

Taxi Fare Prediction Modeling

Introduction

Research Objective

The primary objective of this study is to develop a predictive fare model using the comprehensive NYC yellow taxi trip dataset from 2023. The aim is to identify key fare determinants and leverage advanced machine learning techniques to create a model that can estimate taxi fares accurately. This model will then be adapted to propose a more structured and transparent fare system for Tbilisi, Georgia.

Summary of EDA Insights and Hypotheses

From the exploratory data analysis, several key insights were derived:

- The fare amount, trip distance, trip duration, speed, and tip amount distributions were highly right-skewed, indicating the need for transformations.
- Fare amount had strong positive correlations with trip distance and trip duration.
- Categorical variables such as time of day, day type, season, and holidays significantly impacted fare amounts.
- Proper handling of outliers was crucial to prevent skewing model predictions.

Hypotheses

1. **Trip Distance and Duration:** Longer trips (both in terms of distance and duration) will result in higher fares.
2. **Temporal Factors:** Taxi fares vary significantly by time of day, day of the week, and season, with peak times and special events (e.g., holidays) resulting in higher fares.
3. **Speed:** Higher average speeds may correlate with higher fares due to longer distances covered in shorter times.
4. **Tips:** Higher fares tend to receive higher tips.

The modeling phase will test these hypotheses by evaluating the predictive power of these factors.

Methodology

Data Preparation

The dataset used for modeling has been cleaned and preprocessed based on the insights from the EDA. Key steps include:

- Log transformations to normalize skewed data.
- Handling outliers using Z-score filtering.
- Encoding categorical variables using one-hot encoding.
- Scaling numerical features.
- PCA analysis

Model Selection

To build a robust fare prediction model, we will use both simple and complex machine learning models:

1. Simple Models:

- Linear Regression

2. Complex Models:

- Decision Trees
- Random Forest
- Gradient Boosting Machines (GBM)
- XGBoost

Model Evaluation Metrics

The performance of the models will be evaluated using the following metrics:

- Mean Absolute Error (MAE)
- Mean Squared Error (MSE)
- R-squared (R^2)

These metrics will help in understanding the accuracy and reliability of the models.

```
In [42]: 1 # Import necessary libraries
2 import pandas as pd
3 from sklearn.linear_model import LinearRegression
4 from sklearn.tree import DecisionTreeRegressor
5 from sklearn.ensemble import RandomForestRegressor, GradientBoo
6 import xgboost as xgb
7 from sklearn.metrics import mean_absolute_error, mean_squared_e
8 from sklearn.model_selection import train_test_split
9 import pandas as pd
10 import numpy as np
11 from sklearn.linear_model import LinearRegression
12 from sklearn.tree import DecisionTreeRegressor
13 from sklearn.ensemble import RandomForestRegressor, GradientBoo
14 from sklearn.neighbors import KNeighborsRegressor
15 from sklearn.svm import SVR
16 import xgboost as xgb
17 from sklearn.metrics import mean_absolute_error, mean_squared_e
18 from sklearn.model_selection import train_test_split
19 from sklearn.preprocessing import OneHotEncoder, StandardScaler
20 from sklearn.compose import ColumnTransformer
21 from sklearn.model_selection import GridSearchCV
22 from sklearn.feature_selection import SelectKBest, f_regression
23 from sklearn.model_selection import train_test_split, Randomize
24
25
26
```

```
In [9]: 1 # Load the prepared dataset
2 df = pd.read_parquet('modeling_dataset_v.0.parquet')
```

```
In [10]: 1 df.head()
```

Out[10]:

	fare_amount	trip_duration	trip_distance	total_amount	tip_amount	speed_mph
19276873	19.1	20.366667	2.32	30.72	5.12	6.834697
25505385	17.7	18.933333	1.90	25.20	3.50	6.021127
23056552	8.6	8.950000	1.10	12.10	2.00	7.374302
22171336	10.0	8.966667	1.00	16.00	2.00	6.691450
25241061	7.9	5.783333	1.07	14.88	2.98	11.100865

