

Soul Bike Sharing Prediction Report

Created By Alice Fu

2024/12/8



CONTENT



- Executive Summary
- Introduction
- Methodology
- Results
 - Visualization – Charts
 - Dashboard
- Discussion
 - Findings & Implications
- Conclusion
- Appendix

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 number:

- 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 weather's significant impact on bike rentals.

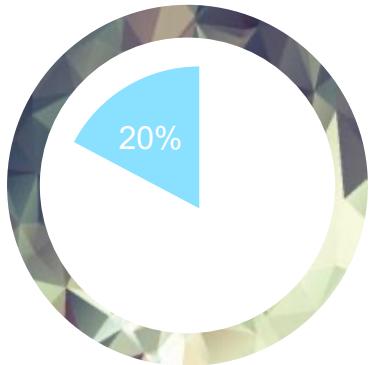
❖ Dataset Overview:

The dataset encompasses weather information (including temperature, humidity, windspeed, etc.), hourly bike rental counts, and date details for the Capital bike share system from 2017 to 2018.





INTRODUCTION



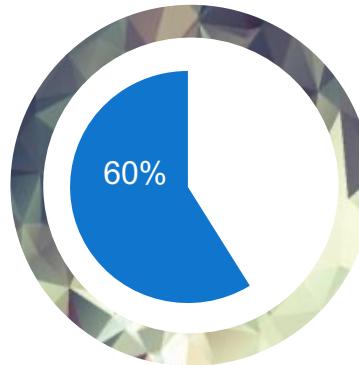
Problems Defining

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



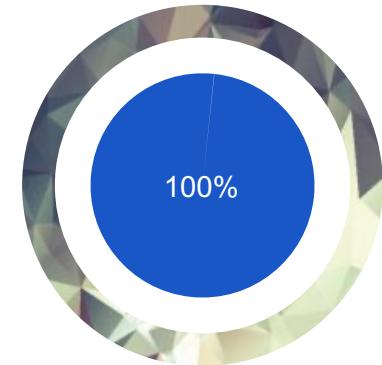
Solid Background

The Global Bike Sharing Cities Dataset is an HTML table on the Wikipedia page List of bicycle-sharing systems



Supportive Tool

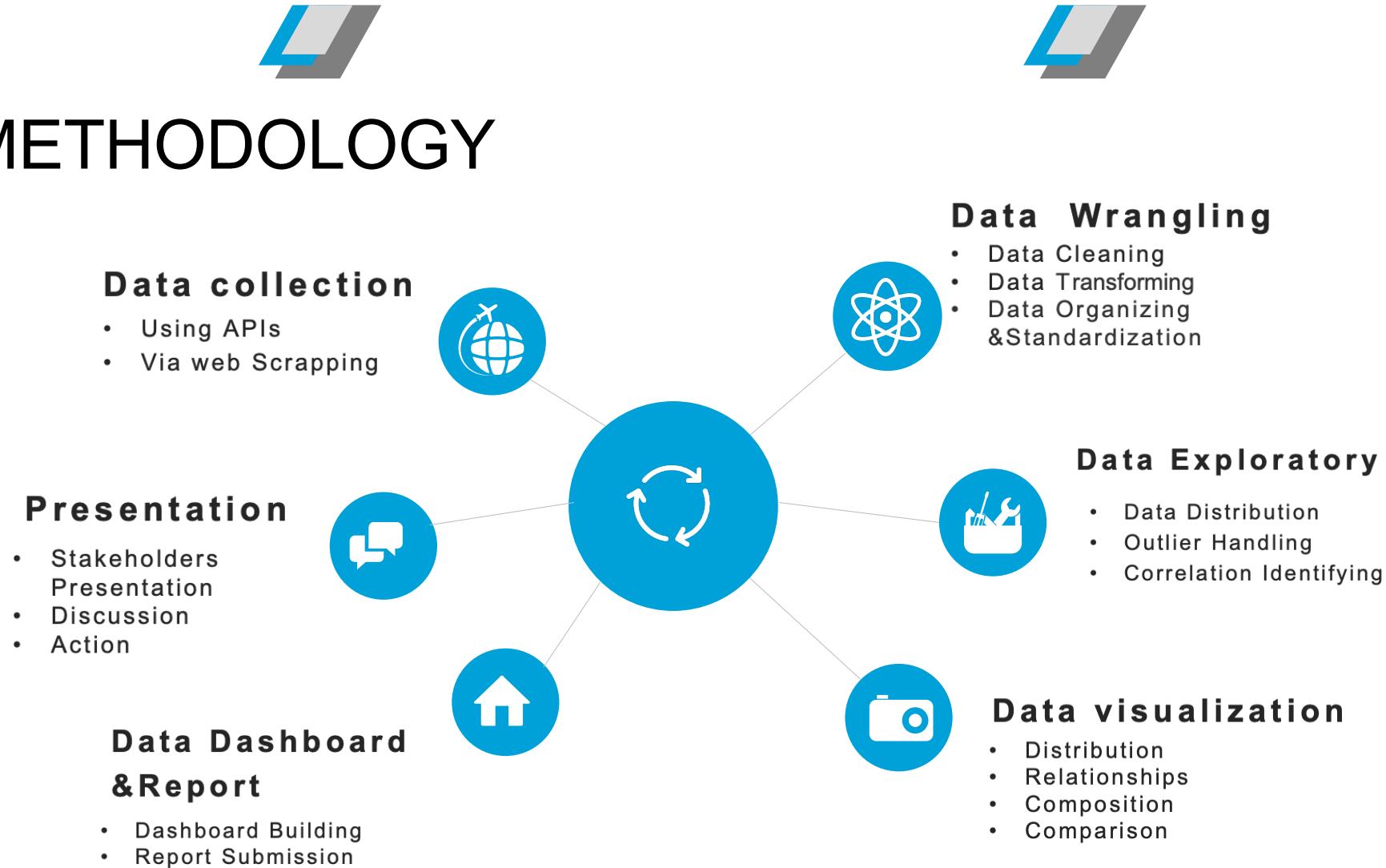
The Open Weather API allows users to access current and forecasted weather data for any location including over 200,000 cities.



