# Final Project
# DSC 323: Data Analysis and Regression

Seoul Bike Sharing Dataset

# Problem Statement

Currently Rental bikes are introduced in many urban cities for the enhancement of mobility comfort. It is important to make the rental bike available and accessible to the public at the right time as it lessens the waiting time. Eventually, providing the city with a stable supply of rental bikes becomes a major concern. The crucial part is the prediction of bike count required at each hour for the stable supply of rental bikes.

# Overview Of the Data

- There are 14 Different Variables with 8760 Rows of data.

- There are 3 categorical columns and 11 numerical columns.

- No Null Values

## Data Description

- Date : year-month-day
- Rented Bike count - Count of bikes rented at each hour
- Hour - Hour of the day
- Temperature-Temperature in Celsius
- Humidity - %
- Windspeed - m/s
- Visibility - 10m
- Dew point temperature - Celsius
- Solar radiation - MJ/m2
- Rainfall - mm
- Snowfall - cm
- Seasons - Winter, Spring, Summer, Autumn
- Holiday - Holiday/No holiday

- Functional Day - NoFunc(Non Functional Hours), Fun(Functional hours)

# Importing Data

```
/*Importing the Dataset*/
title 'Importing Data: Seoul Bike Sharing Dataset';

data SeoulBikeSharing;
   infile 'SeoulBikeSharing.csv' delimiter = ',' firstobs = 2 missover;
   input Date $ Rented_Bike_Count Hour Temperature Humidity Wind_speed Visibility De
run;

/* Print the entire dataset */
proc print data=SeoulBikeSharing;
run;
```

**Importing Data: Seoul Bike Sharing Dataset**

| Obs | Date | Rented_Bike_Count | Hour | Temperature | Humidity | Wind_speed | Visibility | Dew_point_temperature | Solar_Radiation | Rainfall | Snowfall | Seasons | Holiday | Functioning_Day |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 01/12/20 | 254 | 0 | -5.2 | 37 | 2.2 | 2000 | -17.6 | 0.00 | 0.0 | 0 | Winter | No Holid | Yes |
| 2 | 01/12/20 | 204 | 1 | -5.5 | 38 | 0.8 | 2000 | -17.6 | 0.00 | 0.0 | 0 | Winter | No Holid | Yes |
| 3 | 01/12/20 | 173 | 2 | -6.0 | 39 | 1.0 | 2000 | -17.7 | 0.00 | 0.0 | 0 | Winter | No Holid | Yes |
| 4 | 01/12/20 | 107 | 3 | -6.2 | 40 | 0.9 | 2000 | -17.6 | 0.00 | 0.0 | 0 | Winter | No Holid | Yes |
| 5 | 01/12/20 | 78 | 4 | -6.0 | 36 | 2.3 | 2000 | -18.6 | 0.00 | 0.0 | 0 | Winter | No Holid | Yes |
| 6 | 01/12/20 | 100 | 5 | -6.4 | 37 | 1.5 | 2000 | -18.7 | 0.00 | 0.0 | 0 | Winter | No Holid | Yes |
| 7 | 01/12/20 | 181 | 6 | -6.6 | 35 | 1.3 | 2000 | -19.5 | 0.00 | 0.0 | 0 | Winter | No Holid | Yes |
| 8 | 01/12/20 | 460 | 7 | -7.4 | 38 | 0.9 | 2000 | -19.3 | 0.00 | 0.0 | 0 | Winter | No Holid | Yes |
| 9 | 01/12/20 | 930 | 8 | -7.6 | 37 | 1.1 | 2000 | -19.8 | 0.01 | 0.0 | 0 | Winter | No Holid | Yes |
| 10 | 01/12/20 | 490 | 9 | -6.5 | 27 | 0.5 | 1928 | -22.4 | 0.23 | 0.0 | 0 | Winter | No Holid | Yes |
| 11 | 01/12/20 | 339 | 10 | -3.5 | 24 | 1.2 | 1996 | -21.2 | 0.65 | 0.0 | 0 | Winter | No Holid | Yes |
| 12 | 01/12/20 | 360 | 11 | -0.5 | 21 | 1.3 | 1936 | -20.2 | 0.94 | 0.0 | 0 | Winter | No Holid | Yes |

# Dummy Variables

- **Snowfall – cm**
- **Seasons - Winter, Spring, Summer, Autumn**
- **Holiday - Holiday/No holiday**

