
Predicting Rental Demand for Public Bike Sharing

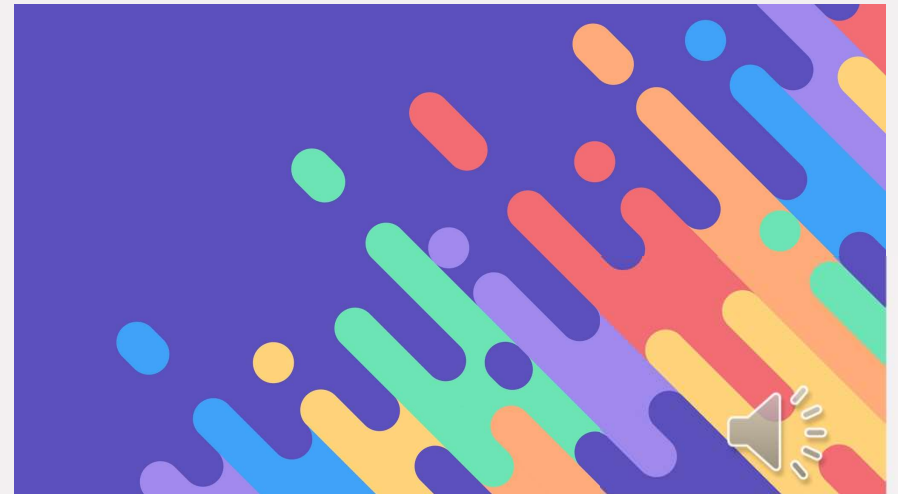
- A Case Study of the Seoul Bike Sharing System

ANA630 Project Paper Presentation

By

Magnus Aghe

M.S. Data Science Program at National University



ABSTRACT

Public bike sharing is becoming an increasing transportation feature of major cities across the world with over 3,000 cities now offering the service as at 2022 (O'Sullivan, 2022).

Because of the high demand for this service, there is a need to make the rental bike available and accessible to the public at the right time thereby reducing waiting time.

This research paper attempts to predict the hourly rental demand for public bikes using the Seoul Bike Sharing System by evaluating various machine learning regression models namely Linear Regression, Polynomial Regression, Decision Tree Regression, Random Forest Regression, Bagging Regressor, and Stacking Regressor.



INTRODUCTION



Seoul's public bike sharing system, called **Ttareungyi**, was set up in 2015, 15 years after the city's first public bicycle rental service started. Inspired by Canadian city, Montreal's modern bike rental called Bixi, this new generation bike sharing system took off.

Ttareungyi is an unmanned rental system that can conveniently be used anywhere, anytime, by anyone. The system was designed to improve citizens' health by enabling the use of bicycles in daily life, achieve a national vision of "Low Carbon, Green Growth", and reduce CO₂ emissions by using bicycles as an alternative form of transportation (Seoul Metropolitan Government).



PROBLEM STATEMENT

With over 25,000 bicycles in more than 2,000 rental bike stations located near substations and in heavily crowded areas, the demand for rental of Seoul Bike has soared, registering hundreds of millions of trips since launching in December 2015. Availability and accessibility of rental bikes to the public at the right time whilst reducing waiting times, has become a major concern.

It is necessary to be able to predict accurately the number of prospective renters of public bikes per time in any given area within the city where Seoul Public Bike operates. Knowing the bike count required each hour will help optimal allocation of bikes to locations where they are needed most.



OBJECTIVE

The objective of this study is to build a predictive model that would predict the required bike count at each hour, for the stable supply of rental bikes to meet demand.

The Seoul Bike Sharing Demand dataset is used in this study. This dataset covers 12 months of data from December 1, 2017, to November 30, 2018, and was obtained from the UC Irvine machine learning data repository website. The dataset has 8760 rows and 14 attributes.

