
Bike Sharing Demand

Phudish Prateepamornkul
Department of Statistics
University of California, Berkeley
Berkeley, CA 94704
phudish_p@berkeley.edu

Abstract

Since the rise of the demand of the bike-sharing scheme especially in Seoul. Where it was reported that in 2021, the daily rental number amounted to 87.8 thousand, up from about 65 thousand in the previous year. Additionally, since the start of the service in 2015, the rate has increase by more than twice every year [1]?. Therefore, it is essential to make the rental bike available and accessible to the public. The purpose of this paper is therefore to estimate the demand of the bike rental hourly in Seoul. The research involved training six unique statistical regression models including 1.Linear Regression 2. Support Vector Machine 3. Decision Tree 4. Gradient Boosting Machine 5. eXtreme Gradient Boosting and 6. Neural Network on a training dataset. Additionally, the study utilized the grid search technique for identifying optimal hyper parameters through cross-validation. These optimal hyper parameters were then applied to assess the models' performance on a separate test dataset. The most successful outcome was observed with the Gradient Boosting Machine, which has a root mean square training error of 37.38 and an R^2 value of 0.99 in the train data set, a testing error of 157.23, and an R^2 value of 0.94 in the test data set.

1 Background

Bike sharing has become increasingly popular in urban cities as a convenient, affordable and environmentally friendly mode of transportation. It operates as a shared service accessible to individuals by allowing user to borrow a bike from one station and return it to another station.

Currently, Seoul has introduced Ttareungi which is the public bicycle sharing system that can be used by anywhere, anyone and any time. This application has been one of the most popular mode for transportation especially during COVID where people are trying to avoid the crowded of people in public transport such as train and buses. It has been reported that the number of Ttareungi users increased by 24% year on year to 2.78 million in 2020 [2]. Additionally, the bicycle sharing program in Seoul has produces more than 1 millions rides a month with more than 1500 bicycle stations and 20000 bicycles operating in the city [2]. Hence, it is very crucial and important for us to be able to predict the rented bike count accurately to meet the people demand.

2 Data Preparation

2.1 Information of the data

The data set use is publicly available at UC Irvine Machine Learning Repository website <http://archive.ics.uci.edu/dataset/560/seoul+bike+sharing+demand>. The dataset consists of hourly rented bike count one year from 2017 December to 2018 November from Seoul which contains 8760 rows and 14 columns including Date, Rented Bike Count, Hour,

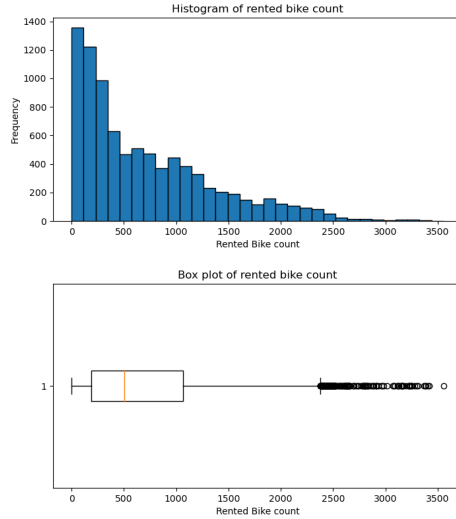


Figure 1: (a) Histogram of Rented Bike count (b) Box plot of Rented Bike count

Temperature($^{\circ}\text{C}$), Humidity(%), Wind speed(m/s), Visibility(10m), Dew point temperature($^{\circ}\text{C}$), Solar Radiation(MJ/m 2), Rainfall(mm), Snowfall(cm), Seasons, Holiday, Functioning Day. Where Seasons(Winter, Spring, Summer, Autumn), Holiday(Yes, No), Functioning Day(Yes, No) are categorical data while the rest are continuous data.

2.2 Exploratory Data Analysis

We want to explore our data first before we do any modelling. Figure 1 shows the Histogram and Box plot of the rented bike count. Clearly, we see that the data is right skewed and there are a few outlier show in the box plot. Therefore, before we put our data into the model we should do some transformation to our target variable first. Figure 2 shows the Histogram of different numerical value in our data. This time we see that most of the features are skew. So we should apply standardization to our features.

2.3 Data preparation and engineering

We started by transforming our response variable first since it is right skew we will apply square root to each of the response so that the distribution look more like normal distribution as shown in Figure 3. Additionally, we should try to utilize most of our features first. We created two more features which are 1. Month (look at the Month of the corresponding Date) 2. Week Type (which is categorical type to look at whether at a given day is it a weekend or weekday). Once we did that then we start to drop our target variable in this case Rented Bike Count column and work with the rest of the features. When dealing with categorical features such as Seasons, Holiday, Functioning Day and Week Type, we use one-hot encoding. Which creates new columns for each category and fill a value with 1 if the record belongs to that category and 0 otherwise. After handling the categorical data we focused on our numerical features. We applied standardization to our numerical features by subtract the mean of the features from each value and divided by the standard deviation. Note that we did not apply standardization to the month. We then split our dataset into training and testing where 25% of the data is use for testing. We have 6570 rows for training data and 2190 rows for testing data.

