



Predict Hourly Rented Bike Count using Basic Linear Regression Models

Estimated time needed: **90** minutes

Lab Overview:

Now that you have performed exploratory analysis on the bike sharing demand dataset and obtained some insights on the attributes, it's time to build predictive models to predict the hourly rented bike count using related weather and date information.

In this lab, you will be asked to use `tidymodels` to build some baseline linear regression models:

- **TASK: Split data into training and testing datasets**
- **TASK: Build a linear regression model using only the weather variables**
- **TASK: Build a linear regression model using both weather and date variables**
- **TASK: Evaluate the models and identify important variables**

Let's start!

The tidyverse and the tidymodels packages can be used to produce high quality statistical and machine learning models. The Tidyverse library is a useful tool that provides various tools for data visualization, data manipulation, and read various datasets into a data frame; our Jupyter notebook platforms have a built-in Tidymodels, Tidyverse and rlang libraries, so we do not need to install these packages prior to loading library. However, if you decide to run this lab on your RStudio Desktop locally on your machine, you can remove the commented lines of code to install these packages before loading.

```
In [ ]: # It may take several minutes to install those libraries in Watson Studio
# install.packages("rlang")
# install.packages("tidymodels")
```

```
In [ ]: library("tidymodels")
library("tidyverse")
library("stringr")
```

The `seoul_bike_sharing_converted_normalized.csv` will be our main dataset which has following variables:

The response variable:

- `RENTED_BIKE_COUNT` - Count of bikes rented at each hour

Weather predictor variables:

- `TEMPERATURE` - Temperature in Celsius
- `HUMIDITY` - Unit is %
- `WIND_SPEED` - Unit is m/s
- `VISIBILITY` - Multiplied by 10m
- `DEW_POINT_TEMPERATURE` - The temperature to which the air would have to cool down in order to reach saturation, unit is Celsius
- `SOLAR_RADIATION` - MJ/m2
- `RAINFALL` - mm
- `SNOWFALL` - cm

Date/time predictor variables:

- `DATE` - Year-month-day
- `HOUR` - Hour of the day
- `FUNCTIONAL_DAY` - NoFunc(Non Functional Hours), Fun(Functional hours)
- `HOLIDAY` - Holiday/No holiday
- `SEASONS` - Winter, Spring, Summer, Autumn

Let's read the dataset as a dataframe first:

```
In [ ]: # Dataset URL
dataset_url <- "https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/bike-sharing-dataset.csv"
bike_sharing_df <- read_csv(dataset_url)
spec(bike_sharing_df)
```

We won't be using the `DATE` column, because 'as is', it basically acts like an data entry index. (However, given more time, we could use the `DATE` column to create a 'day of week' or 'isWeekend' column, which we might expect has an affect on preferred bike rental times.) We also do not need the `FUNCTIONAL_DAY` column because it only has one distinct value remaining (`YES`) after missing value processing.

```
In [ ]: bike_sharing_df <- bike_sharing_df %>%
  select(-DATE, -FUNCTIONING_DAY)
```

TASK: Split training and testing data

First, we need to split the full dataset into training and testing datasets.

The training dataset will be used for fitting regression models, and the testing dataset will be used to evaluate the trained models.

TODO: Use the `initial_split()`, `training()`, and `testing()` functions to generate a training dataset consisting of 75% of the original dataset, and a testing

dataset using the remaining 25%.

```
In [ ]: # Use the `initial_split()`, `training()`, and `testing()` functions to split the
# With seed 1234
set.seed(1234)
# prop = 3/4
# train_data
# test_data
```

TASK: Build a linear regression model using weather variables only

As you could imagine, weather conditions may affect people's bike renting decisions. For example, on a cold and rainy day, you may choose alternate transportation such as a bus or taxi. While on a nice sunny day, you may want to rent a bike for a short-distance travel.

Thus, can we predict a city's bike-sharing demand based on its local weather information? Let's try to build a regression model to do that.

