

# Machine Learning-Based

## Flight Delay Prediction

# intro

## Issue in Air Transportation:

- Flight delays significantly impact consumers and airports.
- Causes include weather, traffic, and scheduling inefficiencies.

## Objective of the Study:

- Examine JFK Airport flight takeoff data.
- Identify primary reasons for flight delays.

## Approach:

- Utilize machine learning models to predict departure delays.
- Compare performance of various regression models.

## Research Goal:

- Determine the most important variables impacting flight departures.
- Assess the effectiveness of different machine learning algorithms.

# Features & Preprocessing

---

01



# Methods & Models

02



# Results

03



# about Dataset

The dataset used for this study is sourced from Kaggle, titled "**Flight Take Off Data – JFK Airport**", created by Deepankur Kansal. It comprises **28,821 rows** and **23 features** related to flight take-off details from JFK Airport. The data is primarily aimed at predicting flight delays, particularly those influenced by weather conditions.



**28,821 rows**  
**23 features**



**Weather and  
flight related data**



**published 4y ago**  
**27.5k views in last month**

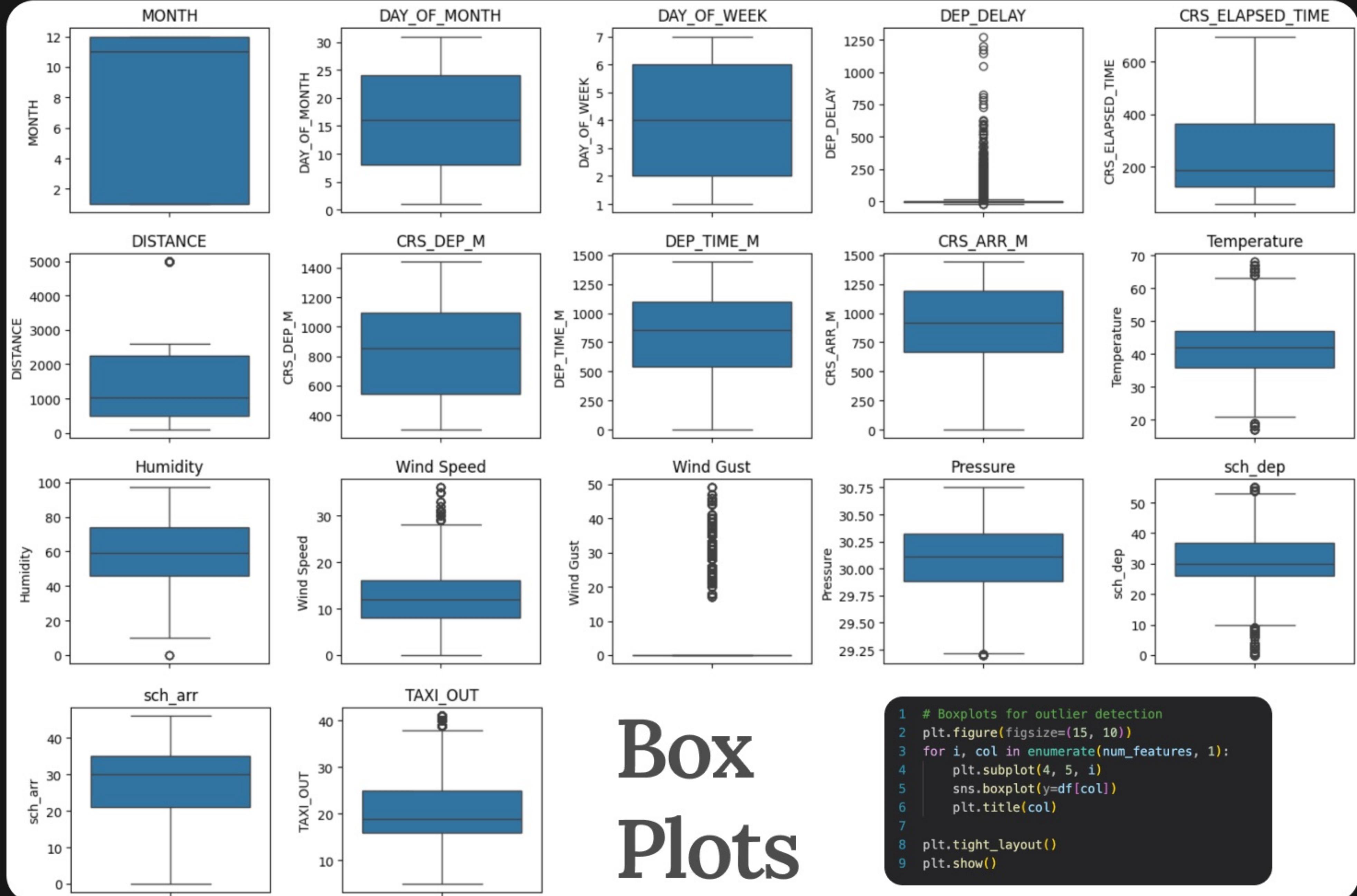


**flights between  
Nov 2019 - Dec 2020**



# Features

**target** feature - **DEP\_DELAY**



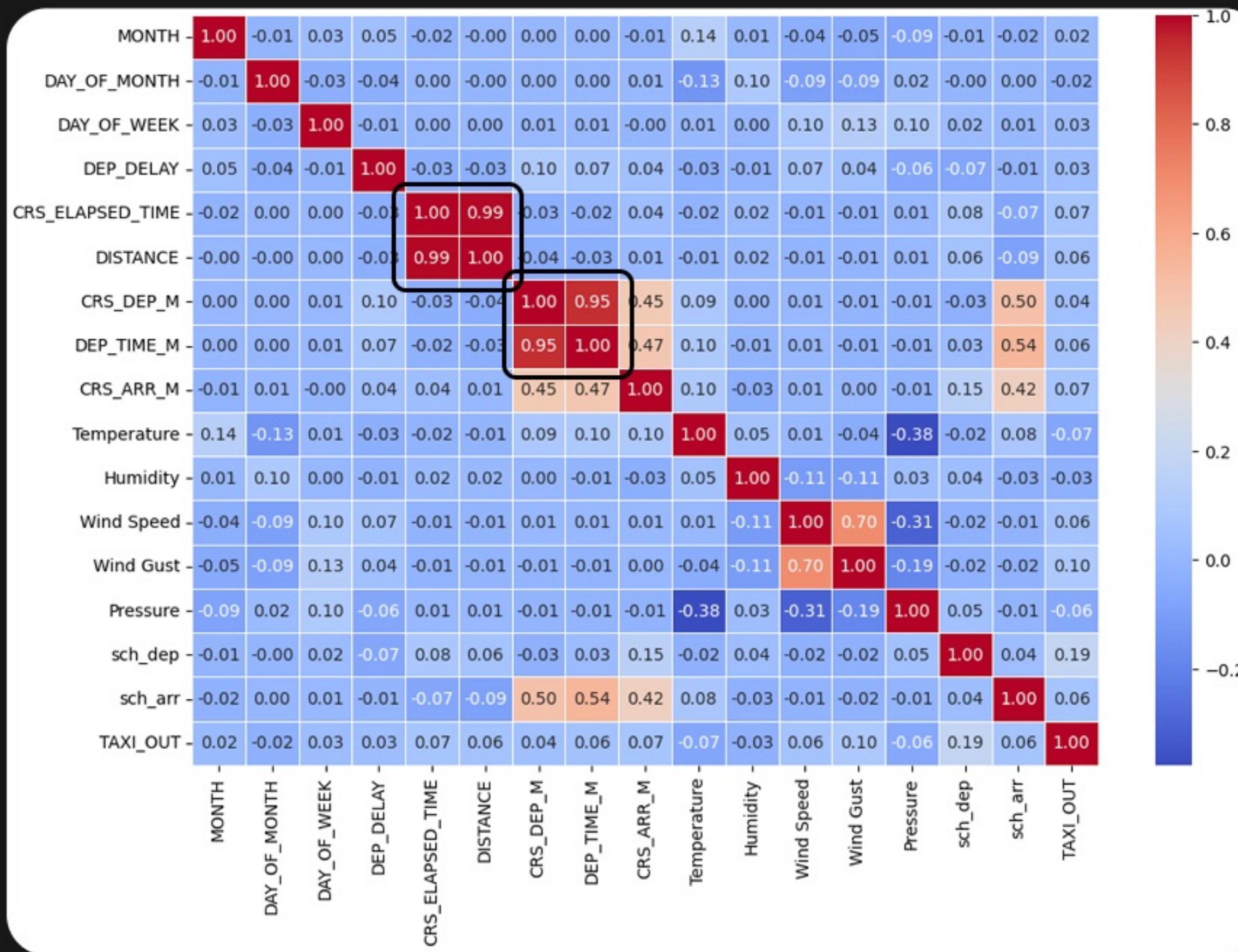
# Box Plots

```

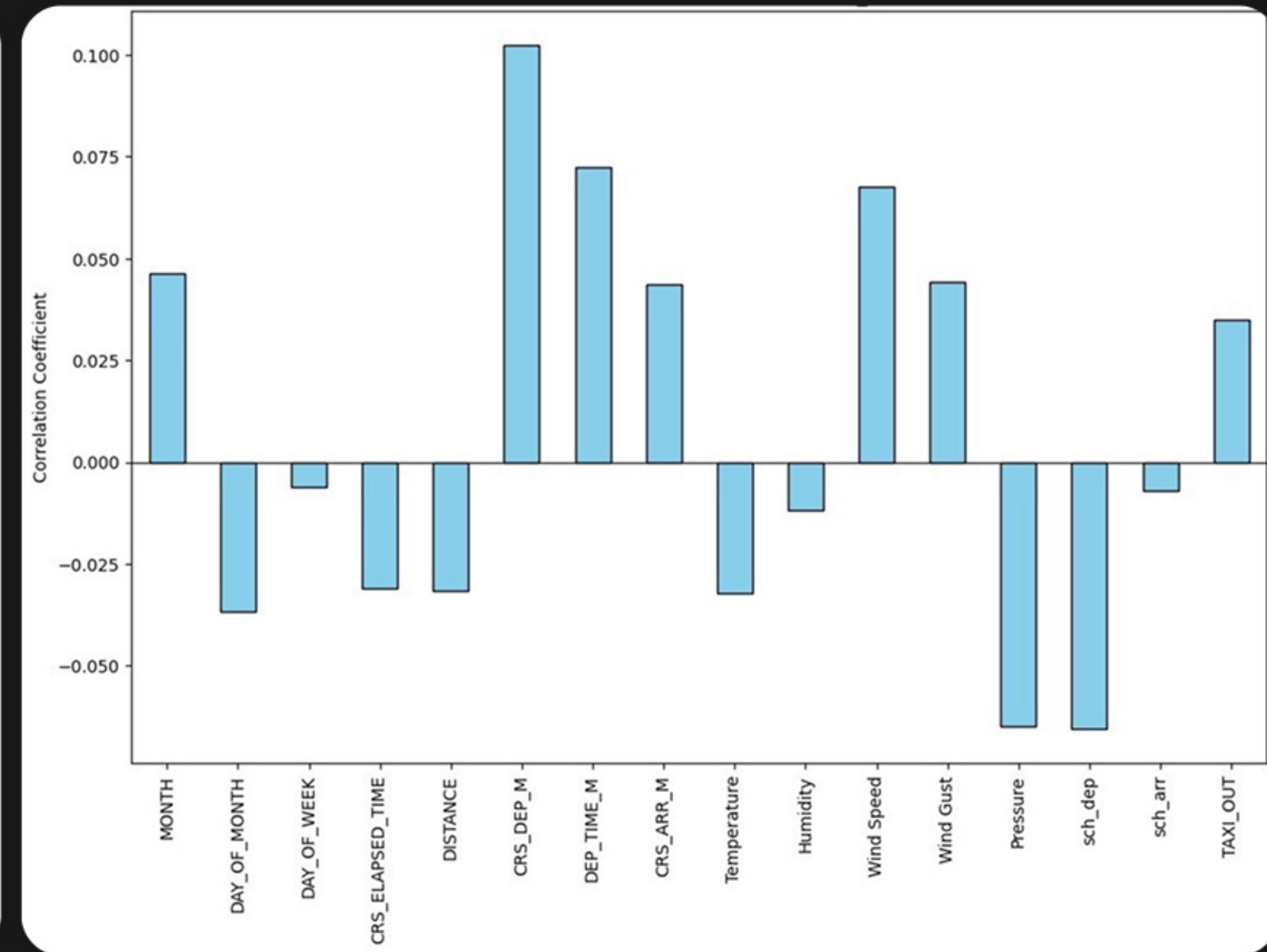
1 # Boxplots for outlier detection
2 plt.figure(figsize=(15, 10))
3 for i, col in enumerate(num_features, 1):
4     plt.subplot(4, 5, i)
5     sns.boxplot(y=df[col])
6     plt.title(col)
7
8 plt.tight_layout()
9 plt.show()