5 rows × 7 columns

```
In [11]: 1 len(df)
```

```
Out[11]: 2805668
```

```
In [12]: 1 df.dtypes
```

```
Out[12]: fare_amount          float64
trip_duration          float64
trip_distance          float64
total_amount           float64
tip_amount             float64
speed_mph              float64
JFK_LGA_Pickup_Fee     float64
General_Airport_Fee    float64
log_fare_amount         float64
log_trip_duration       float64
log_trip_distance       float64
log_tip_amount          float64
is_holiday              int64
payment_type            int64
pickup_time_of_day      category
pickup_season           category
passenger_count_category category
pickup_day_type         category
PUzone                  category
PUborough               category
DOzone                  category
DOborough               category
PCA1                    float64
PCA2                    float64
dtype: object
```

```

In [13]: 1
          2
          3 # Create a combined column for stratification
          4 df['stratify_col'] = df['pickup_season'].astype(str) + '_' + df
          5
          6 # Check the distribution of the combined stratification column
          7 stratify_counts = df['stratify_col'].value_counts()
          8 print(stratify_counts)
          9
         10 # Filter out classes with very few samples
         11 min_class_size = 5 # Define a minimum class size
         12 filtered_classes = stratify_counts[stratify_counts >= min_class
         13 df_filtered = df[df['stratify_col'].isin(filtered_classes)]
         14
         15 # Perform stratified sampling on the filtered dataset
         16 df_sample, _ = train_test_split(df_filtered, test_size=0.9, str
         17
         18 # Drop the stratify column
         19 df_sample = df_sample.drop(columns=['stratify_col'])
         20

```

```

stratify_col
spring_evening_weekday_Manhattan_Manhattan    162096
spring_afternoon_weekday_Manhattan_Manhattan   152673
autumn_evening_weekday_Manhattan_Manhattan      149395
autumn_afternoon_weekday_Manhattan_Manhattan    140089
winter_afternoon_weekday_Manhattan_Manhattan    139463
...
autumn_night_weekend_Unknown_EWR                1
autumn_afternoon_weekday_EWR_EWR                1
autumn_evening_weekday_Unknown_EWR              1
spring_night_weekday_Staten_Island_Queens        1
winter_night_weekday_Unknown_EWR                1
Name: count, Length: 1050, dtype: int64

```

```

In [23]: 1 print(len(df))
          2 print(len(df_sample))

```

```

2805668
280513

```

For the modeling purposes we performed stratified sampling to ensure that all the categorical features remain intact and we have a model that can give us insights based on the categorical variables as well.

```

In [ ]: 1 # Feature selection and encoding
          2 features = df_sample[['PCA1', 'PCA2', 'is_holiday', 'pickup_time_of_day', 'fare_amount']]
          3 target = df_sample['fare_amount']
          4
          5 categorical_features = ['is_holiday', 'pickup_time_of_day', 'pickup_location']
          6 encoder = ColumnTransformer(transformers=[('cat', OneHotEncoder, categorical_features)])

```

```

7
8 X = encoder.fit_transform(features)
9 y = target.values
10
11 # Reduce the number of features using SelectKBest
12 X_new = SelectKBest(f_regression, k=10).fit_transform(X, y)
13
14 # Split data into training, validation, and test sets
15 X_train, X_temp, y_train, y_temp = train_test_split(X_new, y,
16 X_val, X_test, y_val, y_test = train_test_split(X_temp, y_temp,
17
18 # Standardize the data for KNN and SVM
19 scaler = StandardScaler()
20 X_train_scaled = scaler.fit_transform(X_train)
21 X_val_scaled = scaler.transform(X_val)
22 X_test_scaled = scaler.transform(X_test)
23
24 # Linear Regression
25 linear_model = LinearRegression()
26 linear_model.fit(X_train, y_train)
27 y_pred_train_lr = linear_model.predict(X_train)
28 y_pred_val_lr = linear_model.predict(X_val)
29 print("Linear Regression Performance")
30 print("Training MAE:", mean_absolute_error(y_train, y_pred_train_lr))
31 print("Validation MAE:", mean_absolute_error(y_val, y_pred_val_lr))
32 print("Validation R²:", r2_score(y_val, y_pred_val_lr))
33
34 # Decision Tree with Grid Search for parameter tuning
35 param_grid_dt = {'max_depth': [None, 10, 20], 'min_samples_split': [2, 5, 10]}
36 dt_grid_search = GridSearchCV(DecisionTreeRegressor(random_state=42),
37 dt_grid_search.fit(X_train, y_train)
38 best_dt_model = dt_grid_search.best_estimator_
39 y_pred_val_dt = best_dt_model.predict(X_val)
40 print("Decision Tree Performance")
41 print("Validation MAE:", mean_absolute_error(y_val, y_pred_val_dt))
42 print("Validation R²:", r2_score(y_val, y_pred_val_dt))
43
44 # Random Forest with Grid Search for parameter tuning
45 param_grid_rf = {'n_estimators': [10, 50, 100], 'max_depth': [None, 10, 20]}
46 rf_grid_search = GridSearchCV(RandomForestRegressor(random_state=42),
47 rf_grid_search.fit(X_train, y_train)
48 best_rf_model = rf_grid_search.best_estimator_
49 y_pred_val_rf = best_rf_model.predict(X_val)
50 print("Random Forest Performance")
51 print("Validation MAE:", mean_absolute_error(y_val, y_pred_val_rf))
52 print("Validation R²:", r2_score(y_val, y_pred_val_rf))
53
54 # Gradient Boosting with early stopping
55 gb_model = GradientBoostingRegressor(n_estimators=100, validation_fraction=0.1)
56 gb_model.fit(X_train, y_train)
57 y_pred_val_gb = gb_model.predict(X_val)
58 print("Gradient Boosting Performance")
59 print("Validation MAE:", mean_absolute_error(y_val, y_pred_val_gb))