The Goal

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

METHODOLOGY



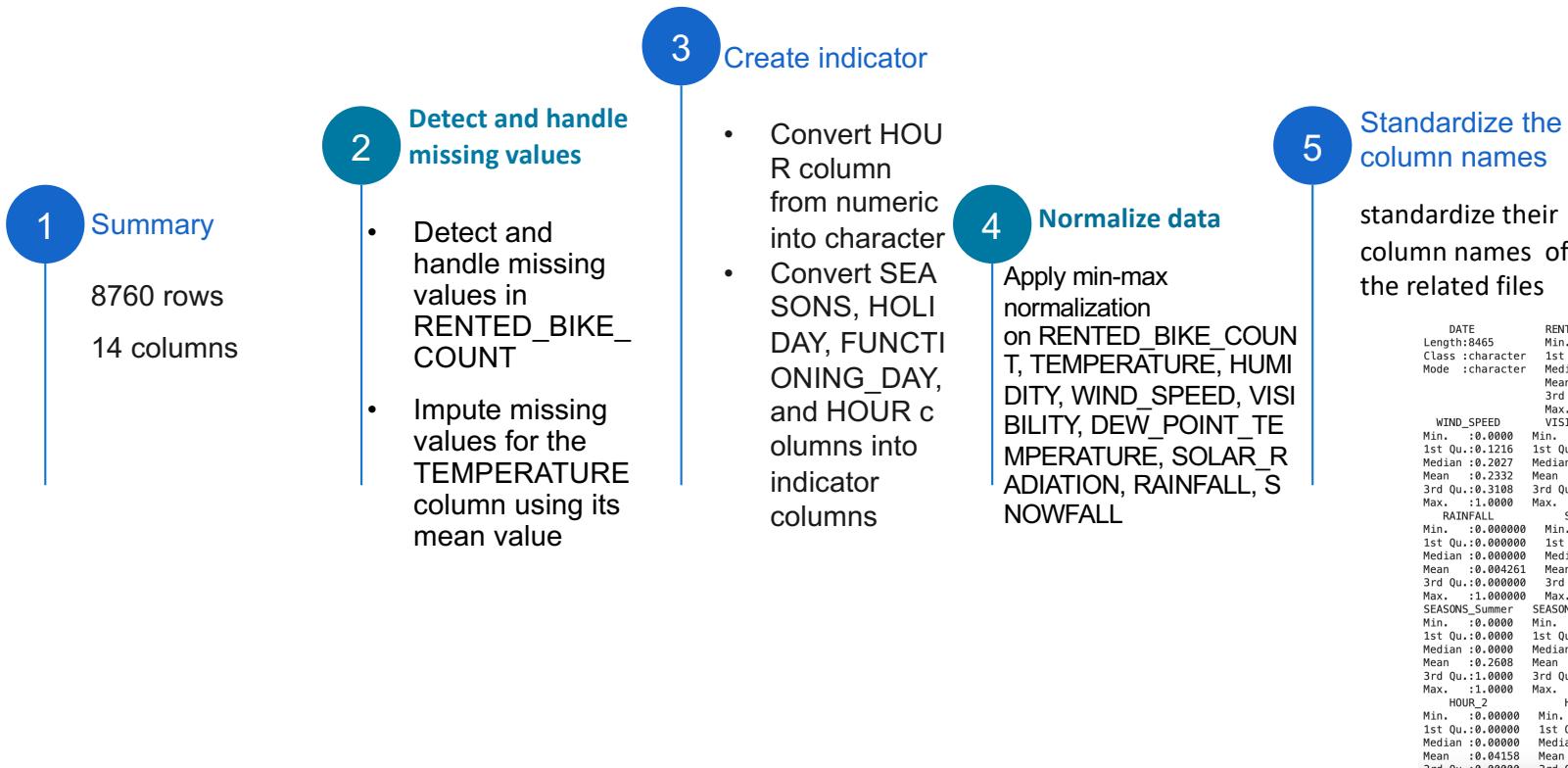


RESULTS



DATA WRANGLING WITH DPLYR

Wrangling the Seoul bike-sharing demand historical dataset using Dplyr





DATA WRANGLING WITH REGULAR EXPRESSIONS

Clean up the bike-sharing systems data using Tidyverse

variable	class
	<chr>
A tibble: 4 × 2	
COUNTRY	character
CITY	character
SYSTEM	character
BICYCLES	character

BICYCLES	CITY
<chr>	<chr>
A spec_tbl_df: 10 × 1	A spec_tbl_df: 10 × 1
4115[22]	Melbourne[12]
310[59]	Brisbane[14][16]
500[72]	Lower Austria[18]
[75]	Namur[19]
180[76]	Brussels[21]
600[77]	Salvador[23]
[78]	Belo Horizonte[24]
initially 800 (later 2500)	João Pessoa[25]
100 (220)	(Pedro de) Toledo[26]
370[114]	Rio de Janeiro[27]

2 Remove undesired reference links

Find any elements in the column containing non-numeric characters

3 Remove reference links

Use the dplyr::mutate() function to apply the remove_ref function to the CITY and SYSTEM columns

```
# A tibble: 480 × 4
  COUNTRY CITY      SYSTEM BICYCLES
  <chr>   <chr>     <chr>   <chr>
  1 Albania Tirana    <NA>    200
  2 Argentina Mendoza <NA>    40
  3 Argentina San Lorenzo, Santa Fe Biciudad 80
  4 Argentina Buenos Aires Serttel Brasil 4000
  5 Argentina Rosario  <NA>    480
  6 Australia Melbourne PBSC & 8D 676
  7 Australia Brisbane 3 Gen. Cyclocity 2000
  8 Australia Melbourne 4 Gen. oBike   1250
  9 Australia Sydney   4 Gen. oBike   1250
  10 Australia Sydney  4 Gen. Ofo    600
# ... with 470 more rows
```

4 Extract the numeric value

Use the mutate() function to apply extract_number on the BICYCLES column

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	NA's
5.0	100.0	335.5	2052.3	1468.2	78000.0	104

1 Type Check

Character Type

5 Summary

Use the summary function to check the descriptive statistics of the numeric BICYCLES column



EDA WITH SQL

Perform exploratory data analysis using SQL queries with the RSQLite R

```
WORLD_CITIES has 26569 rows.
BIKE_SHARING_SYSTEMS has 480 rows.
CITIES_WEATHER_FORECAST has 160 rows.
SEOUL_BIKE_SHARING has 8465 rows.
```

1. 'BIKE_SHARING_SYSTEMS'
2. 'CITIES_WEATHER_FORECAST'
3. 'SEOUL_BIKE_SHARING'
4. 'WORLD_CITIES'

Count_of_Records	<int>
A data.frame: 1 × 1	
1	8465

1 Record Count
8465 rows

Operational Hours

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

Numer_of_hours	<int>
A data.frame: 1 × 1	
1	8465

Start_Date	End_Date
<chr>	<chr>
A data.frame: 1 × 2	
1 01/01/2018	31/12/2017

Weather Outlook

- 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.

SEASONS	HOUR	AVG(RENTED_BIKE_COUNT)		AVG(TEMPERATURE)
		<chr>	<dbl>	
A data.frame: 10 × 4				
1	Summer	18	2135.141	29.38791
2	Autumn	18	1983.333	16.03185
3	Summer	19	1889.250	28.27378
4	Summer	20	1801.924	27.06630
5	Summer	21	1754.065	26.27826
6	Spring	18	1689.311	15.97222
7	Summer	22	1567.870	25.69891
8	Autumn	17	1562.877	17.27778
9	Summer	17	1526.293	30.07691
10	Autumn	19	1515.568	15.06346

BICYCLES	CITY	COUNTRY	LAT	LNG	POPULATION
<dbl>	<chr>	<chr>	<dbl>	<dbl>	<dbl>
1	20000	Seoul	South Korea	37.5833	127.0000
2	20000	Kunshan	China	NA	NA
3	20000	Weifang	China	36.7167	119.1000
4	20000	Xi'an	China	34.2667	108.9000
5	20000	Zhuzhou	China	27.8407	113.1469
6	19165	Shanghai	China	31.1667	121.4667
7	18000	Xuzhou	China	NA	NA
8	16000	Beijing	China	39.9050	116.3914
9	15000	Ningbo	China	29.8750	121.5492

Summary

- Total Bike Count and City Info for Seoul
- Find all city names and coordinates with comparable bike scale to Seoul's bike sharing system

5

Popularity Explore

- Determine which 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.
- Consider the weather over each season

4

Weather Outlook

- 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.

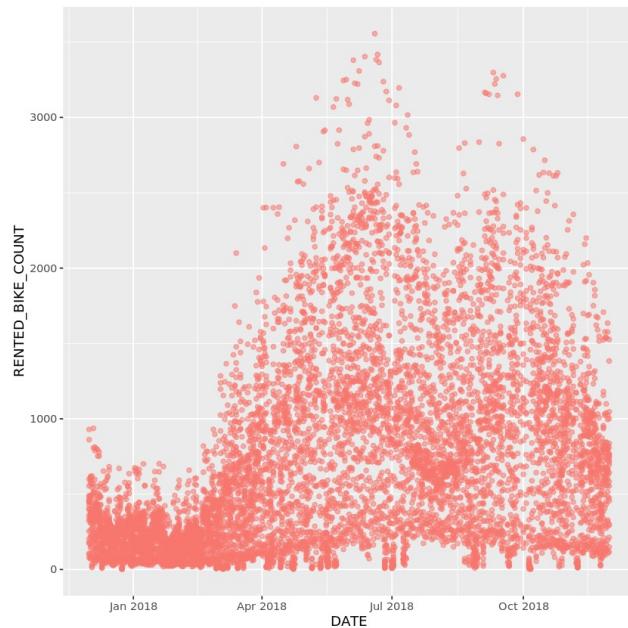
3



EDA WITH VISUALIZATION

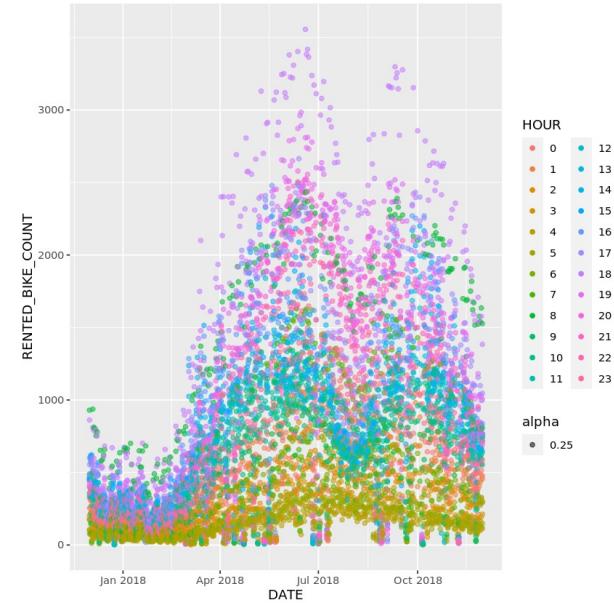
Perform exploratory data analysis using Visualization with tidyverse & ggplot

Create a scatter plot of RENTED_BIKE_COUNT vs DATE



We can see the rented bike count start to increase around **FEB/March** and reach the **max on June** then decrease little bit towards AUG then increase around **SEP** and then start decreasing again towards the end of the year.

Create the same plot of the RENTED_BIKE_COUNT time series but now add HOURS as the colour.



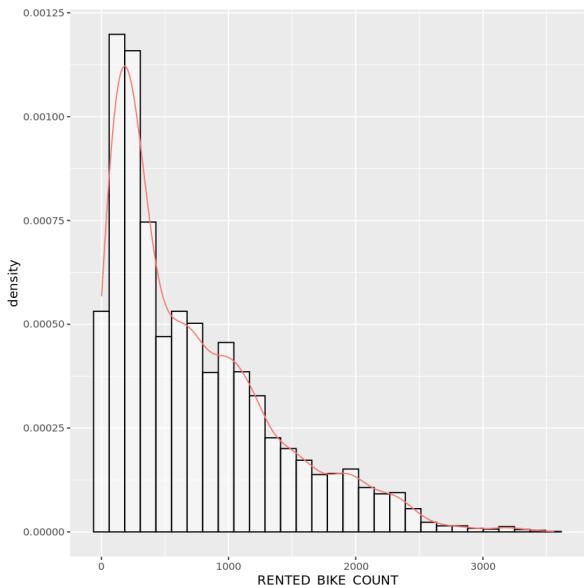
We can see the rented bike count are to low at the dawn and start to increase slowly during the early hours of the morning to reach to max at the evening in 6 or 7 then start decreasing again.