- **Functional Day – NoFunc / Fun**

```
/* Adding Dummy Variables */
title 'Creating Dummy Variables';

data SeoulBikeSharing;
infile 'SeoulBikeSharing.csv' delimiter = ',' firstobs = 2 missover;
input Date $ Rented_Bike_Count Hour Temperature Humidity Wind_speed Visibility Dew_

  Dum_Winter = (Seasons = 'Winter');
  Dum_Spring = (Seasons = 'Spring');
  Dum_Summer = (Seasons = 'Summer');
  IsHoliday = (Holiday = 'Holiday');
  IsFunctionalDay = (Functioning_Day = 'Yes');
  SnowfallDummy = (Snowfall = '1');

  drop Seasons Holiday Functioning_Day Snowfall Snowfall_Category;
run;

/* Print the dataset with the new dummy variables */
proc print data=SeoulBikeSharing;
run;
```

# Dummy Variables

```
/* Adding Dummy Variables */
title 'Creating Dummy Variables';

data SeoulBikeSharing;
infile 'SeoulBikeSharing.csv' delimiter = ',' firstobs = 2 missover;
input Date $ Rented_Bike_Count Hour Temperature Humidity Wind_speed Visibility Dew_

  Dum_Winter = (Seasons = 'Winter');
  Dum_Spring = (Seasons = 'Spring');
  Dum_Summer = (Seasons = 'Summer');
  IsHoliday = (Holiday = 'Holiday');
  IsFunctionalDay = (Functioning_Day = 'Yes');
  SnowfallDummy = (Snowfall = '1');

  drop Seasons Holiday Functioning_Day Snowfall Snowfall_Category;
run;

/* Print the dataset with the new dummy variables */
proc print data=SeoulBikeSharing;
run;
```
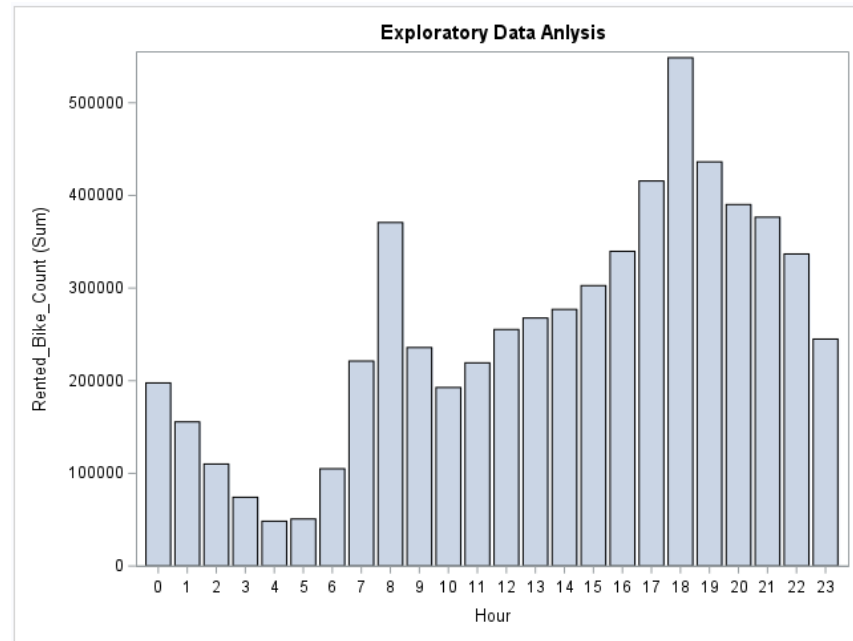
**Creating Dummy Variables**