```

# Correlation matrix of Numerical Features



# Correlation of Other Features with 'DEP\_DELAY'



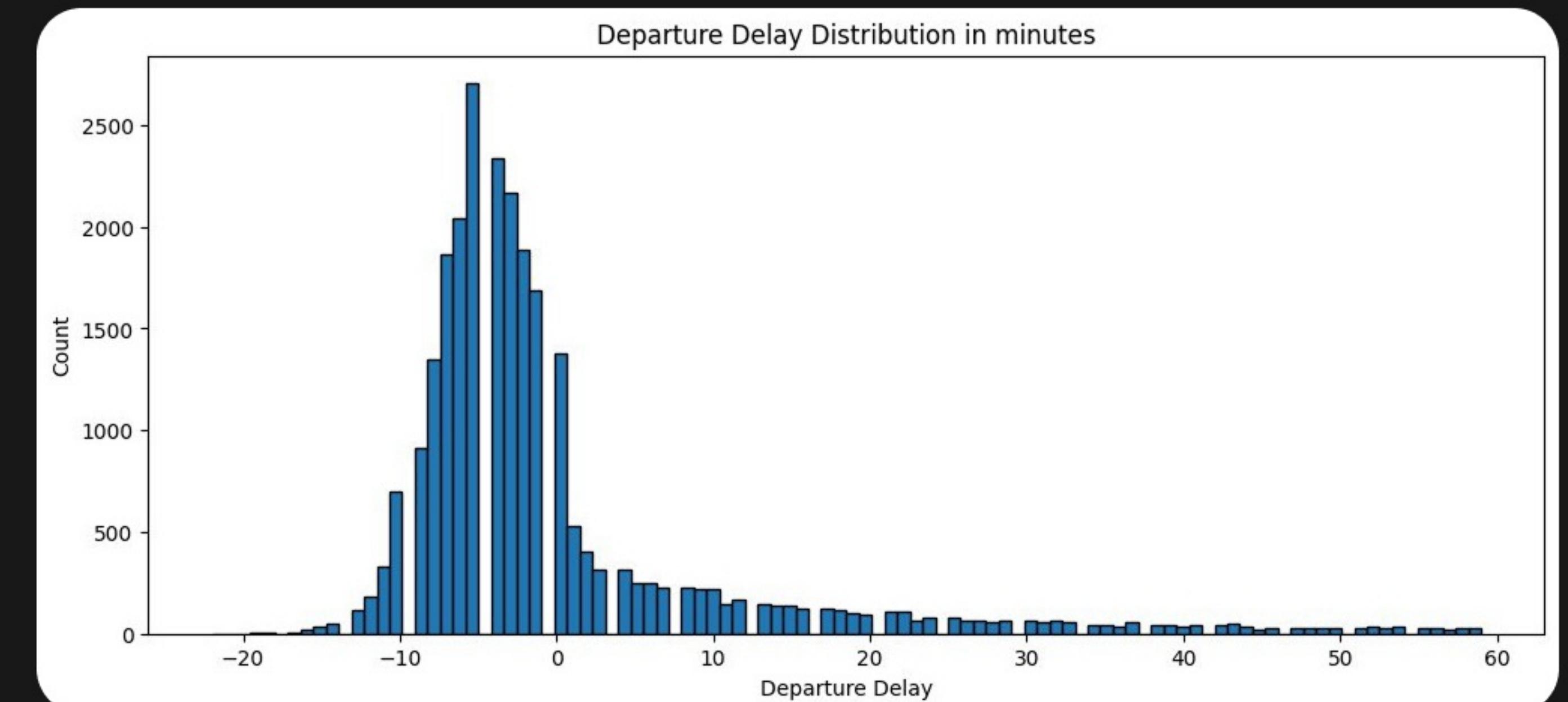
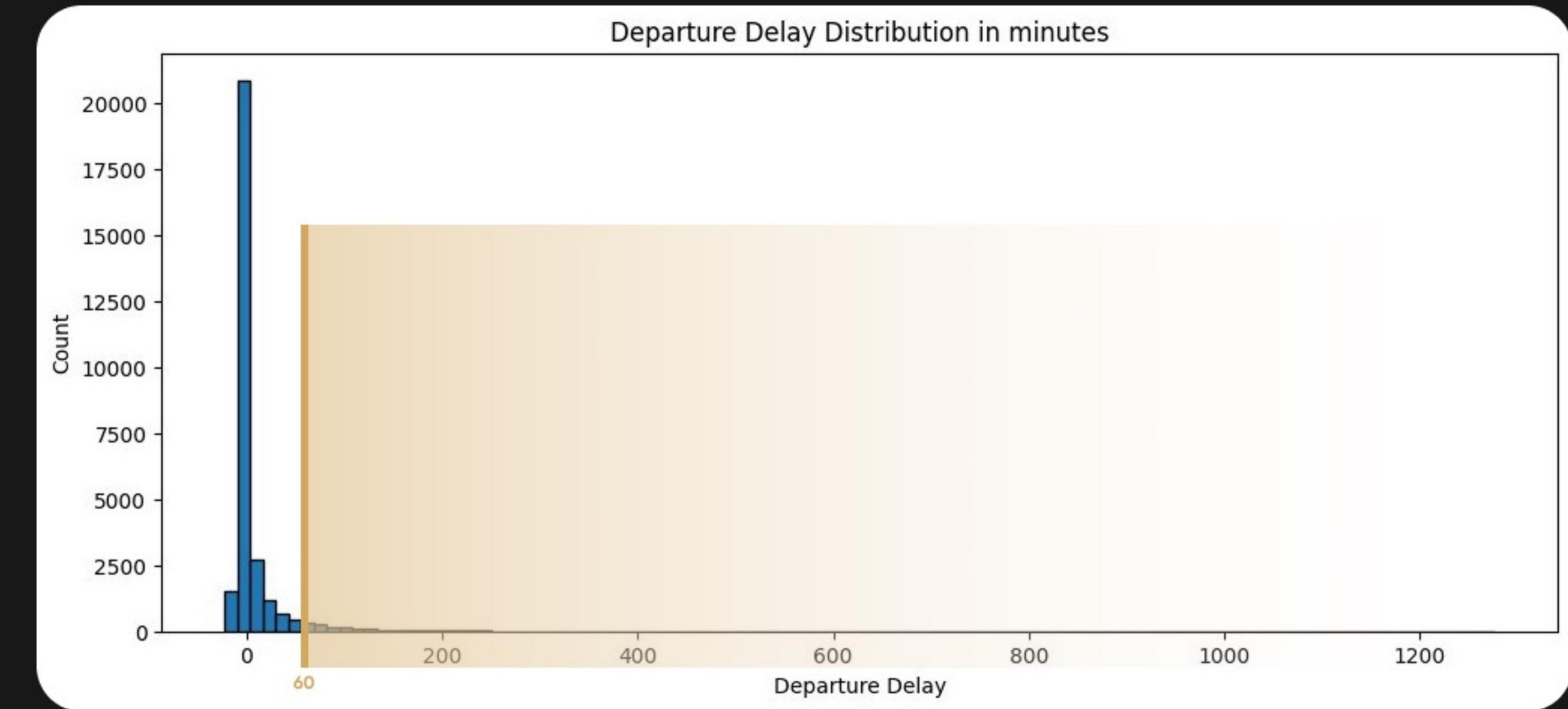
data preprocessing