EDA WITH VISUALIZATION

Perform exploratory data analysis using Visualization with tidyverse & ggplot

Create a histogram overlaid with a kernel density curve



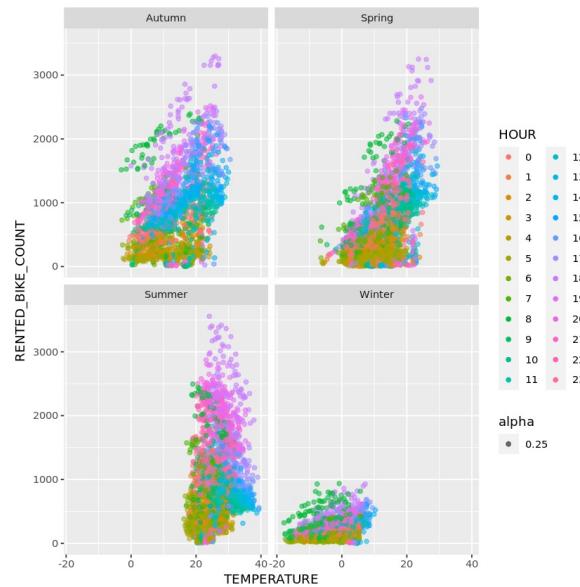
- Histogram Observation: Most times, few bikes are rented; the mode is about 250.
- Modes in Subgroups: Bumps at 700, 900, 1900, and 3200 bikes suggest hidden modes in subgroups.
- Rare Occasions: Occasionally, many more bikes are rented than usual.



EDA WITH VISUALIZATION

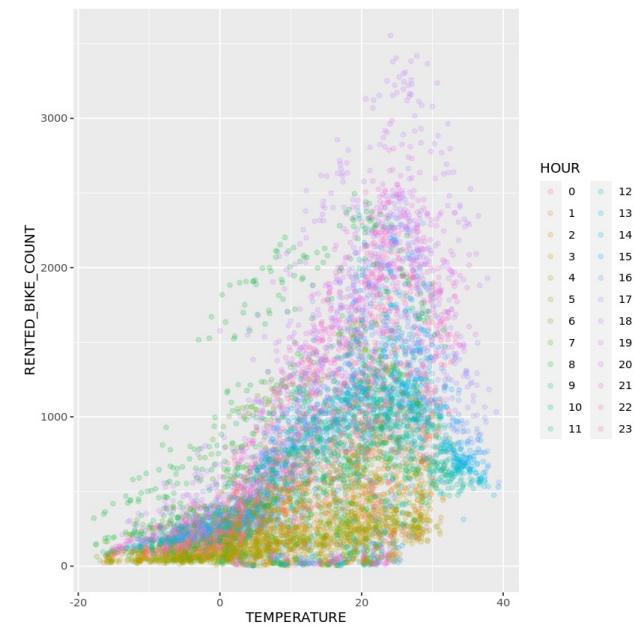
Perform exploratory data analysis using Visualization with tidyverse & ggplot

**Correlation between two variables
(RENTED_BIKE_COUNT and TEMPERATURE by SEASONS)**



Visually, strong correlations are evident as approximately linear patterns.
Autumn & Spring: Similar bike usage patterns with temperatures between 0-20°C; higher usage in warmer weather.
Summer: Consistent usage hours but reduced bike counts in hotter weather.
Winter: Significant drop in bike rentals, with a max of 1000 bikes; peak usage in early morning and evening (6-7 PM).

Create a scatter plot of RENTED_BIKE_COUNT vs TEMPERATURE by Hour as Color



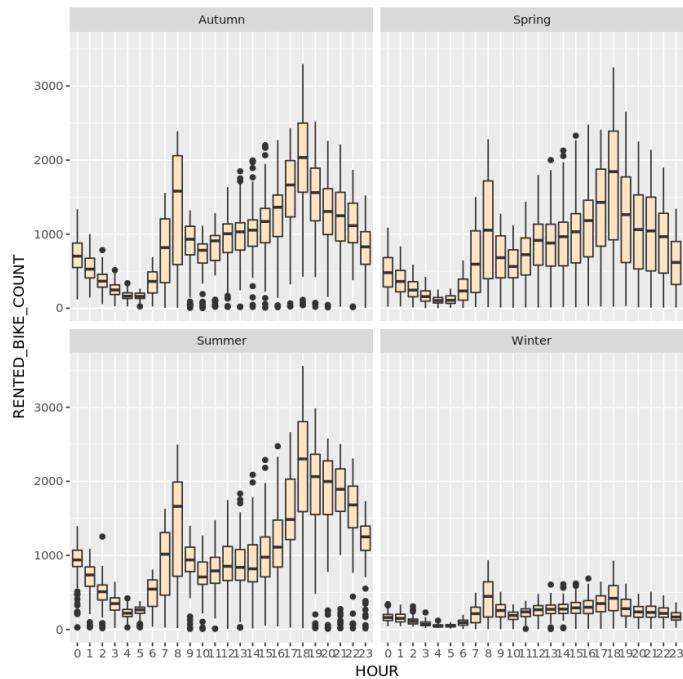
higher usage in warmer weather with peak hour in morning and evening 6-7PM.



EDA WITH VISUALIZATION

Perform exploratory data analysis using Visualization with tidyverse & ggplot

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



Seasonal Variations:

Bike rental counts vary by season but key features remain similar.

Peak Demand:

Peak demand times are consistent across all seasons, at 8 am and 6 pm.

Outliers in Data:

Many outliers in bike count data during different seasons.