TODO: Build a linear regression model called `lm_model_weather` using the following variables:

- `TEMPERATURE` - Temperature in Celsius
- `HUMIDITY` - Unit is %
- `WIND_SPEED` - Unit is m/s
- `VISIBILITY` - Multiplied by 10m
- `DEW_POINT_TEMPERATURE` - The temperature to which the air would have to cool down in order to reach saturation, unit is Celsius
- `SOLAR_RADIATION` - MJ/m²
- `RAINFALL` - mm
- `SNOWFALL` - cm

Define a linear regression model specification.

```
In [ ]: # Use `linear_reg()` with engine `lm` and mode `regression`
```

Fit a model with the response variable `RENTED_BIKE_COUNT` and predictor variables `TEMPERATURE + HUMIDITY + WIND_SPEED + VISIBILITY + DEW_POINT_TEMPERATURE + SOLAR_RADIATION + RAINFALL + SNOWFALL`

```
In [ ]: # Fit the model called `lm_model_weather`
# RENTED_BIKE_COUNT ~ TEMPERATURE + HUMIDITY + WIND_SPEED + VISIBILITY + DEW_POI
```

Print the fit summary for the `lm_model_weather` model.

```
In [ ]: # print(lm_model_weather$fit)
```

You should see the model details such as formula, residuals, and coefficients.

TASK: Build a linear regression model using all variables

In addition to weather, there could be other factors that may affect bike rental demand, such as the time of a day or if today is a holiday or not.

Next, let's build a linear regression model using all variables (weather + date/time) in this task.

TODO: Build a linear regression model called `lm_model_all` using all variables `RENTED_BIKE_COUNT ~ .`.

```
In [ ]: # Fit the model called `lm_model_all`  
# `RENTED_BIKE_COUNT ~ .` means use all other variables except for the response
```

Print the fit summary for `lm_model_all`.

```
In [ ]: # summary(lm_model_all$fit)
```

Now you have built two basic linear regression models with different predictor variables, let's evaluate which model has better performance,

TASK: Model evaluation and identification of important variables

Now that you have built two regression models, `lm_model_weather` and `lm_model_all`, with different predictor variables, you need to compare their performance to see which one is better.

In this project, you will be asked to use very important metrics that are often used in Statistics to determine the performance of a model:

1. R^2 / R-squared
2. Root Mean Squared Error (RMSE)

R-squared

R squared, also known as the coefficient of determination, is a measure to indicate how close the data is to the fitted regression line. The value of R-squared is the percentage of variation of the response variable (y) that is explained by a linear model.

Root Mean Squared Error (RMSE) $RMSE = \sqrt{MSE}$

As you know, the Mean Squared Error measures the average of the squares of errors, where 'error' is the difference between the actual value (y) and the estimated value

(\hat{y}). Another metric that is related to MSE is **Root Mean Squared Error (RMSE)** and is simply the square root of MSE.

We first need to test the `lm_model_weather` and `lm_model_all` models against the test dataset `test_data`, and generate `RENTED_BIKE_COUNT` prediction results.

TODO: Make predictions on the testing dataset using both `lm_model_weather` and `lm_model_all` models

```
In [ ]: # Use predict() function to generate test results for `lm_model_weather` and `lm
# and generate two test results dataframe with a truth column:

# test_results_weather for lm_model_weather model

# test_results_all for lm_model_all
```

NOTE: if you happen to see a warning like : `prediction from a rank-deficient fit may be misleading`, it may be caused by collinearity in the predictor variables.

Collinearity means that one predictor variable can be predicted from other predictor variables to some degree. For example, `RAINFALL` could be predicted by `HUMIDITY`.

But don't worry, you will address `glmnet` models (Lasso and Elastic-Net Regularized Generalized Linear Models) instead of regular `regression` models to solve this issue and further improve the model performance.

Next, let's calculate and print the R-squared and RMSE for the two test results

TODO: Use `rsq()` and `rmse()` functions to calculate R-squared and RMSE metrics for the two test results

```
In [ ]: # rsq_weather <- rsq(...)
# rsq_all <- rsq(...)

# rmse_weather <- rmse(...)
# rmse_all <- rmse(...)
```

From these tables, you should find that the test results from `lm_model_all` are much better. It means that using both weather and datetime variables in the model generates better prediction results.

Since `lm_model_all` has many predictor variables, let's check which predictor variables have larger coefficients. Variables with larger coefficients in the model means they attribute more in the prediction of `RENTED_BIKE_COUNT`. In addition, since all predictor variables are normalized to the same scale, 0 to 1, we thus can compare their coefficients directly.

You could try building another regression model using the non-normalized `seoul_bike_sharing_converted.csv` dataset, and you would find that the coefficients are much different.