# Outlier Handling

Since regression models are sensitive to outliers, 'DEP\_DELAY' values greater than 60 minutes were removed.



28,821 rows →  
27,473 data points



data preprocessing

# Handling Corrupted Variables

'Dew Point' feature contained corrupted values:

- Converted the values to numerical and removed data points with NULL values



27,473 rows →  
25,818 data points

```
1 dataset['Dew Point'].unique()
[12]
...
array(['34', '32', '33', '31', '30', '29', '67', '35', '36', '38', '39',
       '40', '41', '42', '28', '27', '37', '44', '45', '52', '53', '54',
       '56', '57', '59', '58', '50', '46', '43', '48', '49', '47', '21',
       '20', '22', '24', '23', '19', '18', '16', '13', '12', '14', '15',
       '17', '26', '25', '51', '55', '9\x00', '6\x00', '5\x00', '8\x00',
       '10', '4\x00', '7\x00', '3\x00', '2\x00', '11', '0\x00', '-1',
       '1\x00', '-2', '-3'], dtype=object)
```

```
1 dataset['Dew Point'] = pd.to_numeric(dataset['Dew Point'], errors='coerce')
2 dataset.isna().sum()['Dew Point']
[13]
```

... 1654

```
1 dataset = dataset.dropna()
2 dataset.info()
[14]
...
<class 'pandas.core.frame.DataFrame'>
Index: 25818 entries, 0 to 28819
```

data preprocessing

# Encoding Categorical Values

The '**Condition**' and '**Wind**' features had **24** and **18 unique values**, respectively.

- **Converted** these features to **categorical values**

```
1 dataset['Condition'].nunique()
```

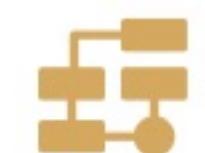
24

```
1 dataset['Wind'].nunique()
```

18

```
1 dataset = pd.get_dummies(dataset, columns=['Condition', 'Wind'])  
2 dataset.info()
```

```
<class 'pandas.core.frame.DataFrame'>  
Index: 25818 entries, 0 to 28819  
Data columns (total 57 columns):
```



17 features →  
57 features



# Features & Preprocessing

01



## Results

03



# Methods & Models

02



# Model Types used



## Linear models

Ordinary Least Squares  
(OLS), Polynomial (Degree = 3)



## Tree models

Decision Tree, Random Forest,  
Gradient Boosting (XGBoost)



## Non-linear models

Support Vector Machines (SVM)  
and k-Nearest Neighbors (k-NN)

Methods

# Settings



**Train 80%  
test 20%**



Linear regression  
models with this split



**Coefficient of  
Determination**



Measures the proportion  
of variance explained by  
the model



**Root Mean  
Squared Error**



Simulate and check  
correctness



**Mean Absolute  
Error**



Represents the average  
absolute difference  
between predicted and  
actual delays

# OLS and 2-Polynomial Regression

```

1 linear_models = {
2     # "LinearRegression w/o Scaling": (make_pipeline(LinearRegression()), {}),
3     "LinearRegression w/ Scaling": (make_pipeline(StandardScaler(), LinearRegression()), {}),
4     "LinearRegression w/ Scaling & Poly": (make_pipeline(StandardScaler(), PolynomialFeatures(), LinearRegression()), {}),
5     # "Ridge w/o Scaling": (make_pipeline(Ridge()), {"ridge_alpha": [0.01, 0.1, 1, 10, 100, 1000]}),
6     # "Lasso w/o Scaling": (make_pipeline(Lasso()), {"lasso_alpha": [0.01, 0.1, 1, 10, 100, 1000]}),
7     "Ridge w/ Scaling & Poly": (make_pipeline(StandardScaler(), PolynomialFeatures(), Ridge()), {"ridge_alpha": [0.01, 0.1, 1, 10, 100, 1000]}),
8     "Lasso w/ Scaling & Poly": (make_pipeline(StandardScaler(), PolynomialFeatures(), Lasso()), {"lasso_alpha": [0.01, 0.1, 1, 10, 100, 1000]})
9 }
10

```

## Linear Models

```
1 best_linear_models = model_tuning(linear_models)
```

```
Best hyperparams for 'LinearRegression w/ Scaling': {}
```

```
Best hyperparams for 'Ridge w/ Scaling': {'ridge_alpha': 1000}
```

```
Best hyperparams for 'Lasso w/ Scaling': {'lasso_alpha': 0.1}
```

## Regularized Linear Regression



To mitigate overfitting and enhance generalization.

