# Statistical Machine Learning Project 2023

**Niclas Svensson**
niclasdan@gmail.com

## Abstract

In this project several machine learning algorithms were trained and tested based on data collected over three years from Washington D.C. The data included parameters such as the snow depth, days of the week, amount of rain etc. Through feature engineering and data analysis the data was cleaned up and deployed to train the algorithms. After tuning and evaluation, the random forest algorithm scored the highest on the chosen metrics, accuracy and f1 score. Based on the findings made in this report, further testing will be done using the random forest algorithm.

## 1 Introduction to the problem

This project focuses on the capital bikeshare system in Washington D.C, where the availability of bikes often does not meet the demand. This problem leads to potential increase in $CO_2$ emissions as people may chose public transport or cars over bicycles.

The group is faced with the challenge to predict whether an increase in available bikes at an certain hour is necessary using different types of data. Ultimately, the goal is to select the best performing machine learning model to put against the test data.

## 2 Data analysis

Prior to commencing model development, we conducted a comprehensive examination of the data utilizing a tool known as ydata_profiling. This analysis was instrumental in gaining insights into the distribution of each data point. Notably, it highlighted anomalies, such as the 'snow' feature containing solely 0 values, rendering it non-contributory to our model's learning process. Further elaboration on this will be provided in the Feature Engineering section.

Moreover, this tool facilitated the identification of crucial data segments and insignificantly impactful ones. Additionally, it offered a visualization of data interconnections, aiding in strategic decisions regarding feature enhancement in subsequent stages.

The Categorical inputs were: hour_of_day, day_of_week, month, holiday, weekday, summertime, snow and increase_stock.
The numerical inputs were: temp, dew, humidity, precip, snowdepth, windspeed, cloudcover and visibility.

### 2.1 Feature engineering

We replaced "month" with *is_winter*, *is_spring*, *is_summer* and *is_fall* to capture a broader weather related impact on bike usage. This was done by mapping each month to it's respective season effectively making it a binary output instead of a numerical.

Introducing *rush_hour* and *night_time* which targets the important times of day with varying bike usage. *Rush_hour* being morning and afternoon and *night_time* being late night and early morning. This binary categorization helps distinguish between peak demand times of the day.

Converting *snowdepth* to a binary *is_there_snow* feature simplifies the model by focusing on the presence of snow rather than the amount. By transforming it into a binary feature we could simply ask: does snow affect bike demand?

Changing the target variable to a binary *increase_stock* will label it under one category making it easier to process the data with our machine learning algorithms because it's a binary output under one category.

Continuous variables such as *temp*, *humidity*, *windspeed*, *cloudcover* and *visibility* were normalized using StandardScaler. This is done to ensure that all these features contribute equally to the final result without any single feature dominating the rest due to it's scale.

In summary we dropped the month column for *is_winter*, *is_spring*, *is_summer* and *is_fall* to define seasons instead. Then we add *rush_hour* and *night_time* from *hour_of_day* to get categorical features for important times during the day. The numerical feature *snowdepth* was also dropped for *is_there_snow* to get a categorical feature instead. Lastly, we also dropped *summertime*, *day_of_week*, *dew*, *precip* and *snow*.

After we managed our features in an appropriate way we started inspecting our data.
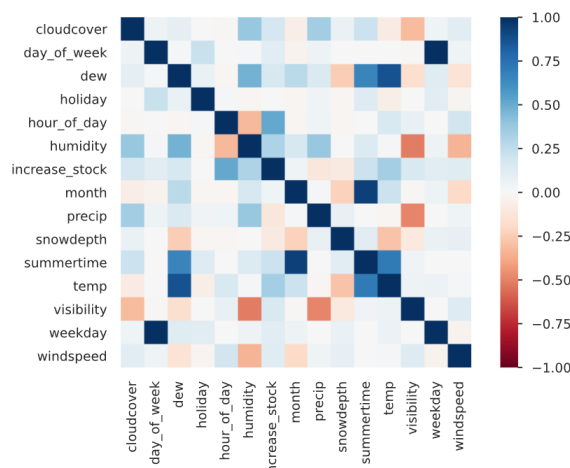


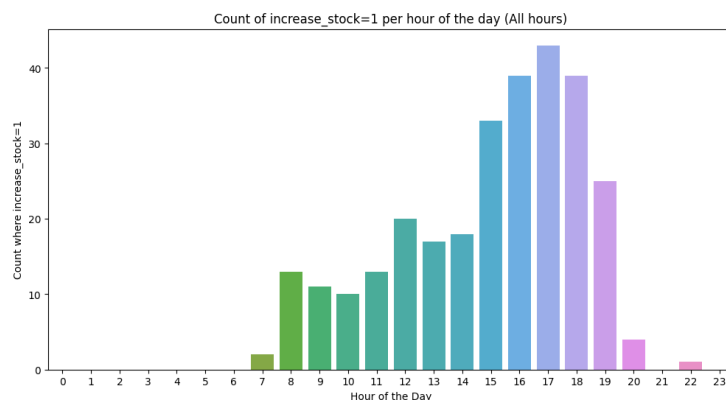Figure 1: Correlations between different input features



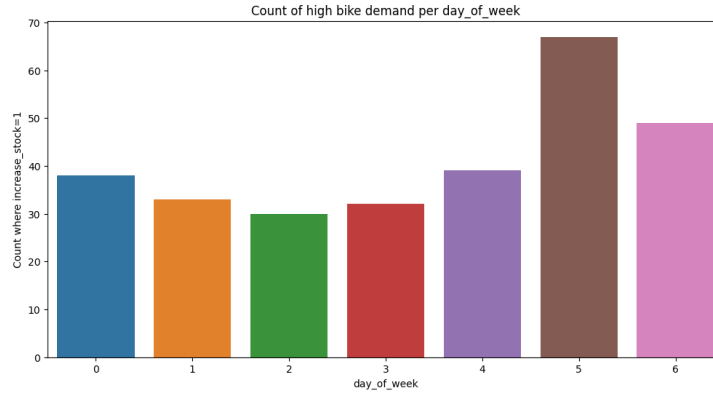Figure 2: Bike demand over different hours of the day

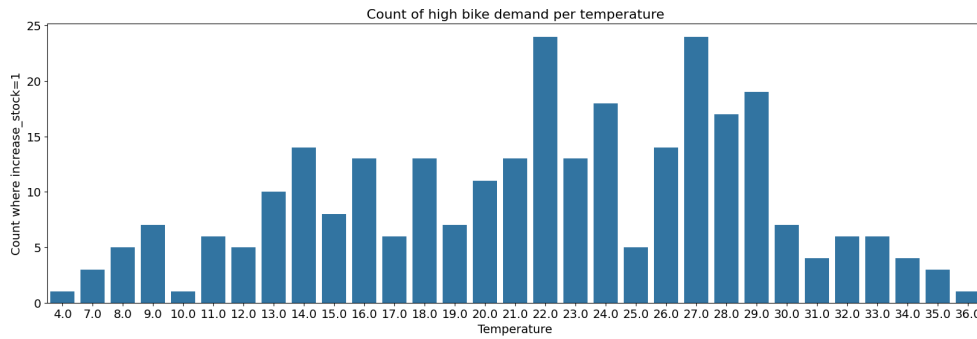Figure 3: Bike demand over different days of the week



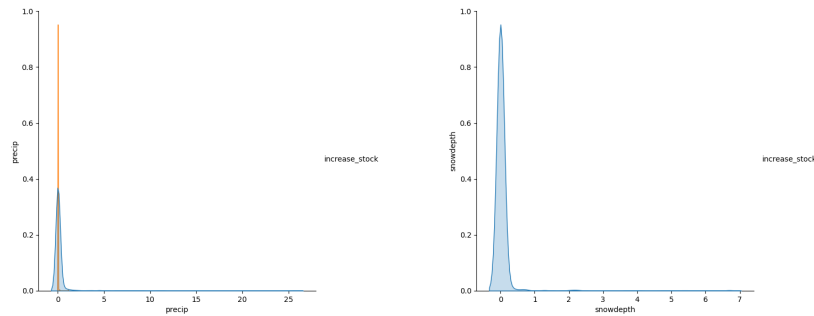Figure 4: Bike demand over different temperatures



Figure 5: Bike Demand depending on Precipitation(left)/Snowdepth(right)

Figure 3 depict the bike demand over a day. Here we can see that independent of the day, around 15 to 19 seem to have the highest bike demand. By inspecting figure 4 we see that weekends have higher demand than weekdays. Based on Figure 5 and 6 we can see weather has an effect on the bike demand. Figure 6 clearly shows that if there is precipitation or if there "is" a snowdepth the bike demand is almost always zero, regardless of the amount of precipitation or magnitude of the snowdepth. Figure 5 displays that the bike demand increases as the weather is above 12 degrees, but decreases rapidly if we surpass 29 degrees.

## 3 Model development

### 3.1 K-nn classifier

The data set is composed of a feature matrix X and a target variable y. The feature matrix X includes all the independent variables while the target variable y represents the binary outcome of high or low bike demand. Each feature in X is then normalized to ensure equal contribution to the result. Without normalization different features might scale differently and influence the model. For the k-nn algorithm we classify data points based on its neighbors. It operates by finding the k nearest neighbors and assigns the most frequent label among these neighbors to the test point. The k value of the k-nn algorithm decides whether it functions more towards bias or variance. A low k might lead to the model reacting to strongly to small fluctuations while a high k might cause the model to miss out on small but important pieces of information. Therefore, identifying the optimal k value involves extensive amount of testing.

$$f(x; k, T) = \arg\max_y \sum_{i=1}^n w(x, X_i; k)\mathbf{1}\{y_i = y\} \tag{1}$$

Eq (1) is the mathematical concept behind the K-nn method. The classification function for a new point x given the number of neighbors k and the training data T. The weight function depends on x, the training data point Xi and the numbers of neighbors k considered. Lastly we have the indicator function which outputs 1 if yi is equal to y and 0 otherwise.

$$d(P, Q) = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2} \tag{2}$$

To calculate the distance between neighbors and input points we use Euclidean distance. In 2 space we use x and y coordinates to define the distance between neighbors and add more parameters if needed. All the parameters are then squared to find the straight line distance between two neighbors.

### 3.2 Discriminant Analysis

Discriminant analysis can be used to predict the labels of the output for new input data. Thus, discriminant analysis aims to estimate $p(y|x)$. The classifier which is obtained after applying the joint probability distribution function, as well as some expressions taken from the Gaussian mixture model (GMM) is:

$$p(y = m|\mathbf{x}_*) = \frac{\hat{\pi}_m N(\mathbf{x}_*|\hat{\mu}_m, \hat{\Sigma}_m s)}{\Sigma_{j=1}^M \hat{\pi}_j N(x_*|\hat{\mu}_j, \hat{\Sigma}_j)} \tag{3}$$