```

```

60 print("Validation R²:", r2_score(y_val, y_pred_val_gb))
61
62 # XGBoost with Grid Search for parameter tuning
63 param_grid_xgb = {'n_estimators': [50, 100], 'max_depth': [3, 6]}
64 xgb_grid_search = GridSearchCV(xgb.XGBRegressor(random_state=42), param_grid_xgb)
65 xgb_grid_search.fit(X_train, y_train)
66 best_xgb_model = xgb_grid_search.best_estimator_
67 y_pred_val_xgb = best_xgb_model.predict(X_val)
68 print("XGBoost Performance")
69 print("Validation MAE:", mean_absolute_error(y_val, y_pred_val_xgb))
70 print("Validation R²:", r2_score(y_val, y_pred_val_xgb))
71
72 # K-Nearest Neighbors
73 knn_model = KNeighborsRegressor(n_neighbors=5, n_jobs=-1)
74 knn_model.fit(X_train_scaled, y_train)
75 y_pred_val_knn = knn_model.predict(X_val_scaled)
76 print("KNN Performance")
77 print("Validation MAE:", mean_absolute_error(y_val, y_pred_val_knn))
78 print("Validation R²:", r2_score(y_val, y_pred_val_knn))
79
80 # Support Vector Machine
81 svm_model = SVR(kernel='rbf')
82 svm_model.fit(X_train_scaled, y_train)
83 y_pred_val_svm = svm_model.predict(X_val_scaled)
84 print("SVM Performance")
85 print("Validation MAE:", mean_absolute_error(y_val, y_pred_val_svm))
86 print("Validation R²:", r2_score(y_val, y_pred_val_svm))
87
88 # Summarize Model Performance
89 models = ['Linear Regression', 'Decision Tree', 'Random Forest']
90 val_mae = [
91     mean_absolute_error(y_val, y_pred_val_lr),
92     mean_absolute_error(y_val, y_pred_val_dt),
93     mean_absolute_error(y_val, y_pred_val_rf),
94     mean_absolute_error(y_val, y_pred_val_gb),
95     mean_absolute_error(y_val, y_pred_val_xgb),
96     mean_absolute_error(y_val, y_pred_val_knn),
97     mean_absolute_error(y_val, y_pred_val_svm)
98 ]
99 val_r2 = [
100     r2_score(y_val, y_pred_val_lr),
101     r2_score(y_val, y_pred_val_dt),
102     r2_score(y_val, y_pred_val_rf),
103     r2_score(y_val, y_pred_val_gb),
104     r2_score(y_val, y_pred_val_xgb),
105     r2_score(y_val, y_pred_val_knn),
106     r2_score(y_val, y_pred_val_svm)
107 ]
108
109 performance_df = pd.DataFrame({
110     'Model': models,
111     'Validation MAE': val_mae,
112     'Validation R²': val_r2

