# Sharing Bike Demand Prediction

Jiatao Yuan
Data Science Initiative

Brown University

GitHub: https://github.com/Francis958/Data1030-Final-Project

December 9, 2021

# Recap

## Intro

- The purpose of this project is to predict the demand for sharing bike per hour
- Regression methods were used
- Good sharing bike demand prediction can ease the traffic congestion and reduce the cost of the company
- Obtained from the UCI Machine Learning repository

## Dataset Recap

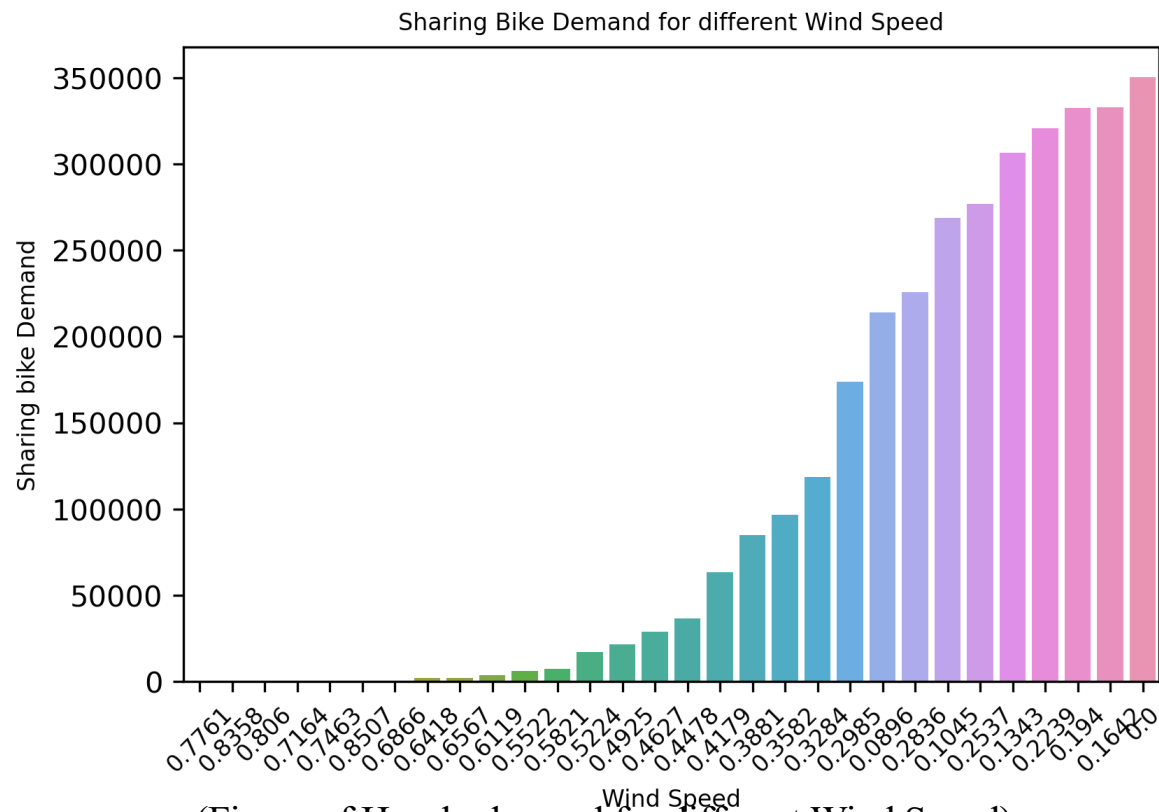| | datetime | season | year | month | hour | holiday | weekday | workingday | weather | temp | atemp | humidity | windspeed | casual | registered | count |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 2011-01-01 | 1 | 2011 | 1 | 0 | 0 | 6 | 0 | 1 | 0.24 | 0.2879 | 0.81 | 0.0 | 3 | 13 | 16 |
| 1 | 2011-01-01 | 1 | 2011 | 1 | 1 | 0 | 6 | 0 | 1 | 0.22 | 0.2727 | 0.80 | 0.0 | 8 | 32 | 40 |
| 2 | 2011-01-01 | 1 | 2011 | 1 | 2 | 0 | 6 | 0 | 1 | 0.22 | 0.2727 | 0.80 | 0.0 | 5 | 27 | 32 |
| 3 | 2011-01-01 | 1 | 2011 | 1 | 3 | 0 | 6 | 0 | 1 | 0.24 | 0.2879 | 0.75 | 0.0 | 3 | 10 | 13 |
| 4 | 2011-01-01 | 1 | 2011 | 1 | 4 | 0 | 6 | 0 | 1 | 0.24 | 0.2879 | 0.75 | 0.0 | 0 | 1 | 1 |

# Preprocessing And Exploratory Data Analysis



(Figure of Correlation Matrix)

- Preprocessing: Casual, Registered and target variable(High correlations good or not?)
- EDA: Wind Speed and target variable(Low correlation bad or not?
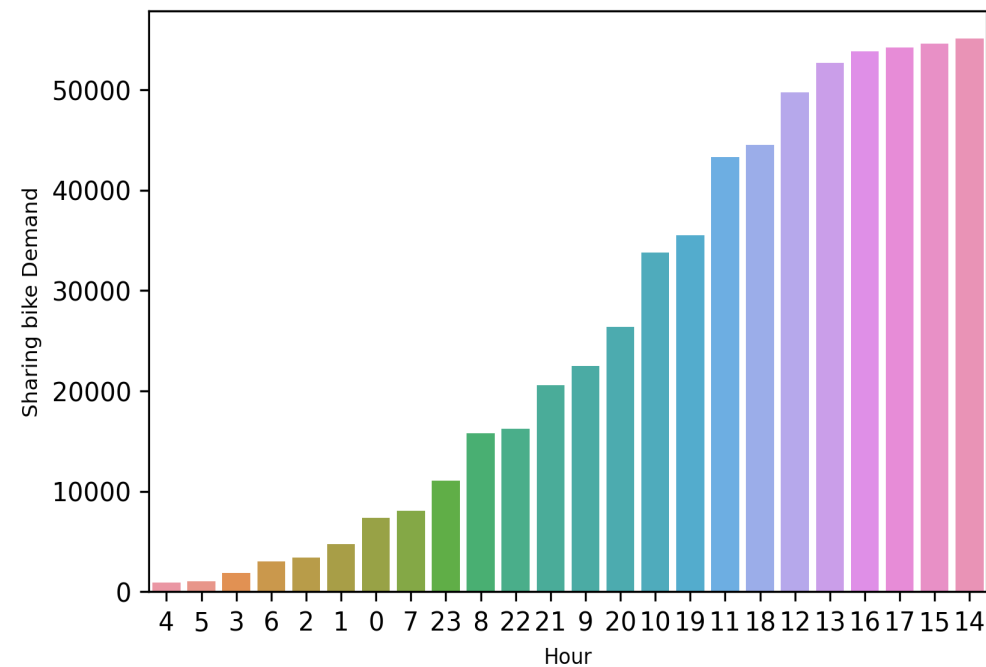- Business insights: Causal and Registered users with different hour slot

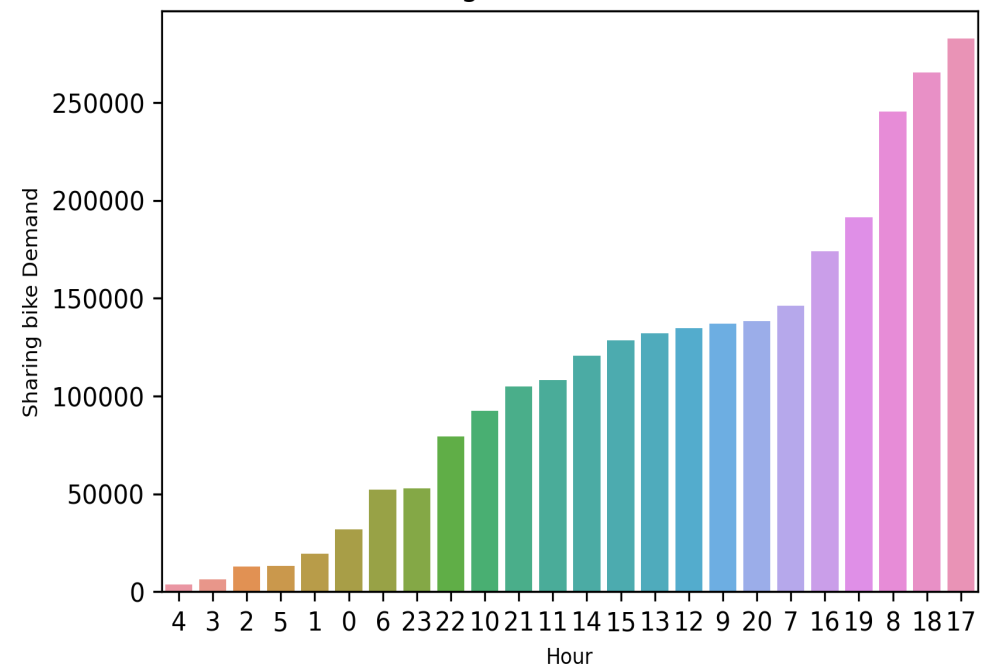# Preprocessing And Exploratory Data Analysis



(Figure of Hourly demand for different Wind Speed)

# Preprocessing And Exploratory Data Analysis



The Demand for the causal users with different hour slot

The Demand for the registered users with different hour slot

# Cross Validation

- Split the data:
    1. Since I add the 6-hour time lags for the dataset, now the data is i.i.d.
    2. I split the data into train, validation, test sets. Test set took up 20% of the whole dataset. Train and validation set took 80% and split into 5-Folds process.

- CV Pipeline
    - $R^2$ is the score
    - Ridge regression: Standard Scaler, One Hot Encoder, drop features of high collinearity, 5 folds, GridSearch and 5 random states
    - Random Forest, XGBoost, GradientBoost: One Hot Encoder, 5 folds, GridSearch and 5 random states

- Hyperparameters Tuning
    - Ridge Regression : L2 regularization term alpha is tuned
    - Random Forest: max_depth and max_features are tuned
    - XGBoost: alpha, max_depth, lambda, learning rate are tuned
    - GradientBoost: learning rate and max_depth are tuned

# Results

- Model Performance

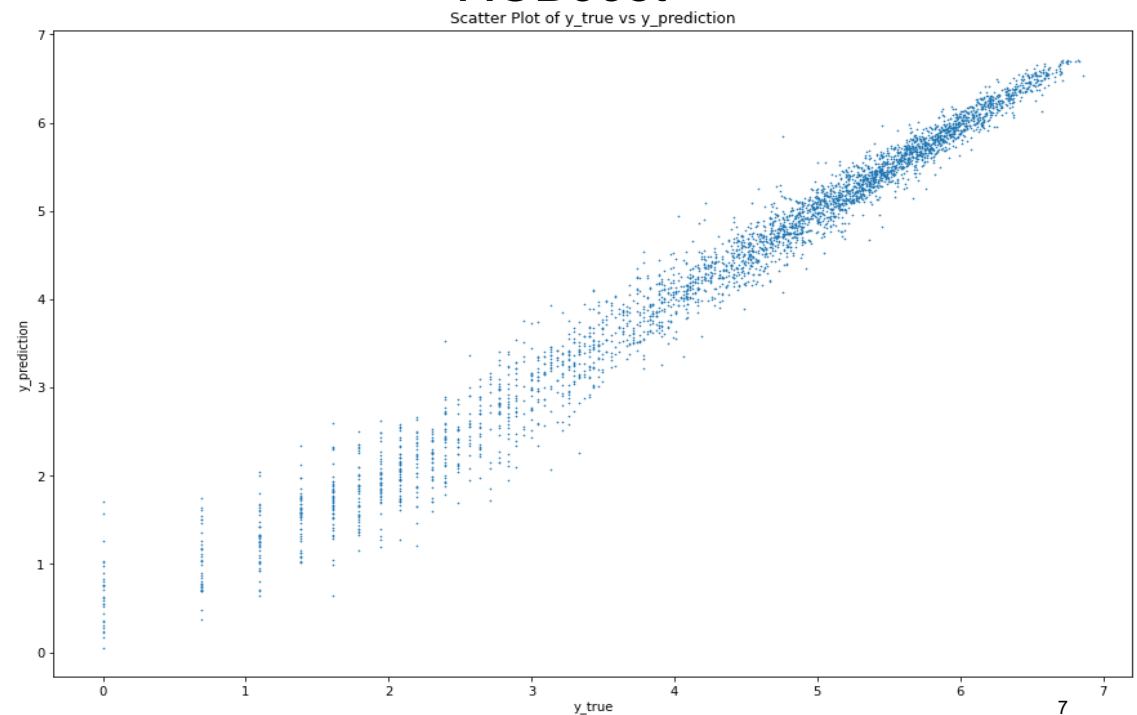| Model | Test Score |
|---|---|
| Ridge Regression | 0.882+ - 0.005 |
| Random Forest | 0.958+- 0.001 |
| GradientBoost | 0.959+-0.001 |
| XGBOOST | 0.940+ - 0.002 |

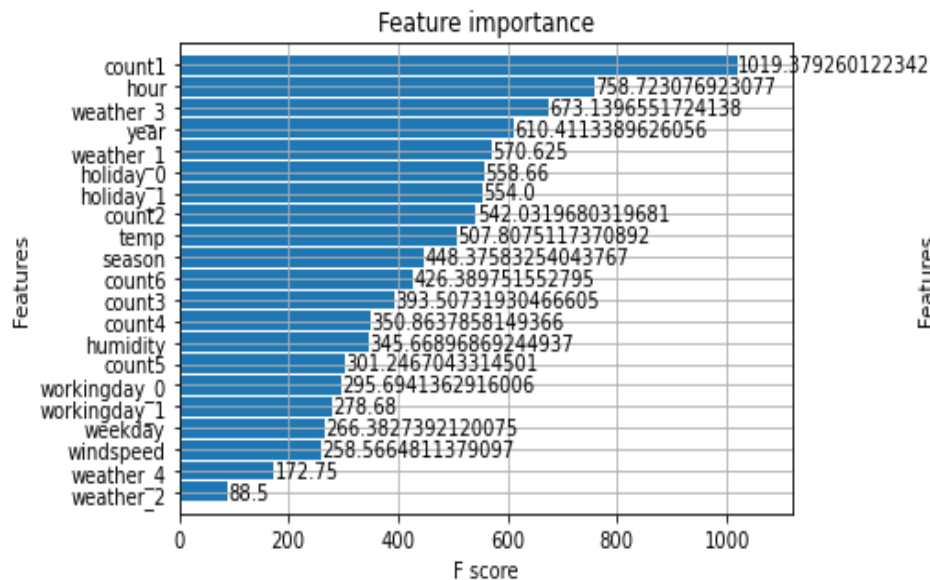- Baseline Model

$$R^2 = 1 - \frac{RSS}{TSS}$$

Become Zero when the RSS equal to TSS

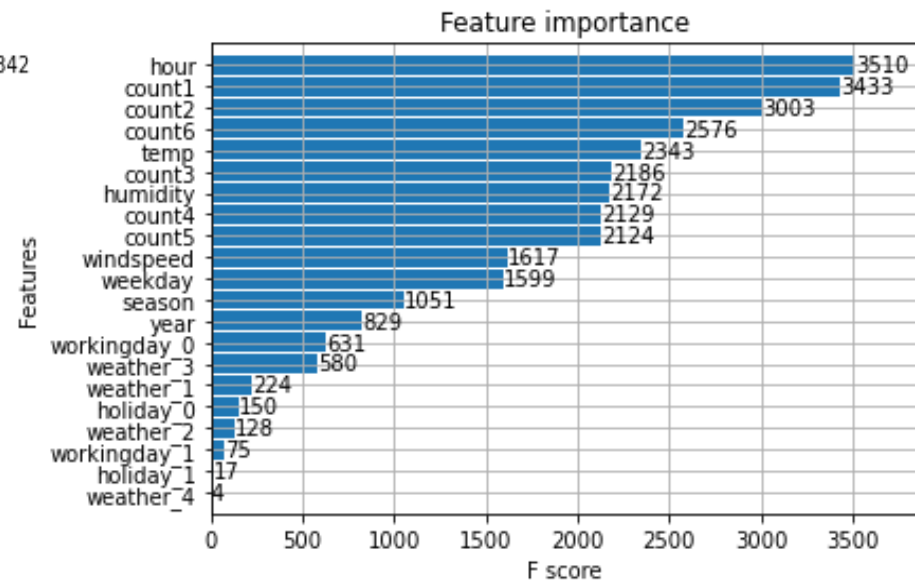- Scatter Plot of y_true vs y_prediction
  - XGBoost


Scatter Plot of y_true vs y_prediction

7

# Global Feature Importance for XGBoost



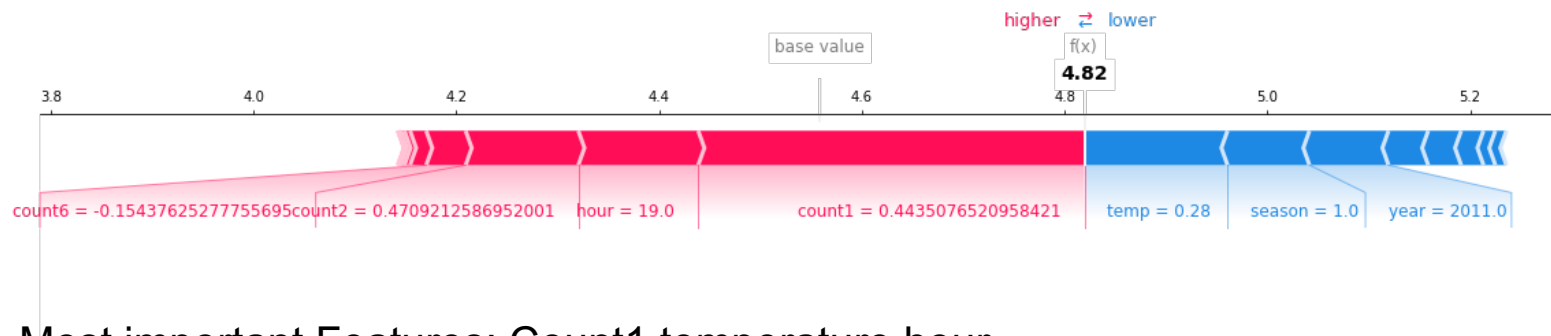Global Feature Importance of Weights

Global Feature Importance of Cover

Most Important Features: Hour, time-lag of 1,2,6 hour(Count1,2,6), weather3(Light Snow, Rain), Holiday, temperatures

# Local Feature Importance for XGBoost

- Data Point 900



Most important Features: Count1,temperature,hour
Least important feature: count6

# Outlook

- For models:
  - Tune parameters more precisely and have a better range
  - Collect more recent data points to make the predictions

- For features:
  - Consider more interactions between features
  - Collect more features including the volume of rainfalls, etc and see their feature importance.

# Questions

Any Questions?