

# Scalable Taxi Fare Prediction

Presented by: Assylbek Khalyk & Abdiakhment Kozhamkulov

---

# Agenda

- Problem Definition
- Data Collection
- Data Exploration (EDA)
- Data Cleaning
- Feature Engineering
- Data Splitting
- Model Selection
- Hyperparameter Tuning
- Model Evaluation
- Interpretation & Conclusion

# Problem Definition

## 01 Regression Problem

target variable = total\_amount

## 02 Input – engineered set of features

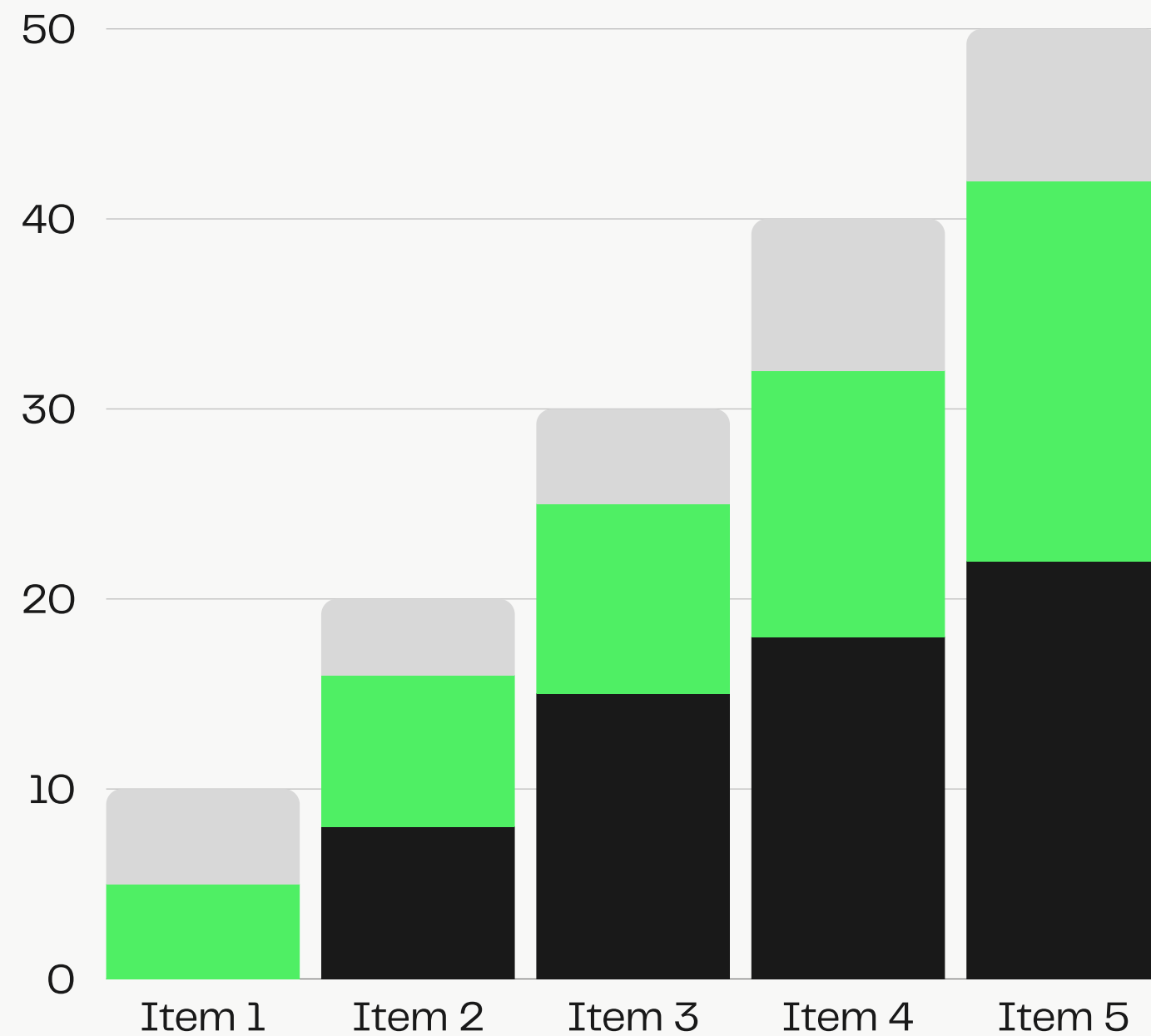
Trip information

Location pairs (PU\_DO\_pair)

Weather conditions

Temporal indicators

# Data Collection



- NYC Taxi & Limousine Commission (TLC) — a government-maintained open data source.
  - ~12 million samples | 17 features
-

---

# Data Cleaning

## Minimizing noise

```
df_cleaned = df_cleaned[(df_cleaned['total_amount'] > 2)
                        & (df_cleaned['passenger_count'] > 0)]
df_cleaned.shape
```

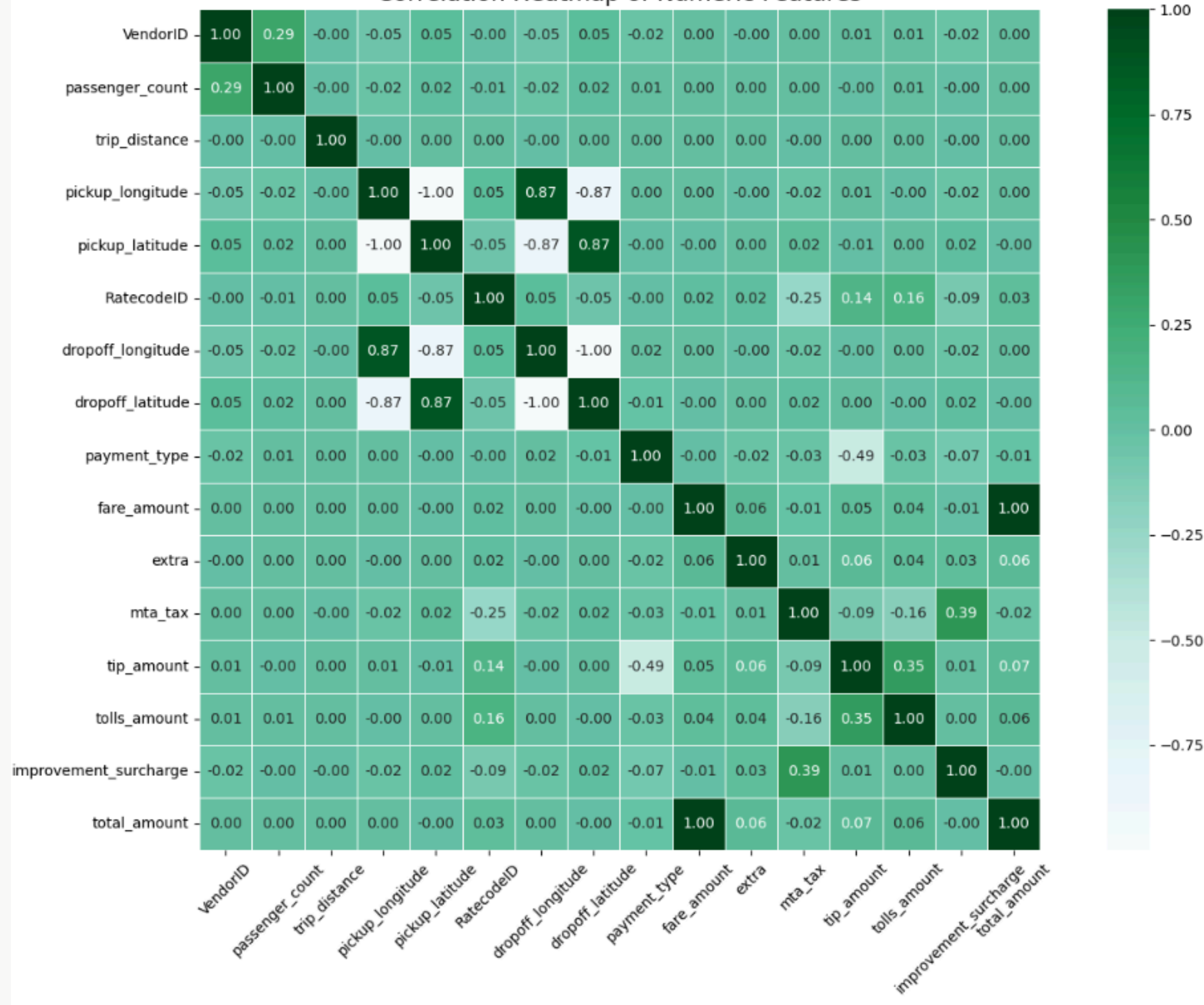
```
df_cleaned = df_cleaned.dropna()
```

```
print(df_cleaned.shape)
df_cleaned.head()
```

```
(11915372, 17)
```

- overall initially clean dataset
-

Correlation Heatmap of Numeric Features



```
# drop unnecessary columns or columns that might cause data Leakage
```

```
df_cleaned = df.drop(columns=[  
    "VendorID",  
    "RatecodeID",  
    "tpep_dropoff_datetime",  
    "fare_amount",  
    "tip_amount",  
    "tolls_amount",  
    "improvement_surcharge",  
    "store_and_fwd_flag",  
    "extra",  
    "mta_tax",  
    "payment_type"  
])
```

# Feature Engineering

## New features

+ haversine\_distance

$$a = \sin^2\left(\frac{\Delta\phi}{2}\right) + \cos(\phi_1) \cdot \cos(\phi_2) \cdot \sin^2\left(\frac{\Delta\lambda}{2}\right)$$

$$c = 2 \cdot \arctan 2\left(\sqrt{a}, \sqrt{1-a}\right)$$

$$\text{distance} = R \cdot c$$

+ is\_airport

+PU\_DO\_pair (mapping lon/lat to city zones using geopandas and lookup table)

+ day of week

+ is\_weekend

+ is\_rush\_hour

+ temperature (Open-Meteo API)

+ wind\_speed

+ precipitation

+ humidity

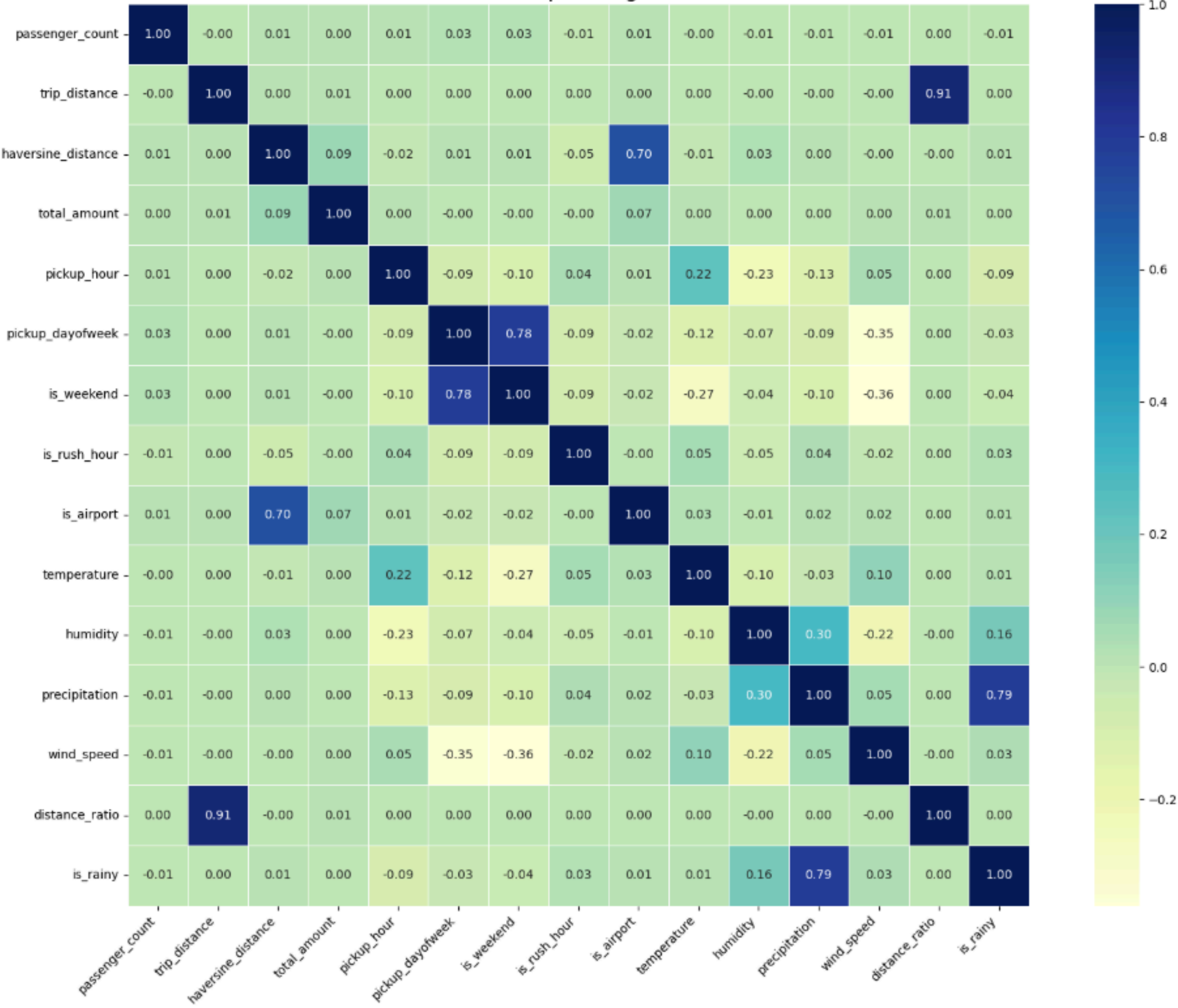
+ distance\_ratio

+ is\_rainy

+ hour\_bin (“morning/midday/evening/night”)



Correlation Heatmap of Engineered Features



---

# Model Selection

**XGBoost**

**Random Forest**

**Ridge**

**HistGBM**

# Hyperparameter tuning

## Random Forest

```
param_grid = {  
    "rf__n_estimators": [100],  
    "rf__max_depth": [10, 15, 20],  
    "rf__min_samples_split": [2, 5],  
    "rf__max_features": ["sqrt"]  
}
```

## XGBoost

```
param_grid = {  
    "xgb__n_estimators": [100, 200],  
    "xgb__max_depth": [5, 7, 10],  
    "xgb__learning_rate": [0.01, 0.1],  
    "xgb__subsample": [0.8, 1.0],  
    "xgb__colsample_bytree": [0.8, 1.0]  
}
```

## Ridge

```
param_grid = {  
    "ridge__alpha": np.logspace(-3, 3, 10) # [0.001, 0.01, ..., 1000]  
}
```

## HistGBM

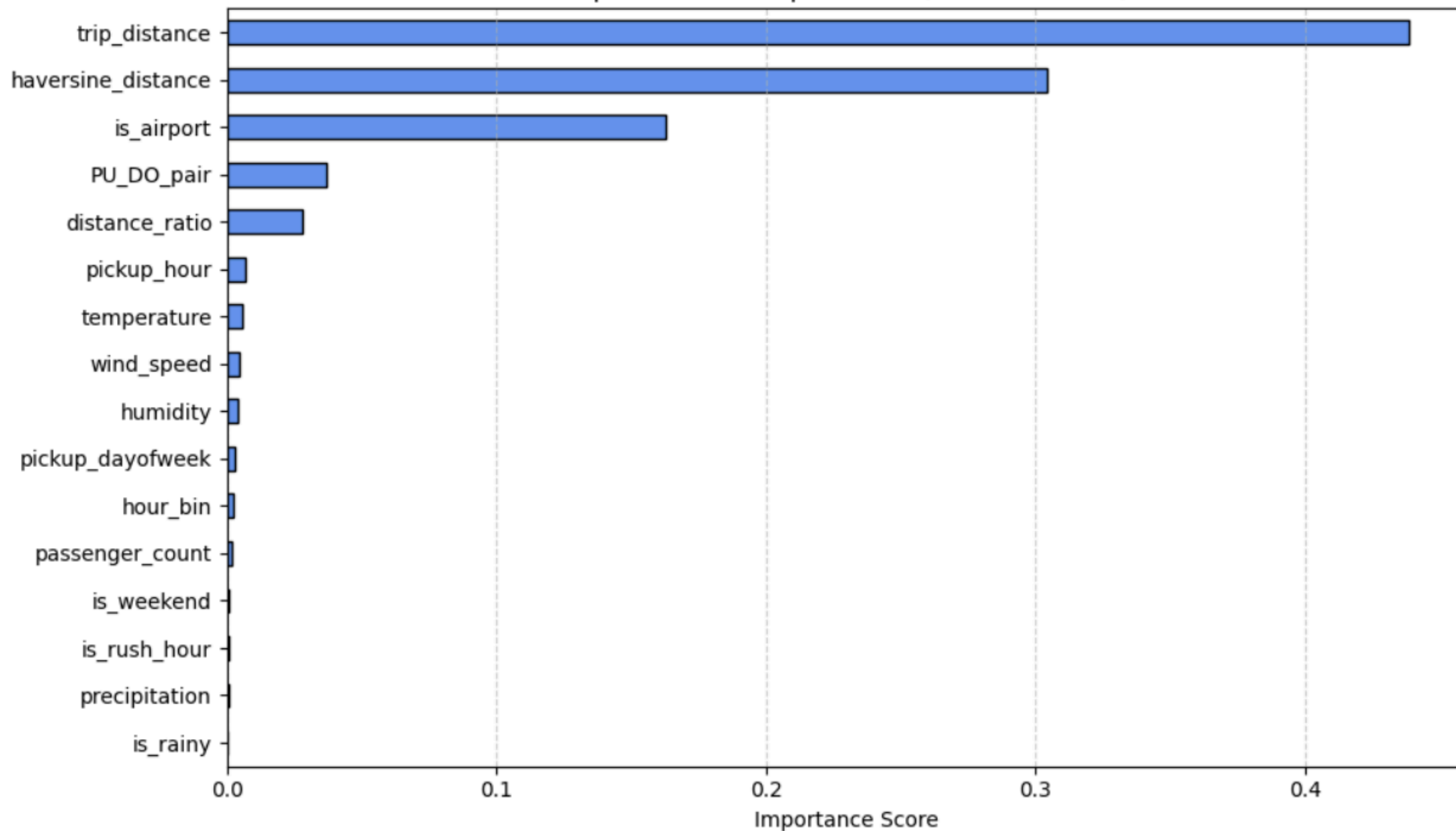
```
param_grid = {  
    "hgb__max_iter": [100, 200],  
    "hgb__max_depth": [5, 10],  
    "hgb__learning_rate": [0.01, 0.1],  
    "hgb__l2_regularization": [0.0, 1.0]  
}
```

---

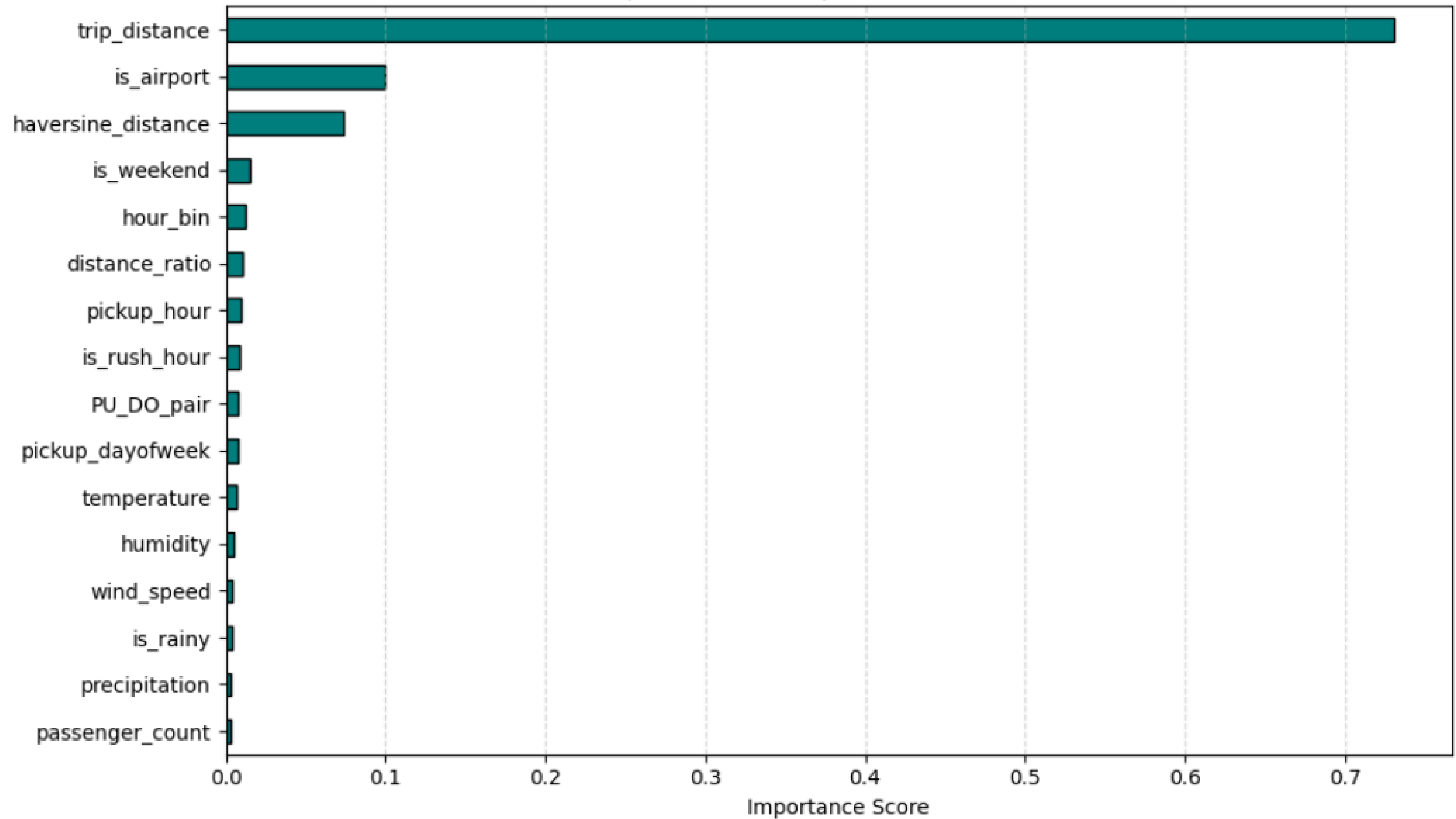
# Model Evaluation

	Random Forest	XGBoost	Ridge	HistGBM
MAE	1.93	1.85	2.33	1.91
R^2	0.91	0.91	0.88	0.91
RMSE	3.66	3.61	4.27	3.64

