



Assignment: Exploratory Data Analysis with tidyverse and ggplot2

Estimated time needed: 60 minutes

Introduction and Objectives

In this Lab, you will use an R notebook to perform exploratory data analysis using tidyverse and the ggplot2 R packages.

You will start by doing some minor data preparation on the SEOUL BIKE SHARING dataset. Then you will generate and explore some statistics from the resulting dataframe and make some observations. Finally, you will generate some informative plots using the ggplot2 library.

Your primary objective is to gather insights from your exploratory analysis. These findings will be part of your story that you will create your final capstone presentation.

Visualization is a very powerful tool for better understanding your data and finding patterns that may exist in it. You can use scatterplots, for example, to display how well two features are correlated with, or similar to each other. When data are highly correlated, it means they vary in similar ways, and so their graphs will look similar (once scaled to a common scale). We can say one variable 'explains' the variation in the other, and that they are 'covariates'. There could be a causal relationship between covariates, meaning that changing one variable has the effect of changing the other, but this need not be the case. Perhaps there is another factor which causes both covariates to respond to variations, or the similarity could be a random coincidence. Either way, the behaviour of one variable can be used to predict the behaviour of the other. The key practical difference is that in the causal case, if we can influence the first variable then we can have a corresponding causal influence on the second. Like turning a light switch on or off to control the light in the room, the state of the switch controls the state of the light bulb. This is an important topic for data science that is beyond our current scope, but we invite you to delve deeper into the subject as you progress in your career.

Other ways visualization can inform your analysis is with spotting outliers and anomalous behaviour in your features. Boxplots are informative in these regards. You can also gain insights about any clear trends and anomalies that may be present in a variable, simply by plotting it directly. For instance, time series and spatial data are particularly interesting

kinds of variables. Outliers can easily consume the range of your plot though, making your data look like a featureless flat line in comparison to these points. So some cleaning, namely outlier removal, may be required to get a clearer picture.

A word of caution: be skeptical about any patterns you find, especially in smaller datasets. In very simple terms, it is true that any two points randomly placed in space always define a unique line; but add a third random point, and it is very unlikely that it will land on that same line. This suggests one of the main advantages of 'big data' - any patterns that emerge in very large datasets are far more likely to persist on unseen data than those found in small datasets.

All right, let's move along and get started with our exploratory analysis!

For reference, we include the Attribute Information for the `seoul_bike_sharing` dataset:

- DATE - format: "2017-12-01"
- RENTED_BIKE_COUNT - Count of bikes rented at each hour
- HOUR - Hour of the day
- TEMPERATURE - Celsius
- HUMIDITY - %
- Windspeed - m/s
- VISIBILITY - 10m
- DEW_POINT_TEMPERATURE - Celsius
- SOLAR_RADIATION - MJ/m2
- RAINFALL - mm
- SNOWFALL - cm
- SEASONS - "Autumn", "Spring", ..
- HOLIDAY - "Holiday", "No holiday"
- FUNCTIONING_DAY - "Yes", "No"

Load the `seoul_bike_sharing` data into a dataframe

Use the following URL to load your dataset.

The dataset is already clean, but you will still need to pay careful attention to data types, especially dates, which you may need to coerce. Also, ensure any categorical variables get typed as factors.

```
seoul_bike_sharing <- "https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBMDeveloperSkillsNetwork-RP0321EN-SkillsNetwork/labs/datasets/seoul_bike_sharing.csv"
```

Task 1 - Load the dataset

Ensure you read `DATE` as type `character` .

Solution 1

```
In [ ]: # provide your solution here
```

Task 2 - Recast `DATE` as a date

Use the format of the data, namely "%d/%m/%Y".

Solution 2

```
In [ ]: # provide your solution here
```

Task 3 - Cast `HOURS` as a categorical variable

Also, coerce its levels to be an ordered sequence. This will ensure your visualizations correctly utilize `HOURS` as a discrete variable with the expected ordering.

Solution 3

```
In [ ]: # provide your solution here
```

Check the structure of the dataframe

```
In [ ]: str(seoul_bike_sharing)
```

Finally, ensure there are no missing values

```
In [ ]: sum(is.na(seoul_bike_sharing))
```

Descriptive Statistics

Now you are all set to take a look at some high level statistics of the `seoul_bike_sharing` dataset.

Task 4 - Dataset Summary

Use the base R `sumamry()` function to describe the `seoul_bike_sharing` dataset.

Solution 4

```
In [ ]: # provide your solution here
```

Some Basic Observations:

- We can see from `DATE` that we have exactly a full year of data.
- No records have zero bike counts.

- Spring and Winter have the same count of records, while autumn has the least and Summer has the most.
- Temperature has a large range, so we might expect it to explain at least some of the variation in bike rentals.
- Precipitation seems to be quite rare, only happening in the fourth quartiles for both `RAINFALL` and `SNOWFALL`.
- The average `WINDSPEED` is very light at only 1.7 m/s, and even the maximum is only a moderate breeze (Google 'Beaufort Wind Scale' to find the different wind descriptions)

By now, you might agree that Exploratory Data Analysis can create more questions than answers. That's okay - you'll have a much deeper understanding and appreciation for your data as a result!

Task 5 - Based on the above stats, calculate how many Holidays there are.

Solution 5:

```
In [ ]: # provide your solution here
```

Task 6 - Calculate the percentage of records that fall on a holiday.