3 Models

Since we are tackling a regression problem, our chosen metrics for evaluating training and testing errors include 1. Root Mean Square Error (RMSE), 2. R^2 score, and 3. Mean Absolute Error (MAE). To optimize each model, we implemented grid search to determine the best possible input parameters. This process involved using 5-fold cross-validation and negative mean squared error as the scoring

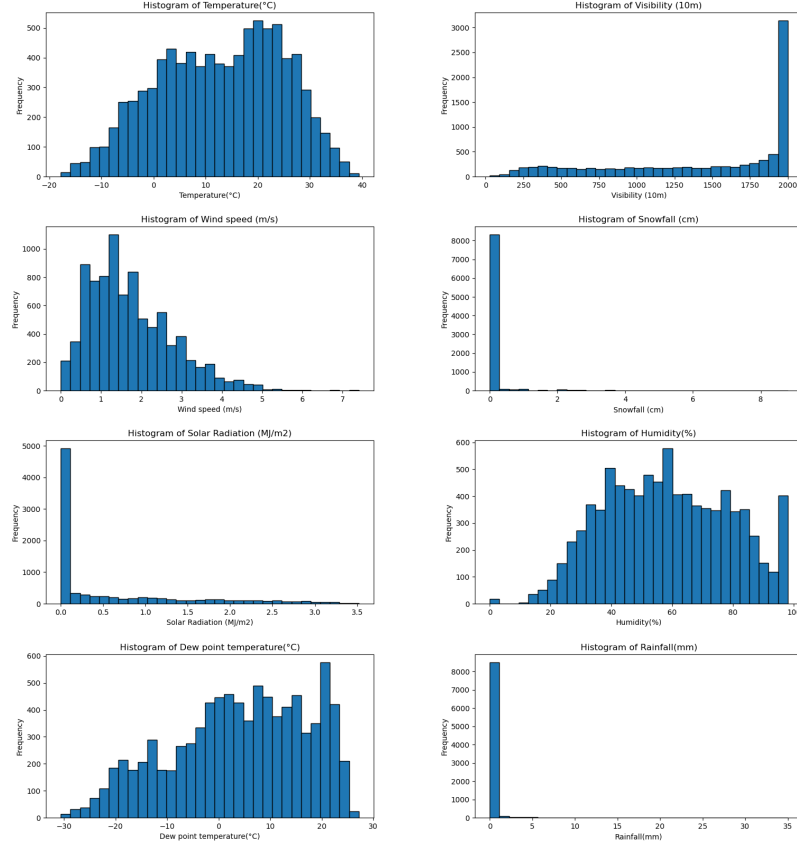


Figure 2: Histogram of numerical features

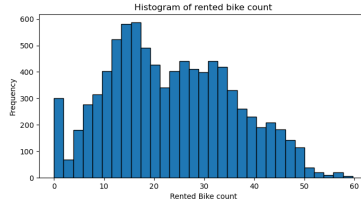


Figure 3: Histogram of square root of response variable

metric. After identifying the optimal parameters, we trained the models on the training set. The way in which we find the grid search is first we tried a few values where each of them is pretty far from each other first. Then after we find the best one in that set of values then we try again using the best one from the previous time but now having each of them to be closer to each other.

3.1 Linear Regression

Linear Regression started by assuming a linear relationship between our target(Rented Bike Count) which we will call it as y and our feature which we will call it as X . We then want to find a vector of β that minimize $(y - X\beta)^2$. Since there is no parameter we have to input we simply fit our data and report the result we got.

3.2 Support Vector Machine (SVM)

SVM is use for regression by transforming data into a higher dimension feature space so that the we can apply linear regression in this new feature space. We have that there are many kernel we can

use to map our data into a higher dimension including linear, polynomial, sigmoid and radial basis function (RBF). In our case we will fix our SVM kernel to be RBF. Additionally, we are going to tune two parameters which are 1. The regularization parameter which tells us how the model would handle errors between predicted values and the actual values in training and 2. Gamma (γ) which is the coefficient in the RBF kernel. After, we did grid search we found that the best parameters are 3200 and 0.5 respectively.

3.3 Decision Tree

Decision Tree regression is a flowchart like tree structure that work by making predictions based on how the previous set of questions were answered. Typically it started by having a root node which branches into possible outcomes. Each of those outcomes also lead to additional node which also branch off into other possibilities. We have that the internal node is the node that symbolizes a choice of our input feature and a leaf node is a node without any child nodes that indicates a numerical value. We are going to tune three parameters which are 1. The maximum depth of the tree 2. The minimum number of samples required to split an internal node and 3. The minimum number of samples required to be at a leaf node. We did grid search and found that the best parameters are 16, 7 and 5 respectively.

3.4 Gradient Boosting Machine (GBM)

Gradient Boosting Machine work by combining a weak learners into a single a stronger learner usually the weak learners are the models that make very few assumptions about the data for example decision trees. This is done by fitting new tree to a residual errors made by existing trees. We then go to the next step by fitting the new tree and then our model will become the first and second tree combine together. We then grow a new tree to fit the residual errors from the combined trees. We are going to tune 5 parameters including 1. number of boosting stages we want to perform 2. the maximum depth of individual regression model 3. the minimum number of samples required to split an internal node 4. the minimum number of samples required at a leaf node and 5. the fraction of samples to be use for fitting the individual base learners. After, grid search we found out that the best parameters are 480, 19, 200, 10 and 0.9 respectively.

3.5 eXtreme Gradient Boosting (XGBoost)

XGBoost applied Gradient Boosting technique like before but the different is that it can use L_1 and L_2 regularization which can help the model to prevent over fitting. We are going to tune 5 parameters including 1. learning rate, 2. maximum tree depth for base learners, 3. L_1 regularization term on weights, 4. L_2 regularization term on weights and 5. number of gradient boosted trees. We found that after grid search the best parameters are 0.1, 45, 15, 27 and 240 respectively.

3.6 Neural Network

We created a simple feed forward Neural Network using ReLU as an activation function. Additionally, we have 3 hidden layers for our model. After grid search we found that the number of node in first, second and third layers are 90, 50 and 40 respectively.

4 Result

We compare 6 models by reporting their corresponding performance using RMSE during cross-validation, and RMSE, R^2 , and MAE on both the training and testing datasets. Note that the RMSE GridSearch (GS) is not the same as the RMSE Training because the RMSE GS is on the cross-validation instead of the whole dataset like the Training error. Additionally, since we transform our response variable by taking the square root which mean that in our prediction for both training and testing datasets we must square our outcome and the response variable first before we can put it in the metric that use to measure the error. However, note that in the RMSE GS we still report the value of error before squaring the prediction and the actual value.

Model	RMSE (GS)	R ² (Train)	MAE (Train)	RMSE (Train)	R ² (Test)	MAE (Test)	RMSE (Test)
Linear Regression	-	0.58	280.94	419.74	0.60	271.92	399.62
SVM	4.49	0.99	4.63	5.28	0.88	127.84	214.85
Decision Tree	4.08	0.96	77.93	132.72	0.87	128.59	225.78
GBM	2.78	0.99	23.84	37.38	0.94	85.44	157.23
XGBoost	2.74	0.99	30.86	51.77	0.94	86.19	157.96
Neural Network	4.68	0.70	219.01	355.99	0.69	218.03	352.79

Table 1: Model Evaluation Metrics (GS: GridSearch, Train: Training, Test: Testing)

5 Discussion

From the table we can see that the best model we achieved is the GBM. This is something that is not surprisingly since GBM has been one of the best machine learning model when it come to tabular data. However, one thing that is surprising is the bad performance of the neural network. I think there are probably two reasons for this one being that we are using tabular data which make neural network face many challenges including mixed feature types of our data and a lack of knowledge of the data structure unlike image and text. Secondly, the size of our dataset is pretty small as it is only 8760 rows only. We then compare our result to this paper [3] where they were able to achieved R^2 of 0.96 , RMSE 117.81 and MAE 79.77 for training data and R^2 of 0.92 , RMSE 172.73 and MAE 109.78 for testing data using GBM. While their XGBoost did second best which achieved R^2 of 0.96 , RMSE 127.63 and MAE 85.21 for training data and R^2 of 0.91 , RMSE 183.80 and MAE 119.59 for testing data. This is very similar result to what we got that GBM and XGBoost have been the best models to achieved the best result. I think that there are a few things that we did differently, one being that we have one more extra feature which is the month we get from the date. Additionally, we also perform square root transformation for our response variable and standardization to most of our numerical columns except for the month which we do not think it is appropriate to do so. When we perform transformation to our response we emphasis our loss on smaller value more than the large value because larger value are reduced more significantly than smaller value so during the loss the smaller value will be emphasis more. This is particular useful since most of our values are pretty small. Lastly, we also have that we did more grid search parameters than what they did in their paper. For example, the paper did grid search for boosting iterations and max tree depth for GBM and XGBoost. However, for our case we did more grid search parameters than what they did which make our result better. To improve our result and model I think that we could try to tune more parameters of the input or we can perform more feature selection by removing certain features and test whether the result improve or not.

References

- [1] Korea Herald. (January 19, 2022). Average daily number of Seoul Bike (public bike rental system) rentals in South Korea from 2015 to 2021 (in 1,000s) [Graph]. In Statista. Retrieved November 27, 2023, from <https://www.statista.com/statistics/997380/south-korea-seoul-bike-daily-rental-number/>
- [2] Stellar Market Research. (n.d.). South Korea Bike Sharing Market: Industry Analysis and Forecast (2023-2029) by Component, Level of Automation, Fuel Type, and Region. Retrieved November 28, 2023, from <https://www.stellarmr.com/report/South-Korea-Bike-Sharing-Market/252>.
- [3] V E, S., Park, J., Cho, Y. (2020). Using data mining techniques for bike sharing demand prediction in metropolitan city. Computer Communications, 153, 353-366. <https://doi.org/10.1016/j.comcom.2020.02.007>

Appendix and Code

All of the dataset and code are available at https://github.com/PhudishTam/STAT254_HW_6.gitcode

```

1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 df = pd.read_csv("SeoulBikeData.csv", encoding='unicode_escape')

```

```

5 plt.figure(figsize=(8,4))
6 plt.hist(df["Rented Bike Count"],bins=31,edgecolor="black")
7 plt.xlabel("Rented Bike count")
8 plt.ylabel("Frequency")
9 plt.title("Histogram of rented bike count")
10 plt.savefig("histogram")
11
12 plt.figure(figsize=(8,4))
13 plt.boxplot(df["Rented Bike Count"],vert=False)
14 plt.xlabel("Rented Bike count")
15 plt.title("Box plot of rented bike count")
16 plt.savefig('boxplot')
17 plt.show()
18
19 plt.figure(figsize=(8,4))
20 plt.hist(df["Temperature( C)"],bins=31,edgecolor="black")
21 plt.xlabel("Temperature( C)")
22 plt.ylabel("Frequency")
23 plt.title("Histogram of Temperature( C)")
24 plt.savefig('histogram_temperature')
25
26 plt.figure(figsize=(8,4))
27 plt.hist(df["Wind speed (m/s)"],bins=31,edgecolor="black")
28 plt.xlabel("Wind speed (m/s)")
29 plt.ylabel("Frequency")
30 plt.title("Histogram of Wind speed (m/s)")
31 plt.savefig('histogram_wind_speed')
32
33 plt.figure(figsize=(8,4))
34 plt.hist(df["Visibility (10m)"],bins=31,edgecolor="black")
35 plt.xlabel("Visibility (10m)")
36 plt.ylabel("Frequency")
37 plt.title("Histogram of Visibility (10m)")
38 plt.savefig('histogram_visibility')
39
40 plt.figure(figsize=(8,4))
41 plt.hist(df["Dew point temperature( C)"],bins=31,edgecolor="black")
42 plt.xlabel("Dew point temperature( C)")
43 plt.ylabel("Frequency")
44 plt.title("Histogram of Dew point temperature( C)")
45 plt.savefig('histogram_dew_point_temperature')
46
47 plt.figure(figsize=(8,4))
48 plt.hist(df["Solar Radiation (MJ/m2)"],bins=31,edgecolor="black")
49 plt.xlabel("Solar Radiation (MJ/m2)")
50 plt.ylabel("Frequency")
51 plt.title("Histogram of Solar Radiation (MJ/m2)")
52 plt.savefig('histogram_solar_radiation')
53
54 plt.figure(figsize=(8,4))
55 plt.hist(df["Rainfall(mm)"],bins=31,edgecolor="black")
56 plt.xlabel("Rainfall(mm)")
57 plt.ylabel("Frequency")
58 plt.title("Histogram of Rainfall(mm)")
59 plt.savefig('histogram_rainfall')
60
61 plt.figure(figsize=(8,4))
62 plt.hist(df["Snowfall (cm)"],bins=31,edgecolor="black")
63 plt.xlabel("Snowfall (cm)")
64 plt.ylabel("Frequency")
65 plt.title("Histogram of Snowfall (cm)")
66 plt.savefig('histogram_snowfall')
67
68 plt.figure(figsize=(8,4))
69 plt.hist(df["Humidity(%)"],bins=31,edgecolor="black")

```

```

70 plt.xlabel("Humidity(%)")
71 plt.ylabel("Frequency")
72 plt.title("Histogram of Humidity(%)")
73 plt.savefig('histogram_humidity')
74
75 df['square'] = np.sqrt(df['Rented Bike Count'].values)
76 plt.figure(figsize=(8,4))
77 plt.hist(df["square"],bins=31,edgecolor="black")
78 plt.xlabel("Rented Bike count")
79 plt.ylabel("Frequency")
80 plt.title("Histogram of rented bike count")
81 plt.savefig("histogram_root")

```

Listing 1: Download all of our plot (image.py)

```

1 import numpy as np
2 import pandas as pd
3 from datetime import datetime
4 from sklearn.preprocessing import StandardScaler
5 from sklearn.metrics import r2_score, mean_squared_error,
    mean_absolute_error
6 from sklearn.model_selection import GridSearchCV
7
8 def preprocess_data(df):
9     df['Month'] = pd.to_datetime(df['Date'], format='%d/%m/%Y').dt.
    month
10    df['Week_Type'] = df['Date'].apply(lambda x: "Weekend" if datetime
    .strptime(x, '%d/%m/%Y').weekday() > 4 else "Weekday")
11
12    X = df.drop(['Date', 'Rented Bike Count'], axis=1)
13    y = df['Rented Bike Count']
14    y = np.sqrt(y)
15
16    cat_columns = ['Seasons', 'Holiday', 'Functioning Day', 'Week_Type
    ,']
17    for col in cat_columns:
18        X[col] = X[col].astype('category')
19
20    X = pd.get_dummies(X, columns=cat_columns, dtype=int)
21    return X, y
22
23 def scale_features(X_train, X_test):
24     numeric_columns = ['Temperature( C )', 'Humidity(%)', 'Dew point
    temperature( C )',
25                        'Solar Radiation (MJ/m2)', 'Rainfall(mm)', '
    Visibility (10m)',
26                        'Snowfall (cm)', 'Wind speed (m/s)']
27     scaler = StandardScaler()
28     X_train[numeric_columns] = scaler.fit_transform(X_train[
    numeric_columns])
29     X_test[numeric_columns] = scaler.transform(X_test[numeric_columns
    ])
30     return X_train, X_test
31
32 def tune_hyperparameters(estimator, param_grid, X_train, y_train):
33     grid_search = GridSearchCV(estimator=estimator, param_grid=
    param_grid, scoring='neg_mean_squared_error', cv=5)
34     grid_search.fit(X_train, y_train)
35     root_mean_square_error_train = -grid_search.best_score_
36     return grid_search.best_params_, root_mean_square_error_train
37
38 def evaluate_model(model, X_train, y_train, X_test, y_test):
39     model.fit(X_train, y_train)
40     y_pred_train_root = model.predict(X_train)
41     y_pred_train = (y_pred_train_root)**2

```

```

42     r2_score_train = r2_score(y_train**2, y_pred_train)
43     mean_absolute_error_score_train = mean_absolute_error(y_train**2,
44     y_pred_train)
44     root_mean_square_error_train = np.sqrt(mean_squared_error(y_train
45     **2, y_pred_train))
45     y_pred_test_root = model.predict(X_test)
46     y_pred_test = (y_pred_test_root)**2
47     r2_score_test = r2_score(y_test**2, y_pred_test)
48     mean_absolute_error_score_test = mean_absolute_error(y_test**2,
49     y_pred_test)
49     mean_squared_error_test = mean_squared_error(y_test**2,
50     y_pred_test)
50     root_mean_square_error_test = np.sqrt(mean_squared_error_test)
51     return (r2_score_train, mean_absolute_error_score_train,
52     root_mean_square_error_train,
53     r2_score_test, mean_absolute_error_score_test,
54     root_mean_square_error_test)

```

Listing 2: Preprocess the data (preprocess.py)

```

1  import numpy as np
2  import pandas as pd
3  import time
4  from datetime import datetime
5  from sklearn.model_selection import train_test_split
6  from sklearn.preprocessing import StandardScaler, OneHotEncoder
7  from sklearn.metrics import r2_score, mean_squared_error,
8  mean_absolute_error
9  from sklearn.linear_model import LinearRegression
10 from preprocess import preprocess_data, scale_features,
11 tune_hyperparameters, evaluate_model
12
13 df = pd.read_csv("SeoulBikeData.csv", encoding='unicode_escape')
14 X, y = preprocess_data(df)
15 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size
16 =0.25, random_state=60)
17 X_train, X_test = scale_features(X_train, X_test)
18 r2_score_train, mean_absolute_error_score_train,
19 root_mean_square_error_train, r2_score_test,
20 mean_absolute_error_score_test, root_mean_square_error_test =
21 evaluate_model(LinearRegression(), X_train, y_train, X_test, y_test)
22
23 print(f"The r2 score on training is : {r2_score_train}")
24 print(f"The mean absolute error training is : {
25     mean_absolute_error_score_train}")
26 print(f"The root mean square error training is : {
27     root_mean_square_error_train}")
28 print(f"The r2 score on testing is : {r2_score_test}")
29 print(f"The mean absolute error testing is : {
30     mean_absolute_error_score_test}")
31 print(f"The root mean square error testing is : {
32     root_mean_square_error_test}")

```

Listing 3: Linear Model (lm.py)

```

1  import numpy as np
2  import pandas as pd
3  import time
4  from datetime import datetime
5  from sklearn.model_selection import train_test_split
6  from sklearn.preprocessing import StandardScaler, OneHotEncoder
7  from sklearn.metrics import r2_score, mean_squared_error,
8  mean_absolute_error
9  from sklearn.svm import SVR
10 from dask.distributed import Client, LocalCluster
11 from dask_ml.model_selection import GridSearchCV

```



```

11 from preprocess import preprocess_data, scale_features,
    tune_hyperparameters, evaluate_model
12
13 if __name__ == "__main__":
14     df = pd.read_csv("SeoulBikeData.csv", encoding='unicode_escape')
15     X, y = preprocess_data(df)
16     X_train, X_test, y_train, y_test = train_test_split(X, y,
17                                                         test_size=0.25, random_state=60)
18     X_train, X_test = scale_features(X_train, X_test)
19     cluster = LocalCluster(n_workers=32, threads_per_worker=3)
20     client = Client(cluster)
21     t0 = time.time()
22     estimator = SVR(kernel = 'rbf')
23     param_test_1 = {'C': [3200, 3500, 3800], 'gamma': [1.5, 1, 0.5]}
24     best_params_svm, best_score_svm = tune_hyperparameters(estimator,
25                                                         param_test_1, X_train, y_train)
26     print(f"Best parameters: {best_params_svm}")
27     print(f"The root mean square error training is: {np.sqrt(
28         best_score_svm)}")
29     best_model = SVR(C= best_params_svm['C'], gamma = best_params_svm["
30         gamma"])
31     r2_score_train, mean_absolute_error_score_train,
32     root_mean_square_error_train, r2_score_test,
33     mean_absolute_error_score_test, root_mean_square_error_test =
34     evaluate_model(best_model, X_train, y_train, X_test, y_test)
35     print(f"Time it took: {time.time() - t0}")
36     print(f"The r2 score on training is : {r2_score_train}")
37     print(f"The mean absolute error training is : {
38         mean_absolute_error_score_train}")
39     print(f"The root mean square error training is : {
40         root_mean_square_error_train}")
41     print(f"The r2 score on testing is : {r2_score_test}")
42     print(f"The mean absolute error testing is : {
43         mean_absolute_error_score_test}")
44     print(f"The root mean square error testing is : {
45         root_mean_square_error_test}")
46     cluster.close()

```

Listing 4: SVM Model (svm.py)

```

1 import numpy as np
2 import pandas as pd
3 import time
4 from datetime import datetime
5 from sklearn.model_selection import train_test_split
6 from sklearn.preprocessing import StandardScaler, OneHotEncoder
7 from sklearn.metrics import r2_score, mean_squared_error,
    mean_absolute_error
8 from sklearn.tree import DecisionTreeRegressor
9 from dask.distributed import Client, LocalCluster
10 from dask_ml.model_selection import GridSearchCV
11 from preprocess import preprocess_data, scale_features,
    tune_hyperparameters, evaluate_model
12
13 if __name__ == "__main__":
14     df = pd.read_csv("SeoulBikeData.csv", encoding='unicode_escape')
15     X, y = preprocess_data(df)
16     X_train, X_test, y_train, y_test = train_test_split(X, y,
17                                                         test_size=0.25, random_state=60)
18     X_train, X_test = scale_features(X_train, X_test)
19     cluster = LocalCluster(n_workers=32, threads_per_worker=3)
20     client = Client(cluster)
21     t0 = time.time()
22     estimator = DecisionTreeRegressor()

```

```

22     param_test_1 = {'max_depth': range(1, 100), 'min_samples_split':
23     range(2,10), 'min_samples_leaf': range(1,9)}
24     best_params, best_score = tune_hyperparameters(estimator,
25     param_test_1, X_train, y_train)
26     print(f"Best parameters: {best_params}")
27     print(f"The The root mean square error training is: {np.sqrt(
28     best_score)}")
29     best_model = DecisionTreeRegressor(max_depth = best_params['
30     max_depth'], min_samples_split = best_params['min_samples_split'],
31     min_samples_leaf = best_params['min_samples_leaf'])
32     r2_score_train, mean_absolute_error_score_train,
33     root_mean_square_error_train, r2_score_test,
34     mean_absolute_error_score_test, root_mean_square_error_test =
35     evaluate_model(best_model, X_train, y_train, X_test, y_test)
36     print(f"Time it took: {time.time() - t0}")
37     print(f"The r2 score on training is : {r2_score_train}")
38     print(f"The mean absolute error training is : {
39     mean_absolute_error_score_train}")
40     print(f"The root mean square error training is : {
41     root_mean_square_error_train}")
42     print(f"The r2 score on testing is : {r2_score_test}")
43     print(f"The mean absolute error testing is : {
44     mean_absolute_error_score_test}")
45     print(f"The root mean square error testing is : {
46     root_mean_square_error_test}")
47     cluster.close()

```

Listing 5: Decision Tree Model (decision_tree.py)

```

1  import numpy as np
2  import pandas as pd
3  import time
4  from datetime import datetime
5  from sklearn.model_selection import train_test_split
6  from sklearn.preprocessing import StandardScaler, OneHotEncoder
7  from sklearn.metrics import r2_score, mean_squared_error,
8  mean_absolute_error
9  from sklearn.ensemble import GradientBoostingRegressor
10 from dask.distributed import Client, LocalCluster
11 from dask_ml.model_selection import GridSearchCV
12 from preprocess import preprocess_data, scale_features,
13 tune_hyperparameters, evaluate_model
14
15 if __name__ == "__main__":
16     df = pd.read_csv("SeoulBikeData.csv", encoding='unicode_escape')
17     X, y = preprocess_data(df)
18     X_train, X_test, y_train, y_test = train_test_split(X, y,
19     test_size=0.25, random_state=60)
20     X_train, X_test = scale_features(X_train, X_test)
21     cluster = LocalCluster(n_workers=32, threads_per_worker=3)
22     client = Client(cluster)
23     t0 = time.time()
24     estimator = GradientBoostingRegressor(learning_rate=0.1,
25     min_samples_split=500, min_samples_leaf=50, max_depth=8, subsample
26     =0.8, random_state=10)
27     param_test1 = {'n_estimators': range(20,500, 20)}
28     best_params_n_estimators, best_score_n_estimators =
29     tune_hyperparameters(estimator, param_test1, X_train, y_train)
30     print(f"Best parameters: {best_params_n_estimators}")
31     print(f"The root mean square error training is: {np.sqrt(
32     best_score_n_estimators)}")
33     estimator = GradientBoostingRegressor(learning_rate=0.1,
34     n_estimators=best_params_n_estimators['n_estimators'],
35     min_samples_leaf = 50, subsample=0.8, random_state=10)

```

```

27     param_test2 = {'max_depth': range(5, 20, 2), 'min_samples_split':
range(200, 1500, 200)}
28     best_params_max_depth_min_samples_split,
best_score_max_depth_min_samples_split = tune_hyperparameters(
estimator, param_test2, X_train, y_train)
29     print(f"Best parameters: {best_params_max_depth_min_samples_split}")
30     print(f"The root mean square error training is: {np.sqrt(
best_score_max_depth_min_samples_split)}")
31     estimator = GradientBoostingRegressor(n_estimators=
best_params_n_estimators['n_estimators'], max_depth =
best_params_max_depth_min_samples_split["max_depth"],
min_samples_split = best_params_max_depth_min_samples_split["
min_samples_split"], random_state=10)
32     param_test3 = {'min_samples_leaf': range(10, 70, 5), 'subsample':
[0.6, 0.7, 0.75, 0.8, 0.85, 0.9]}
33     best_params_min_samples_leaf_subsample,
best_score_min_samples_leaf_subsample = tune_hyperparameters(
estimator, param_test3, X_train, y_train)
34     # feature_importances = best_model.feature_importances_
35     # features_and_importances = zip(X_train.columns,
feature_importances)
36     # sorted_features_and_importances = sorted(
features_and_importances, key=lambda x: x[1], reverse=True)
37     # print("Feature Importances:")
38     # for feature, importance in sorted_features_and_importances:
39     #     print(f"{feature}: {importance}")
40     print(f"Best parameters: {best_params_min_samples_leaf_subsample}")
41     print(f"The The root mean square error training is: {np.sqrt(
best_score_min_samples_leaf_subsample)}")
42     best_model = GradientBoostingRegressor(n_estimators=
best_params_n_estimators['n_estimators'], max_depth =
best_params_max_depth_min_samples_split["max_depth"],
min_samples_split = best_params_max_depth_min_samples_split["
min_samples_split"], min_samples_leaf=
best_params_min_samples_leaf_subsample["min_samples_leaf"],
43     subsample = best_params_min_samples_leaf_subsample['subsample'],
random_state=10)
44     r2_score_train, mean_absolute_error_score_train,
root_mean_square_error_train, r2_score_test,
mean_absolute_error_score_test, root_mean_square_error_test =
evaluate_model(best_model, X_train, y_train, X_test, y_test)
45     print(f"Time it took: {time.time() - t0}")
46     print(f"The r2 score on training is : {r2_score_train}")
47     print(f"The mean absolute error training is : {
mean_absolute_error_score_train}")
48     print(f"The root mean square error training is : {
root_mean_square_error_train}")
49     print(f"The r2 score on testing is : {r2_score_test}")
50     print(f"The mean absolute error testing is : {
mean_absolute_error_score_test}")
51     print(f"The root mean square error testing is : {
root_mean_square_error_test}")
52     cluster.close()
53
54     cluster.close()

```

Listing 6: GBM Model (gbm.py)

```

1 import numpy as np
2 import pandas as pd
3 import time
4 from datetime import datetime
5 from sklearn.model_selection import train_test_split

```

```

6 from sklearn.preprocessing import StandardScaler, OneHotEncoder
7 from sklearn.metrics import r2_score, mean_squared_error,
  mean_absolute_error
8 import xgboost as xgb
9 from dask.distributed import Client, LocalCluster
10 from dask_ml.model_selection import GridSearchCV
11 from preprocess import preprocess_data, scale_features,
  tune_hyperparameters, evaluate_model
12
13 if __name__ == "__main__":
14     df = pd.read_csv("SeoulBikeData.csv", encoding='unicode_escape')
15     X, y = preprocess_data(df)
16     X_train, X_test, y_train, y_test = train_test_split(X, y,
17     test_size=0.25, random_state=60)
18     X_train, X_test = scale_features(X_train, X_test)
19     cluster = LocalCluster(n_workers=32, threads_per_worker=3)
20     client = Client(cluster)
21     t0 = time.time()
22     estimator = xgb.XGBRegressor()
23     param_test_1 = {'learning_rate': [1,0.1], 'max_depth'
24     : [45,50,55,60], 'n_estimators': [200,220,240], 'reg_alpha':
25     [15,16,17,18], 'reg_lambda': [25,27,29,31]}
26     best_params, best_score = tune_hyperparameters(estimator,
27     param_test_1, X_train, y_train)
28     print(f"Best parameters: {best_params}")
29     print(f"The root mean square error training is: {np.sqrt(
30     best_score)}")
31     best_model = xgb.XGBRegressor(learning_rate = best_params['
32     learning_rate'], max_depth = best_params['max_depth'], n_estimators
33     = best_params['n_estimators'], reg_alpha = best_params['reg_alpha']
34     , reg_lambda = best_params['reg_lambda'] )
35     r2_score_train, mean_absolute_error_score_train,
36     root_mean_square_error_train, r2_score_test,
37     mean_absolute_error_score_test, root_mean_square_error_test =
38     evaluate_model(best_model, X_train, y_train, X_test, y_test)
39     print(f"Time it took: {time.time() - t0}")
40     print(f"The r2 score on training is : {r2_score_train}")
41     print(f"The mean absolute error training is : {
42     mean_absolute_error_score_train}")
43     print(f"The root mean square error training is : {
44     root_mean_square_error_train}")
45     print(f"The r2 score on testing is : {r2_score_test}")
46     print(f"The mean absolute error testing is : {
47     mean_absolute_error_score_test}")
48     print(f"The root mean square error testing is : {
49     root_mean_square_error_test}")
50     cluster.close()