### Ridge (L2 Regularization)

Introduces a penalty on large coefficients, improving stability and reducing variance



### Lasso (L1 Regularization)

Encourages sparsity by shrinking some coefficients to zero, effectively performing feature selection



# Capture nonlinear relationships

```

11 tree_ensemble_models = {
12     "DecisionTree": (DecisionTreeRegressor(), {"max_depth": [3, 5, 10, 15]}),
13     "RandomForest": (RandomForestRegressor(), {"n_estimators": [50, 100, 200], "max_depth": [3, 5, 10, 15]}),
14     "BaggingRegressor": (
15         BaggingRegressor(estimator=DecisionTreeRegressor(), n_estimators=50),
16         {"n_estimators": [50, 100, 200], "estimator__max_depth": [3, 5, 10, 15]}
17     ),
18     "GradientBoosting": (GradientBoostingRegressor(), {"n_estimators": [50, 100, 200], "learning_rate": [0.01, 0.1, 0.2]}),
19     "AdaBoost": (AdaBoostRegressor(), {"n_estimators": [50, 100, 200], "learning_rate": [0.01, 0.1, 0.2]}),
20 }

```

## Tree-based and Ensemble Models

```

1 best_tree_ensemble_models = model_tuning(tree_ensemble_models)

Best hyperparams for 'DecisionTree': {'max_depth': 5}

Best hyperparams for 'RandomForest': {'max_depth': 15, 'n_estimators': 200}

Best hyperparams for 'BaggingRegressor': {'estimator__max_depth': 15, 'n_estimators': 200}

Best hyperparams for 'GradientBoosting': {'learning_rate': 0.2, 'n_estimators': 200}

Best hyperparams for 'AdaBoost': {'learning_rate': 0.01, 'n_estimators': 50}

```

### Decision Tree Regression

A non-parametric model that recursively splits data based on feature values

### Random Forest Regression

An ensemble learning method that reduces overfitting by averaging multiple decision trees

### Gradient Boosting Regression

building decision trees in a multi-stage fashion, correcting errors of the previous stages using gradient descent

### Bagging Regression

ensemble learning method combining multiple models trained on different subsets of the dataset

# Other models

```

22 other_models = {
23     # "Linear SVR w/o Scaling": (make_pipeline(StandardScaler(), SVR(kernel="linear")), {"svr_C": [0.1, 1, 10, 100, 1000]}),
24     "Linear SVR w/ Scaling": (make_pipeline(StandardScaler(), SVR(kernel="linear")), {"svr_C": [0.1, 1, 10, 100, 1000]}),
25     # "RBF SVR w/o Scaling": (make_pipeline(StandardScaler(), SVR(kernel="rbf")), {"svr_C": [0.1, 1, 10, 100, 1000]}),
26     "RBF SVR w/ Scaling": (make_pipeline(StandardScaler(), SVR(kernel="rbf")), {"svr_C": [0.1, 1, 10, 100, 1000]}),
27     # "KNN w/o Scaling": (make_pipeline(StandardScaler(), KNeighborsRegressor()), {"kneighborsregressor_n_neighbors": [3, 5, 7, 9]}),
28     "KNN w/ Scaling": (make_pipeline(StandardScaler(), KNeighborsRegressor()), {"kneighborsregressor_n_neighbors": [3, 5, 7, 9]}),
29 }

```

Test Metrics for 'Linear SVR w/ Scaling' with params {'svr\_C': 1}  
MSE: 1.42e+02 RMSE: 11.9 MAE: 6.43 R2: -0.0554

Test Metrics for 'RBF SVR w/ Scaling' with params {'svr\_C': 100}  
MSE: 1.33e+02 RMSE: 11.5 MAE: 6.48 R2: 0.00978

Test Metrics for 'KNN w/o Scaling' with params {'kneighborsregressor\_n\_neighbors': 9}  
MSE: 1.32e+02 RMSE: 11.5 MAE: 7.47 R2: 0.0161

Test Metrics for 'KNN w/ Scaling' with params {'kneighborsregressor\_n\_neighbors': 9}  
MSE: 1.32e+02 RMSE: 11.5 MAE: 7.47 R2: 0.0161