The parameters to be trained are $\theta = \{\mu_m, \Sigma_m, \pi_m\}_{m=1}^M$. where $\mu_m$, $\Sigma_m$ and $\pi_m$ is the mean vector, covariance matrix and the marginal class probabilities for each class. The mean vector of each class is given by:

$$\hat{\mu}_m = \frac{1}{n_m} \sum_{i:y_i=m} x_i \tag{4}$$

And the covariance matrix for the different classes are estimated as:

$$\hat{\sum}_m = \frac{1}{n_m}(x_i - \hat{\mu}_m)(x_i - \hat{\mu}_m)^T \tag{5}$$

The marginal class probabilities for each class is $\hat{\pi_m} = \frac{n_m}{n}$. The 'hard' predictions ($\hat{y}_*$) can be obtained, by choosing the class which has the most probable prediction with:

$$\hat{y}_* = \arg\max_m p(y = m|\mathbf{x}_*) \tag{6}$$

As the decision boundaries of this classifier takes on a quadratic form, this classifier approach is denoted as Quadratic Discriminant Analysis (QDA). However, if the assumption is made that the covariance matrix is the same for all of the classes, the Linear Discriminant Analysis (LDA) is obtained instead. (1) LDA and QDA both assumes that the dataset is drawn from a Gaussian distribution. The assumption of all the covariance matrices being the same results in a dimensionality reduction of the dataset, where LDA projects the input data onto a linear subspace which is used to

classify the output labels. As QDA does no assumption about the covariance between classes in the dataset, the parameter it handles is increased. Because of this QDA has an amount of parameters which is quadratic in regards of the amount of input. This results in QDA having its quadratic decision boundary.(2)

### 3.3 Tree Based Method

Tree based methods are supervised learning techniques, that are adaptable for both classification and regression problems. They can be used for both linear and non-linear problems.

#### 3.3.1 Decision Trees

Often referred to as rule-based models, they divide the input space into multiple region, based on the "rules" imposed on certain features. These ordered rules can be represented as a tree structured graph, hence the name. With this model, we are looking for the prediction $\hat{y}(x_*)$, such as:

$$\hat{y}(\mathbf{x}_*) = \sum_{l=1}^{L} \hat{y}_l I\{\mathbf{x}_* \in R_l\} \tag{7}$$

where L denotes the number of regions, $R_l$ is the l-th region and $\hat{y}_l$ is the prediction for the region by the model. $I\{x_* \in R_l\}$ is 1 if the point is in the given region, 0 otherwise. By finding suitable $\hat{y}$ values and R regions, we can teach our model to classify data. We use recursive binary splitting to do this, a greedy algorithm that tries to minimize the cost function after each split. Two commonly used cost functions are the squared error or an optimization problem defined as:

$$argmin_{nextsplit}(n_1 Q_1 + n_2 Q_2) \tag{8}$$

Here $n_i$ is the number of samples in the i-th region, while $Q$ is a function that measures the quality of the split, often called the criterion. In our solution we used the *Gini-index*, which is defined as:

$$Q_l = \sum_{m=1}^{M} \hat{\pi}_{lm}(1 - \hat{\pi}_{lm}) \tag{9}$$

where $\hat{\pi}_{lm}$ is the ratio of training observations within the l-th region that are part of the m-th class.

By defining this metrics, we can start building the classification tree. To do this, we divide the input space into two regions, based on one of the selected features from the input and the loss measuring function. Once we have the two regions, defined by the cut point s:

$$R_1(i, s) = \{\mathbf{x}|x_i < s\} \quad and \quad R_2(i, s) = \{\mathbf{x}|x_i > s\} \tag{10}$$

We repeat this step on the regions, until a certain depth, or some kind of other criterion is reached.

Once we are done building the tree, we have the prediction for this regions defined, as the majority class of the points in that region:

$$\hat{y}_1 = MajorityVote\{y_i : \mathbf{x_i} \in R_1\} \quad and \quad \hat{y}_2 = MajorityVote\{y_i : \mathbf{x_i} \in R_2\} \tag{11}$$

### 3.4 Random Forest

Ensemble techniques create a single, more accurate model using several base ones. Random forests have decision trees, as their base models. We bootstrap the data to create multiple data sets and different decision trees are then trained on these. Each tree is utilizing only a subset of inputs. This further injections randomness into the structure helps prevent any dominant set of inputs from consistently winning in every tree, making this more effective than bagging. Through this technique, we successfully reduce the correlation between different trees, resulting in a significant reduction in variance when averaging out their votes. During training we can grow each tree in parallel and control their attributes such as height, splits, number of trees, with hyperparameters. Each tree is trained on a bootstrapped dataset using selected subsets of features, so it is a computationally heavy model, however it produces better results than single decision trees.

### 3.5   Logistic Regression

Logistic regression offers a method to model conditional class probabilities. It's essentially an adaptation of the linear regression model tailored to address classification problems rather than regression. It serves as a statistical approach for binary classification tasks, estimating the probability of a binary outcome based on input features.

#### 3.5.1   Mathematical model

The aim is to find a function $g(x)$ that estimates the likelihood of the positive class. Initially, we employ the linear regression model, expressed as

$$z = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \ldots + \theta_p x_p = \theta^T x \tag{12}$$

which maps $x$ to $z$. Here, $z$ spans the entire real number space, but we require a function that outputs within the range $[0, 1]$. Logistic regression addresses this by compressing $z$ towards the interval $[0, 1]$ using the logistic function $h(z) = \frac{e^z}{1+e^z}$. This transformation yields:

$$g(x) = \frac{e^{\theta^T x}}{1 + e^{\theta^T x}} \tag{13}$$

This revised function $g(x)$ ensures values between 0 and 1, allowing for interpretation as probabilities.

For logistic regression, training the model involves finding the optimal parameters that minimize the logistic loss.

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^{m} \left[ y^{(i)} \log(h_\theta(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_\theta(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^{n} \theta_j^2 \tag{14}$$

For solving the problem `sklearn.linear_model.LogisticRegression`, the `liblinear` solver has been used. The `liblinear` solver is a numerical optimization algorithm designed for logistic regression and linear support vector machines (SVMs). It employs an iterative numerical optimization technique, typically a coordinate descent algorithm.

### 3.6   Hyperparameters

The tuning of the various models was conducted with GridSearchCV. For logistic regression model we covered parameters such as the regularization strength, penalty type, solver algorithms, maximum iterations etc. The hyperparameter tuned for LDA was the solver, resulting in the best solver being svd. The tuned hyperparameter for QDA was the regularization parameter, which regularizes the estimates per-class covariance by transforming S2, which corresponds to the scaling attributes of the different classes. Random forest was more difficult to tune as it tended to overfit the data when using few trees. The alteration of the number of trees, their depth and number of samples per tree showed to have the most impact. For k-NN the tuned hyperparameter was the amount of neighbors (k) to include for the test value. The optimal k was found using for-loops.

### 3.7   Evaluation and performance

For the naive model we used (`sklearn.dummy.DummyClassifier`) that serves as a simple baseline to compare against our other more complex models. It generates predictions uniformly at random from the list of unique classes, so each class has equal probability.

The data set was divided into training and validation sets employing a 20-fold cross-validation strategy (`sklearn.model_selection.KFold`). This method splits the data set into subsets, trains the model on various combinations, and computes the metrics. The mean value of each metric, across validation sets provided an estimate of the model's performance. We chose to evaluate our models based on accuracy and F1 score. The accuracy would give us a peek into the general performance of the model. Precision and recall have their own reasons behind them, higher precision would mean higher customer satisfaction and less emission, while higher recall yields higher profits for the operating company. We decided on evaluating the model's F1 score, that takes both of them into account.

After training the models on the training set using k-fold cross validation, we evaluated them on the test set as well. We achieved the following results



Figure 6: Accuracy and F1 scores of the models

It is easy to see that why we ended up choosing the random forest model, it performed the best including accuracy and F1 score as well before the logistic regression one.

## 4   Conclusions

The final model was evaluated as the best based on accuracy and F1 score, using a 20-fold cross validation strategy and further tested on a separate test set. This evaluation of each model provided insight into it's strengths and limitations for the bike-sharing system.

Furthermore, the project takes a methodical approach to machine learning in a real-world scenario. The project delved into feature engineering, implementation of different models and optimizing the hyper parameters for the best performance. This methodical and comprehensive strategy not only addresses the immediate need of the bike-sharing system but also lays a foundation for future enhancements that machine learning algorithms can perform.

# References

[1] Lindholm, Andreas, Niklas Wahlström, Fredrik Lindsten, & Thomas B. Schön. (2022). *Machine Learning - A First Course for Engineers and Scientists*. Cambridge University Press.

[2] Scikit-learn, *Linear and Quadratic Discriminant Analysis (2023)*, last accessed 20 december 2023, `https://scikit-learn.org/stable/modules/lda_qda.html`

# A Appendix

## A.1 Main code

```python
#importing all neccesary libraries
import pandas as pd
import matplotlib.pyplot as plt
import numpy as np
import seaborn as sns
import sklearn.discriminant_analysis as skl_da
import sklearn.preprocessing as skl_pre
import sklearn.linear_model as skl_lm
import sklearn.discriminant_analysis as skl_da
import sklearn.neighbors as skl_nb
import sklearn.model_selection as skl_ms
from sklearn.dummy import DummyClassifier
from sklearn.metrics import accuracy_score
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import f1_score
from sklearn.ensemble import BaggingClassifier, RandomForestClassifier
from sklearn.model_selection import GridSearchCV

np.random.seed(1)

# loading the data
data = pd.read_csv('/content/training_data.csv', dtype={'ID':
    str}).dropna().reset_index(drop=True)

data['increase_stock'] = data['increase_stock'].map({'low_bike_demand': 0,
    'high_bike_demand': 1})

# Plotting
data_snow = data[['snowdepth', 'increase_stock']]
plot_snow = sns.pairplot(data_snow, hue = 'increase_stock', height = 6);

data_precip = data[['precip', 'increase_stock']]
plot_precip = sns.pairplot(data_precip, hue = 'increase_stock', height = 6);

#do not modify this seed value, shuffle the train data in an another part of
    the code, so we all have the same test set

# Categorizing different variables
data['is_winter'] = data['month'].apply(lambda x: 1 if x in [12, 1, 2] else 0)
data['is_spring'] = data['month'].apply(lambda x: 1 if x in [3, 4, 5] else 0)
data['is_summer'] = data['month'].apply(lambda x: 1 if x in [6, 7, 8] else 0)
data['is_fall'] = data['month'].apply(lambda x: 1 if x in [9, 10, 11] else 0)
data['rush_hour'] = data['hour_of_day'].apply(lambda x: 1 if (7 <= x <= 9) or
    (17 <= x <= 19) else 0)
data['night_time'] = data['hour_of_day'].apply(lambda x: 1 if (20 <= x) or (7
    >= x) else 0)
data['is_there_snow'] = data['snowdepth'].apply(lambda x: 1 if (0 < x) else 0)
data['increase_stock'] = data['increase_stock'].map({'low_bike_demand': 0,
    'high_bike_demand': 1})

# dropping variables from the dataframe
```

```
46  columns_to_drop = ['snowdepth','summertime','day_of_week','month','dew',
    ↪  'precip', 'snow']
47  data.drop(columns=columns_to_drop, inplace=True)
48
49  #normalizing numerical coloumns
50  normalized_cols = ['temp', 'humidity', 'windspeed', 'cloudcover',
    ↪  'visibility']
51  dummy_cols = ['hour_of_day', 'holiday', 'weekday']
52
53  scaler = StandardScaler()
54  data[normalized_cols] = scaler.fit_transform(data[normalized_cols])
55
56  # creating data set with all the data to train the final chosen method
57  x_final = data.drop(columns = ['increase_stock'])
58  y_final = data['increase_stock']
59
60  seed_value = 42
61  #Splitting the data into training set and a testing set
62  data, test_data = skl_ms.train_test_split(data, test_size=0.2,
    ↪  random_state=seed_value, shuffle = False)
63
64  # seperating labels and input data
65  x_training = data.drop(columns = ['increase_stock'])
66  y_training = data['increase_stock']
67  x_testing = test_data.drop(columns = ['increase_stock'])
68  y_testing = test_data['increase_stock']
69
70  # assigning models
71  LDA = skl_da.LinearDiscriminantAnalysis()
72  QDA = skl_da.QuadraticDiscriminantAnalysis()
73
74  # splitting data into training set and validation set
75  x_train, x_val, y_train, y_val = skl_ms.train_test_split(x_training,
    ↪  y_training, train_size = 0.80, random_state = 1 )
76
77  # fit the basic models
78  LDA.fit(x_train, y_train)
79  QDA.fit(x_train, y_train)
80
81  # accuracy on training data
82  LDA_score_train = LDA.score(x_train, y_train) # 0.869
83  QDA_score_train = QDA.score(x_train, y_train) # 0.196
84
85  # accuracy on validation data
86  LDA_score_val = LDA.score(x_val, y_val) # 0.859
87  QDA_score_val = QDA.score(x_val, y_val) # 0.164
88
89  # Tuning LDA, grid search for best solver
90  # tuning done on the left out validation set
91  # increasing accuracy to 0.898
92  param_grid = {'solver': ['svd', 'lsqr', 'eigen']}
93  grid_search = skl_ms.GridSearchCV(LDA, param_grid, cv=5)
94  grid_search.fit(x_val, y_val)
95
96  best_paramss = grid_search.best_params_
97
98  best_LDA = grid_search.best_estimator_
99
100 y_pred = best_LDA.predict(x_val)
101 accuracy = accuracy_score(y_val, y_pred) # LDA tuned result = 0.848
102
103 # Tuning QDA, grid search to optimize the regularization parameter
104 # tuning done on the left out validation set
105 # increasing accuracy to 0.887
```

```python
106  params = [{'reg_param': [0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1.0]}]
107  cv = skl_ms.RepeatedStratifiedKFold(n_splits = 10, n_repeats = 3, random_state
     ↪  = 1)
108  search = skl_ms.GridSearchCV(QDA, params, cv = cv)
109  search.fit(x_val, y_val)
110
111
112  best_params = search.best_params_
113  best_QDA = search.best_estimator_
114
115  y_pred = best_QDA.predict(x_val)
116  acc = accuracy_score(y_val, y_pred)
117
118  print('Accuracy on Train data')
119  print(f'LDA: {LDA_score_train:.3f}')