```

```
113 }  
114  
115 print(performance_df)
```

Preparation For Modeling

```
In [29]: 1 # Feature selection and encoding
          2 features = df_sample[['PCA1', 'PCA2', 'is_holiday', 'pickup_time']
          3 target = df_sample['fare_amount']
```

```
In [30]: 1 categorical_features = ['is_holiday', 'pickup_time_of_day', 'pi
2 encoder = ColumnTransformer(transformers=[('cat', OneHotEncoder
```

```
In [32]: 1 X = encoder.fit_transform(features)
2 y = target.values
3
4 # Reduce the number of features using SelectKBest
5 X_new = SelectKBest(f_regression, k=10).fit_transform(X, y)
6
7 # Split data into training, validation, and test sets
8 X_train, X_temp, y_train, y_temp = train_test_split(X_new, y, t
9 X_val, X_test, y_val, y_test = train_test_split(X_temp, y_temp,
10
```

```
In [35]: 1 # Standardize the data for KNN and SVM
2 scaler = StandardScaler(with_mean=False)
3 X_train_scaled = scaler.fit_transform(X_train)
4 X_val_scaled = scaler.transform(X_val)
5 X_test_scaled = scaler.transform(X_test)
6
```

Baseline Model


```
In [36]: 1 # Baseline Model: Linear Regression
2 linear_model = LinearRegression()
3 linear_model.fit(X_train, y_train)
4 y_pred_train = linear_model.predict(X_train)
5 y_pred_val = linear_model.predict(X_val)
6
7 print("Linear Regression Performance")
8 print("Training MAE:", mean_absolute_error(y_train, y_pred_train))
9 print("Validation MAE:", mean_absolute_error(y_val, y_pred_val))
10 print("Validation R²:", r2_score(y_val, y_pred_val))
11
```

Linear Regression Performance
Training MAE: 2.697683248247949
Validation MAE: 2.6931883885075365
Validation R²: 0.944835876489503

Linear Regression Model

Performance Metrics:

- **Training MAE:** 2.70
- **Validation MAE:** 2.69
- **Validation R²:** 0.9448

Interpretation:

1. Mean Absolute Error (MAE):

- The average absolute error in fare prediction is approximately 2.70 on the training data and 2.69 on the validation data. The closeness of these values indicates that the model is not overfitting and generalizes well to unseen data.

2. R-squared (R²):

- The R² value of 0.9448 indicates that approximately 94.48% of the variance in taxi fares can be explained by the model. This high R² value suggests that the linear regression model fits the data well.

Conclusion:

- The linear regression model performs well with a high R² value and low MAE, indicating strong predictive capability for the given features. However, it is a simple model and may not capture complex, non-linear relationships in the data.

Complex Models

```
In [39]: 1 # Decision Tree with Grid Search for parameter tuning
2 param_grid_dt = {'max_depth': [None, 10, 20], 'min_samples_split': [2, 5, 10]}
3 dt_grid_search = GridSearchCV(DecisionTreeRegressor(random_state=42), param_grid_dt, cv=5)
4 dt_grid_search.fit(X_train, y_train)
5 best_dt_model = dt_grid_search.best_estimator_
6 y_pred_val_dt = best_dt_model.predict(X_val)
7 print("Decision Tree Performance")
8 print("Validation MAE:", mean_absolute_error(y_val, y_pred_val_dt))
9 print("Validation R²:", r2_score(y_val, y_pred_val_dt))
10
11
```

Decision Tree Performance
Validation MAE: 2.3716119571271497
Validation R²: 0.9545403292306762

Decision Tree Model

Performance Metrics:

- **Validation MAE:** 2.37
- **Validation R²:** 0.9545

Interpretation:

1. Mean Absolute Error (MAE):

- The average absolute error in fare prediction is approximately \$2.37 on the validation data. This is slightly lower than the MAE of the linear regression model, indicating better predictive accuracy.

2. R-squared (R²):

- The R² value of 0.9545 indicates that approximately 95.45% of the variance in taxi fares can be explained by the decision tree model. This is higher than the R² value for the linear regression model, suggesting that the decision tree captures more variability in the data.

Conclusion:

- The decision tree model performs better than the linear regression model in terms of both MAE and R². It captures more complex relationships between the features and the target variable. However, decision trees can be prone to overfitting, so it is essential to validate the model on unseen data and potentially prune the tree or limit its depth to prevent overfitting.

```
In [44]: 1 # Gradient Boosting with early stopping
2 gb_model = GradientBoostingRegressor(n_estimators=100, validation_
3 gb_model.fit(X_train, y_train)
4 y_pred_val_gb = gb_model.predict(X_val)
5 print("Gradient Boosting Performance")
6 print("Validation MAE:", mean_absolute_error(y_val, y_pred_val_
7 print("Validation R²:", r2_score(y_val, y_pred_val_gb))
8
```

Gradient Boosting Performance
 Validation MAE: 2.41214002650482
 Validation R²: 0.9550262508155019

```
In [45]: 1 # XGBoost with Grid Search for parameter tuning
2 param_grid_xgb = {'n_estimators': [50, 100], 'max_depth': [3, 6
3 xgb_grid_search = GridSearchCV(xgb.XGBRegressor(random_state=42
4 xgb_grid_search.fit(X_train, y_train)
5 best_xgb_model = xgb_grid_search.best_estimator_
6 y_pred_val_xgb = best_xgb_model.predict(X_val)
7 print("XGBoost Performance")
8 print("Validation MAE:", mean_absolute_error(y_val, y_pred_val_
9 print("Validation R²:", r2_score(y_val, y_pred_val_xgb))
```

XGBoost Performance
 Validation MAE: 2.3678238791443933
 Validation R²: 0.9528451666259213

```
In [ ]: 1 # Random Forest with Randomized Search for parameter tuning usi
2 param_dist_rf = {'n_estimators': [10, 50], 'max_depth': [10, 20
3 X_train_sample, _, y_train_sample, _ = train_test_split(X_train
4 rf_random_search = RandomizedSearchCV(RandomForestRegressor(ran
5 rf_random_search.fit(X_train_sample, y_train_sample)
6 best_rf_model = rf_random_search.best_estimator_
7 final_rf_model = RandomForestRegressor(**best_rf_model.get_para
8 final_rf_model.fit(X_train, y_train)
9 y_pred_val_rf = final_rf_model.predict(X_val)
10 print("Random Forest Performance")
11 print("Validation MAE:", mean_absolute_error(y_val, y_pred_val_
12 print("Validation R²:", r2_score(y_val, y_pred_val_rf))
13
```

```
In [47]: 1 # K-Nearest Neighbors
2 knn_model = KNeighborsRegressor(n_neighbors=5, n_jobs=-1)
3 knn_model.fit(X_train_scaled, y_train)
4 y_pred_val_knn = knn_model.predict(X_val_scaled)
5 print("KNN Performance")
6 print("Validation MAE:", mean_absolute_error(y_val, y_pred_val_
7 print("Validation R²:", r2_score(y_val, y_pred_val_knn))
8
```


 KeyboardInterrupt Traceback (most recent call)

```

KeyboardInterrupt                                Traceback (most recent call last)
/var/folders/4d/8tkcz58x0md0_v7fj12j6ht80000gn/T/ipykernel_2648/1713529712.py in <module>
      2 knn_model = KNeighborsRegressor(n_neighbors=5, n_jobs=-1)
      3 knn_model.fit(X_train_scaled, y_train)
----> 4 y_pred_val_knn = knn_model.predict(X_val_scaled)
      5 print("KNN Performance")
      6 print("Validation MAE:", mean_absolute_error(y_val,
y_pred_val_knn))

~/opt/anaconda3/lib/python3.9/site-packages/sklearn/neighbors/_regression.py in predict(self, X)
    206         X = check_array(X, accept_sparse='csr')
    207
--> 208         neigh_dist, neigh_ind = self.kneighbors(X)
    209
    210         weights = _get_weights(neigh_dist, self.weights)

~/opt/anaconda3/lib/python3.9/site-packages/sklearn/neighbors/_base.py in kneighbors(self, X, n_neighbors, return_distance)
    703         kws = self.effective_metric_params_
    704
--> 705         chunked_results = list(pairwise_distances_chunked(
    706             X, self._fit_X, reduce_func=reduce_func,
    707             metric=self.effective_metric_, n_jobs=n_jobs,
bs,

~/opt/anaconda3/lib/python3.9/site-packages/sklearn/metrics/pairwise.py in pairwise_distances_chunked(X, Y, reduce_func, metric, n_jobs, working_memory, **kws)
   1631         if reduce_func is not None:
   1632             chunk_size = D_chunk.shape[0]
-> 1633             D_chunk = reduce_func(D_chunk, sl.start)
   1634             _check_chunk_size(D_chunk, chunk_size)
   1635         yield D_chunk

~/opt/anaconda3/lib/python3.9/site-packages/sklearn/neighbors/_base.py in _kneighbors_reduce_func(self, dist, start, n_neighbors, return_distance)
    580         """
    581         sample_range = np.arange(dist.shape[0])[:, None]
--> 582         neigh_ind = np.argpartition(dist, n_neighbors - 1,
axis=1)
    583         neigh_ind = neigh_ind[:, :n_neighbors]
    584         # argpartition doesn't guarantee sorted order, so
we sort again

~/opt/anaconda3/lib/python3.9/site-packages/numpy/core/fromnumeric.py in argpartition(a, kth, axis, kind, order)
    856
    857         """
--> 858         return wrapfunc(a, 'argpartition', kth, axis=axis,