Support  
Vector  
Regression

Finds a hyperplane that best fits the data while minimizing errors and maximizing margins using support points

k-Nearest  
Neighbors  
Regression

Predicts delays based on the average outcome of the k most similar observations

# PORTATIONS TO BE

## Preprocessing approaches

President Trump calls for tripling of ICE force; riots continue

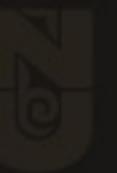
01

Models & Results

02

ed  
s  
is has set in  
most countries  
in previous,  
as to finding  
it by trying  
more aggressive  
agents,  
is another  
and address  
Post Office  
U.S. Postal  
service

LIVE NOW:  
PRESIDENT  
ADDRESSES  
THE NATION



R S T I



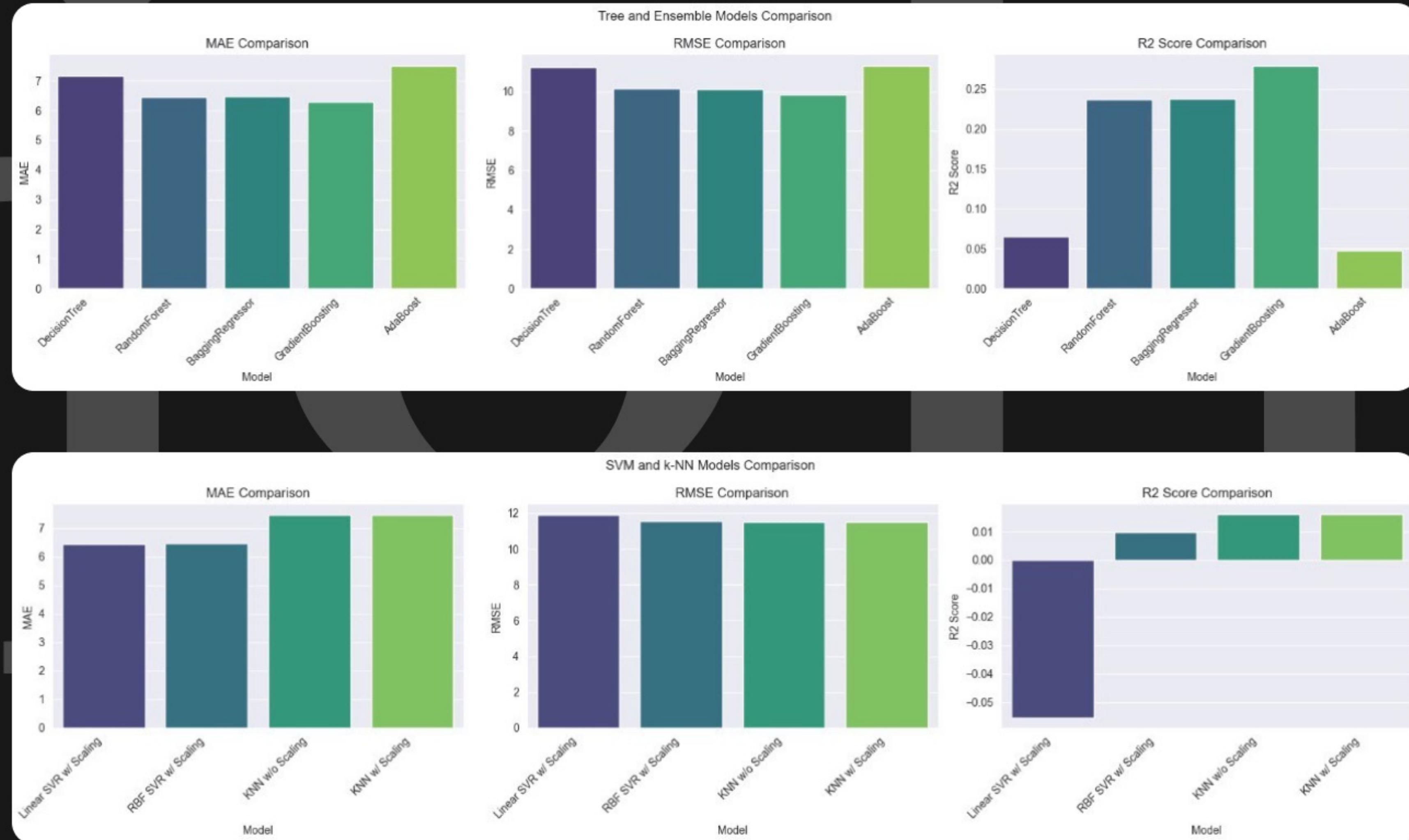
## Results

03



model results

# Results



model results

# Results

Model	Our R <sup>2</sup> Score	Research Paper R <sup>2</sup> Score
Linear Regression	0.0395	0.0214
Ridge	0.0546	0.0212
Lasso	0.062	0.0083
Decision Tree	0.0651	-0.9243
Random Forest	0.238	-0.1021
Bagging	0.243	-
Gradient Boosting	0.278	-
AdaBoost	0.0464	-
Linear SVM	-0.0554	-0.0297
RBF SVM	0.00978	-
K-Nearest Neighbors	0.0161	-