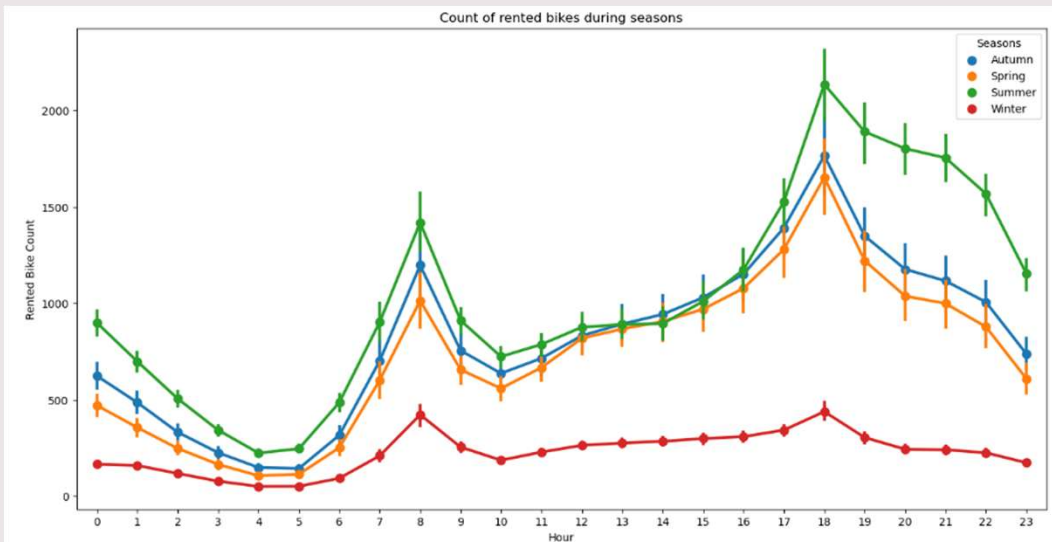
SUMMARY STATISTICS

The mean Rented Bike Count per hour is 704.6 with a minimum of 0, a maximum of 3566 and a median of 504.5. The mean temperature is 12.88 °C, with a minimum of 0, a maximum of 39.4 °C and a median of 13.7 °C. Solar radiation has a mean of 0.57 MJ/m², with a minimum of 0, a maximum of 3.52 MJ/m² and a median of 0.01 MJ/m².

	Rented Bike Count	Hour	Temperature	Humidity	Wind speed	Visibility	Dew point temperature	Radiation	Rainfall	Snowfall
count	8760.00	8760.00	8760.00	8760.00	8760.00	8760.00	8760.00	8760.00	8760.00	8760.00
mean	704.60	11.50	12.88	58.23	1.72	1436.83	4.07	0.57	0.15	0.08
std	645.00	6.92	11.94	20.36	1.04	608.30	13.06	0.87	1.13	0.44
min	0.00	0.00	-17.80	0.00	0.00	27.00	-30.60	0.00	0.00	0.00
25%	191.00	5.75	3.50	42.00	0.90	940.00	-4.70	0.00	0.00	0.00
50%	504.50	11.50	13.70	57.00	1.50	1698.00	5.10	0.01	0.00	0.00
75%	1065.25	17.25	22.50	74.00	2.30	2000.00	14.80	0.93	0.00	0.00
max	3556.00	23.00	39.40	98.00	7.40	2000.00	27.20	3.52	35.00	8.80



