# Bike-Sharing Demand

Shril Patel

March 3, 2025

# Outline

# Executive Summary

❖ Rising Importance of Bike Sharing:

➤ Over the past few decades, bike sharing has become increasingly significant as more people seek healthier and more livable cities where such activities are readily accessible.

❖ Prediction of bike rental:

➤ We discovered that a polynomial model with more terms and interactions, achieved the best performance.

➤ Key factors include temperature, rainfall, humidity, peak hour, and weekdays, indicating weathers significant impact on bike rentals.

❖ Dataset Overview:

➤ The Dataset encompasses weather information (including temperature, humidity, windspeed, etc.). Hourly bike rentals counts, and date details for the capital bike share system from 2017 to 2018.

# Introduction

❖ Problem Defining:

➢ It is important for each of these cities to provide a reliable supply of rental bikes to optimize accessibility at all times.

❖ Solid Background:

➢ The global bike sharing cities dataset is an HTML table on the Wikipedia page list of bicycle-sharing systems.

❖ Supportive Tool:

➢ The open Weather API allows users to access current and forecasted weather data for any location including over 200,00 cities.

❖ The Goal:

➢ Minimizing program costs, including bike supply to meet demand, is important. Predicting hourly bike needs based on weather helps optimize supply.
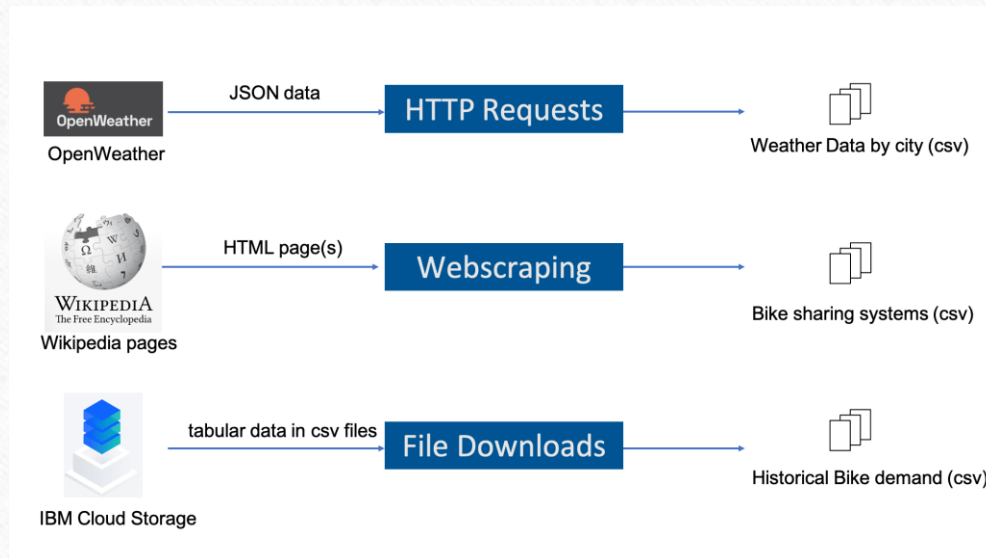
# Methodology