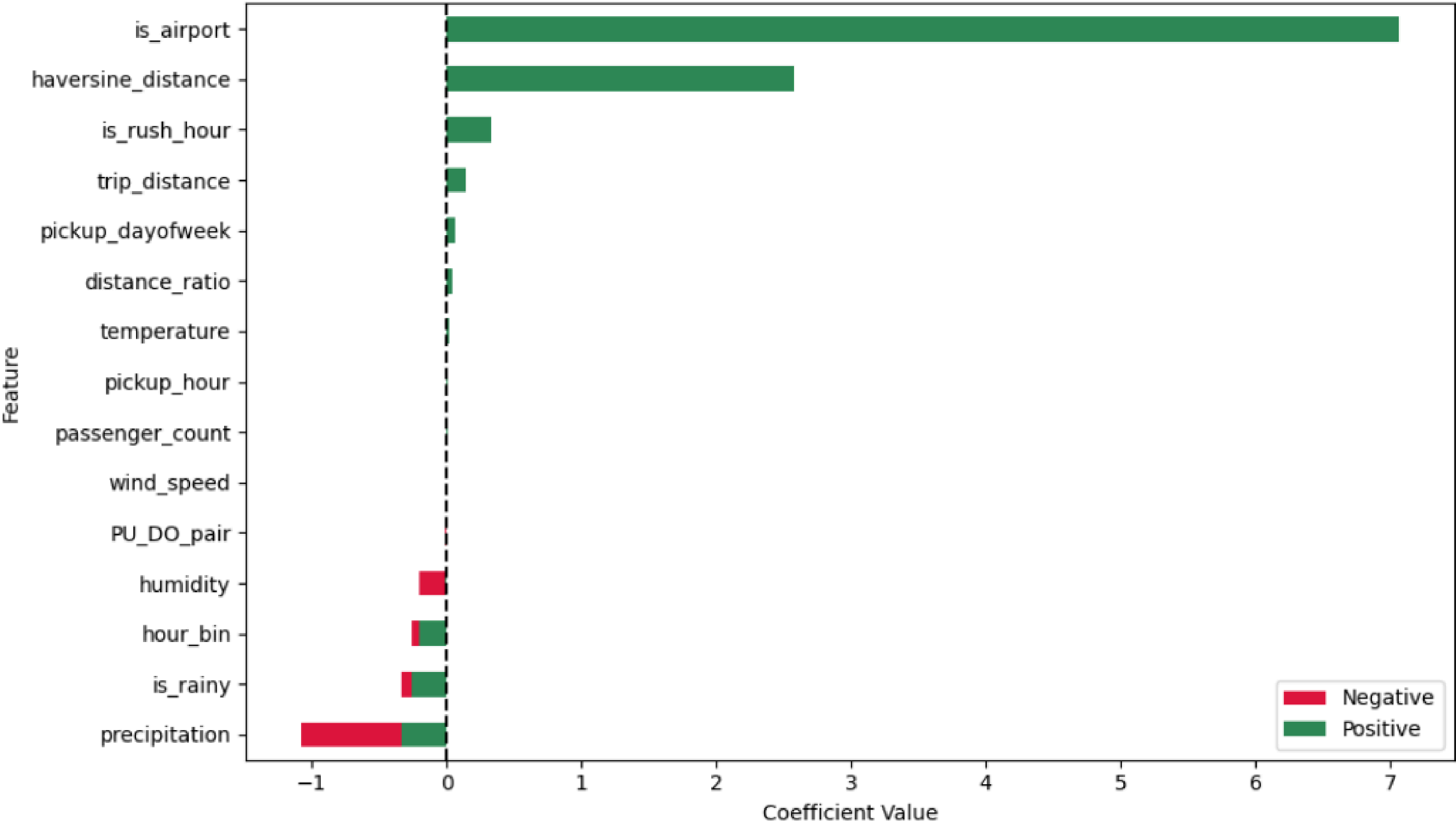
Top Feature Importances (Random Forest)



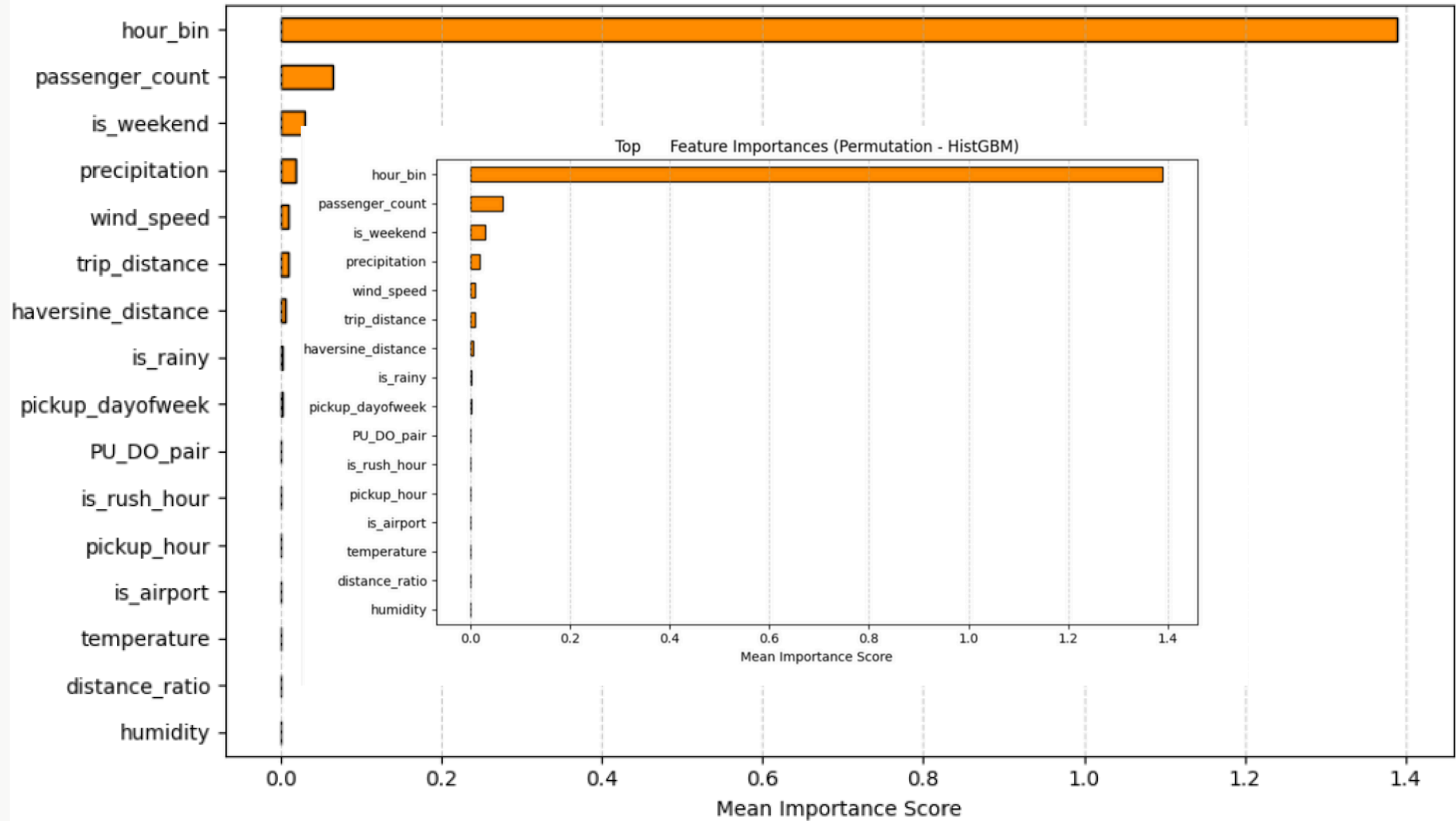
Top Feature Importances (XGBoost)



Top Influential Features (Ridge Coefficients)



Top Feature Importances (Permutation - HistGBM)





# Limitations & Conclusion

switched from 12million samples → 1million  
computational power constraints  
long training time  
subsampling the data

Project shows potential for scalable fare  
prediction in real-world settings

# Q&A



Best Params: {'rf\_\_max\_depth': 20, 'rf\_\_max\_features': 'sqrt',  
'rf\_\_min\_samples\_split': 5, 'rf\_\_n\_estimators': 100}  
MAE: 1.9339617291844997  
RMSE: 3.669324583550365  
R<sup>2</sup>: 0.9156658972962635

Fitting 5 folds for each of 48 candidates, totalling 240 fits  
Best Params: {'xgb\_\_colsample\_bytree': 0.8,  
'xgb\_\_learning\_rate': 0.1, 'xgb\_\_max\_depth': 10,  
'xgb\_\_n\_estimators': 200, 'xgb\_\_subsample': 1.0}  
MAE: 1.8611311739613294  
RMSE: 3.6188917960831843  
R<sup>2</sup>: 0.9179682147080562

Fitting 5 folds for each of 10 candidates, totalling 50 fits  
Best alpha: {'ridge\_\_alpha': np.float64(1000.0)}  
MAE: 2.336484721459339  
RMSE: 4.277048386395624  
R<sup>2</sup>: 0.8854172359770055

Fitting 5 folds for each of 16 candidates, totalling 80 fits  
Best Params: {'hgb\_\_l2\_regularization': 1.0, 'hgb\_\_learning\_rate': 0.1,  
'hgb\_\_max\_depth': 10, 'hgb\_\_max\_iter': 200}  
MAE: 1.9156819838390302  
RMSE: 3.6456114998803493  
R<sup>2</sup>: 0.9167523968111276