```

```

kind=kind, order=order)
859
860

~/opt/anaconda3/lib/python3.9/site-packages/numpy/core/fromnumeric
.py in _wrapfunc(obj, method, *args, **kwds)
    57
    58     try:
---> 59         return bound(*args, **kwds)
    60     except TypeError:
    61         # A TypeError occurs if the object does have such
a method in its

```

KeyboardInterrupt:

```

In [ ]: 1 # Support Vector Machine
2 svm_model = SVR(kernel='rbf')
3 svm_model.fit(X_train_scaled, y_train)
4 y_pred_val_svm = svm_model.predict(X_val_scaled)
5 print("SVM Performance")
6 print("Validation MAE:", mean_absolute_error(y_val, y_pred_val_
7 print("Validation R²:", r2_score(y_val, y_pred_val_svm))
8

```

```

In [ ]: 1 # Neural Network (MLP)
2 model = Sequential()
3 model.add(Dense(64, input_dim=X_train.shape[1], activation='rel
4 model.add(Dense(32, activation='relu'))
5 model.add(Dense(1))
6
7 # Compile the model
8 model.compile(optimizer='adam', loss='mean_squared_error')
9
10 # Train the model
11 history = model.fit(X_train, y_train, epochs=50, batch_size=32,
12
13 # Predict on training and validation data
14 y_pred_train_nn = model.predict(X_train)
15 y_pred_val_nn = model.predict(X_val)
16
17 print("Neural Network Performance")
18 print("Training MAE:", mean_absolute_error(y_train, y_pred_trai
19 print("Validation MAE:", mean_absolute_error(y_val, y_pred_val_
20 print("Validation R²:", r2_score(y_val, y_pred_val_nn))
21

```

```
In [ ]: 1 models = ['Linear Regression', 'Decision Tree', 'Random Forest']
2 train_mae = [
3     mean_absolute_error(y_train, linear_model.predict(X_train)),
4     mean_absolute_error(y_train, dt_model.predict(X_train)),
5     mean_absolute_error(y_train, rf_model.predict(X_train)),
6     mean_absolute_error(y_train, gb_model.predict(X_train)),
7     mean_absolute_error(y_train, xgb_model.predict(X_train)),
8     mean_absolute_error(y_train, knn_model.predict(X_train_scaled)),
9     mean_absolute_error(y_train, svm_model.predict(X_train_scaled))
10 ]
11
12 val_mae = [
13     mean_absolute_error(y_val, linear_model.predict(X_val)),
14     mean_absolute_error(y_val, dt_model.predict(X_val)),
15     mean_absolute_error(y_val, rf_model.predict(X_val)),
16     mean_absolute_error(y_val, gb_model.predict(X_val)),
17     mean_absolute_error(y_val, xgb_model.predict(X_val)),
18     mean_absolute_error(y_val, knn_model.predict(X_val_scaled)),
19     mean_absolute_error(y_val, svm_model.predict(X_val_scaled))
20 ]
21
22 val_r2 = [
23     r2_score(y_val, linear_model.predict(X_val)),
24     r2_score(y_val, dt_model.predict(X_val)),
25     r2_score(y_val, rf_model.predict(X_val)),
26     r2_score(y_val, gb_model.predict(X_val)),
27     r2_score(y_val, xgb_model.predict(X_val)),
28     r2_score(y_val, knn_model.predict(X_val_scaled)),
29     r2_score(y_val, svm_model.predict(X_val_scaled)),
30 ]
31
32
33 performance_df = pd.DataFrame({
34     'Model': models,
35     'Training MAE': train_mae,
36     'Validation MAE': val_mae,
37     'Validation R²': val_r2
38 })
39
40 print(performance_df)
41
```