❖ Perform Data Collection

❖ Perform Data Wrangling

➢ With Regular Expressions

➢ With dplyr

❖ Perform exploratory data analysis (EDA) using SQL and visualization

❖ Perform predictive analysis using regression models

➢ How to build baseline model

➢ How to improve the baseline model
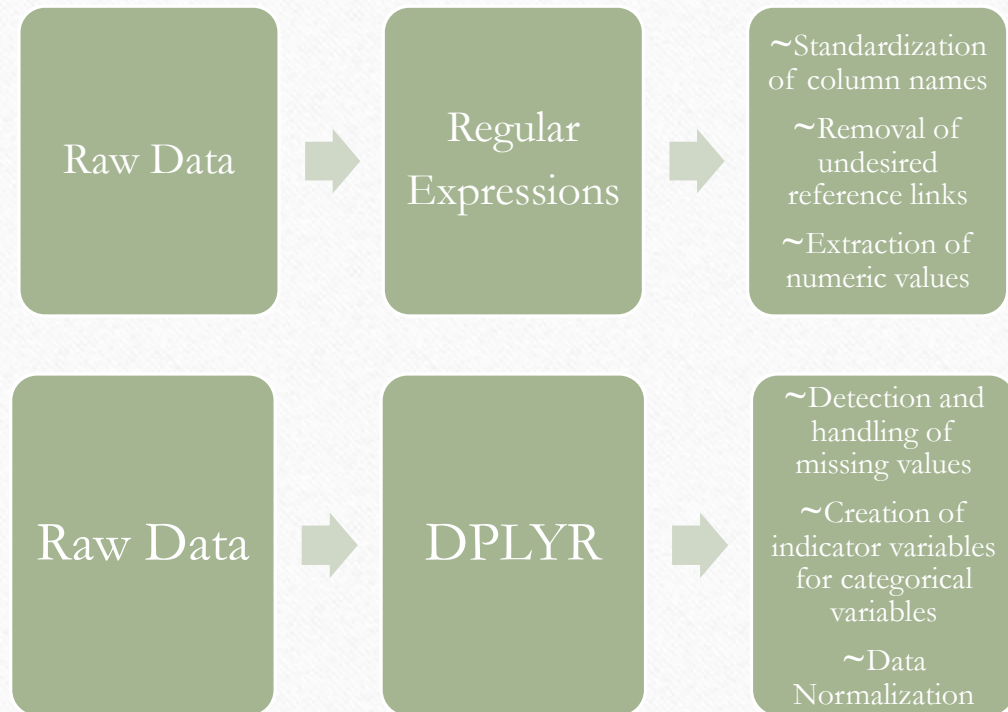
❖ Build a R Shiny dashboard

# Methodology

# Data Collection



- 5-day weather forecast for a list of cities from OpenWeather API.

- Global Bike Sharing Systems Dataset which is an HTML table on the Wikipedia page list of bike-sharing systems: URL. It lists active bicycle-sharing systems around the world.

-  World Cities Data which contains information such as name, latitude, and longitude, about major cities around the world.

- The Seoul Sharing Demand Data set which contains weather information (Temperature, Humidity, Visibility, Dewpoint, Solar radiation, Snowfall, Rainfall), and the number of bikes rented per hour and date

# Data Wrangling

❖ This stage involved cleaning the data by checking for missing values, mis-formatted values and/or unexpected noises.

❖ The first step involved the use of the R package 'string' and regular expression language to standardize column names, remove unwanted reference links from tables and extract numeric values from rows.

❖ The next step involved the use of the 'dplyr' package to detect and handle missing value in the data tables, create indicator (dummy) variable for categorical variables, and normalize the data using min-max normalization.

| Raw Data | ➡ | Regular Expressions | ➡ | ~Standardization of column names<br><br>~Removal of undesired reference links<br><br>~Extraction of numeric values |

| Raw Data | ➡ | DPLYR | ➡ | ~Detection and handling of missing values<br><br>~Creation of indicator variables for categorical variables<br><br>~Data Normalization |

# EDA with SQL

❖ Exploratory Data Analysis was performed on the datasets using Structured Query Language (SQL). The queries performed to;

- ➤ Determine how many records are in the seoul_bike_sharing dataset.

- ➤ Determine how many hours had non-zero rented bike count.

- ➤ Query the weather forecast for Seoul over the next 3 hours.

- ➤ Find which seasons are included in the Seoul bike sharing dataset.

- ➤ Find the first and last dates in the Seoul bike sharing dataset.

- ➤ Determine the date and hour had the most bike rentals.

- ➤ Determine the average hourly temperature and the average number of bike rentals per hour over each season. List the top ten results by average bike count.

- ➤ Find the average hourly bike count during each season.

- ➤ Determine the average TEMPERATURE, HUMIDITY, WIND_SPEED, VISIBILITY, DEW_POINT_TEMPERATURE, SOLAR_RADIATION, RAINFALL, and SNOWFALL per season.

- ➤ Determine the Total Bike Count and City Info for Seoul.

- ➤ Find all city names and coordinates with comparable bike scale to Seoul's sharing systems.

# EDA with Data Visualization

❖ Exploratory Data Analysis was also performed to visually inspect dataset using the ggplot2 library. The ggplot2 library was used;

➤ Create a scatter plot of 'RENTED_BIKE_COUNT' vs 'DATE'

➤ Create the same plot of the 'RENTED_BIKE_COUNT' time series, but now add 'HOURS' as the color.