First let's print all coefficients:

```
In [ ]: lm_model_all$fit$coefficients
```

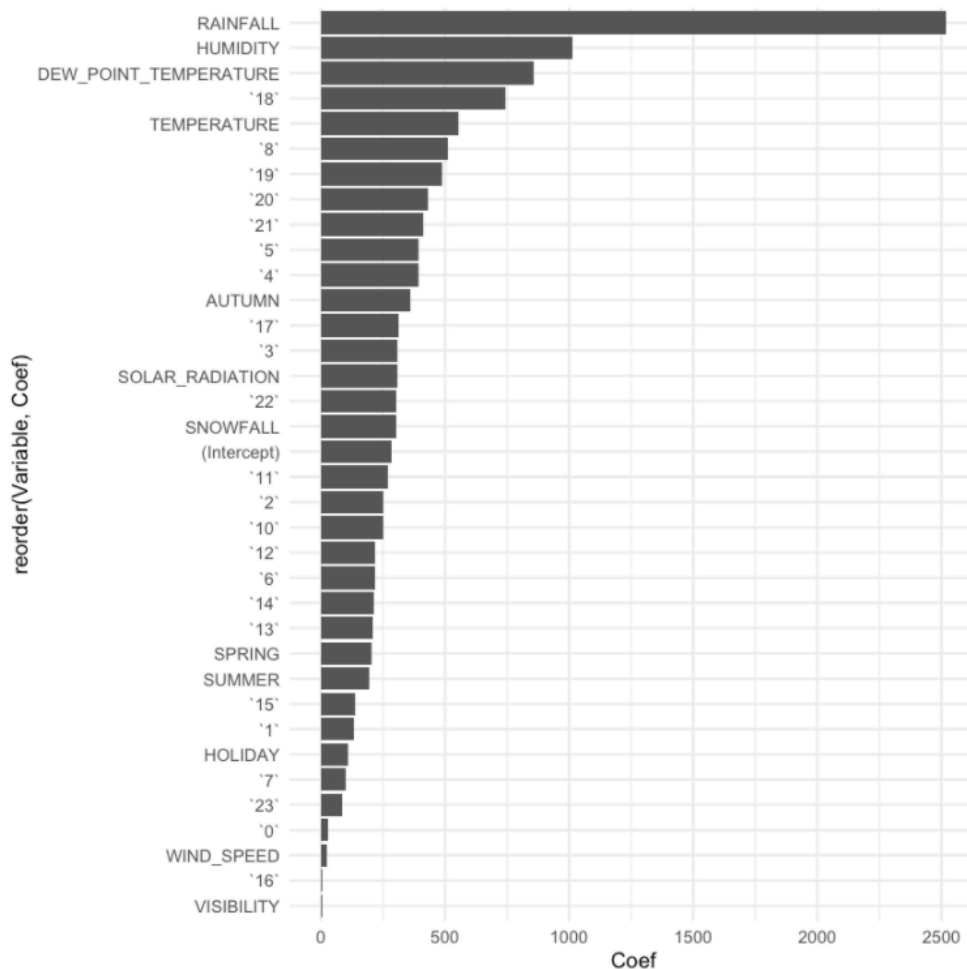
hmm, it's not very clear to compare the coefficients from a long and unsorted list. Next, you need to sort and visualize them using a bar chart

TODO: Sort the coefficient list in descending order and visualize the result using `ggplot` and `geom_bar`

```
In [ ]: # Sort coefficient list
```

```
In [ ]: # Visualize the list using ggplot and geom_bar
```

You should see a sorted coefficient bar chart like the following example:



Mark down these 'top-ranked variables by coefficient', which will be used for model refinements in the next labs.

Note that here the main reason we use absolute value is to easily identify important variables, i.e. variables with large magnitudes, no matter it's negative or positive. If we want to interpret the model then it's better to separate the positive and negative coefficients.

Next Steps

Great! Now you have built a baseline linear regression model to predict hourly bike rent count, with reasonably good performance. In the next lab, you will be refining the baseline model to improve its performance.

Authors

[Yan Luo](#)

Other Contributors

Jeff Grossman

Change Log

Date (YYYY-MM-DD)	Version	Changed By	Change Description
2021-04-08	1.0	Yan	Initial version created

© IBM Corporation 2021. All rights reserved.