120  print(f'QDA: {QDA_score_train:.3f}')
121  print('Val data')
122  print(f'LDA: {LDA_score_val:.3f}')
123  print(f'QDA: {QDA_score_val:.3f}')
124  print('Grid Search alteration')
125  print(f'Grid Search result on LDA {accuracy:.3f}')
126  print(f'Grid search result on QDA {acc:.3f}');
127
128  #train the tuned models
129  best_LDA.fit(x_training, y_training)
130  best_QDA.fit(x_training, y_training)
131
132  accuraccy_rates=[]
133  # Specified parameters
134  params = {
135      'C': 1.0,
136      'class_weight': None,
137      'fit_intercept': True,
138      'intercept_scaling': 1,
139      'max_iter': 100,
140      'penalty': 'l1',
141      'solver': 'liblinear',
142      'tol': 0.001
143  }
144
145  #Cross validation
146
147  n_fold=10
148  cv=skl_ms.KFold(n_splits=n_fold, random_state=2, shuffle=True)
149
150  for train_index, val_index in cv.split(x_training):
151    X_train, X_val = x_training.iloc[train_index], x_training.iloc[val_index]
152    y_train, y_val = y_training.iloc[train_index], y_training.iloc[val_index]
153
154    model=skl_lm.LogisticRegression(**params)
155    model.fit(X_train, y_train)
156    prediction = model.predict(X_val)
157    accuracy=np.mean(prediction==y_val)
158    accuraccy_rates.append(accuracy)
159    print('accuracy rate: ' + str(accuracy))
160
161  print('Mean accuracy rate: ' + str(np.mean(accuraccy_rates)))
162
163  # Define hyperparameters grid
164  param_grid = {
165      'C': [0.1, 1.0, 10.0, 30.0],
166      'penalty': ['l1', 'l2'],
167      'solver': ['liblinear', 'saga'],  # Different solvers
168      'max_iter': [100, 200, 300],  # Different max_iter values
```

```python
169         'class_weight': [None, 'balanced'],  # Different class weight options
170         'fit_intercept': [True, False],  # Whether to fit intercept
171         'tol': [1e-3, 1e-4],  # Tolerance for stopping criteria
172         'intercept_scaling': [1, 2]  # Scaling for intercept (relevant for
        ↪   'liblinear')
173         }
174
175 # Initialize GridSearchCV
176 grid_search =
    ↪   GridSearchCV(estimator=skl_lm.LogisticRegression(solver='liblinear'),
177                             param_grid=param_grid,
178                              cv=cv)
179
180 # Fit the model
181 grid_search.fit(X_val, y_val)
182
183 # Get the best parameters
184 best_params = grid_search.best_params_
185 print("Best Parameters:", best_params)
186
187 # Get the best model
188 best_LG = grid_search.best_estimator_
189 # Fit the best LG model
190 best_LG.fit(x_training, y_training)
191 ## Split the data into training(80%) and validation sets(20%)
192 x_train, x_val, y_train, y_val = skl_ms.train_test_split(x_training,
    ↪   y_training, test_size=0.2, random_state=1)
193
194 # Implement and evaluate k-NN Classifier with k=2
195 knn = skl_nb.KNeighborsClassifier(n_neighbors=2).fit(x_train, y_train)
196 print("Test value: k = 2")
197 print(f"Accuracy: {accuracy_score(y_val, knn.predict(x_val))}\n")
198
199 # Experiment with different values of k to find the optimal one
200 misclassification = []
201 k_values = range(1, 50)  # 50 - 200 is optimal
202
203 for k in k_values:
204     knn = skl_nb.KNeighborsClassifier(n_neighbors=k)
205     knn.fit(x_train, y_train)
206     y_pred = knn.predict(x_val)
207     misclassification.append(np.mean(y_pred != y_val))
208
209 # Plotting misclassification rate vs. k
210 plt.plot(k_values, misclassification, marker='.')
211 plt.xlabel('Number of Neighbors (k)')
212 plt.ylabel('Misclassification Rate')
213 plt.title('k-NN Varying number of Neighbors')
214 plt.show()
215
216 # Identify the optimal k value and retrain the model
217 optimal_k = k_values[misclassification.index(min(misclassification))]
218 print("For loop")
219 print(f"Optimal value of k: {optimal_k}")
220
221 # Retraining k-NN classifier with the optimal number of neighbors
222 knn_optimal = skl_nb.KNeighborsClassifier(n_neighbors=optimal_k)
223 knn_optimal.fit(x_train, y_train)
224 y_pred_optimal = knn_optimal.predict(x_val)
225 accuracy_optimal = accuracy_score(y_val, y_pred_optimal)
226 print(f"Optimized Accuracy with k = {optimal_k}: {accuracy_optimal}\n")
227
228 # Find the optimal k using GridSearchCV
```

```python
229  grid_search = GridSearchCV(skl_nb.KNeighborsClassifier(), {'n_neighbors':
     ↪ range(1, 50)}, cv=5, scoring='accuracy').fit(x_train, y_train)
230  print("Grid search")
231  print(f"Best parameters: {grid_search.best_params_}")
232  print(f"Best model accuracy: {accuracy_score(y_val,
     ↪ grid_search.best_estimator_.predict(x_val))}\n")
233
234  # Find the optimal k using RandomizedSearchCV, 1 - 200 seems to be getting
     ↪ best result
235  random_search = skl_ms.RandomizedSearchCV(skl_nb.KNeighborsClassifier(),
     ↪ {'n_neighbors': range(1, 200)}, cv=5, scoring='accuracy',
     ↪ random_state=1).fit(x_train, y_train)
236  print("Random search")
237  print(f"Best parameters: {random_search.best_params_}")
238  print(f"Best model accuracy: {accuracy_score(y_val,
     ↪ random_search.best_estimator_.predict(x_val))}")