➤ Create a histogram overlaid with a kernel density curve.

➤ Create a scatter plot to visualize the correlation between 'RENTED_BIKE_COUNT' by 'SEASONS'

➤ Create a display of four boxplots of 'RENTED_BIKE_COUNT' vs. 'HOUR' grouped by 'SEASONS'.

# Predictive Analysis

❖ This stage involved building and refining a regression model to predict the number of bikes that would be rented at each hour of the day by factoring in weather and non-weather conditions.

❖ The cleaned and processed data is divided into two sets: Train set (for model creation and refinement via polynomial regression, interaction, and regularization) and Test set (for Model Evaluation). The constructed models were evaluated using RMSE and RSQ as our assessment metrics, and the regression algorithm was chosen based on the lowest RMSE and highest RSQ values.

Input Data — Train/Test Split — Train Data — Polynomial Regression Interaction Terms Elastic Net Regulation — Test Data — Model Evaluation (RSQ &RMSE) — Best Model Selection

# Build a R Shiny Dashboard

❖ The results of the predictive linear regression model were combined with an interactive dashboard created using the shiny package in R. This dashboard contained;

➢ A basic max bike prediction overview map.

➢ A static temperature trend line.

➢ An interactive bike-sharing demand prediction trend line.

➢ A static humidity and bike-sharing demand prediction correlation plot.

# Results

- ❖ Exploratory Data Analysis Results
- ❖ Predictive Analysis Results
- ❖ A dashboard demo in screenshots

# EDA with SQL

# Busiest Bike Rental Times

❖ On June 19, 2018, at 6:00pm, the highest number of bike rentals in Seoul was recorded.



Task 6 - Subquery - 'all-time high'

determine which date and hour had the most bike rentals.

Solution 6

```
[10]: # provide your solution here
dbGetQuery(conn, "SELECT DATE, HOUR, RENTED_BIKE_COUNT as Maximum_COUNT FROM seoul_bike_sharing
WHERE RENTED_BIKE_COUNT = (SELECT MAX(RENTED_BIKE_COUNT) FROM seoul_bike_sharing)")
```

A data.frame: 1 × 3

| DATE | HOUR | Maximum_COUNT |
| --- | --- | --- |
| <chr> | <dbl> | <dbl> |
| 19/06/2018 | 18 | 3556 |

# Hourly Popularity and Temperature by Seasons

❖ The result of the query showed that more bikes were rented during Summer, Autumn or Spring, between the hours of 5pm and 10pm and at temperatures between 15 degrees Celsius and 30 degrees Celsius which are typical Autumn, Summer and Spring temperatures.

## Task 7 - Hourly popularity and temperature by season ¶

Determine the average hourly temperature and the average number of bike rentals per hour over each season. List the top ten results by average bike count.

### Solution 7

```
# provide your solution here
dbGetQuery(conn, "SELECT SEASONS, HOUR, AVG(RENTED_BIKE_COUNT), AVG(TEMPERATURE) FROM seoul_bike_sharing GROUP
```

A data.frame: 10 × 4

| SEASONS | HOUR | AVG(RENTED_BIKE_COUNT) | AVG(TEMPERATURE) |
|---------|------|------------------------|------------------|
| <chr>   | <dbl> | <dbl>                 | <dbl>            |
| Summer  | 18   | 2135.141               | 29.38791         |
| Autumn  | 18   | 1983.333               | 16.03185         |
| Summer  | 19   | 1889.250               | 28.27378         |
| Summer  | 20   | 1801.924               | 27.06630         |
| Summer  | 21   | 1754.065               | 26.27826         |
| Spring  | 18   | 1689.311               | 15.97222         |
| Summer  | 22   | 1567.870               | 25.69891         |
| Autumn  | 17   | 1562.877               | 17.27778         |
| Summer  | 17   | 1526.293               | 30.07691         |
| Autumn  | 19   | 1515.568               | 15.06346         |

# Rental Seasonality

❖ The result of the query showed that on average more bikes were rented during Summer and fewer bikes were rented during Winter



## Task 8 - Rental Seasonality

**Find the average hourly bike count during each season.**

Also include the minimum, maximum, and standard deviation of the hourly bike count for each season.

> Hint : Use the SQRT(AVG(col*col) - AVG(col)*AVG(col) ) function where col refers to your column name for finding the standard deviation

## Solution 8

```
# provide your solution here

dbGetQuery(conn, " SELECT SEASONS, AVG(RENTED_BIKE_COUNT) as AVG_S_COUNT, MIN(RENTED_BIKE_COUNT) as M
group by (SEASONS)
Order by AVG_S_COUNT desc")

# No Standard deviation function available in SQLlite so I will not include it.
```

A data.frame: 4 × 4

| SEASONS | AVG_S_COUNT | MIN_S_COUNT | MAX_S_COUNT |
|---------|-------------|-------------|-------------|
| <chr> | <dbl> | <dbl> | <dbl> |
| Summer | 1034.0734 | 9 | 3556 |
| Autumn | 924.1105 | 2 | 3298 |
| Spring | 746.2542 | 2 | 3251 |
| Winter | 225.5412 | 3 | 937 |

Let's explore a bit and see what might be the most significant contributing factors in terms of the provided data.

# Weather Seasonality

❖ The results showed that the temperature, rainfall, humidity, snowfall, etc. readings were very correlated with the season of the year. For example, there was a higher average temperature and rainfall in the summer, while snowfall was confined to the winter season. It also showed a correlation between the number of rented bikes and the season of the year.



**Task 9 - Weather Seasonality**

Consider the weather over each season. On average, what were the TEMPERATURE, HUMIDITY, WIND_SPEED, VISIBILITY, DEW_POINT_TEMPERATURE, SOLAR_RADIATION, RAINFALL, and SNOWFALL per season?

Include the average bike count as well , and rank the results by average bike count so you can see if it is correlated with the weather at all.

**Solution 9**

```
# provide your solution here
dbGetQuery(conn," SELECT SEASONS, AVG(RENTED_BIKE_COUNT) as AVG_S_COUNT, AVG(TEMPERATURE) as AVG_S_TEMP, AVG(HUMIDITY) as AVG_S_HUMIDITY, AVG(WIND_SPEED) as AV
group by (SEASONS)
Order by AVG_S_COUNT desc")
```

A data.frame: 4 × 10

| SEASONS | AVG_S_COUNT | AVG_S_TEMP | AVG_S_HUMIDITY | AVG_WIND_SPEED | AVG_VISIBILITY | AVG_DEW_POINT | AVG_SOLAR_RADIATION | AVG_RAINFALL | AVG_SNOWFALL |
|---------|-------------|------------|----------------|----------------|----------------|---------------|---------------------|--------------|--------------|
| <chr> | <dbl> | <dbl> | <dbl> | <dbl> | <dbl> | <dbl> | <dbl> | <dbl> | <dbl> |
| Summer | 1034.0734 | 26.587711 | 64.98143 | 1.609420 | 1501.745 | 18.750136 | 0.7612545 | 0.25348732 | 0.00000000 |
| Autumn | 924.1105 | 13.821580 | 59.04491 | 1.492101 | 1558.174 | 5.150594 | 0.5227827 | 0.11765617 | 0.06350026 |
| Spring | 746.2542 | 13.021685 | 58.75833 | 1.857778 | 1240.912 | 4.091389 | 0.6803009 | 0.18694444 | 0.00000000 |
| Winter | 225.5412 | -2.540463 | 49.74491 | 1.922685 | 1445.987 | -12.416667 | 0.2981806 | 0.03282407 | 0.24750000 |

# Bike-Sharing Info in Seoul

❖ The result of the query showed that there were 20,000 bikes available for rent in the city of Seoul.



Task 10 - Total Bike Count and City Info for Seoul

Use an implicit join across the WORLD_CITIES and the BIKE_SHARING_SYSTEMS tables to determine the total number of bikes avaialble in Seoul, plus the following city information about Seoul: CITY, COUNTRY, LAT, LON, POPULATION, in a single view.

Notice that in this case, the CITY column will work for the WORLD_CITIES table, but in general you would have to use the CITY_ASCII column.

Solution 10

```
# provide your solution here
dbGetQuery(conn, "SELECT B.BICYCLES, B.CITY, B.COUNTRY, W.LAT, W.LNG, W.POPULATION  FROM BIKE_
LEFT JOIN WORLD_CITIES AS W ON B.CITY = W.CITY_ASCII
WHERE B.CITY = 'Seoul'")
```

A data.frame: 1 × 6

| BICYCLES | CITY | COUNTRY | LAT | LNG | POPULATION |
|---|---|---|---|---|---|
| <dbl> | <chr> | <chr> | <dbl> | <dbl> | <dbl> |
| 20000 | Seoul | South Korea | 37.5833 | 127 | 21794000 |

# Cities Similar to Seoul

❖ The result of the query showed that six cities in China namely; Zhuzhou, Xi'an, Weifang, Shanghai, Beijing and Ningbo had similar number of rental bikes like Seoul.



**Task 11 - Find all city names and coordinates with comparable bike scale to Seoul's bike sharing system**

Find all cities with total bike counts between 15000 and 20000. Return the city and country names, plus the coordinates (LAT, LNG), population, and number of bicycles for each city.

Later we will ask you to visualize these similar cities on leaflet, with some weather data.

**Solution 11**

```
# provide your solution here
dbGetQuery(conn, "SELECT B.BICYCLES, B.CITY, B.COUNTRY, W.LAT, W.LNG, W.POPULATION  FROM BIKE_
LEFT JOIN WORLD_CITIES AS W ON B.CITY = W.CITY_ASCII
WHERE B.CITY = 'Seoul' OR B.BICYCLES BETWEEN 15000 AND 20000
order by B.BICYCLES desc")

dbListTables(conn)
```

A data.frame: 9 × 6

| BICYCLES | CITY | COUNTRY | LAT | LNG | POPULATION |
|---|---|---|---|---|---|
| <dbl> | <chr> | <chr> | <dbl> | <dbl> | <dbl> |
| 20000 | Kunshan | China | NA | NA | NA |
| 20000 | Weifang | China | 36.7167 | 119.1000 | 9373000 |
| 20000 | Xi'an | China | 34.2667 | 108.9000 | 7135000 |
| 20000 | Zhuzhou | China | 27.8407 | 113.1469 | 3855609 |
| 20000 | Seoul | South Korea | 37.5833 | 127.0000 | 21794000 |
| 19165 | Shanghai | China | 31.1667 | 121.4667 | 22120000 |
| 18000 | Xuzhou | China | NA | NA | NA |
| 16000 | Beijing | China | 39.9050 | 116.3914 | 19433000 |
| 15000 | Ningbo | China | 29.8750 | 121.5492 | 7639000 |

# EDA with Visualization

# Bike Rental Vs. Date

❖ The number of bicycle rentals starts to rise in the spring and peaks in the summer months of June and July. We also observe a slight decline in August as the season draws to a close, followed by another high in September. The number of hired bikes also gradually declines as the winter months approach as we approach the conclusion of the year

# Bike Rental Vs. Datetime

❖ The largest hourly bike rentals are observed between the hours of 6pm and 7pm. In the mornings around 7 a.m., we also observe a moderate to high volume of bike rentals. Generally, more bike rentals occur during the day, and this pattern continues throughout the year.

# Bike Rental Histogram

❖ It also demonstrates that rental bike counts ranging from 0 to 500 occur most frequently in the sample. The 'mode,' or most often rented number of bikes, is about 250. Because the density curve is tilted to the left, the mean is smaller than the median.

# Daily Total Rainfall & Snowfall

❖ Snowfall primarily occurs during winter days. Rainfall occurs almost year-round, with its main peaks in June, July, and September

# Predictive Analysis

# Ranked Coefficients

❖ Rainfall, humidity, the 6 PM temperature, and overall temperature are the strongest predictors of bike rentals.

❖ Rainfall and humidity make biking less comfortable and practical.

❖ Temperature, especially at 6 PM—a key time for commuting and leisure—encourages biking when conditions are favorable. Together, these factors highlight the influence of weather on outdoor activity preferences.

# Model Evaluation

❖ The Root Mean Squared Error (RMSE) and R Squared (RSQ) metrics were used to assess the performance of each model. The bar charts on the right show the RMSE and RSQ performance of all models.



Comparing RSME AND RSQ Across Models

# Q-Q Plot of the Best Model

❖ Q-Q plot of the best model's test results vs the truths

# DashBoard

# Global Bike Demand Prediction

Map of predicted bike demand for the cities of Seoul, Suzhou, London, New York and Paris.

# Dashboard City View: Seoul

The following screenshot displays Seoul's temperature forecast and bike demand forecast. However, the impact of humidity on predicted bike demand is missing due to technical reasons.

# Dashboard City View: London

The following screenshot displays London's temperature forecast, bike demand forecast, and the impact of humidity on predicted bike demand.

# Conclusion

❖ The summer months of June and July saw the biggest demand for bike sharing.

❖ During the day, the peak demand for rental bikes is between 6 and 7 p.m.

❖ During the winter, there is virtually little demand for bicycles.

❖ The quantity of bikes hired throughout each hour of any given day is heavily influenced by temperature and humidity.

# Appendix

❖ Web Scrapping

❖ OpenWeather  API

❖ Data Wrangling – Regular Expression

❖ EDA – SQL

❖ EDA with Visualiztion

# Appendix: Web Scrapping

# Appendix: OpenWeather API

# Appendix: Data Wrangling – Regular Expressions

# Appendix: EDA – SQL

```
1    #Install RSQLite package
2    install.packages("https://cran.r-project.org/src/contrib/Archive/RSQLite/RSQLite_0.10.0.tar.gz", repos = NULL, type = "source", dependencies = TRUE)
3    library("RSQLite")
4    setwd("C:/Users/Mr... /Desktop/IBM CERTIFICATE/9 Capstone Project")
5
6    library(tidyverse)
7    library(stringr)   # For string manipulation
8    library(readr)     # For reading CSV files (read_csv())
9    library(dplyr)
10   library(ggplot2)
11
12   #Establish Connection
13   conn <- dbConnect(RSQLite::SQLite(),"RDB.sqlite")
14
15   #Read Datasets
16   SEOUL_BIKE_SHARING <- read_csv("seoul_bike_sharing.csv")
17   CITIES_WEATHER_FORECAST <- read_csv("cities_weather_forecast.csv")
18   BIKE_SHARING_SYSTEMS <- read_csv("bike_sharing_systems.csv")
19   WORLD_CITIES <- read_csv("world_cities.csv")
20
21   #Load Tables
22   dbWriteTable(conn, "SEOUL_BIKE_SHARING", SEOUL_BIKE_SHARING, overwrite=TRUE, header = TRUE)
23   dbWriteTable(conn, "CITIES_WEATHER_FORECAST", CITIES_WEATHER_FORECAST, overwrite=TRUE, header = TRUE)
24   dbWriteTable(conn, "BIKE_SHARING_SYSTEMS", BIKE_SHARING_SYSTEMS, overwrite=TRUE, header = TRUE)
25   dbWriteTable(conn, "WORLD_CITIES", WORLD_CITIES, overwrite=TRUE, header = TRUE)
26   dbListTables(conn)
27
28
29   #### TASK 1: Determine how many records are in the seoul_bike_sharing dataset.
30   dbGetQuery(conn, 'SELECT COUNT(*) AS Records FROM SEOUL_BIKE_SHARING')
31
32   #### TASK 2: Determine how many hours had non-zero rented bike count.
33   dbGetQuery(conn, "SELECT count(HOUR) as Numer_of_hours FROM SEOUL_BIKE_SHARING
34   WHERE RENTED_BIKE_COUNT > 0")
35
36   #### TASK 3: Query the the weather forecast for Seoul over the next 3 hours.
37   #Recall that the records in the CITIES_WEATHER_FORECAST dataset are 3 hours apart, so we just need the first record from the query
38
39   dbGetQuery(conn, "SELECT * FROM CITIES_WEATHER_FORECAST
40   WHERE CITY = 'Seoul'
41   Limit 1")
42
43   #### TASK 4: Find which seasons are included in the seoul bike sharing dataset.
44   dbGetQuery(conn, "SELECT distinct SEASONS as Seasons
45   FROM SEOUL_BIKE_SHARING")
46
47   #### TASK 5: Find the first and last dates in the Seoul Bike Sharing dataset.
48   dbGetQuery(conn, "SELECT MIN(DATE) as Start_Date, MAX(DATE) as End_Date
49   FROM SEOUL_BIKE_SHARING")
50
51   #### TASK 6: Determine which date and hour had the most bike rentals.
52   dbGetQuery(conn, "
53   SELECT DATE, HOUR, RENTED_BIKE_COUNT AS Maximum_COUNT
54   FROM SEOUL_BIKE_SHARING
55   WHERE RENTED_BIKE_COUNT = (SELECT MAX(RENTED_BIKE_COUNT) FROM SEOUL_BIKE_SHARING)")
56
57   #### TASK 7: Determine the average hourly temperature and the average number of bike rentals per hour over each season. List the top ten results by average bike count.
58   dbGetQuery(conn, "SELECT SEASONS, HOUR, AVG(RENTED_BIKE_COUNT), AVG(TEMPERATURE)
59                FROM SEOUL_BIKE_SHARING
60                GROUP BY SEASONS, HOUR
61                ORDER BY AVG(RENTED_BIKE_COUNT) DESC
62                LIMIT 10")
```

```
64   #### TASK 8: Find the average hourly bike count during each season.
65   #Include:  min, max, and sd of the hourly bike count for each season
66   dbGetQuery(conn, "
67   SELECT
68       SEASONS,
69       AVG(RENTED_BIKE_COUNT) AS Avg_Bike_Count,
70       MIN(RENTED_BIKE_COUNT) AS Min_Bike_Count,
71       MAX(RENTED_BIKE_COUNT) AS Max_Bike_Count,
72       SQRT(AVG(RENTED_BIKE_COUNT * RENTED_BIKE_COUNT) - AVG(RENTED_BIKE_COUNT) * AVG(RENTED_BIKE_COUNT)) AS Std_Dev_Bike_Count
73   FROM SEOUL_BIKE_SHARING
74   GROUP BY SEASONS")
75
76   #### TASK 9: Consider the weather over each season.
77   #On avg, what were the TEMPERATURE, HUMIDITY, WIND_SPEED, VISIBILITY, DEW_POINT_TEMPERATURE, SOLAR_RADIATION, RAINFALL, and SNOWFALL per season?
78   #Include the average bike count as well , and rank the results by average bike count
79   dbGetQuery(conn, "
80   SELECT
81       SEASONS,
82       AVG(TEMPERATURE) AS Avg_Temperature,
83       AVG(HUMIDITY) AS Avg_Humidity,
84       AVG(WIND_SPEED) AS Avg_WindSpeed,
85       AVG(VISIBILITY) AS Avg_Visibility,
86       AVG(DEW_POINT_TEMPERATURE) AS Avg_DewPointTemp,
87       AVG(SOLAR_RADIATION) AS Avg_SolarRad,
88       AVG(RAINFALL) AS Avg_Rainfall,
89       AVG(SNOWFALL) AS Avg_Snowfall,
90       AVG(RENTED_BIKE_COUNT) AS Avg_RentedBikes
91   FROM SEOUL_BIKE_SHARING
92   GROUP BY SEASONS
93   ORDER BY AVG(RENTED_BIKE_COUNT) DESC")
94
95   #### TASK 10: Use an implicit join across the WORLD_CITIES and the BIKE_SHARING_SYSTEMS tables to determine the total number of bikes avaialble in Seoul, plus the following c
96   dbGetQuery(conn, "SELECT B.BICYCLES, B.CITY, B.COUNTRY, W.LAT, W.LNG, W.POPULATION
97   FROM BIKE_SHARING_SYSTEMS AS B
98   LEFT JOIN WORLD_CITIES AS W
99   ON B.CITY = W.CITY_ASCII
100  WHERE B.CITY = 'Seoul'")
101
102  #TASK 11: Find all cities with total bike counts between 15000 and 20000.
103  #Return the city and country names, plus the coordinates (LAT, LNG), population, and number of bicycles for each city.
104
105  dbGetQuery(conn, "
106  SELECT
107      B.BICYCLES,
108      B.CITY,
109      B.COUNTRY,
110      W.LAT,
111      W.LNG,
112      W.POPULATION
113  FROM BIKE_SHARING_SYSTEMS AS B
114  LEFT JOIN WORLD_CITIES AS W
115      ON B.CITY = W.CITY_ASCII
116  WHERE B.CITY = 'Seoul'
117      OR (B.BICYCLES BETWEEN 15000 AND 20000)
118  ORDER BY B.BICYCLES DESC
119  ")
120
121  dbListTables(conn)
122
123  # Once you're done with the connection, close it
124  dbDisconnect(conn)
```

# EDA with Visualizations