Solution 6

```
In [ ]: # provide your solution here
```

Task 7 - Given there is exactly a full year of data, determine how many records we expect to have.

Solution 7

```
In [ ]: # provide your solution here
```

Task 8 - Given the observations for the 'FUNCTIONING_DAY' how many records must there be?

Solution 8

```
In [ ]: # provide your solution here
```

Drilling Down

Let's calculate some seasonally aggregated measures to help build some more context.

Task 9 - Load the dplyr package, group the data by SEASONS, and use the summarize() function to calculate the seasonal total rainfall and snowfall.

Solution 9

```
In [ ]: # provide your solution here
```

Wow, that seems like a lot of snow.

Now that you have some ideas about what sorts of questions can be answered through descriptive statistics, let's start visualizing the data.

Data Visualization

Let's take a closer look at our main variable of interest, namely, RENTED_BIKE_COUNT . Think of this variable as the key *measure* or *dependent variable* in your analysis.

Indeed, it is a measured quantity, and we expect it to depend on factors such as the expected weather.

Evidently, if the immediate or forecasted weather is harsh or unpleasant, many people could choose to use alternate transit or simply wait for better weather rather than rent a bike.

On the other hand, many people may be inspired to ride under pleasant expected weather conditions.

The weather is largely influenced by the time of day and the seasons, so these are also factors.

The time of day, the day of week, and Holidays all matter because they control commuting schedules.

Finer granularity data such as a unique ID for each bike and/or rider, when and where each bike was rented, or even finer - a history of when and where each bike was used or idle - would be interesting as well.

Load the ggplot2 package so we can generate some data visualizations.

```
In [ ]: # provide your solution here
```

Our variable of interest is a time series, so why not start by taking a look at it in its natural form?

Task 10 - Create a scatter plot of RENTED_BIKE_COUNT vs DATE .

Tune the opacity using the alpha parameter such that the points don't obscure each other too much.

Solution 10

```
In [ ]: # provide your solution here
```

Ungraded Task: We can see some patterns emerging here.

Describe them and keep your findings for your presentation in the final project.

Solution

provide your solution here

Using colour

Let's see if we can enhance some of these features by incorporating colour. Given our observations so far, `HOURS` is a great candidate for this task.

Task 11 - Create the same plot of the `RENTED_BIKE_COUNT` time series, but now add `HOURS` as the colour.

Solution 11

```
In [ ]: # provide your solution here
```

Ungraded Task: The trends are much more clear now.

Describe them and keep your findings for your presentation in the final project.

Solution

provide your solution here

Distributions

Task 12 - Create a histogram overlaid with a kernel density curve

Normalize the histogram so the y axis represents 'density'. This can be done by setting `y=..density..` in the aesthetics of the histogram.

► [Click here for a hint](#)

► [Click here for another hint](#)

Solution 12

```
In [ ]: # provide your solution here
```

Ungraded Task: Describe the main features you see in your plot.

Consider what its shape tells you, and keep your findings for your presentation in the final project.

► [Click here for a solution](#)

Correlation between two variables (scatter plot)

Task 13 - Use a scatter plot to visualize the correlation between `RENTED_BIKE_COUNT` and `TEMPERATURE` by `SEASONS`.

Start with `RENTED_BIKE_COUNT` vs. `TEMPERATURE`, then generate four plots corresponding to the `SEASONS` by adding a `facet_wrap()` layer. Also, make use of colour and opacity to emphasize any patterns that emerge. Use `HOURL` as the color.

Solution 13

```
In [ ]: # provide your solution here
```

Ungraded Task: Describe the patterns you see.

What do these patterns imply about the relationships between these variables? Keep your findings for your presentation in the final project.

► [Click here for a solution](#)

Comparing this plot to the same plot below, but without grouping by `SEASONS`, shows how important seasonality is in explaining bike rental counts.

```
In [ ]: ggplot(seoul_bike_sharing) +  
  geom_point(aes(x=TEMPERATURE, y=RENTED_BIKE_COUNT, colour=HOURL), alpha=1/5)
```

Outliers (boxplot)

Task 14 - Create a display of four boxplots of `RENTED_BIKE_COUNT` vs. `HOURL` grouped by `SEASONS`.

Use `facet_wrap` to generate four plots corresponding to the seasons.

Solution 14

```
In [ ]: # provide your solution here
```

Ungraded Task: Compare and contrast the key features of these boxplots between seasons.

At this point, a story should be taking shape. Again, keep your findings for your presentation in the final project.

► [Click here for a solution](#)

Task 15 - Group the data by `DATE`, and use the `summarize()` function to calculate the daily total rainfall and snowfall.

Also, go ahead and plot the results if you wish.

Solution 15

```
In [ ]: # provide your solution here
```

Task 16 - Determine how many days had snowfall.

Solution 16

```
In [ ]: # provide your solution here
```

There are many more visualizations we could have chosen to cover here, but the important thing was that you deepen your understanding of the dataset.

I hope we succeeded in that endeavour!

(Keep going, you are getting closer to the finish line with each step you take. :-))

Further Reading

[1] 'Causal Model' (2021) *Wikipedia*. Available at ["https://en.wikipedia.org/wiki/Causal_model"](https://en.wikipedia.org/wiki/Causal_model) (Accessed: 22 April 2021).

Author(s)

Jeff Grossman

Contributor(s)

Yan Luo, Rav Ahuja

© IBM Corporation 2021. All rights reserved.

In []: