# SEOUL BIKE SHARING

Analysis on demand estimation

Brocco Mattia – 2044714

Magliani Jacopo – 2040912

# OBJECTIVES

- Discover patterns in the **bike sharing** system of Seoul through data exploration and modelling

- Make use of the available features to derive a **regression model** to estimate the count of rented bikes

- **Produce insights** on the shared public transport for further development

**8760**

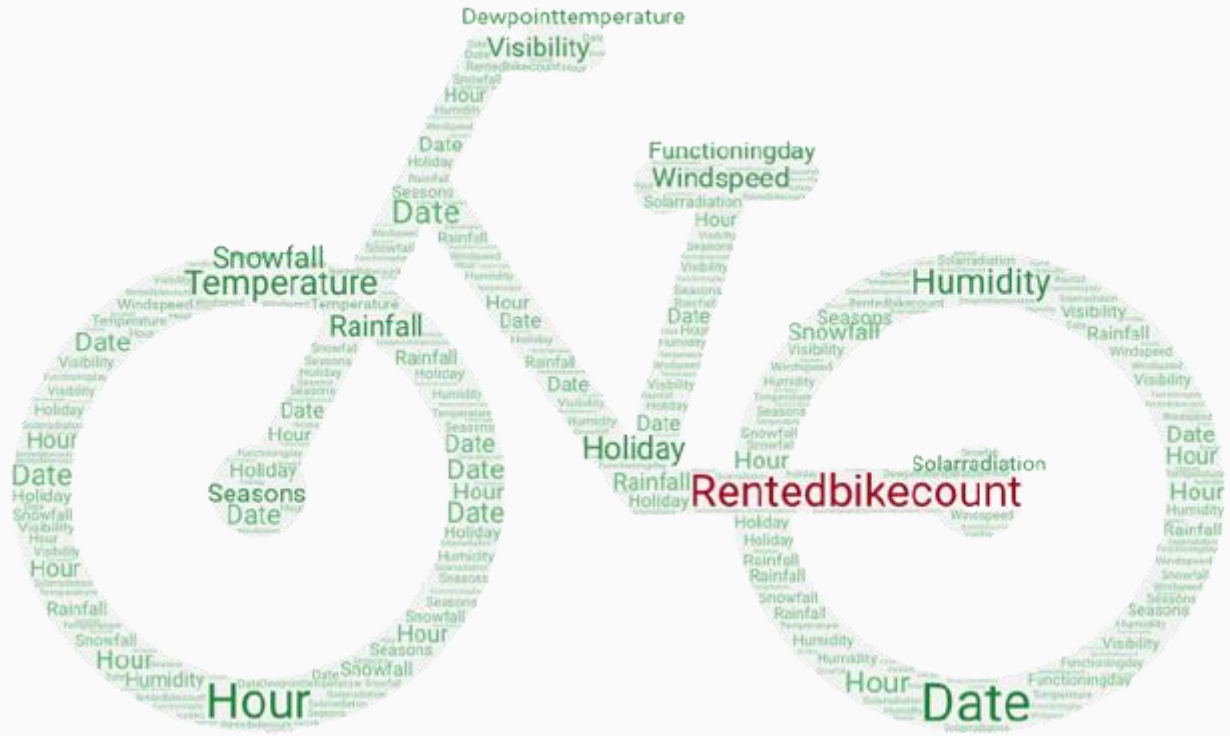the number of observations

**14**
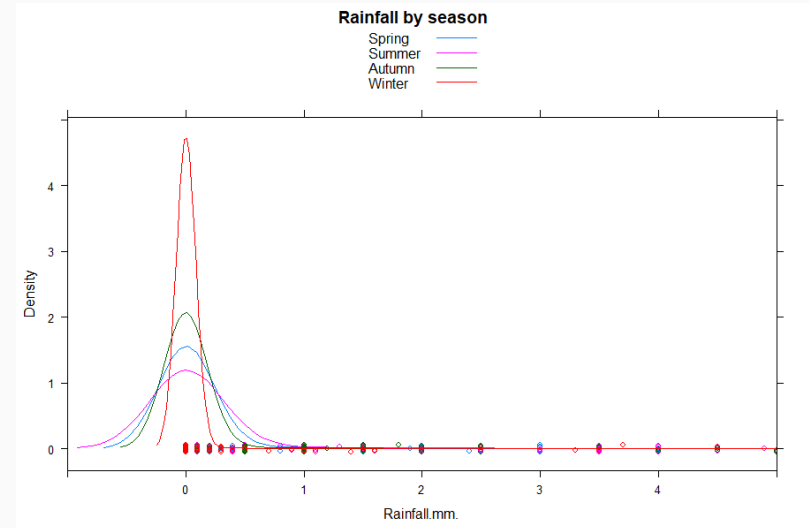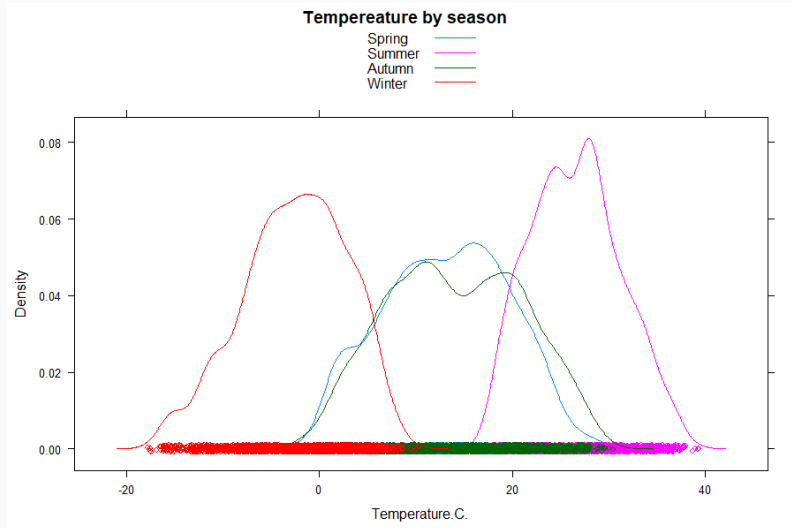
the number of variables

| 01/12 | ... | 30/11 |
|-------|-----|-------|
| 2017  |     | 2018  |

the time interval
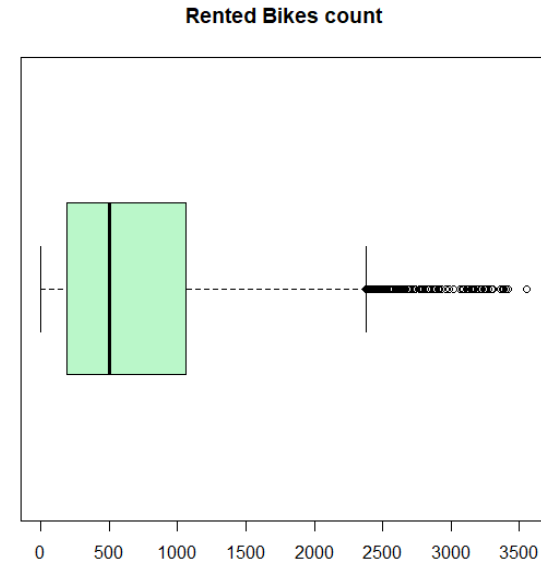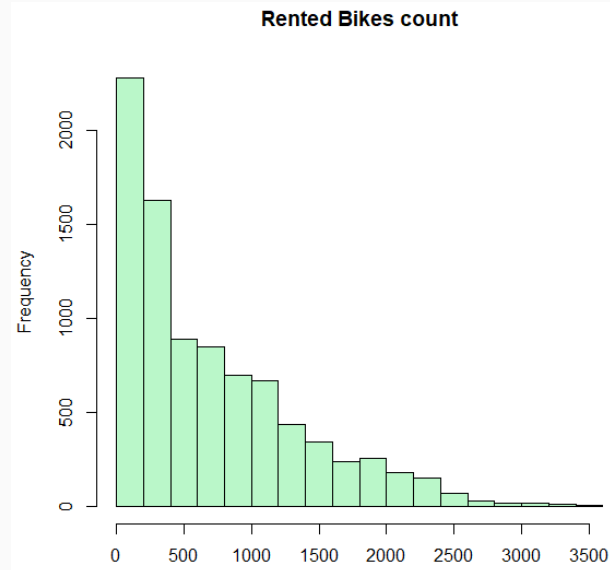
1