Usage Patterns:

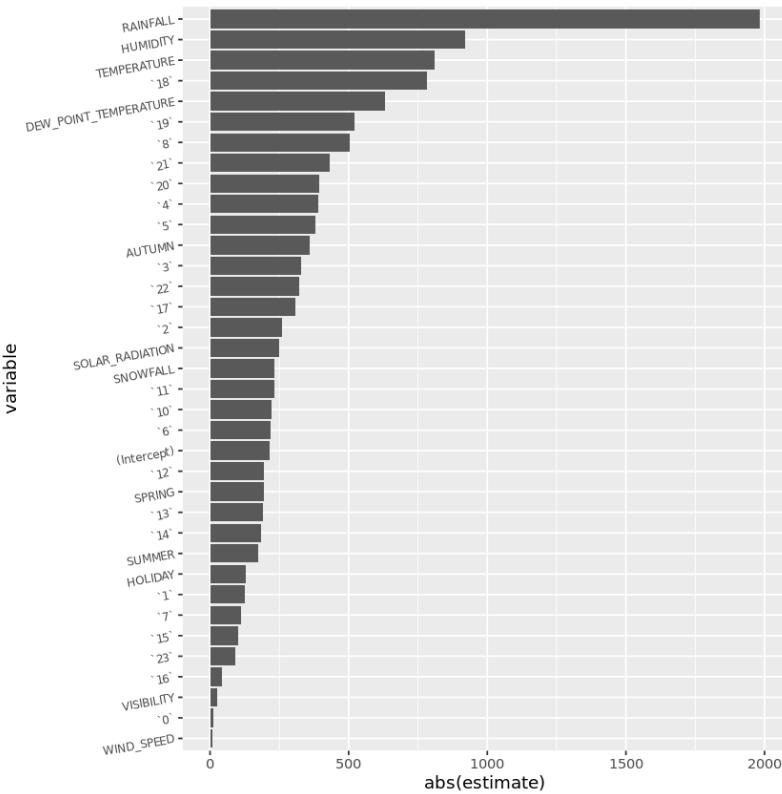
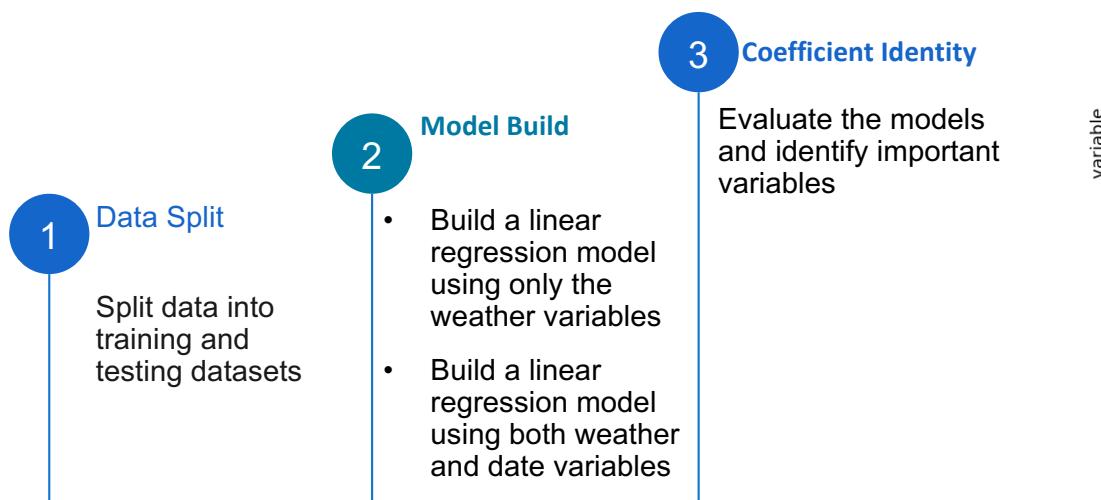
People generally use bikes at similar times in different seasons, with slight variations in counts.

Winter Drop:

Significant drop in bike rentals during Winter.



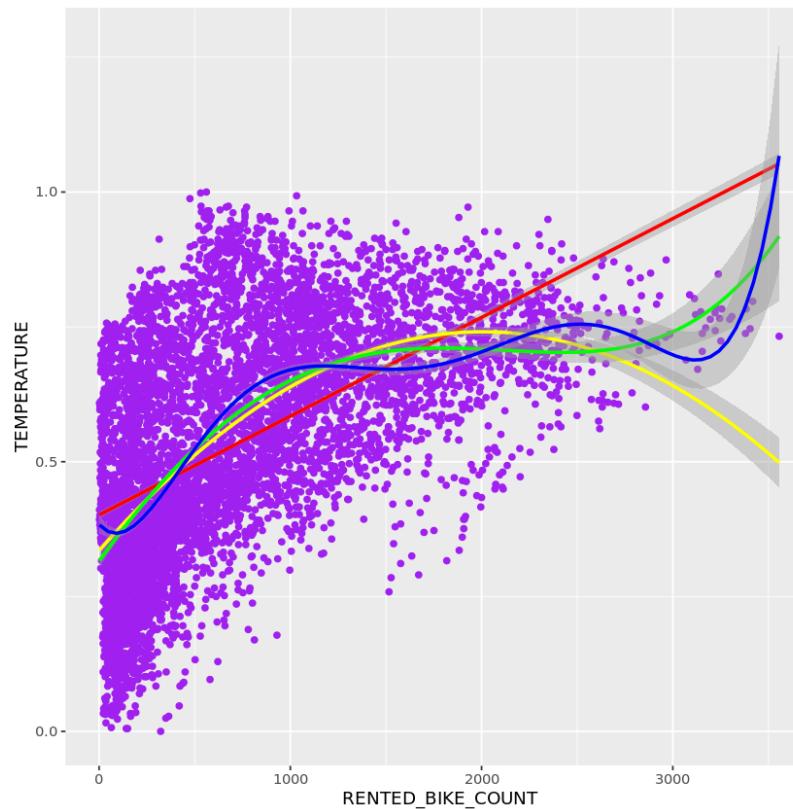
Predict Hourly Rented Bike Count using Basic Linear Regression Models





Refine the Baseline Regression Models

Correlation between RENTED_BIKE_COUNT and TEMPERATURE with the higher order polynomial fits