TABLE I: Comparison of R<sup>2</sup> Scores Between Our Study and Research Paper after pre-processing

Model	Our R <sup>2</sup> Score	Research Paper R <sup>2</sup> Score
Linear Regression	0.7523	-
Ridge	0.7548	-
Lasso	0.7421	-
Decision Tree	0.7210	0.8986
Random Forest	0.8327	0.9629
Gradient Boosting	0.8456	-
XGBoosting	0.9018	0.9864
AdaBoost	0.8012	-
Polynomial SVM	0.04	-
K-Nearest Neighbors	0.6894	-

TABLE II: Comparison of R<sup>2</sup> Scores Between Our Study and Research Paper with all features

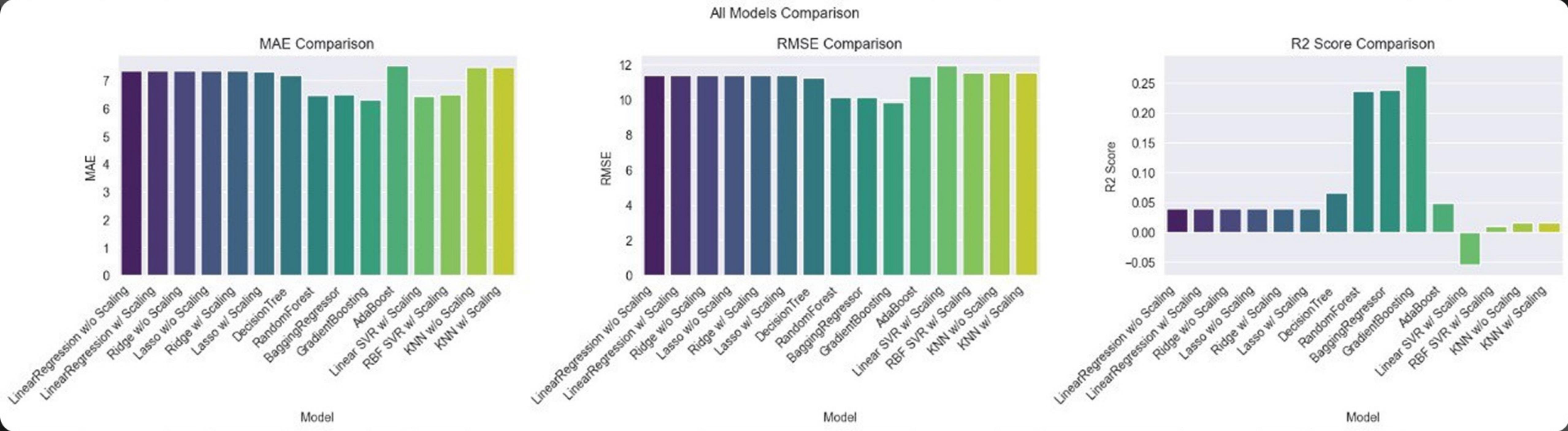
Trees total (Decision tree fastest, random forest slowest, ~17m, xgboost fast): 25m 4.5s

Linear total (all fast): 17.3s  
kNN total: 17.1s

SVM total (linear kernel slowest, ~70m, poly kernel fast, RBF kernel mid, 15-20m): 95m 7.5s

model results

# Results



model results

# Results

Model	Our R <sup>2</sup> Score	Research Paper R <sup>2</sup> Score
Linear Regression	0.0395	0.0214
Ridge	0.0546	0.0212
Lasso	0.062	0.0083
Decision Tree	0.0651	-0.9243
Random Forest	0.238	-0.1021
Bagging	0.243	-
Gradient Boosting	0.278	-
AdaBoost	0.0464	-
Linear SVM	-0.0554	-0.0297
RBF SVM	0.00978	-
K-Nearest Neighbors	0.0161	-

TABLE I: Comparison of R<sup>2</sup> Scores Between Our Study and Research Paper after pre-processing

Model	Our R <sup>2</sup> Score	Research Paper R <sup>2</sup> Score
Linear Regression	0.7523	-
Ridge	0.7548	-
Lasso	0.7421	-
Decision Tree	0.7210	0.8986
Random Forest	0.8327	0.9629
Gradient Boosting	0.8456	-
XGBoosting	0.9018	0.9864
AdaBoost	0.8012	-
Polynomial SVM	0.04	-
K-Nearest Neighbors	0.6894	-

TABLE II: Comparison of R<sup>2</sup> Scores Between Our Study and Research Paper with all features

Trees total (Decision tree fastest, random forest slowest, ~17m, xgboost fast): 25m 4.5s

Linear total (all fast): 17.3s  
kNN total: 17.1s

SVM total (linear kernel slowest, ~70m, poly kernel fast, RBF kernel mid, 15-20m): 95m 7.5s

with all the features

# Results