239  best_kNN = grid_search.best_estimator_
240
241  #training the best kNN model
242  best_kNN.fit(x_training, y_training)
243
244  # seperating the dataset into a training set and a validation set
245  x_train, x_val, y_train, y_val = skl_ms.train_test_split(x_training,
     ↪ y_training, train_size = 0.80, random_state = 1 )
246
247  #fitting the basic model
248  random_forest = RandomForestClassifier()
249  random_forest.fit(x_train, y_train)
250
251  # accuracy on training data
252  random_forest_score_train = random_forest.score(x_train, y_train) # 0.869
253  random_forest_score_val = random_forest.score(x_val, y_val)
254
255  print('Without tuning')
256  print(f'Random forest accuracy {random_forest_score_train} on training data')
257  print(f'Random forest accuracy {random_forest_score_val} on validation data')
258
259  # Hyperparameters to evaluate over (tuning)
260  param_grid = {
261      'n_estimators': [25, 50, 100, 150],
262      'max_features': ['sqrt', 'log2', None],
263      'max_depth': [3, 6, 9],
264      'max_leaf_nodes': [3, 6, 9],
265  }
266
267  grid_search = GridSearchCV(RandomForestClassifier(),
268                            param_grid=param_grid)
269  grid_search.fit(x_val, y_val)
270  print(grid_search.best_estimator_)
271
272  # Training the best Random Forest model
273  best_RF = RandomForestClassifier(max_depth = 9, max_features = None,
     ↪ max_leaf_nodes = 6, n_estimators = 25)
274  best_RF.fit(x_training, y_training)
275  y_pred = best_RF.predict(x_val)
276  accuracy = accuracy_score(y_val, y_pred)
277  print(f' Accuracy of tuned random forest classifier {accuracy}')
278
279  # A dummy classifier for comparison of models.
280  model = DummyClassifier(strategy = 'uniform')
281  model.fit(x_training, y_training)
282  Dummy = model
283
284  #final test on hold out dataset
```

```
285
286  models = []
287  models.append(Dummy)
288  models.append(best_LG)
289  models.append(best_LDA)
290  models.append(best_QDA)
291  models.append(best_kNN)
292  models.append(best_RF)
293  test_accuracy = []
294  test_fscore = []
295
296  for m in range(np.shape(models)[0]):
297          model = models[m]
298          pred = model.predict(x_testing)
299          test_accuracy.append(np.mean(pred == y_testing))
300          test_fscore.append(f1_score(y_testing, pred))
301
302
303  # Plotting test accuracy
304  plt.figure(figsize=(10, 5))
305  plt.subplot(1, 2, 1)
306  plt.plot(('Dummy', 'LG', 'Lda', 'Qda', 'K-nn', 'Random Forest'), test_accuracy,
       ↪  marker='o', linestyle='-', color='skyblue')
307  plt.xlabel('Models')
308  plt.ylabel('Test Accuracy')
309  plt.title('Test Accuracy of Different Models')
310  plt.xticks(rotation=45)
311  plt.grid(True)
312
313  # Plotting F1 score
314  plt.subplot(1, 2, 2)
315  plt.plot(('Dummy', 'LG', 'Lda', 'Qda', 'K-nn', 'Random Forest'), test_fscore,
       ↪  marker='o', linestyle='-', color='salmon')
316  plt.xlabel('Models')
317  plt.ylabel('Test F1 Score')
318  plt.title('Test F1 Score of Different Models')
319  plt.xticks(rotation=45)
320  plt.grid(True)
321
322  plt.tight_layout()
323  plt.show()
324
325  print(f'Mean accuracy per model: {test_accuracy}')
326  print(f'Mean F score per model: {test_fscore}')
327
328  # to be able to download results
329  from google.colab import files
330
331
332  # training chosen model on the whole training dataset
333  best_RF.fit(x_final, y_final)
334
335  # loading the data
336  data = pd.read_csv('/content/test_data.csv', dtype={'ID':
       ↪  str}).dropna().reset_index(drop=True)
337
338  #do not modify this seed value, shuffle the train data in an another part of
       ↪  the code, so we all have the same test set
339
340  # Categorizing different variables
341  data['is_winter'] = data['month'].apply(lambda x: 1 if x in [12, 1, 2] else 0)
342  data['is_spring'] = data['month'].apply(lambda x: 1 if x in [3, 4, 5] else 0)
343  data['is_summer'] = data['month'].apply(lambda x: 1 if x in [6, 7, 8] else 0)
344  data['is_fall'] = data['month'].apply(lambda x: 1 if x in [9, 10, 11] else 0)
```

```python
345  data['rush_hour'] = data['hour_of_day'].apply(lambda x: 1 if (7 <= x <= 9) or
     ↪  (17 <= x <= 19) else 0)
346  data['night_time'] = data['hour_of_day'].apply(lambda x: 1 if (20 <= x) or (7
     ↪  >= x) else 0)
347  data['is_there_snow'] = data['snowdepth'].apply(lambda x: 1 if (0 < x) else 0)
348
349  # dropping variables from the dataframe
350  columns_to_drop = ['snowdepth','summertime','day_of_week','month','dew',
     ↪  'precip', 'snow']
351  data.drop(columns=columns_to_drop, inplace=True)
352
353  #normalizing numerical coloumns
354  normalized_cols = ['temp', 'humidity', 'windspeed', 'cloudcover',
     ↪  'visibility']
355  dummy_cols = ['hour_of_day', 'holiday', 'weekday']
356
357  scaler = StandardScaler()
358  data[normalized_cols] = scaler.fit_transform(data[normalized_cols])
359
360  y_pred = best_RF.predict(data)
361  np.savetxt("predictions.csv", y_pred, delimiter=",")
362
363  files.download("predictions.csv")
```