Atmospheric variables show trends not too far from what we expect from a continental climate throughout the different seasons (e.g., higher temperatures in Summer, lower rainfall in Winter).
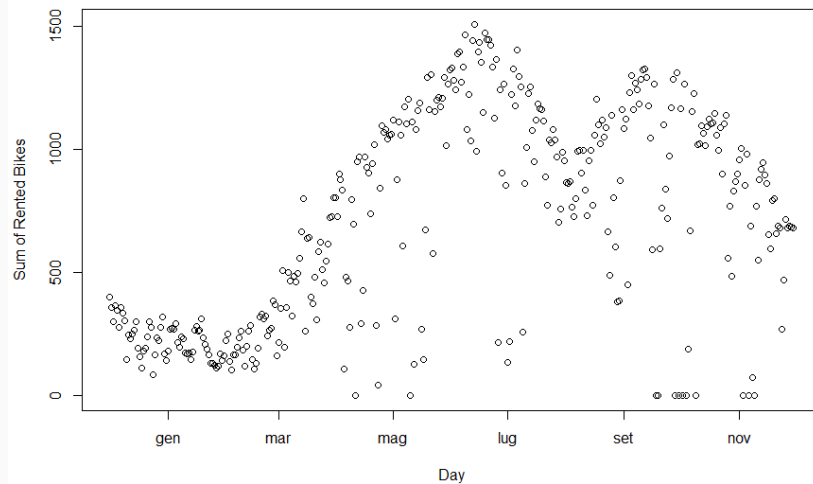
The target variable shows a left-skewed distribution, far from a normal distribution.



Rented Bikes count
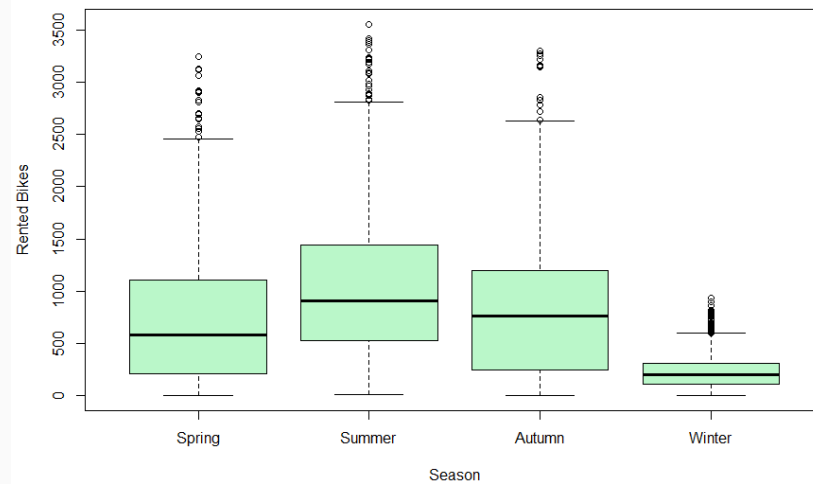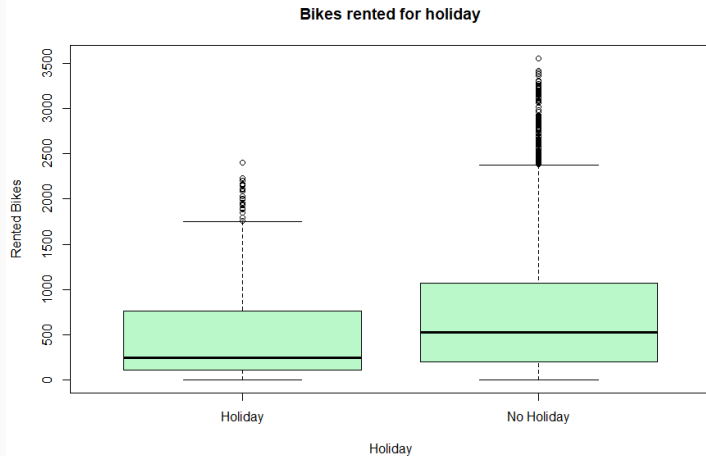
Bikes rented for each day
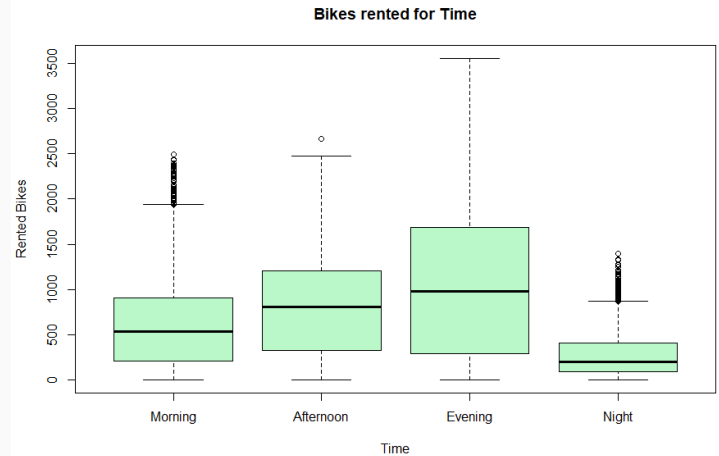


Bikes rented for season

- 18 holiday days

- Non-Holiday show an **average of rented bikes greater** than that of Holidays days. This may be also due to the small number of Holiday.

- The influence of a **specific holiday** can't be understood as the time interval considered is of 1 year.

Consideration on the target variable with respect to an additional **factor variable computed from the Hour** in a day.
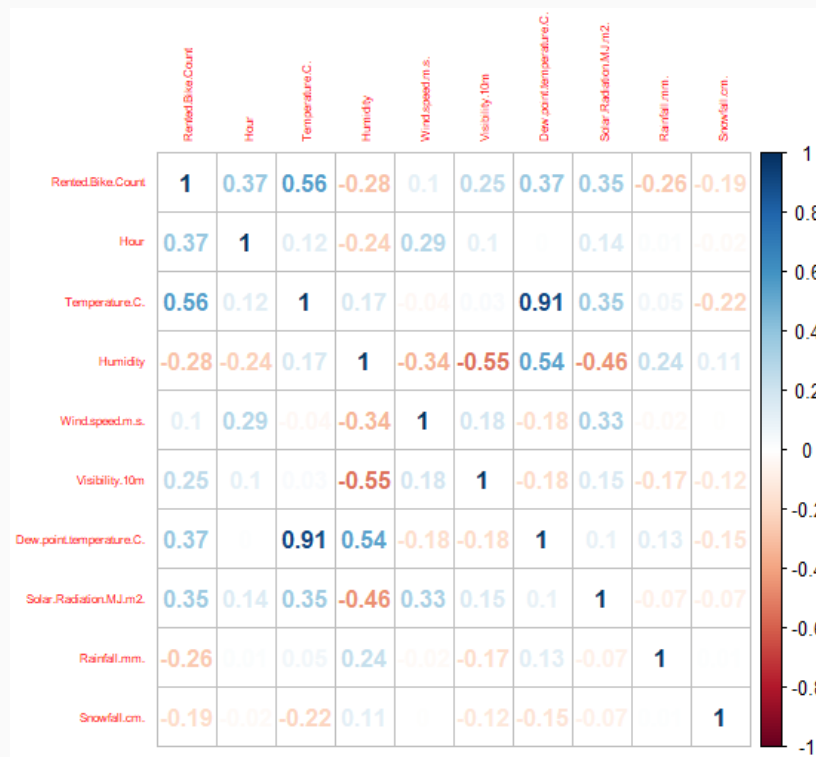


Bikes rented for holiday



Bikes rented for Time

6

Removed examples with **Functioning Day** flag **NO**

Removed variables:
Date, Time and Functioning Day

Target variable: applied with log to make it more similar to normal distribution

There is a strong natural correlation between Temperature (°C) and Dew Point Temperature (°C)

Adjusted $R^2$: **0.6058**
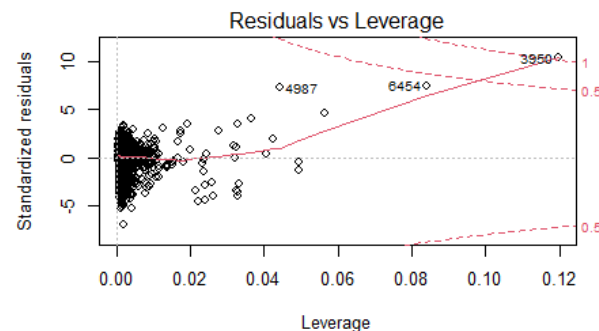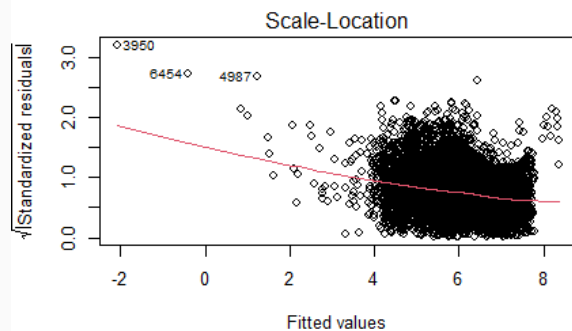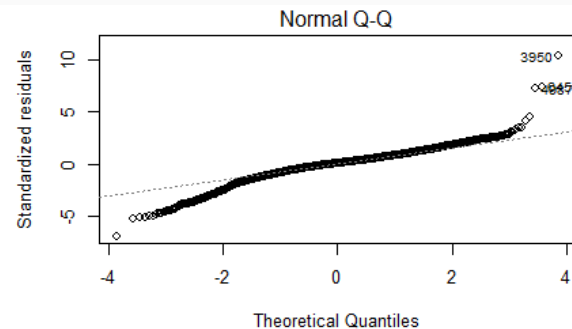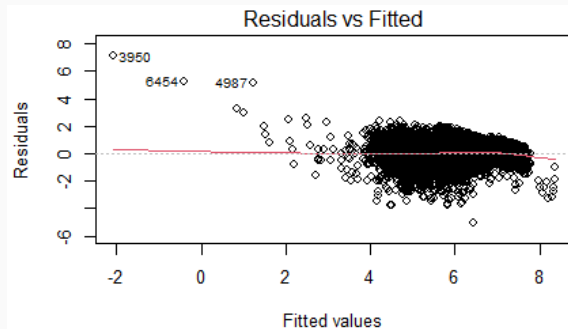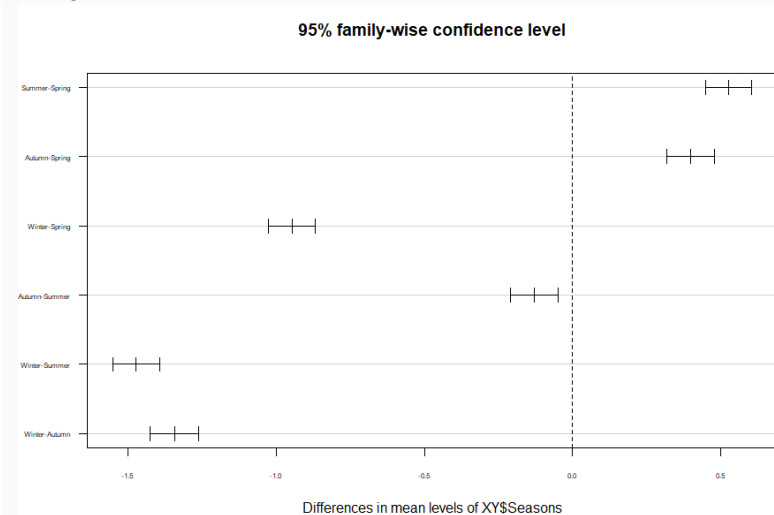
Variables whose estimate are **not significantly different from zero**:
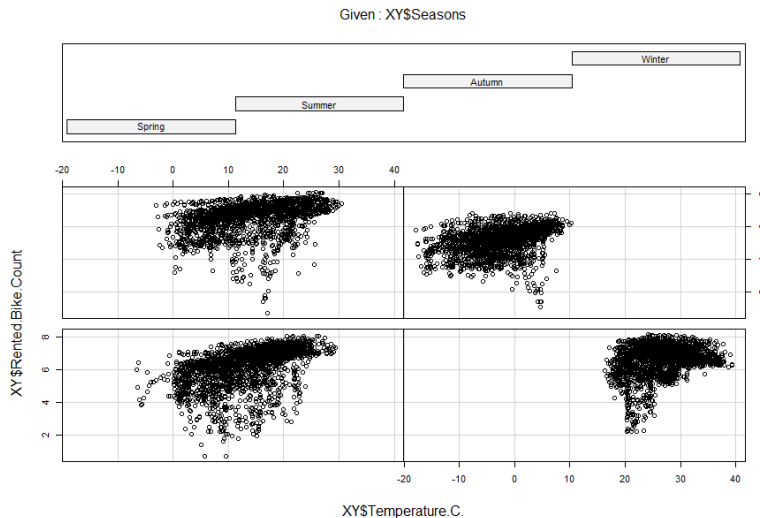- Wind Speed
- Visibility
- Solar radiation
- Snowfall
- Season = Summer

ANOVA on levels of Seasons

1. Rejected H0 of Bartlett's test on equal variances
2. Performed Fisher's test anyway and rejected the null hypothesis of equal means between levels
3. Performed Tukey's test on pairwise differences in means and rejected all null hypothesis on equality

# CHANGES IN THE MODEL

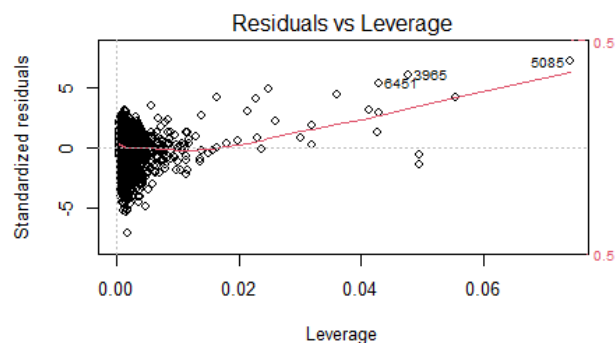GVIF of Dew Point Temperature > 10 provides evidence to remove the variable
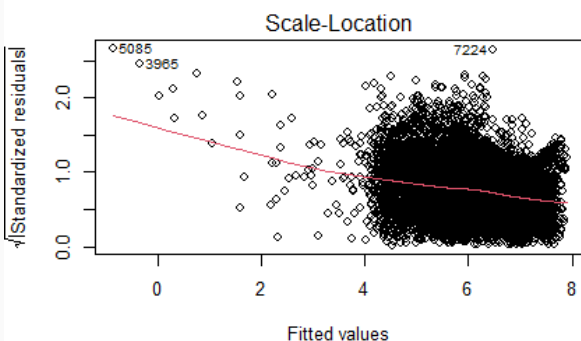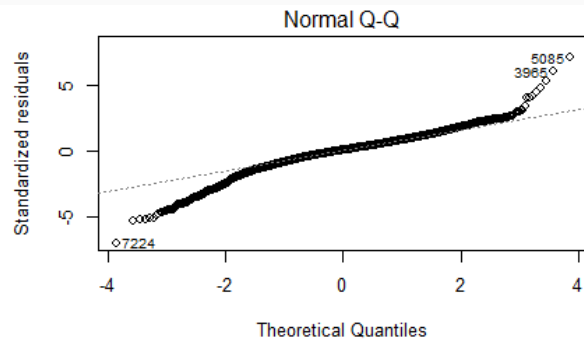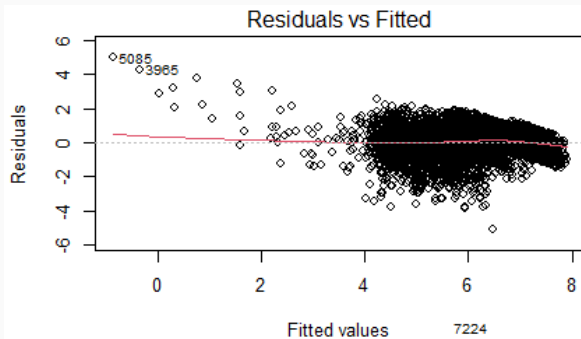
3 outliers removed

Adjusted $R^2$: **0.6125**

**Collinearity** problems **removed** (GVIF now all around 1)

ANOVA table shows **all regressors are significant**



12

# CHANGES IN THE MODEL

Binary factor variables for:
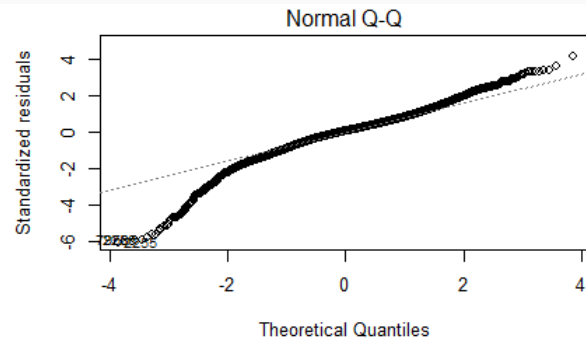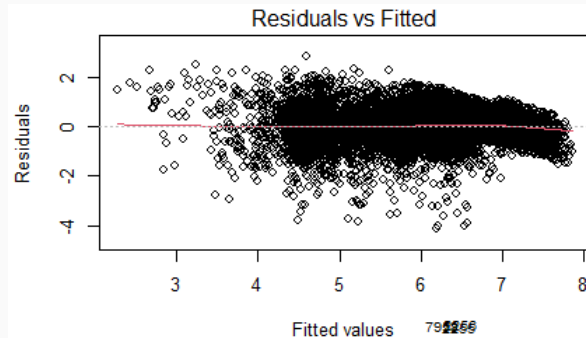  Snowfall
  Rainfall

Adjusted R$^2$: **0.6531**

No changes in collinearity and against the assumptions of the linear regression.

Are all variables gathered necessary?

- **Best Subset Selection**
- Backward step-wise selection
- Forward step-wise selection

Turns out the explained variance does not decrease when the following are **removed**:
- Visibility
- Wind Speed
- Snowfall



15

## METHODS

Best Subset, Forward step-wise and Backward step-wise show
**identical behavior**

## CROSS-VALIDATION

We considered the output of Best Subset Selection, and confirmed it through
**10-fold cross-validation**

## TEST RESULTS

ANOVA (Chi-squared) test on best selection against full model
with **p-value 0.9693**

## ADJ. R²

Same of the full model,
**Adj. R² = 0.6532**

# SHRINKAGE METHODS

## LASSO



Training MSE of 10-fold CV Lasso

Lambda: **0.1**

| Training MSE | **0.5095** |
|---|---|
| Test MSE | **0.5164** |
| $R^2$ | **0.6182** |
| Adj. $R^2$ | **0.6177** |

## RIDGE



Training MSE of 10-fold CV Ridge

Lambda: **0.25**

| Training MSE | **0.4793** |
|---|---|
| Test MSE | **0.4887** |
| $R^2$ | **0.6387** |
| Adj. $R^2$ | **0.6382** |

Number of non-zero components against regularization intensity

Given results on Variable Selection and Shrinkage, the **combined use** of the Best Subset variables and Ridge was tried, but with no improvement.

In this case **lambda** was even **lower** due to a simpler model.

FINAL CHOICE
Non-regularized model with only a subset of variables that estimates the log of the response variable.

**Comparison of true and predicted (log) values**



Log of Rented.Bike.count

Index

True Values
Predicted

# PERFORMANCE RECAP

| | |
|---|---|
| **0.6058** | (1) **All variables**, only «functioning days», **log** of Rented Bike Count |
| **0.6125** | (2) Outliers **removed**, excluded Dew Point Temp. due to **collinearity** |
| **0.6531** | (3) Transformation of **Snowfall** & **Rainfall** in binary factors |
| **0.6532** | (4) **Best Subset Selection**: excluded three variables |
| **0.6177** | (5) **LASSO** shrinkage from model described in **(3)** |
| **0.6382** | (6) **RIDGE** shrinkage from model described in **(3)** |
| **0.6497** | (7) **RIDGE** shrinkage from model described in **(4)** |

# CONCLUSION

## VARIABLES

The **exclusion** of Dew Point Temperature, Wind Speed, Visibility, Snowfall and the use of a binary factor for Rainfall.

## MODELS

The best performing model comes from the **selection of the best subset** of variables **without** the use of any **regularization**.

## FURTHER WORK

Atmospheric variables may not be enough. Variables related to **road traffic**, presence of **events**, data from general **public transportation**, etc. may help in estimating Rented Bike Count.

# Q&A