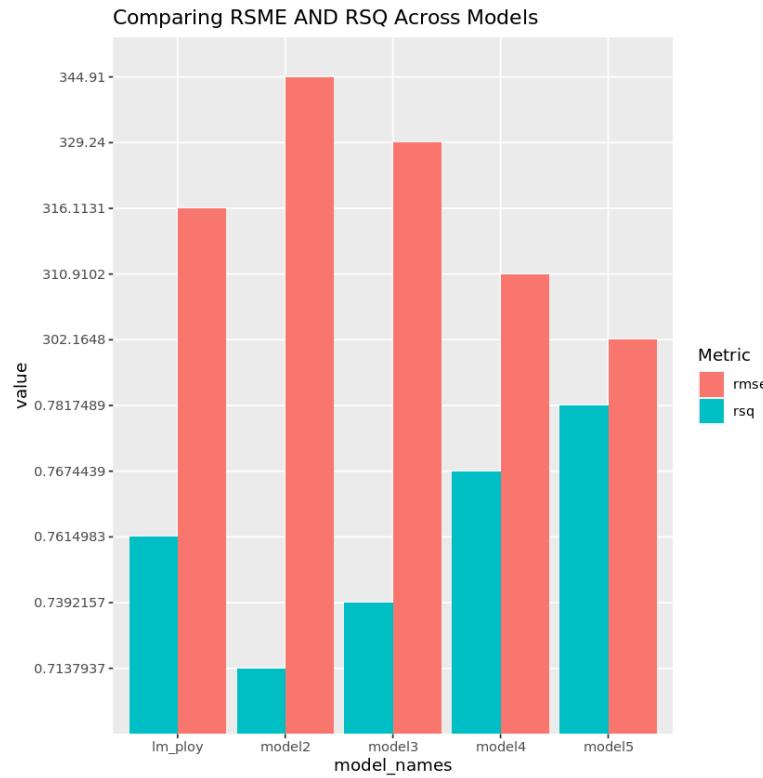
- 1 Improve Model
 - Add polynomial terms
 - Add interactions terms
 - Add regularizations terms

- 2 Best Model Selection
 - Experiment to search for improved models

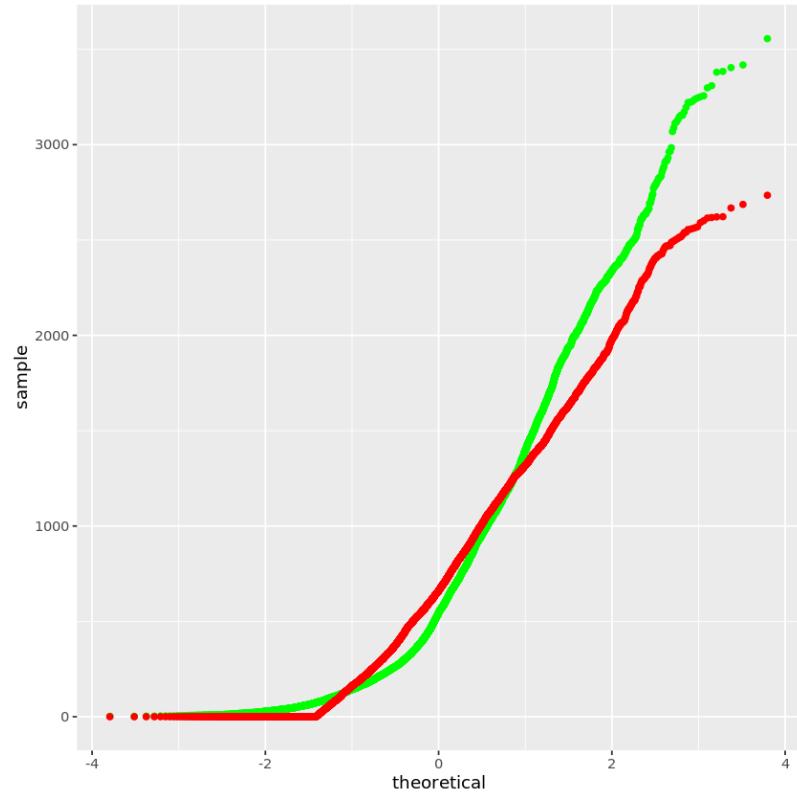


Refine the Baseline Regression Models

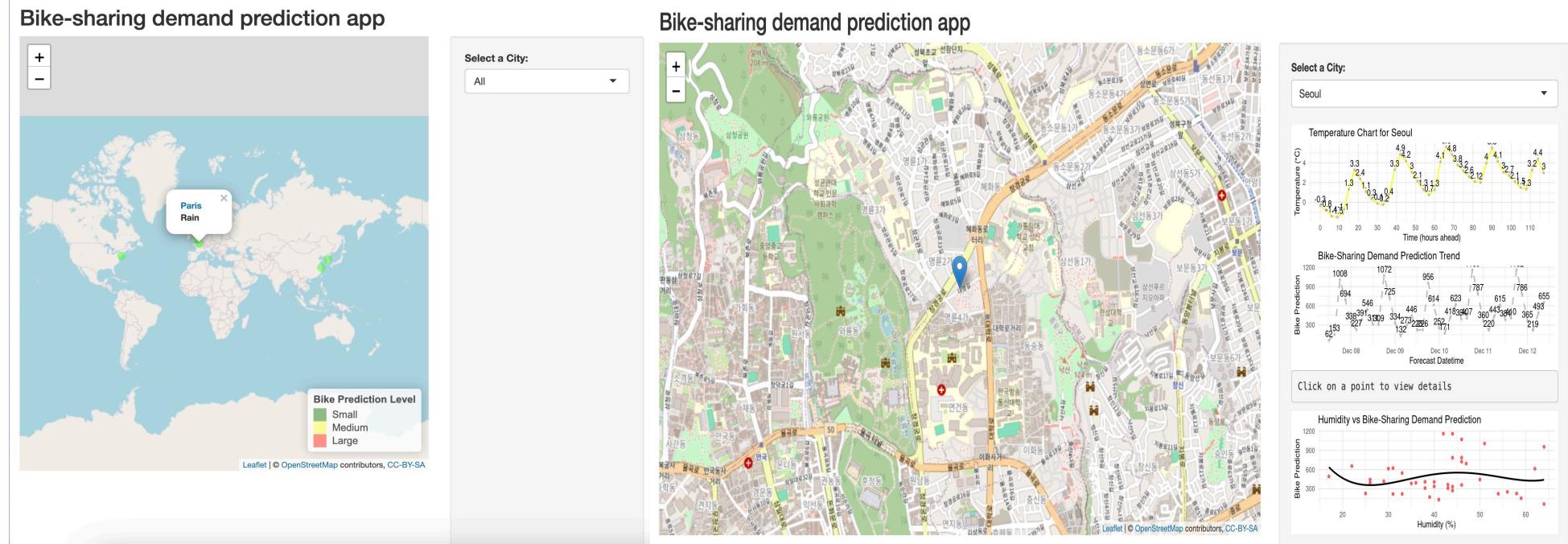
Model 5 Reported the best performed model in terms of rmse and rsq



Plotting the distribution difference between the predictions generated by your best model vs the true values on test dataset



R Shiny dashboard





Discussion

FINDINGS

Weekday vs. Weekend: Bike rental count is higher during weekdays than weekends.

