively. From Table 1 it seems that there
 42 is no significant difference in availability need based on `holiday`. There does, however, seem to be
 43 higher bike demand when `weekday` is 0, i.e. it is the weekend, as seen in the difference in percentage
 44 in Table 2.

Table 1: Percentage distribution for Holidays

Holiday	Increase Stock	Percentage
0	Low Bike Demand	81.97%
	High Bike Demand	18.03%
1	Low Bike Demand	83.02%
	High Bike Demand	16.98%

Table 2: Percentage distribution for Weekdays

Weekday	Increase Stock	Percentage
0	Low Bike Demand	75.00%
	High Bike Demand	25.00%
1	Low Bike Demand	84.86%
	High Bike Demand	15.14%

45 The weather features trend against `increase_stock` can be seen in Figure 2.

46 For the temperature, the demand for bicycles is greater between 15-25 degrees Celsius. When the
 47 temperature is lower, between 8-22 degrees Celsius, the demand is generally lower.

48 Increased humidity, around 55-82, aligns with a lower demand, lower humidity, between 37-62,
 49 aligns with a higher demand.

50 The value distribution for windspeed, cloudcover, and visibility seem to be similar between high and
 51 low bike demand, indicating a lack of inference ability.

52 When the demand is high there is generally low or close to no precipitation. For low demand, there is
 53 generally moderate precipitation. However, there are days with no precipitation where the demand is
 54 still low. The same behavior can be seen for snowdepth.

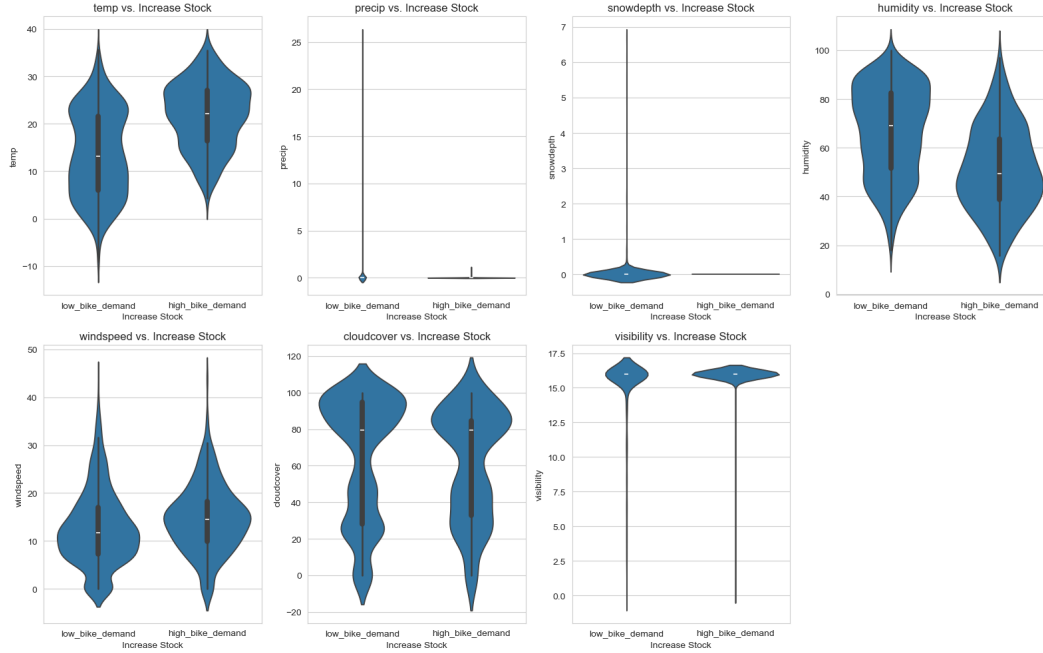


Figure 2: Weather features plotted against increase_stock.

2.1 Pre-processing

Following the EDA, a unified pre-processing was conducted, the code for which can be seen in Appendix B. Due to inconsistencies in the original feature, the pre-processing re-codes the feature `summertime` so that it is set to 1 if the month is between March and November. Furthermore, cyclical encoding of the time variables `day_of_week`, `hour_of_day`, and `month` was conducted through the function `cyclical_encoding` to catch the cyclical nature of the features. `increase_stock`'s values were converted to 1 if `high_bike_demand` and 0 if `low_bike_demand`. Then, the features `precip`, `snowdepth`, and `visibility` were converted to binary values where the values were set to 1 if the values were not their most common values which was 0, 0, and 16 respectively. Lastly, the features `holiday`, `snow` and the original time features were dropped. The pre-processing produced a new csv-file with data that was used for all models.

3 Model development

After the data had been analysed and processed the machine learning models were developed. Four models: logistic regression, linear discriminant analysis, K-nearest neighbours, and random forest and a naive classifier were developed.

First, all model developments used stratified 5-fold cross validation for training and validation because of the small dataset. Stratified cross validation was chosen because it keeps the proportion of classes for each fold which is useful for an imbalanced dataset such as the one used in the project. Further, to allow for hyperparameter tuning the data was initially split for all models into a training and test set. The training set was used in the cross validation and used for hyperparameter tuning while the test set was only used when evaluating the final model. The data was split into 80% training and 20% test data using `random_state=42` between all models to guarantee consistency across models and reproducibility.

3.1 Evaluation metrics

To evaluate model performance accuracy, recall, precision and F1-score were used. F1-score was used as the cross validation metric. Due to the imbalanced dataset, accuracy was not single-handedly used to evaluate performance. Recall was seen as important due to the District Department of

Transportation’s goal of identifying if bike availability should be increased. This makes it slightly more important to not miss cases when there actually is a need to increase availability. However, precision is also important to take into account since the department would not want to increase availability without there actually being demand for it. Thus, the F1-score, which averages precision and recall, is viewed as a useful metric when evaluating performance.

3.2 Naive classification model

A naive classifier was developed to bench mark the performance of the other models. The classifier was designed to only predict low_bike_demand and had an accuracy of roughly 0.82 due to the imbalanced dataset. See Appendix C for the naive classifier code.

3.3 Logistic regression

3.3.1 Mathematical description

Logistic Regression is a supervised binary classification problem that models the probability of an instance belonging to a particular class (Lindholm et al. [1]). The sigmoid function transforms the linear combination of input features into a probability bounded between 0 and 1.

$$z = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_p x_p = \theta^T \mathbf{x} \quad (1)$$

$$p(y = 1|\mathbf{x}) = g(\mathbf{x}) = \frac{e^{\theta^T \mathbf{x}}}{1 + e^{\theta^T \mathbf{x}}} \quad (2)$$

For binary classification, the decision is made based on a set threshold which is typically 0.5.

$$\hat{y} = \begin{cases} 1, & \text{if } g(\mathbf{x}) \geq \tau \\ 0, & \text{if } g(\mathbf{x}) < \tau \end{cases} \quad (3)$$

The model is trained by finding the optimal parameters θ that maximises the likelihood of correctly classifying the training data.

$$\ell(\theta) = \sum_{i=1}^n [y_i \log g(\mathbf{x}_i) + (1 - y_i) \log(1 - g(\mathbf{x}_i))] \quad (4)$$

3.3.2 Method application

Logistic Regression was applied to the dataset after processing per subsection 2.1 and the common modeling choices. See Appendix D for the implementation code. Class imbalance was considered when splitting the data, with a balanced version tested, but the final model without balancing performed the best. The 18% minority class, while imbalanced, was perhaps not low enough (<10%) possibly explaining why balancing did not improve performance (van den Goorbergh et al. [3]). Grid search cross-validation, though computationally expensive, was preferred over random search to obtain the optimal hyperparameters (Sukamto et al. [2]). Regularisation strength (C) was tested at 0.1, 0.5, 1, 5, 10 and 20 along with liblinear and saga solvers and L1 (Lasso) and L2 (Ridge) penalties to prevent overfitting. The final model, optimised for F1-score to handle class imbalance, was trained with C = 20, solver = liblinear and penalty = L2. To further refine the classification decision, different threshold values were tested with 0.4 maximising the F1-score.

3.4 Linear discriminant analysis

3.4.1 Mathematical description

Linear discriminant analysis (LDA) can be expressed as Equation 5. This model follows from the generative model for classification and assuming a prior according to Equation 6, and that the conditional probability density of x for class m is a multivariate normal density with the mean

described by Equation 7 and a covariance matrix, that is further assumed to be same between all classes, described by Equation 8.

$$\delta_m(x) = x^T \Sigma^{-1} \hat{\mu}_m - \frac{1}{2} \hat{\mu}_m^T \Sigma^{-1} \hat{\mu}_m + \log \hat{\pi}_m \quad (5)$$

Where:

$$\hat{\pi}_m = \frac{n_m}{n}, \quad m = 1, \dots, M \quad (6)$$

$$\hat{\mu}_m = \frac{1}{n_m} \sum_{i:y_i=m} x_i, \quad m = 1, \dots, M \quad (7)$$

$$\hat{\Sigma} = \frac{1}{n - M} \sum_{m=1}^M \sum_{i:y_i=m} (x_i - \hat{\mu}_m)(x_i - \hat{\mu}_m)^T \quad (8)$$

3.4.2 Method application

LDA was applied to the pre-processed data and used the common modelling choices, see Appendix E for the implementation code. Further, grid search was used to tune the LDA hyperparameters. Within each cross validation run the numerical weather and time features were scaled on the training set which was then used to transform the validation set. This approach was used to avoid data leakage from scaling on all training data before using cross validation. humidity and cloudcover were min-max normalized because of their defined ranges, while temp, dew, windspeed, and the time features were z-transform standardized. The lsqr solver and auto shrinkage parameters were found to be best from the hyperparameter tuning. Then, the classification threshold was tested by using the tuned model with different thresholds and choosing the best one. A threshold of 0.36 was found to perform best. Lastly, the final model was trained on all the training data and evaluated on the, until now, unseen test data.

3.5 K-nearest neighbours

3.5.1 Mathematical description

For a classification task, the k-NN method determines a new data point's class based on a majority vote of the k closest data points' classes (Lindholm et al., 2022). To calculate the distance between points, different distance metrics can be used to adapt the model to the characteristics of the input data. The Minkowski metric is a generalized expression which can simulate both Euclidean and Manhattan distance. The Minkowski distance, d between two points x_* and x_i with m input variables is given as:

$$d = \left(\sum_{j=1}^m |x_{ij} - x_{*j}|^p \right)^{\frac{1}{p}} \quad (9)$$

where ($p = 1$) gives the Manhattan distance, ($p = 2$) the Euclidean distance, and larger p -values render greater influence for greater distances in the input space. To avert overfitting and suppress sensitivity to noise, the method can aggregate output from multiple data points. However, using a large k might result in underfitting where patterns that characterise the different classes are lost.

3.5.2 Method application

An important aspect of using the k-NN method is to normalise all input variables to be of the similar scale (Lindholm et al. [1]). Utilising the preprocessed data, the input variables humidity and cloudcover were min-max-scaled since their upper and lower bounds were known. The variables temp, dew, and windspeed were scaled using Z-transform since they were assumed to have normally distributed values.

To evaluate the performance of different model variations the GridSearchCV function was used to test all combinations of different hyperparameters. The distance metric parameter was set to

151 Minkowski distance with another parameter p ranging between $[1, 5]$, including Manhattan ($p = 1$)
 152 and Euclidean ($p = 2$) distance. The parameter `n_neighbours`, which determines the k -value, was
 153 set to range between $[1, 50]$. Additionally, the parameter `weights`, was added with the options
 154 `uniform`, which treats all neighbours' influence equally, and `distance`, which weighs neighbours'
 155 influence based on their distance to the new data point.

156 Manual examinations of different input variables were tested to achieve the greatest key metric
 157 scores, ultimately dropping `cloudcover`, `windspeed`, `dew`, and `snowdepth` from the input space.
 158 The function `GridSearchCV` then performed a 5-fold cross validation on the training set for all 490
 159 combinations of hyperparameters resulting in 2450 fits. The best hyperparameters found were based
 160 on the averaged F1-score of each model fit's cross validation. Lastly, the model was evaluated on
 161 these parameters with the unseen test dataset to achieve a final performance report. See Appendix F
 162 for the implementation code for k -NN.

163 3.6 Random forest

164 3.6.1 Mathematical description

165 Tree-based methods in supervised machine learning builds on the concept of iteratively splitting data
 166 into smaller subsets based on specific decision rules. This is done by partitioning the input space into
 167 distinct regions, aiming to maximize class separation based on specific criterias (Lindholm et al. [1]).
 168 In this task the Gini impurity was used as the splitting criterion as described in Equation 10, where p_i
 169 is the proportion of class i instances in region R and C is the number of classes.

$$G(R) = 1 - \sum_{i=1}^c p_i^2 \quad (10)$$

170 In the formula, a lower $G(R)$ indicates higher purity, guiding the algorithm to split on the feature
 171 that minimizes impurity the most. By repeatedly splitting the data based on this guidance, tree-based
 172 models create clear decision rules that help classify new data points while also recognizing patterns
 173 and relationships between different features. However, the performance of (simple) classification
 174 and regression trees (CARTs) tend to overfit the training data, leading to high variance and reduced
 175 generalization performance. To address this, Random forest was chosen as the preferred tree-based
 176 method, as it combines bagging (bootstrap aggregation) and random feature selection when creating
 177 multiple decision trees and combining their predictions, reducing overfitting and increasing stability
 178 (Lindholm et al. [1]). Thus, this ensemble approach effectively reduces variance while maintaining
 179 low bias.

180 3.6.2 Method application

181 Random forest was applied to the pre-processed data and followed the common modeling choices.
 182 Given the imbalance in the dataset, where high bike demand was less frequent, class weighting was
 183 set to "balanced" when splitting the data to adjust for this.

184 Hyperparameter tuning was performed using grid search cross-validation over 16 parameter com-
 185 binations, totaling 80 fits. The tested parameters included `n_estimators` (100, 300), `max_depth` (15,
 186 20), `min_samples_split` (5, 10) and `min_samples_leaf` (2, 4), ensuring a balance between model
 187 complexity and generalization. `max_features = "sqrt"` was used to enhance model diversity by
 188 reducing feature correlation. The final model, optimized for F1-score (0.7018), was trained with
 189 `n_estimators = 300`, `max_depth = 20`, `min_samples_split = 5`, `min_samples_leaf = 4` and `class_weight`
 190 `= "balanced"`. See Appendix G for the implementation code for Random forest.

191 3.7 Performance evaluation

192 The naive classifier has an accuracy of roughly 0.82 when only predicting `low_bike_demand`. The
 193 models developed within this project all perform better than the naive classifier measured by accuracy,
 194 as seen in Figure 3. This indicates that the models have learned actual patterns from the training data
 195 and were able to generalize to unseen data better than a naive classifier. The highest accuracy was
 196 attained by the random forest model, which correctly classified 94% of the test data.

Furthermore, all models were better at classifying the class 0 (low_bike_demand), as seen by the higher F1-score. This can be explained by the unbalanced dataset making it more likely that the true value actually is 0 and there being more data corresponding to the classification 0 to train on. The F1-score was seen as a useful metric for this project and evaluating on the macro average F1-score could thus be fruitful. The metric averages the F1-score between the two classes equally, providing an unweighted estimate of the F1-score for both classes. It being unweighted is seen as favorable in this case due to class 1 prediction being slightly more important since the District Department of Transportation wanted to determine when to increase bike availability.

	precision	recall	f1-score	support		precision	recall	f1-score	support
0	0.94	0.92	0.93	262	0	0.92	1.00	0.96	262
1	0.66	0.74	0.70	58	1	1.00	0.62	0.77	58
accuracy			0.88	320	accuracy			0.93	320
macro avg	0.80	0.83	0.81	320	macro avg	0.96	0.81	0.86	320
weighted avg	0.89	0.88	0.89	320	weighted avg	0.94	0.93	0.92	320

(a) Logistic regression classification report.

(b) Linear discriminant analysis classification report.

	precision	recall	f1-score	support		precision	recall	f1-score	support
0	0.92	0.96	0.94	262	0	0.97	0.96	0.96	262
1	0.79	0.64	0.70	58	1	0.82	0.84	0.83	58
accuracy			0.90	320	accuracy			0.94	320
macro avg	0.86	0.80	0.82	320	macro avg	0.89	0.90	0.90	320
weighted avg	0.90	0.90	0.90	320	weighted avg	0.94	0.94	0.94	320

(c) K-nearest neighbour classification report.

(d) Random forest classification report.

Figure 3: Performance of all developed models shown through respective classification reports.

3.8 Model selection

The random forest model was chosen as the final model due to it outperforming every other model approach on all metrics. Due to the imbalance in the data, the occurrences of high bike demand become more important to capture. However, the bike rental service might not want to increase their stock on false occasions. Therefore, the F1-score proves as a significant evaluation metric since it provides an estimation of the model's ability to catch high bike demand instances, while avoiding false alarms. Since the F1-score is a balance between the precision and recall scores, imbalanced performance on these evaluation metrics can undermine the F1-score's actual significance. Therefore, due to the random forest model's balanced precision and recall scores, the F1-score proves the model's reliability to predict true instances of high bike demand while also maintaining very high performance in predicting low bike demand as well.

4 Conclusion

This project aimed to develop a classification model to predict whether an increase in bike availability would be necessary based on various temporal and weather-related features. Out of the four classification models: Logistic regression, Linear discriminant analysis, K-nearest neighbours and Random forest. The model best suited for this task was found to be Random forest, achieving a macro-average F1-score of 0.90. A key challenge during the project was dealing with the class imbalance, risking bias toward the majority class. To address this, different methods were tested, including class weighting, which improved Random forest but was less effective for Logistic regression. Future related work to this project could explore other classification methods such as boosting or using neural networks to offer potential further performance. To conclude, this project demonstrates that machine learning can enhance bike demand forecasting, while varying performances of the different classification models highlight the importance of careful model selection and hyper parameter tuning, tailored to the given task.

229 **References**

- 230 [1] A. Lindholm, N. Wahlström, F. Lindsten, and T. B. Schön. *Machine Learning A First Course for*
231 *Engineers and Scientists*. Cambridge University Press, 2022.
- 232 [2] Sukamto, Hadiyanto, and Kurnianingsih. Knn optimization using grid search algorithm for
233 preeclampsia imbalance class. *E3S Web of Conferences*, 448:2057—2067, 2023. doi: <https://doi.org/10.1051/e3sconf/202344802057>.
234
- 235 [3] R. van den Goorbergh, M. van Smeden, D. Timmerman, and B. Van Calster. The harm of
236 class imbalance corrections for risk prediction models: illustration and simulation using logistic
237 regression. *Journal of the American Medical Informatics Association: JAMIA*, 29(9):1525—
238 1534, 2022. doi: <https://doi.org/10.1093/jamia/ocac093>.

239 A Appendix - Data analysis code

```

240 1 # Import packages
241 2
242 3 import pandas as pd
243 4 import numpy as np
244 5 import matplotlib.pyplot as plt
245 6 import seaborn as sns
246 7
247 8 import sklearn.preprocessing as skl_pre
248 9 import sklearn.linear_model as skl_lm
249 10 import sklearn.discriminant_analysis as skl_da
250 11 import sklearn.neighbors as skl_nb
251 12 import sklearn.model_selection as skl_ms
252 13
253 14 #from IPython.display import set_matplotlib_formats
254 15 #set_matplotlib_formats('png')
255 16 from IPython.core.pylabtools import figsize
256 17 figsize(10, 6) # Width and height
257 18 plt.style.use('seaborn-white')
258 19
259 20 # Import training data
260 21 data = pd.read_csv('training_data_vt2025.csv')
261 22
262 23 print(data.head())
263 24
264 25 print(data.info())
265 26
266 27 print(data.describe().T)
267 28
268 29 # Plot distributions of each feature
269 30 data.hist(figsize=(12, 8), bins=30, edgecolor="black")
270 31 plt.suptitle("Feature Distributions", fontsize=16)
271 32 plt.show()
272 33
273 34 # Plot boxplots for each feature
274 35 plt.figure(figsize=(12, 6))
275 36 data.plot(kind="box", subplots=True, layout=(5, 4), figsize=(14, 8),
276 37         sharex=False, sharey=False)
277 37 plt.suptitle("Box Plots for Outlier Detection", fontsize=16)
278 38 plt.show()
279 39
280 40 # Plot correlation heatmap
281 41 df_encoded = data.copy()
282 42 df_encoded["increase_stock"] = df_encoded["increase_stock"].astype('
283 43     category').cat.codes # Encode target as numeric
284 44 corr = df_encoded.corr()
285 45
286 46 plt.figure(figsize=(10, 8))
287 47 sns.heatmap(corr, annot=True, cmap="coolwarm", fmt=".2f")
288 48 plt.title("Feature Correlation Heatmap")
289 49 plt.show()
290 50
291 51 # Plot time features against increase_stock
292 52
293 53 # Convert target variable to string if needed
294 54 data['increase_stock'] = data['increase_stock'].astype(str)
295 55
296 56 # Define categorical features
297 57 categorical_features = ['month', 'hour_of_day', 'day_of_week']
298 58
299 59 # Create a figure with subplots (1 row, 3 columns)
300 60 fig, axes = plt.subplots(1, 3, figsize=(18, 5)) # Adjust size as
301 61     needed
302

```

```

30361
30462 # Loop through features and plot each on a separate subplot
30563 for ax, feature in zip(axes, categorical_features):
30664     data.groupby(feature)['increase_stock'].value_counts().unstack().
307     plot(kind='bar', ax=ax)
30865     ax.set_title(f'{feature} vs. Increase Stock')
30966     ax.set_xlabel(feature)
31067     ax.set_ylabel('Count')
31168     ax.legend(title='Increase Stock Category')
31269
31370 # Adjust layout and display the plots
31471 plt.tight_layout()
31572 plt.show()
31673
31774
31875 # Plot holiday and weekday against increase_stock
31976
32077 # Define categorical features
32178 categorical_features = ['holiday', 'weekday']
32279
32380 # Create a figure with subplots (1 row, 3 columns)
32481 fig, axes = plt.subplots(1, 2, figsize=(12, 4)) # Adjust size as
325     needed
32682
32783 # Loop through features and plot each on a separate subplot
32884 for ax, feature in zip(axes, categorical_features):
32985     data.groupby(feature)['increase_stock'].value_counts().unstack().
330     plot(kind='bar', ax=ax)
33186     ax.set_title(f'{feature} vs. Increase Stock')
33287     ax.set_xlabel(feature)
33388     ax.set_ylabel('Count')
33489     ax.legend(title='Increase Stock Category')
33590
33691 # Adjust layout and display the plots
33792 plt.tight_layout()
33893 plt.show()
33994
34095
34196 # Calculate percentage distribution of increase_stock for holidays
34297 holiday_counts = data.groupby("holiday")["increase_stock"].
343     value_counts(normalize=True) * 100
34498
34599 # Calculate percentage distribution of increase_stock for weekdays
34600 weekday_counts = data.groupby("weekday")["increase_stock"].
347     value_counts(normalize=True) * 100
34801
34902 # Display the results
35003 print("Percentage Distribution of Increase Stock for Holidays:")
35104 print(holiday_counts, "\n")
35205
35306 print("Percentage Distribution of Increase Stock for Weekdays:")
35407 print(weekday_counts)
35508
35609
35710 # Plot numerical features against increase_stock
35811
35912 # Define the numerical features to compare against increase_stock
36013 numeric_features = ['temp', 'precip', 'snowdepth', 'humidity', '
361     windspeed', 'cloudcover', 'visibility']
36214
36315 # Set plot style
36416 sns.set_style("whitegrid")
36517
36618 # Create subplots for better visualization
36719 fig, axes = plt.subplots(nrows=2, ncols=4, figsize=(16, 10))

```

```

36820 axes = axes.flatten()
36921
37022 # Plot each numeric feature against increase_stock
37123 for i, feature in enumerate(numeric_features):
37224     if i < len(axes): # Avoid index errors if fewer than 8 features
37325         sns.violinplot(x="increase_stock", y=feature, data=data, ax=
374         axes[i])
37526         axes[i].set_title(f"{feature} vs. Increase Stock")
37627         axes[i].set_xlabel("Increase Stock")
37728         axes[i].set_ylabel(feature)
37829
37930 # Remove empty subplot (since we have 7 plots but a 2x4 grid)
38031 fig.delaxes(axes[-1])
38132
38233 # Adjust layout
38334 plt.tight_layout()
38435 plt.show()

```

385 B Appendix - Pre-processing code

```
386 1 # Import packages
387 2 import pandas as pd
388 3 import numpy as np
389 4
390 5 # Import data
391 6 data = pd.read_csv('training_data_vt2025.csv')
392 7
393 8 def cyclical_encoding(df, column, period):
394 9     df[column + '_sin'] = np.round(np.sin(2 * np.pi * df[column] /
395     period), 6)
396 10     df[column + '_cos'] = np.round(np.cos(2 * np.pi * df[column] /
397     period), 6)
398 11     df.drop(columns=[column], inplace=True) # Remove the original
399     column
400 12     return df
401 13
402 14 def pre_processing(data):
403 15     # Make copy of dataset
404 16     data_processed = data.copy()
405 17
406 18     # Create new summertime feature
407 19     data_processed['is_summer'] = ((data_processed['month'] >= 3) & (
408     data_processed['month'] <= 11)).astype(int)
409 20
410 21     # Normalize calendar data using cosine encoding
411 22     data_processed = cyclical_encoding(data_processed, 'day_of_week',
412     7)
413 23     data_processed = cyclical_encoding(data_processed, 'hour_of_day',
414     24)
415 24     data_processed = cyclical_encoding(data_processed, 'month', 12)
416 25
417 26     # Give target feature numerical values
418 27     data_processed['increase_stock'] = data_processed['increase_stock']
419     .replace({'high_bike_demand': 1, 'low_bike_demand': 0})
420 28
421 29     # Create binary category of features
422 30     data_processed['is_raining'] = (data_processed['precip'] != 0).
423     astype(int)
424 31     data_processed['is_snowing'] = (data_processed['snowdepth'] != 0).
425     astype(int)
426 32     data_processed['is_visible'] = (data_processed['visibility'] !=
427     16).astype(int)
428 33
429 34     # Drop columns
430 35     data_processed = data_processed.drop(columns=['holiday', 'snow', '
431     snowdepth', 'precip', 'visibility', 'summertime'])
432 36
433 37     return data_processed
434 38
435 39 new_data= pre_processing(data)
436 40
437 41 # Save the processed data
438 42 new_data.to_csv("preprocessed_data_1.csv", index=False)
439 43
440 44 print("Preprocessing complete. File saved as preprocessed_data_1.csv")
```

441 C Appendix - Naive classifier code

```
442 1 import pandas as pd
443 2 from sklearn.dummy import DummyClassifier
444 3 from pathlib import Path
445 4 from sklearn.metrics import accuracy_score, f1_score, roc_auc_score
446 5 from sklearn.model_selection import train_test_split
447 6
448 7 # Read the preprocessed data
449 8
450 9 # Construct the full path to the CSV file
451 10 csv_file_path_pre_processed = Path.cwd().parent / 'preprocessed_data_1
452     .csv'
453 11
454 12 # Read the CSV file using pandas
455 13 data = pd.read_csv(csv_file_path_pre_processed)
456 14
457 15 # Split the data into input values, X, and output value, y
458 16 X = data.drop(columns=['increase_stock'])
459 17 y = data['increase_stock']
460 18
461 19 # Split Data into Train & Test Sets (80% Train, 20% Test)
462 20 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size
463     =0.2, stratify=y, random_state=42)
464 21
465 22 # Create a Dummy Classifier (always predicts the majority class)
466 23 dummy_clf = DummyClassifier(strategy="most_frequent")
467 24 dummy_clf.fit(X_train, y_train)
468 25
469 26 # Make predictions
470 27 y_dummy_pred = dummy_clf.predict(X_test)
471 28
472 29 # Evaluate performance
473 30 dummy_accuracy = accuracy_score(y_test, y_dummy_pred)
474 31 dummy_f1 = f1_score(y_test, y_dummy_pred, zero_division=1) # Avoid
475     division errors
476 32 dummy_roc_auc = roc_auc_score(y_test, y_dummy_pred)
477 33
478 34 # Print Results
479 35 print("Dummy Classifier Performance (Majority Class Strategy):")
480 36 print(f"Accuracy: {dummy_accuracy:.4f}")
481 37 print(f"F1-Score: {dummy_f1:.4f}")
482 38 print(f"ROC-AUC Score: {dummy_roc_auc:.4f}")
```