```

Listing 7: XGB Model (xgb.py)

```

1 import numpy as np
2 import pandas as pd
3 import time
4 from sklearn.model_selection import train_test_split
5 from sklearn.metrics import r2_score, mean_squared_error,
  mean_absolute_error
6 import torch
7 from dask_ml.model_selection import GridSearchCV
8 from preprocess import preprocess_data, scale_features, evaluate_model
9 from torch import nn
10 from skorch import NeuralNetRegressor
11
12 df = pd.read_csv("SeoulBikeData.csv", encoding='unicode_escape')
13 X, y = preprocess_data(df)

```

```

14 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size
    =0.25, random_state=60)
15 X_train, X_test = scale_features(X_train, X_test)
16 device = "cuda" if torch.cuda.is_available() else "cpu"
17 X_train = torch.tensor(X_train.values, dtype=torch.float32)
18 X_test = torch.tensor(X_test.values, dtype= torch.float32)
19 y_train = torch.tensor(y_train.values, dtype= torch.float32).unsqueeze(
    dim = 1)
20 y_test = torch.tensor(y_test.values, dtype= torch.float32).unsqueeze(
    dim = 1)
21 X_train = X_train.to(device)
22 X_test = X_test.to(device)
23 y_train = y_train.to(device)
24 y_test = y_test.to(device)
25
26 class NeuralNetwork(nn.Module):
27     def __init__(self, input_features = 20, out_features = 1,
        hidden_units_1 = 10, hidden_units_2 = 10, hidden_units_3 = 10):
28         super().__init__()
29         self.linear_layer_stack = nn.Sequential(
30             nn.Linear(in_features=input_features, out_features=
        hidden_units_1),
31             nn.ReLU(),
32             nn.Linear(in_features=hidden_units_1, out_features=
        hidden_units_2),
33             nn.ReLU(),
34             nn.Linear(in_features=hidden_units_2, out_features=
        hidden_units_3),
35             nn.ReLU(),
36             nn.Linear(in_features=hidden_units_3, out_features=
        out_features)
37         )
38     def forward(self, x):
39         return self.linear_layer_stack(x)
40
41
42 t0 = time.time()
43 net = NeuralNetRegressor(
44     NeuralNetwork,
45     criterion=torch.nn.MSELoss,
46     optimizer=torch.optim.Adam,
47     max_epochs=10,
48     device=device,
49     verbose=0
50 )
51 params = {
52     'module__hidden_units_1': [50,60,70,90],
53     'module__hidden_units_2': [30,40,50,70],
54     'module__hidden_units_3': [10,20,30,40],
55 }
56 best_model = GridSearchCV(net, params, cv=5, scoring='
    neg_mean_squared_error', n_jobs = -1)
57 X_train = X_train.cpu()
58 y_train = y_train.cpu()
59 X_test = X_test.cpu()
60 y_test = y_test.cpu()
61 r2_score_train, mean_absolute_error_score_train,
    root_mean_square_error_train, r2_score_test,
    mean_absolute_error_score_test, root_mean_square_error_test =
    evaluate_model(best_model, X_train, y_train, X_test, y_test)
62 print(f"Time it took: {time.time() - t0}")
63 print(f"The best parameters: {best_model.best_params_}")
64 print(f"The root mean sqare error training is : {np.sqrt(-best_model.
    best_score_)}")
65 print(f"The r2 score on training is : {r2_score_train}")

```

```
66 print(f"The mean absolute error training is : {  
    mean_absolute_error_score_train}")  
67 print(f"The root mean square error training is : {  
    root_mean_square_error_train}")  
68 print(f"The r2 score on testing is : {r2_score_test}")  
69 print(f"The mean absolute error testing is : {  
    mean_absolute_error_score_test}")  
70 print(f"The root mean square error testing is : {  
    root_mean_square_error_test}")
```

Listing 8: Neural Network Model (neural_network.py)