## A.2  Data visualization

```python
1   import pandas as pd
2   import matplotlib.pyplot as plt
3   import numpy as np
4   from sklearn.dummy import DummyClassifier
5   from sklearn.metrics import accuracy_score
6   import sklearn.discriminant_analysis as skl_da
7   import sklearn.preprocessing as skl_pre
8   import sklearn.linear_model as skl_lm
9   import sklearn.discriminant_analysis as skl_da
10  import sklearn.neighbors as skl_nb
11  import sklearn.model_selection as skl_ms
12  from sklearn.preprocessing import StandardScaler
13  from sklearn.metrics import f1_score
14  import seaborn as sns
15
16  data = pd.read_csv('/content/training_data.csv', dtype={'ID':
     ↪  str}).dropna().reset_index(drop=True)
17
18  # Map categorical values to numerical values
19  data['increase_stock'] = data['increase_stock'].map({'low_bike_demand': 0,
     ↪  'high_bike_demand': 1})
20
21
22  data.shape
23  display(data)
24  data.info()
25  data.describe()
26
27  # Count occurrences of increase_stock=1 for each hour_of_day
28  hourly_counts_all = data[data['increase_stock'] ==
     ↪  1]['hour_of_day'].value_counts().sort_index()
29
30  # Create a Series containing counts of increase_stock=1 for all hours
     ↪  (including those with count 0)
31  all_hours_counts = hourly_counts_all.reindex(range(24), fill_value=0)
32
```

```python
33  # Create the bar plot
34  plt.figure(figsize=(12, 6))
35  sns.barplot(x=all_hours_counts.index, y=all_hours_counts.values)
36  plt.title('Count of increase_stock=1 per hour of the day (All hours)')
37  plt.xlabel('Hour of the Day')
38  plt.ylabel('Count where increase_stock=1')
39  plt.show()
40
41  # Filter the DataFrame for increase_stock=1
42  increase_stock_1 = data[data['increase_stock'] == 1]
43
44  # Drop rows with missing values in the 'day_of_week' column
45  increase_stock_1 = increase_stock_1.dropna(subset=['day_of_week'])
46
47  # Count occurrences of increase_stock=1 for each day_of_week
48  day_of_week_counts =
    ↪  increase_stock_1['day_of_week'].value_counts().sort_index()
49
50  # Create the bar plot
51  plt.figure(figsize=(12, 6))
52  sns.barplot(x=day_of_week_counts.index, y=day_of_week_counts.values)
53  plt.title('Count of high bike demand per day_of_week')
54  plt.xlabel('day_of_week')
55  plt.ylabel('Count where increase_stock=1')
56  plt.show()
57
58  # Filter the DataFrame for increase_stock=1
59  increase_stock_1 = data[data['increase_stock'] == 1]
60
61  # Count occurrences of increase_stock=1 for each temp
62  temp_counts = increase_stock_1['temp'].value_counts().sort_index()
63
64  # Find the range of temp values in the dataset
65  temp_range = range(int(data['temp'].min()), int(data['temp'].max()) + 1)
66
67  # Fill in missing temp counts with zeros
68  for temp in temp_range:
69      if temp not in temp_counts.index:
70          temp_counts[temp] = 0
71
72  # Sort the index after adding the missing values
73  temp_counts = temp_counts.sort_index()
74
75  # Create the bar plot
76  plt.figure(figsize=(60, 6))
77  sns.barplot(x=temp_counts.index, y=temp_counts.values)
78  plt.title('Count of high bike demand per temperature')
79  plt.xlabel('Temperature')
80  plt.ylabel('Count where increase_stock=1')
81  plt.show()
82
83  # Filter the DataFrame for increase_stock=1
84  increase_stock_1 = data[data['increase_stock'] == 1]
85
86  # Count occurrences of increase_stock=1 for each dew
87  dew_counts = increase_stock_1['dew'].value_counts().sort_index()
88
89  # Find the range of dew values in the dataset
90  dew_range = range(int(data['dew'].min()), int(data['dew'].max()) + 1)
91
92  # Fill in missing dew counts with zeros
93  for dew in dew_range:
94      if dew not in dew_counts.index:
95          dew_counts[dew] = 0
```

```
96
97   # Sort the index after adding the missing values
98   dew_counts = dew_counts.sort_index()
99
100  # Create the bar plot
101  plt.figure(figsize=(60, 6))
102  sns.barplot(x=dew_counts.index, y=dew_counts.values)
103  plt.title('Count of high bike demand per dew')
104  plt.xlabel('Dew')
105  plt.ylabel('Count where increase_stock=1')
106  plt.show()
107
108  # Filter the DataFrame for increase_stock=1
109  increase_stock_1 = data[data['increase_stock'] == 1]
110
111  # Count occurrences of increase_stock=1 for each humidity
112  humidity_counts = increase_stock_1['humidity'].value_counts().sort_index()
113
114  # Find the range of humidity values in the dataset
115  humidity_range = range(int(data['humidity'].min()), int(data['humidity'].max())
     ↪  + 1)
116
117  # Fill in missing humidity counts with zeros
118  for humidity in humidity_range:
119      if humidity not in humidity_counts.index:
120          humidity_counts[humidity] = 0
121
122  # Sort the index after adding the missing values
123  humidity_counts = humidity_counts.sort_index()
124
125  # Create the bar plot
126  plt.figure(figsize=(200, 6))
127  sns.barplot(x=humidity_counts.index, y=humidity_counts.values)
128  plt.title('Count of high bike demand per humidity')
129  plt.xlabel('Humidity')
130  plt.ylabel('Count where increase_stock=1')
131  plt.show()
132
133  # Filter the DataFrame for increase_stock=1
134  increase_stock_1 = data[data['increase_stock'] == 1]
135
136  # Count occurrences of increase_stock=1 for each windspeed
137  windspeed_counts = increase_stock_1['windspeed'].value_counts().sort_index()
138
139  # Find the range of windspeed values in the dataset
140  windspeed_range = range(int(data['windspeed'].min()),
     ↪  int(data['windspeed'].max()) + 1)
141
142  # Fill in missing windspeed counts with zeros
143  for windspeed in windspeed_range:
144      if windspeed not in windspeed_counts.index:
145          windspeed_counts[windspeed] = 0
146
147  # Sort the index after adding the missing values
148  windspeed_counts = windspeed_counts.sort_index()
149
150  # Create the bar plot
151  plt.figure(figsize=(60, 6))
152  sns.barplot(x=windspeed_counts.index, y=windspeed_counts.values)
153  plt.title('Count of high bike demand per windspeed')
154  plt.xlabel('Windspeed')
155  plt.ylabel('Count where increase_stock=1')
156  plt.show()
157
```

```python
158  # Define the number of bins and create bins for visibility
159  num_bins = 10   # You can adjust the number of bins as needed
160  visibility_bins = pd.cut(data['visibility'], bins=num_bins)
161
162  # Filter the DataFrame for increase_stock=1
163  increase_stock_1 = data[data['increase_stock'] == 1]
164
165  # Count occurrences of increase_stock=1 for each visibility bin
166  visibility_counts =
     ↪  increase_stock_1.groupby(visibility_bins)['increase_stock'].count()
167
168  # Create the bar plot
169  plt.figure(figsize=(12, 6))
170  visibility_counts.plot(kind='bar')
171  plt.title('Count of increase_stock=1 for each visibility bin')
172  plt.xlabel('Visibility Bins')
173  plt.ylabel('Count where increase_stock=1')
174  plt.xticks(rotation=45)   # Rotate x-axis labels for better readability
175  plt.show()
176
177  # Define the number of bins and create bins for snowdepth
178  num_bins = 100   # You can adjust the number of bins as needed
179  snowdepth_bins = pd.cut(data['snowdepth'], bins=num_bins)
180
181  # Filter the DataFrame for increase_stock=1
182  increase_stock_1 = data[data['increase_stock'] == 1]
183
184  # Count occurrences of increase_stock=1 for each snowdepth bin
185  snowdepth_counts =
     ↪  increase_stock_1.groupby(snowdepth_bins)['increase_stock'].count()
186
187  # Create the bar plot
188  plt.figure(figsize=(60, 6))
189  snowdepth_counts.plot(kind='bar')
190  plt.title('Count of increase_stock=1 for each snowdepth bin')
191  plt.xlabel('Snowdepth Bins')
192  plt.ylabel('Count where increase_stock=1')
193  plt.xticks(rotation=45)   # Rotate x-axis labels for better readability
194  plt.show()
195
196  # Define the number of bins and create bins for cloudcover
197  num_bins = 30   # You can adjust the number of bins as needed
198  cloudcover_bins = pd.cut(data['cloudcover'], bins=num_bins)
199
200  # Filter the DataFrame for increase_stock=1
201  increase_stock_1 = data[data['increase_stock'] == 1]
202
203  # Count occurrences of increase_stock=1 for each cloudcover bin
204  cloudcover_counts =
     ↪  increase_stock_1.groupby(cloudcover_bins)['increase_stock'].count()
205
206  # Create the bar plot
207  plt.figure(figsize=(12, 6))
208  cloudcover_counts.plot(kind='bar')
209  plt.title('Count of increase_stock=1 for each cloudcover bin')
210  plt.xlabel('Cloudcover Bins')
211  plt.ylabel('Count where increase_stock=1')
212  plt.xticks(rotation=45)   # Rotate x-axis labels for better readability
213  plt.show()
214
215  !pip install ydata_profiling
216  !pip install lida==0.0.10 kaleido python-multipart uvicorn lmx==0.0.15a0
     ↪  tensorflow-probability==0.22.0
217
```

```python
218  from ydata_profiling import ProfileReport
219
220
221  profile = ProfileReport(data,title="Bike sharing data report")
222  profile.to_notebook_iframe()
223  profile.to_file("eda.html")
```