VISUALIZING THE DATA CATEGORICAL VARIABLES

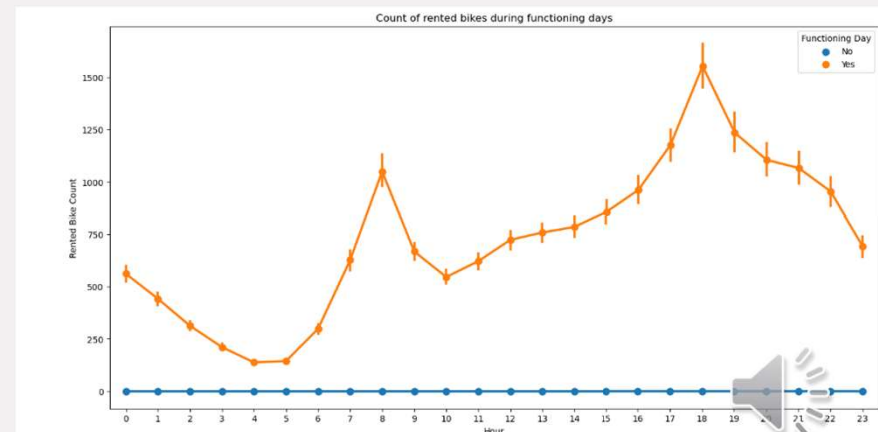
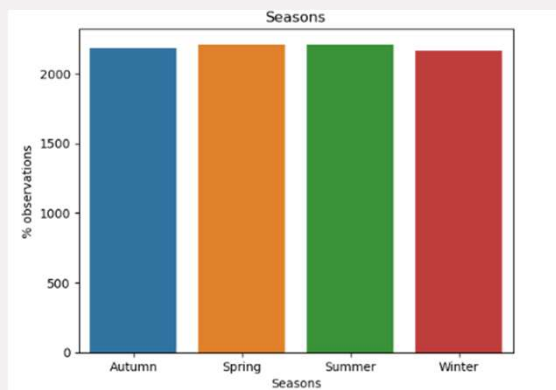


Across all four seasons we observe peak bike rental demands at 18:00 hours for any given day, with the demand much higher during the summer and lowest across peaks during the winter. Spring, Summer and Autumn show similar graph pattern unlike Winter.



There seems to be more demand for rented bikes during non-holidays than during holidays. Also, there appears to be negligible demand for rented bikes across non-functioning days. Almost all demand is experienced during functioning days.

The spread of number of observations in the dataset is uniform across seasons. There were slightly more Spring and Summer days than Autumn and Winter.

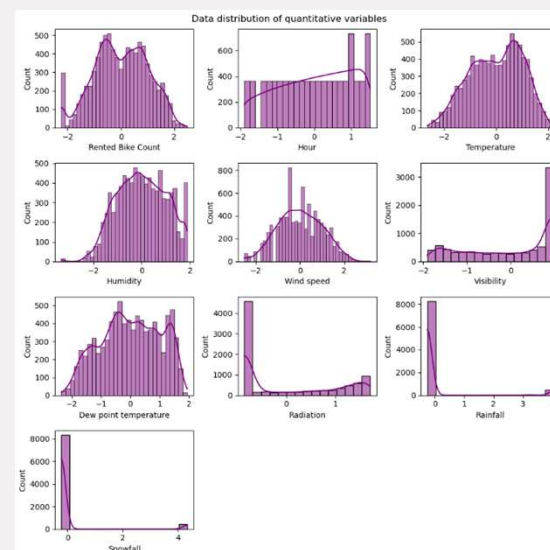
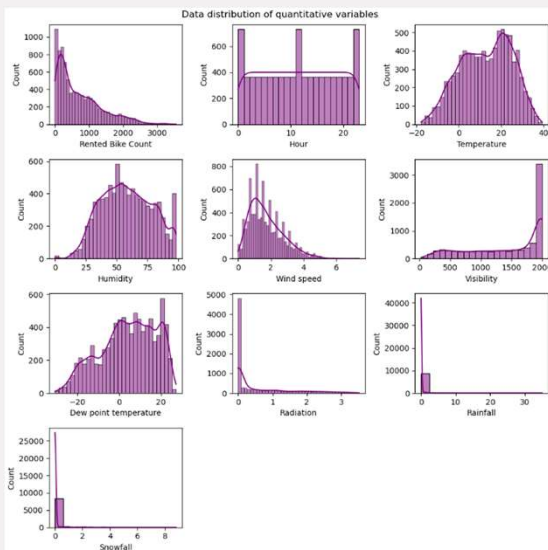