Peak Hours: Rental bike counts peak at 8 AM and 6-7 PM, with demand gradually increasing from 5 AM to 8 AM, then dipping, and rising again until 6-7 PM.



IMPLICATIONS

Temperature & Wind: People prefer renting bikes at moderate to high temperatures and even with light winds, suggesting a need for comfortable weather conditions.

Seasonal Trends: Highest bike rentals in Autumn and Summer, lowest in Winter, indicating seasonal preferences.

Weather Conditions: Bike rentals are highest on clear days and lowest on snowy and rainy days, impacting rental decisions.



Conclusion

- These observations suggest that bike rentals are influenced by various factors **including time of day, weather conditions, hours and seasons.**
- Understanding these patterns can help bike rental companies optimize their services and better meet customer demand.
- **Bike rentals peak during morning and evening commutes.** Bike rental services should place more bikes at popular stations and increase redistribution frequency during these periods.
- They may consider **increasing the availability of bikes** during peak hours and seasons, and **adjusting prices based on weather conditions to attract more customers.** This ensures smoother commutes, enhances user satisfaction, and encourages continued bike usage.



Appendix

https://github.com/AliceInBRisbane/Seoul_Bike_Sharing_Demands_Analysis_with_R/blob/main/collection-data-by-apicall.ipynb

First import `httr` library. Run `install.packages("httr")` prior to loading the package only if you are running this lab locally on RStudio on your system.

```
:  
  
install.packages("httr") # for Http requests  
  
library(httr) # for Http requests  
  
Updating HTML index of packages in '.Library'  
Making 'packages.html' ... done
```

The API base URL to get current weather is <https://api.openweathermap.org/data/2.5/weather>

```
: # URL for Current Weather API  
current_weather_url <- 'https://api.openweathermap.org/data/2.5/weather'
```

Next, let's create a list to hold URL parameters for current weather API

```
: # need to be replaced by your real API key  
your_api_key <- "d858c9d0c145defdd6a7fe642bb3ef6"  
# Input 'q' is the city name  
# Input 'appid' is your API KEY,  
# Input 'units' are preferred units such as Metric or Imperial  
current_query <- list(q = "Seoul", appid = your_api_key, units="metric")
```

It contains very detailed weather data about the city of `Seoul`. Feel free to try other cities as well. We need to convert the named list to a data frame so that we can use data frame operations to process the data. Below is a simple example, which you may implement your own way to convert it to a data frame.

```
# Create some empty vectors to hold data temporarily  
weather <- c()  
visibility <- c()  
temp <- c()  
temp_min <- c()  
temp_max <- c()  
pressure <- c()  
humidity <- c()  
wind_speed <- c()  
wind_deg <- c()
```

Now assign the values in the `json_result` list into different vectors

```
# weather is also a list with one element, its main element indicates the weather status such as clear or rain  
weather <- c(weather, json_result$weather[[1]]$main)  
# Get Visibility  
visibility <- c(visibility, json_result$visibility)  
# Get current temperature  
temp <- c(temp, json_result$main$temp)  
# Get min temperature  
temp_min <- c(temp_min, json_result$main$temp_min)  
# Get max temperature  
temp_max <- c(temp_max, json_result$main$temp_max)  
# Get pressure  
pressure <- c(pressure, json_result$main$pressure)  
# Get humidity
```



Appendix

https://github.com/AliceInBRisbane/Seoul_Bike_Sharing_Demands_Analysis_with_R/blob/main/collectio n-data-by-webscraping.ipynb

TASK: Extract bike sharing systems HTML table from a Wiki page and convert it into a data frame

TODO: Get the root HTML node

```
url <- "https://en.wikipedia.org/wiki/List_of_bicycle-sharing_systems"
# Get the root HTML node by calling the `read_html()` method with URL
html_nodes(root_node, "table")

xml_nodeset(4)
[1] <table class="wikitable sortable sticky-header" style="background:#f8f9fa ...
[2] <table class="nowraplinks mw-collapsible autocollapse navbox-inner" style ...
[3] <table class="nowraplinks navbox-subgroup" style="border-spacing:0"><tbody ...
[4] <table class="nowraplinks navbox-subgroup" style="border-spacing:0"><tbody ...

Note that this HTML page at least contains three child <table> nodes under the root HTML node. So, you will need to use
html_nodes(root_node, "table") function to get all its child <table> nodes:

<html>
  <table>(table1)</table>
  <table>(table2)</table>
  <table>(table3)</table>
  ...
</html>
```

```
table_nodes <- html_nodes(root_node, "table")
```

```
# Convert the bike-sharing system table into a dataframe
bike_sharing_df <- html_table(table_nodes[[1]], fill = TRUE)
```

Summarize the bike sharing system data frame

```
# Summarize the dataframe
summary(bike_sharing_df)
```

Country	Country	City / Region	Name
Length:885	Length:885	Length:885	Length:885
Class :character	Class :character	Class :character	Class :character
Mode :character	Mode :character	Mode :character	Mode :character
System	Operator	Launched	Discontinued
Length:885	Length:885	Length:885	Length:885
Class :character	Class :character	Class :character	Class :character
Mode :character	Mode :character	Mode :character	Mode :character

Export the data frame as a csv file called `raw_bike_sharing_systems.csv`

```
# Export the dataframe into a csv file
write.csv(bike_sharing_df, file = "raw_bike_sharing_systems.csv", row.names = FALSE)
```

For more details about webscraping with `rvest`, please refer to the previous webscraping notebook here:



Appendix

https://github.com/AliceInBRisbane/Seoul_Bike_Sharing_Demands_Analysis_with_R/blob/main/lab-jupyter-data-wrangling-with-dplyr.ipynb

```
# Drop rows with `RENTED_BIKE_COUNT` column == NA
bike_sharing_df <- bike_sharing_df[!is.na(bike_sharing_df$RENTED_BIKE_COUNT), ]  
  
# Print the dataset dimension again after those rows are dropped
print(dim(bike_sharing_df))
```

```
[1] 8465   14
```

Now that you have handled missing values in the `RENTED_BIKE_COUNT` variable, let's continue processing missing values for the `TEMPERATURE` column.