483 D Appendix - Logistic regression code

```
484 1 import pandas as pd
485 2 import numpy as np
486 3 from sklearn.model_selection import StratifiedKFold, train_test_split,
487     GridSearchCV
488 4 from sklearn.linear_model import LogisticRegression
489 5 from sklearn.metrics import classification_report, f1_score,
490     roc_auc_score, accuracy_score, precision_score, recall_score
491 6
492 7 # Load Preprocessed data
493 8 file_path = "preprocessed_data_1.csv"
494 9 df = pd.read_csv(file_path)
495 10
496 11 # Drop increase stock
497 12 X = df.drop(columns=['increase_stock'])
498 13 y = df['increase_stock']
499 14
500 15 # Split data into train and test
501 16 X_train, X_test, y_train, y_test = train_test_split(
502 17     X, y, test_size=0.2, stratify=y, random_state=42
503 18 )
504 19
505 20 # Stratified k-fold
506 21 k_folds = 5
507 22 skf = StratifiedKFold(n_splits=k_folds, shuffle=True, random_state=42)
508 23
509 24 # Hyperparameter for gridsearch
510 25 param_grid = [
511 26     {'C': [0.1, 0.5, 1, 5, 10, 20], 'penalty': ['l1'], 'solver': ['
512     saga'], 'max_iter': [500]},
513 27     {'C': [0.1, 0.5, 1, 5, 10, 20], 'penalty': ['l2'], 'solver': ['
514     liblinear'], 'max_iter': [500]}
515 28 ]
516 29
517 30 # Logistic regression model
518 31 logistic_model = LogisticRegression(random_state=42)
519 32
520 33 # Grid search cross-validation
521 34 grid_search = GridSearchCV(
522 35     logistic_model, param_grid, cv=skf, scoring='f1', n_jobs=-1,
523     verbose=1
524 36 )
525 37
526 38 grid_search.fit(X_train, y_train)
527 39
528 40 # Print the best parameters and f1 score
529 41 best_params = grid_search.best_params_
530 42 best_f1 = grid_search.best_score_
531 43
532 44 print("\n Best Hyperparameters Found:")
533 45 print(best_params)
534 46 print(f"Best F1 Score from Cross-Validation: {best_f1:.4f}")
535 47
536 48 # Train the final Model with best parameters
537 49 final_model = LogisticRegression(**best_params, random_state=42)
538 50 final_model.fit(X_train, y_train)
539 51
540 52 y_prob_test = final_model.predict_proba(X_test)[: , 1]
541 53
542 54 # Threshold tuning
543 55 thresholds = np.arange(0.4, 0.6, 0.8)
544 56 f1_scores = []
545 57
546 58 for threshold in thresholds:
```

```

54759     y_pred = (y_prob_test >= threshold).astype(int)
54860     f1_scores.append(f1_score(y_test, y_pred))
54961
55062     best_threshold = thresholds[np.argmax(f1_scores)]
55163     best_f1_final = max(f1_scores)
55264
55365     print(f"\n Best Threshold Found: {best_threshold:.2f}")
55466     print(f"Best F1-Score at Best Threshold: {best_f1_final:.4f}")
55567
55668     # Best threshold
55769     y_pred_final = (y_prob_test >= best_threshold).astype(int)
55870
55971     # Evaluate final model
56072     final_accuracy = accuracy_score(y_test, y_pred_final)
56173     final_precision = precision_score(y_test, y_pred_final)
56274     final_recall = recall_score(y_test, y_pred_final)
56375     final_f1 = f1_score(y_test, y_pred_final)
56476     final_roc_auc = roc_auc_score(y_test, y_prob_test)
56577
56678     print("\n Final Model Evaluation:")
56779     print(f"Accuracy: {final_accuracy:.4f}")
56880     print(f"Precision: {final_precision:.4f}")
56981     print(f"Recall: {final_recall:.4f}")
57082     print(f"F1-Score: {final_f1:.4f}")
57183     print(f"ROC-AUC Score: {final_roc_auc:.4f}")
57284
57385     # Classification report
57486     print("\n Classification Report:\n")
57587     print(classification_report(y_test, y_pred_final))

```

576 E Appendix - Linear discriminant analysis code

```
577 1 import numpy as np
578 2 import pandas as pd
579 3 from pathlib import Path
580 4 from sklearn.model_selection import StratifiedKFold, GridSearchCV,
581     train_test_split
582 5 from sklearn.metrics import classification_report, f1_score,
583     roc_auc_score, accuracy_score, precision_score, recall_score
584 6 from sklearn.preprocessing import StandardScaler, MinMaxScaler,
585     RobustScaler
586 7 from sklearn.compose import ColumnTransformer
587 8 from sklearn.pipeline import Pipeline
588 9 from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
589 0 from sklearn.metrics import confusion_matrix
590 1
591 12 # Read the preprocessed data
592 13
593 14 # Construct the full path to the CSV file
594 15 csv_file_path_pre_processed = Path.cwd().parent / 'preprocessed_data_1
595     .csv'
596 16 csv_file_path_raw = Path.cwd().parent / 'training_data_vt2025.csv'
597 17
598 18 # Read the CSV file using pandas
599 19 pre_processed_data = pd.read_csv(csv_file_path_pre_processed)
600 20 raw_data = pd.read_csv(csv_file_path_raw)
601 21
602 22 # Split the data into input values, X, and output value, y
603 23 X = pre_processed_data.drop(columns=['increase_stock'])
604 24 y = pre_processed_data['increase_stock']
605 25
606 26 # Set aside test data that is kept unseen until final model evaluation
607 27 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size
608     =0.2, stratify=y, random_state=42)
609 28
610 29 # Define column indices for different scalers
611 30 standardize_cols = ['temp', 'dew', 'windspeed', 'day_of_week_sin', '
612     day_of_week_cos', 'hour_of_day_sin', 'hour_of_day_cos', 'month_sin
613     ', 'month_cos']
614 31 normalize_cols = ['humidity', 'cloudcover']
615 32
616 33 # Define the column transformer with different scalers
617 34 preprocessor = ColumnTransformer([
618 35     ('std_scaler', StandardScaler(), standardize_cols),
619 36     ('minmax_scaler', MinMaxScaler(), normalize_cols)
620 37 ], remainder='passthrough')
621 38
622 39 # Define Stratified K-Fold
623 40 skf = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
624 41
625 42 # Create a pipeline: Scaling + LDA
626 43 pipeline = Pipeline([
627 44     ('preprocessor', preprocessor), # Apply different scalers
628 45     ('lda', LinearDiscriminantAnalysis())
629 46 ])
630 47
631 48 param_grid = [
632 49     {'lda__solver': ['svd']}, # SVD does not take shrinkage
633 50     {'lda__solver': ['lsqr', 'eigen'], 'lda__shrinkage': ['auto', None]
634     } # Only these solvers use shrinkage
635 51 ]
636 52
637 53 # Perform Grid Search with Cross-Validation
638 54 grid_search = GridSearchCV(pipeline, param_grid, cv=skf, scoring='f1',
639     n_jobs=-1, verbose=1)
```



```

64055 grid_search.fit(X_train, y_train)
64156
64257 # Get the best parameters, but remove the 'lda_' prefix
64358 best_params = {key.replace("lda_", ""): value for key, value in
644         grid_search.best_params_.items()}
64559
64660 print("\n Best Hyperparameters Found:")
64761 print(best_params)
64862 best_f1 = grid_search.best_score_
64963 print(f"Best F1 Score from Cross-Validation: {best_f1:.4f}")
65064
65165 final_lda_model = LinearDiscriminantAnalysis(**best_params)
65266 final_lda_model.fit(X_train, y_train)
65367
65468 y_prob_test = final_lda_model.predict_proba(X_test)[:,-1]
65569
65670 # Threshold tuning
65771 thresholds = np.linspace(0,1,101)
65872 f1_scores = []
65973
66074 for r in thresholds:
66175     y_pred = (y_prob_test >= r).astype(int)
66276     f1_scores.append(f1_score(y_test, y_pred))
66377
66478 best_threshold = thresholds[np.argmax(f1_scores)]
66579 best_f1_final = max(f1_scores)
66680
66781 print(f"\n Best Threshold Found: {best_threshold:.2f}")
66882 print(f"Best F1-Score at Best Threshold: {best_f1_final:.4f}")
66983
67084 # Best threshold
67185 y_pred_final = (y_test - y_prob_test >= best_threshold).astype(int)
67286
67387 # Evaluate final model
67488 final_accuracy = accuracy_score(y_test, y_pred_final)
67589 final_precision = precision_score(y_test, y_pred_final)
67690 final_recall = recall_score(y_test, y_pred_final)
67791 final_f1 = f1_score(y_test, y_pred_final)
67892 final_roc_auc = roc_auc_score(y_test, y_pred_final)
67993
68094 print("\n Final Model Evaluation:")
68195 print(f"Accuracy: {final_accuracy:.4f}")
68296 print(f"Precision: {final_precision:.4f}")
68397 print(f"Recall: {final_recall:.4f}")
68498 print(f"F1-Score: {final_f1:.4f}")
68599 print(f"ROC-AUC Score: {final_roc_auc:.4f}")
68600
68701 # Classification report
68802 print("\n Classification report:\n")
68903 print(classification_report(y_test, y_pred_final))
69004
69105 # Confusion Matrix
69206 cm = confusion_matrix(y_test, y_pred_final)
69307
69408 # Visualize cm using a heatmap
69509 import seaborn as sns
69610 import matplotlib.pyplot as plt
69711
69812 plt.figure(figsize=(6, 4))
69913 sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', xticklabels=['Pred:
700     0', 'Pred: 1'], yticklabels=['True: 0', 'True: 1'])
70114 plt.title('Confusion Matrix')
70215 plt.xlabel('Predicted')
70316 plt.ylabel('True')
70417 plt.show()

```

705 F Appendix - K-nearest neighbor code

```
706 1 import pandas as pd
707 2 import numpy as np
708 3 import matplotlib.pyplot as plt
709 4 from sklearn.neighbors import KNeighborsClassifier
710 5 from sklearn.model_selection import KFold, cross_val_predict,
711     StratifiedKFold
712 6 from sklearn.model_selection import train_test_split, GridSearchCV
713 7 from sklearn.metrics import classification_report
714 8 from sklearn.preprocessing import MinMaxScaler, StandardScaler
715 9 from collections import defaultdict
716 10 from sklearn.metrics import accuracy_score, classification_report,
717     confusion_matrix, f1_score, precision_score, recall_score,
718     roc_auc_score
719 11 from sklearn.preprocessing import LabelEncoder
720 12
721 13 # Load data
722 14 bikedata = pd.read_csv("training_data_vt2025.csv")
723 15
724 16
725 17 # Normalization scaler
726 18 n_scaler = MinMaxScaler()
727 19 bikedata["humidity"] = n_scaler.fit_transform(bikedata[["humidity"]])
728 20 bikedata["cloudcover"] = n_scaler.fit_transform(bikedata[["cloudcover"]])
729 21
730 22
731 23 # Standardization
732 24 s_scaler = StandardScaler()
733 25 bikedata["temp"] = s_scaler.fit_transform(bikedata[["temp"]])
734 26 bikedata["dew"] = s_scaler.fit_transform(bikedata[["dew"]])
735 27 bikedata["windspeed"] = s_scaler.fit_transform(bikedata[["windspeed"]])
736 28
737 29
738 30 # Binary transformation
739 31 bikedata['precip'] = bikedata['precip'].apply(lambda x: 1 if x > 0
740     else 0)
741 32 bikedata['snow'] = bikedata['snowdepth'].apply(lambda x: 1 if x > 0
742     else 0)
743 33 bikedata['visibility'] = bikedata['visibility'].apply(lambda x: 1 if x
744     >= 16 else 0)
745 34
746 35 # Sine & Cosine encoding
747 36 bikedata['hour_sin'] = np.sin(2 * np.pi * bikedata['hour_of_day'] /
748     24)
749 37 bikedata['hour_cos'] = np.cos(2 * np.pi * bikedata['hour_of_day'] /
750     24)
751 38 bikedata['day_sin'] = np.sin(2 * np.pi * bikedata['day_of_week'] / 7)
752 39 bikedata['day_cos'] = np.cos(2 * np.pi * bikedata['day_of_week'] / 7)
753 40 bikedata['month_sin'] = np.sin(2 * np.pi * bikedata['month'] / 12)
754 41 bikedata['month_cos'] = np.cos(2 * np.pi * bikedata['month'] / 12)
755 42
756 43 #bikedata = pd.read_csv("fully_processed_data1.csv")
757 44
758 45 X = bikedata.drop(columns=['increase_stock', 'hour_of_day', '
759     day_of_week', 'month', 'snowdepth', 'cloudcover', 'windspeed', '
760     dew'])
761 46 y = bikedata['increase_stock']
762 47
763 48
764 49 X_train, X_test, y_train, y_test = train_test_split(
765 50     X, y, test_size=0.2, stratify=y, random_state=42)
766 51
767 52 # Encode the target variable
768 53 mapping = {'low_bike_demand': 0, 'high_bike_demand': 1}
769 54 y_train_encoded = [mapping[label] for label in y_train]
```

```

76952 y_test_encoded = [mapping[label] for label in y_test]
77053 '''
77154 label_encoder = LabelEncoder()
77255 y_train_encoded = label_encoder.fit_transform(y_train)
77356 y_test_encoded = label_encoder.transform(y_test)
77457 '''
77558
77659
77760 kf = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
77861 k_values = list(range(1, 50))
77962
78063 parameter_grid = {
78164     'n_neighbors' : (k_values),
78265     'weights' : ['uniform', 'distance'],
78366     'metric' : ['minkowski'],
78467     'p' : [1, 2, 3, 4, 5]
78568 }
78669
78770 KNN_model = KNeighborsClassifier()
78871
78972 grid_search = GridSearchCV(KNN_model, parameter_grid, cv = kf, scoring
790    = 'f1', verbose = 1)
79173
79274 grid_search.fit(X_train, y_train_encoded)
79375
79476 best_parameters = grid_search.best_params_
79577 best_f1 = grid_search.best_score_
79678
79779 print("\nBest Hyperparameters Found:")
79880 print(best_parameters)
79981 print("\n Best F1 Score from Cross-Validation:")
80082 print(best_f1)
80183
80284
80385
80486 # Train final model with best hyperparameters
80587 Final_KNN_model = KNeighborsClassifier(**best_parameters)
80688 Final_KNN_model.fit(X_train, y_train_encoded)
80789
80890 y_pred = Final_KNN_model.predict(X_test)
80991
81092 # Evaluate model performance
81193 print("Classification Report:\n")
81294 print(classification_report(y_test_encoded, y_pred))
81395 print("\n Confusion Matrix:\n")
81496 print(confusion_matrix(y_test_encoded, y_pred))

```

815 G Appendix - Random forest code

```

816 1 import numpy as np
817 2 import pandas as pd
818 3 from sklearn.model_selection import train_test_split, StratifiedKFold,
819   GridSearchCV
820 4 from sklearn.ensemble import RandomForestClassifier
821 5 from sklearn.metrics import accuracy_score, classification_report,
822   confusion_matrix, f1_score, precision_score, recall_score,
823   roc_auc_score
824 6
825 7 # Load dataset
826 8 data = pd.read_csv("preprocessed_data_1.csv")
827 9
828 10 # Separate features (X) and target (y)
829 11 X = data.drop(columns=['increase_stock'])
830 12 y = data['increase_stock']
831 13
832 14 # Split data (80% train, 20% test)
833 15 X_train, X_test, y_train, y_test = train_test_split(
834 16     X, y, test_size=0.2, stratify=y, random_state=42)
835 17
836 18 # Decide amount of folds and do stratified K-fold
837 19 k_folds = 5
838 20 kfold = StratifiedKFold(n_splits=k_folds, shuffle=True, random_state
839   =42) # shuffle ensures the data is randomized before splitting.
840 21
841 22 # Defines the hyperparameters to test in GridSearchCV
842 23 param_grid = {
843 24     'n_estimators': [100, 300],          # Number of trees (more for
844     'max_depth': [15, 20],              # Maximum depth of each tree
845 25     'min_samples_split': [5, 10],       # Minimum samples required to split
846 26     'min_samples_leaf': [2, 4],         # Minimum samples per leaf node (
847     'max_features': ['sqrt'],           # Each tree only uses square root
848 27     'class_weight': ['balanced']       # Adjusts weights to handle
849     'max_depth': [15, 20],              # Maximum depth of each tree
850 28     'min_samples_split': [5, 10],       # Minimum samples required to split
851     'min_samples_leaf': [2, 4],         # Minimum samples per leaf node (
852 29     'max_features': ['sqrt'],           # Each tree only uses square root
853     'class_weight': ['balanced']       # Adjusts weights to handle
854 30 }
855 31
856 32 # Initialize the random forest classifier
857 33 rf_model = RandomForestClassifier(random_state=42)
858 34
859 35 # Grid search with cross-validation to find the best hyperparameters
860 36 grid_search = GridSearchCV(
861 37     estimator = rf_model,              # Random Forest as the base model
862 38     param_grid = param_grid,           # Tests all combinations of
863     hyperparameters
864 39     cv = kfold,                        # Stratified K-Fold cross-
865     validation (5 folds)
866 40     scoring = 'f1',                    # Optimizes the F1-score (useful
867     for imbalanced datasets)
868 41     n_jobs = -1,                       # Use all available CPU cores
869 42     verbose = 1                        # Displays progress updates
870 43 )
871 44
872 45 # Train model with grid_search
873 46 grid_search.fit(X_train, y_train)
874 47
875 48 # Get best hyperparameters & best F1 score
876 49 best_params = grid_search.best_params_
877 50 best_f1 = grid_search.best_score_
878 51

```

```

87952 # Print best hyperparameters and F1 Score
88053 print("\nBest Hyperparameters Found:")
88154 print(best_params)
88255 print("\nBest F1 Score from Cross-Validation:")
88356 print(best_f1)
88457
88558 # Train final model with best hyperparameters
88659 best_rf = RandomForestClassifier(**best_params, random_state=42)
88760 best_rf.fit(X_train, y_train)
88861
88962 # Make predictions on test data
89063 y_pred = best_rf.predict(X_test)
89164
89265 # Evaluate model performance
89366 print("Classification Report:\n")
89467 print(classification_report(y_test, y_pred)) # Confusion Matrix: Top:
895      TP, FP & Bottom: FN, TN
89668 print("\n Confusion Matrix:\n")
89769 print(confusion_matrix(y_test, y_pred))
89870
89971 # Evaluating final model performance
90072
90173 f1 = f1_score(y_test, y_pred) # Measures the balance
902      between precision & recall (good for imbalanced data)
90374 accuracy = accuracy_score(y_test, y_pred) # Overall correctness of
904      predictions
90575 precision = precision_score(y_test, y_pred) # Proportion of positive
906      predictions that were actually correct
90776 recall = recall_score(y_test, y_pred) # Proportion of actual
908      positive cases correctly identified
90977 roc_auc = roc_auc_score(y_test, y_pred) # Measures models
910      ability to distinguish between classes (higher = better)
91178
91279 # Print results
91380 print("\nSummary of Model Performance:")
91481 print(f"- F1 Score: {f1:.4f} ")
91582 print(f"- Accuracy: {accuracy:.4f} ")
91683 print(f"- Precision: {precision:.4f} ")
91784 print(f"- Recall: {recall:.4f} ")
91885 print(f"- AUC-ROC: {roc_auc:.4f} ")

```