VISUALIZING THE DATA

QUANTITATIVE VARIABLES

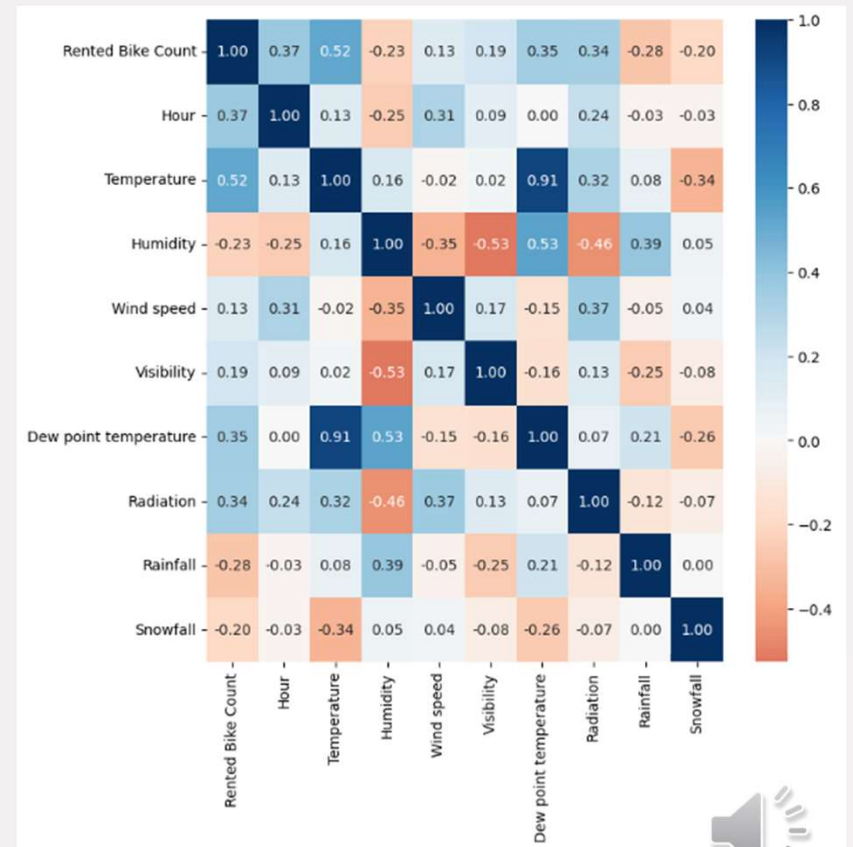
From the plotted histograms of the numeric attributes below, we can see that both positively and negatively skewed distributions exist in the dataset. Since most of the predictor variables are not normally distributed, we have to transform the variables.

Upon applying power transformation, the distribution plots now look more symmetrical as seen below.



CORRELATION ANALYSIS

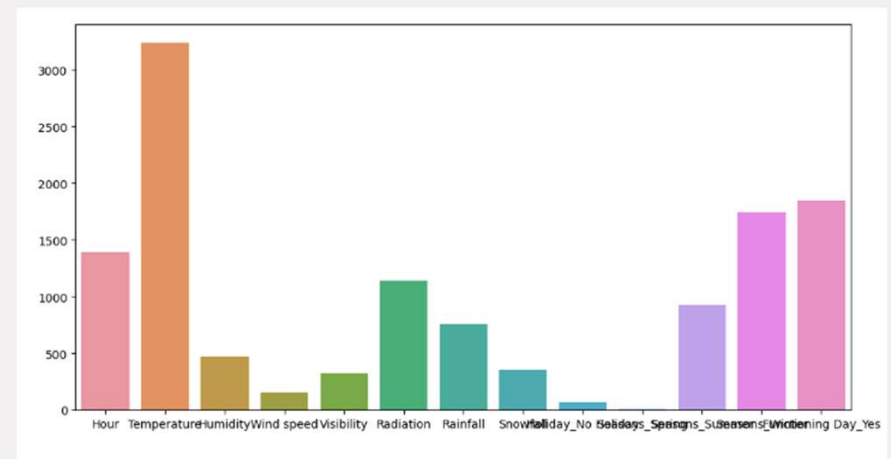
We see that Temperature is the most correlated with Rented Bike Count. Hour, Dew point temperature and Radiation are also quite strongly correlated with Rented Bike Count. Amongst the predictor variables, we see high collinearity between Temperature and Dew point temperature (0.91). This suggests we might remove Dew point temperature from the list of predictor variables.