Unlike the `RENTED_BIKE_COUNT` variable, `TEMPERATURE` is not a response variable. However, it is still an important predictor variable - as you could imagine, there may be a positive correlation between `TEMPERATURE` and `RENTED_BIKE_COUNT`. For example, in winter time with lower temperatures, people may not want to ride a bike, while in summer with nicer weather, they are more likely to rent a bike.

How do we handle missing values for `TEMPERATURE`? We could simply remove the rows but it's better to impute them because `TEMPERATURE` should be relatively easy and reliable to estimate statistically.

Let's first take a look at the missing values in the `TEMPERATURE` column.

```
bike_sharing_df %>%
  filter(is.na(TEMPERATURE))
```

```
library(dplyr)  
  
summer_temp_mean <- bike_sharing_df %>%
  filter(SEASONS == "Summer") %>%
  summarise(mean_temp = mean(TEMPERATURE, na.rm = TRUE)) %>%
  pull(mean_temp)  
  
# Impute missing values for TEMPERATURE column with summer average temperature
bike_sharing_df <- bike_sharing_df %>%
  mutate(TEMPERATURE = ifelse(is.na(TEMPERATURE) & SEASONS == "Summer", summer_temp_mean, TEMPERATURE))  
  
# Print the summary of the dataset again to make sure no missing values in all columns
missing_values_per_column <- colSums(is.na(bike_sharing_df))
print(missing_values_per_column)  
  
if (all(missing_values_per_column == 0)) {
  cat("No missing values in any column of the dataset.\n")
} else {
  cat("There are missing values in the dataset. Check the output above for details.\n")
}  
  
print(summary(bike_sharing_df))
```

Appendix

https://github.com/AliceInBRisbane/Seoul_Bike_Sharing_Demands_Analysis_with_R/blob/main/lab-jupyter-data-wrangling-with-regex.ipynb

```
# Select the four columns
sub_bike_sharing_df <- bike_sharing_df %>% select(COUNTRY, CITY, SYSTEM, BICYCLES)
```

Let's see the types of the selected columns

```
sub_bike_sharing_df %>%
  summarize_all(class) %>%
  gather(variable, class)
```

variable	class
<chr>	<chr>
tibble: 4 x 2	
COUNTRY	character
CITY	character
SYSTEM	character
BICYCLES	character

They are all interpreted as character columns, but we expect the `BICYCLES` column to be of numeric type. Let's see why it wasn't loaded as a numeric column - possibly some entries contain characters. Let's create a simple function called `find_character` to check that.

```
# grep searches a string for non-digital characters, and returns TRUE or FALSE
# if it finds any non-digital characters, then the bicycle column is not purely numeric
find_character <- function(strings) grep("[^0-9]", strings)
```

```
]# Define a 'reference link' character class,
#[A-z0-9] means at least one character
#[` and `] means the character is wrapped by [], such as for [12] or [abc]
ref_pattern <- "\[[A-z0-9]+\]"
find_reference_pattern <- function(strings) grepl(ref_pattern, strings)
```

```
]# Check whether the COUNTRY column has any reference links
sub_bike_sharing_df %>%
  select(COUNTRY) %>%
  filter(find_reference_pattern(COUNTRY)) %>%
  slice(0:10)
```

COUNTRY
<chr>
A spec_tbl_df: 0 x 1

Ok, looks like the `COUNTRY` column is clean. Let's check the `CITY` column.

```
]# Check whether the CITY column has any reference links
sub_bike_sharing_df %>%
  select(CITY) %>%
  filter(find_reference_pattern(CITY)) %>%
  slice(0:10)
```

CITY
<chr>
A spec_tbl_df: 10 x 1



Appendix



https://github.com/AliceInBRisbane/Seoul_Bike_Sharing_Demands_Analysis_with_R/blob/main/Lab-SQL-EDA_SQLite.ipynb

```
# provide your solution here
dbGetQuery(conn, "SELECT count(HOUR) as Numer_of_hours FROM seoul_bike_sharing
WHERE RENTEDBIKECOUNT > 0")

Numer_of_hours
<int>
data.frame: 1 x 1
1 8465
```

Task 3 - Weather Outlook

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

Recall that the records in the CITIES_WEATHER_FORECAST dataset are 3 hours apart, so we just need the first record from the query.

Solution 3

```
# provide your solution here
dbGetQuery(conn, "SELECT * FROM CITIES_WEATHER_FORECAST
WHERE CITY = 'Seoul'
Limit 1")

CITY WEATHER VISIBILITY TEMP TEMP_MIN TEMP_MAX PRESSURE HUMIDITY WIND_SPEED WIND_DEG SEASON
<chr> <chr> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <chr>
```

Solution 6

```
# provide your solution here
dbGetQuery(conn, "SELECT DATE, HOUR, RENTEDBIKECOUNT as Maximum_COUNT FROM seoul_bike_sharing
WHERE RENTEDBIKECOUNT = (SELECT MAX(RENTEDBIKECOUNT) FROM seoul_bike_sharing)")

DATE HOUR Maximum_COUNT
<chr> <dbl> <dbl>
data.frame: 1 x 3
1 19/06/2018 18 3556
```

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(RENTEDBIKECOUNT), AVG(TEMPERATURE) FROM seoul_bike_sharing GROUP BY
SEASONS HOUR AVG(RENTEDBIKECOUNT) AVG(TEMPERATURE)
<chr> <dbl> <dbl> <dbl>
```



SCOURCE



Weather & Bike-Sharing



OpenWeather APIs Calls -Current & Prediction of Weather
https://home.openweathermap.org/users/sign_up

Web scrape a Global Bike-Sharing Systems Wiki Page
https://en.wikipedia.org/wiki/List_of_bicycle-sharing_systems



THANK YOU

2024/12/8