DATA PREPARATION

FEATURE SELECTION & IMPORTANCE

From the feature importance graph of our predictor variables, Temperature is the most important predictor at 26.02% then Functioning Day (14.88%), Season's Winter (14.02%), Hour (11.2%), Solar Radiation (9.19%), Season's Summer (7.48%), Rainfall (6.08%), Humidity (3.79%). These 8 variables account for 92.66% of the total feature importance. In total, we have 13 predictor variables to be used in our predictive model.



REGRESSION MODELS BUILDING

1. LINEAR REGRESSION
2. POLYNOMIAL REGRESSION
3. DECISION TREES
4. RANDOM FOREST REGRESSION
5. BAGGING REGRESSOR
6. STACKING REGRESSION

First, we split our dataset into two parts 70% training set, and 30% test set. Then we normalize or scale our data. We then build models using Linear Regression, Polynomial Regression, Decision Tree Regression, Random Forest Regression, Bagging Regressor, and Stacking Regressor.

Performance measures for regression used in this study include R-square, Mean Square Error, and K-Cross Validation. Cross-validation is a technique in which we train our model; using the subset of the dataset and then evaluate using the complementary subset of the dataset.



RESULTS: MODEL EVALUATION

	R-Square	MSE	CV Accuracy	CV std
Stacking Regressor	89.56%	10.40%	89.88%	0.49%
Bagging Regressor	89.38%	10.57%	89.73%	0.54%
Random forest Regression	89.02%	10.94%	88.90%	0.69%
Decision Tree Regression	81.36%	18.56%	80.60%	1.38%
Polynomial Regression	80.16%	19.76%	70.99%	0.49%
Linear Regression	70.67%	29.22%	70.99%	0.49%

From lowest to highest in terms of prediction accuracy, the linear regression model had an R-square of 70.67%, a mean square error of 29.22%, cross validation accuracy of 70.99% and CV standard deviation of 0.49%. Polynomial regression model had an R-square of 80.16%, MSE of 19.76%, CV accuracy of 70.99% and a CV deviation of 0.49%. The decision tree regression model had an R-square of 81.36%, MSE of 18.56%, cross validation accuracy of 80.60% and a CV deviation of 1.38%.

Random forest regression model had an R-square of 89.02%, MSE of 10.94%, CV accuracy of 88.90% and a CV deviation of 0.69%. Bagging regressor model had a better R-square of 89.38%, MSE of 10.57%, CV accuracy of 89.73% and a CV deviation of 0.54%. Stacking regressor, the best model, had an R-square of 89.56%, MSE of 10.40%, CV accuracy of 89.88% and a CV deviation of 0.49%.

CONCLUSION



This study focused on predicting the hourly rental demand for public bikes using the Seoul Bike Sharing System Dataset. Six machine learning regression methods were applied namely, Linear Regression, Polynomial Regression, Decision Tree Regression, Random Forest Regression, Bagging Regressor, and Stacking Regressor.

In my findings, stacking regressor, bagging regressor and random forest regression were the three models that stood out with the best R-square and CV accuracy scores, as well as the least mean squared error and least CV standard deviation.

Stacking regressor topped with the overall best model with r-square of 89.56%, CV accuracy of 89.88%, MSE of 10.4% and CV standard deviation of 0.49%.



REFERENCES



UC Irvine Machine Learning Repository (2020). Seoul Bike Sharing Demand.
<http://archive.ics.uci.edu/dataset/560/seoul+bike+sharing+demand>

Seoul Metropolitan Government (2021). 2019 Safe convenient people-centered Seoul Transportation., pp. 31-34.
<https://english.seoul.go.kr/wp-content/uploads/2021/05/2019-Safe-convenient-people-centered-Seoul-Transportation.pdf>

Seoul Metropolitan Government (Undated). Seoul Public Bike.
<https://english.seoul.go.kr/service/movement/seoul-public-bike/1-seoul-public-bike/>

O'Sullivan, F., (2022). The Radical Roots of Bikesharing.
https://getpocket.com/explore/item/the-radical-roots-of-bikesharing?utm_source=pocket-newtab

Brianna, X., (2019). How to Use the Seoul Bike Rental System
<https://10mag.com/seoul-bike-rental/>



THANK YOU!