| Count | Hour | Temperature | Humidity | Wind_speed | Visibility | Dew_point_temperature | Solar_Radiation | Rainfall | Dum_Winter | Dum_Spring | Dum_Summer | IsHoliday | IsFunctionalDay | SnowfallDummy |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 254 | 0 | -5.2 | 37 | 2.2 | 2000 | -17.6 | 0.00 | 0.0 | 1 | 0 | 0 | 0 | 1 | 0 |
| 204 | 1 | -5.5 | 38 | 0.8 | 2000 | -17.6 | 0.00 | 0.0 | 1 | 0 | 0 | 0 | 1 | 0 |
| 173 | 2 | -6.0 | 39 | 1.0 | 2000 | -17.7 | 0.00 | 0.0 | 1 | 0 | 0 | 0 | 1 | 0 |
| 107 | 3 | -6.2 | 40 | 0.9 | 2000 | -17.6 | 0.00 | 0.0 | 1 | 0 | 0 | 0 | 1 | 0 |
| 78 | 4 | -6.0 | 36 | 2.3 | 2000 | -18.6 | 0.00 | 0.0 | 1 | 0 | 0 | 0 | 1 | 0 |
| 100 | 5 | -6.4 | 37 | 1.5 | 2000 | -18.7 | 0.00 | 0.0 | 1 | 0 | 0 | 0 | 1 | 0 |
| 181 | 6 | -6.6 | 35 | 1.3 | 2000 | -19.5 | 0.00 | 0.0 | 1 | 0 | 0 | 0 | 1 | 0 |
| 460 | 7 | -7.4 | 38 | 0.9 | 2000 | -19.3 | 0.00 | 0.0 | 1 | 0 | 0 | 0 | 1 | 0 |
| 930 | 8 | -7.6 | 37 | 1.1 | 2000 | -19.8 | 0.01 | 0.0 | 1 | 0 | 0 | 0 | 1 | 0 |
| 490 | 9 | -6.5 | 27 | 0.5 | 1928 | -22.4 | 0.23 | 0.0 | 1 | 0 | 0 | 0 | 1 | 0 |
| 339 | 10 | -3.5 | 24 | 1.2 | 1996 | -21.2 | 0.65 | 0.0 | 1 | 0 | 0 | 0 | 1 | 0 |

# Exploratory Data Analysis

# Histogram
## Bike Share – Hour



Exploratory Data Anlysis
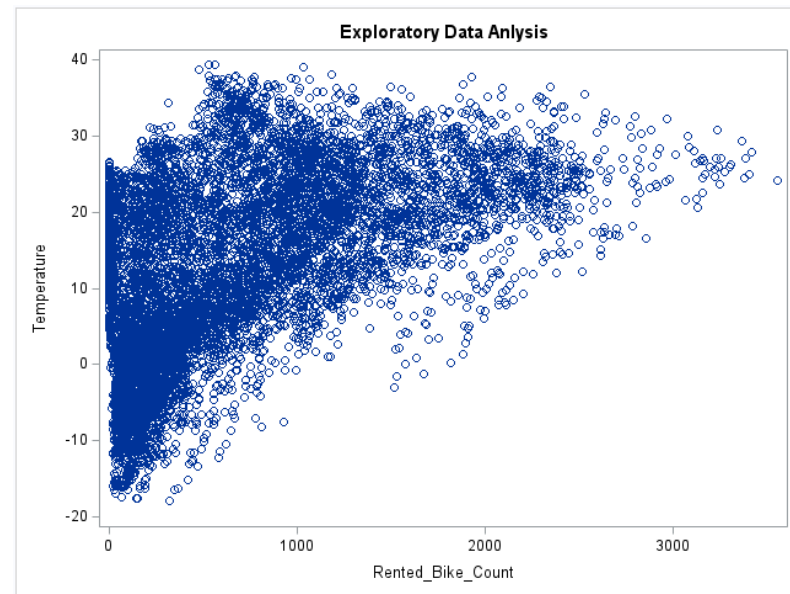
```
/* Create a histogram of Rented Bike Counts by Hour */
title 'Exploratory Data Anlysis';
proc sgplot data=SeoulBikeSharing;
  vbar Hour / response=Rented_Bike_Count;
run;
```

# Scatterplot
Bike Share – Temperature



```
/* Rented Bike vs Temperature */
proc sgplot data=SeoulBikeSharing;
   scatter x=Rented_Bike_Count y=Temperature;
run;
```

# Boxplot
Snowfall vs No-Snowfall

```
/* Calculate the total number of bikes rented when SnowfallDummy is 0 vs 1 */
proc means data=SeoulBikeSharing sum;
  class SnowfallDummy;
  var Rented_Bike_Count;
run;

/* Sort the data by 'SnowfallDummy' variable */
PROC SORT data=SeoulBikeSharing;
BY SnowfallDummy;
RUN;

/* Create the boxplot */
PROC BOXPLOT data=SeoulBikeSharing;
PLOT Rented_Bike_Count*SnowfallDummy;
RUN;
```

**Exploratory Data Anlysis**

**The MEANS Procedure**

| Analysis Variable : Rented_Bike_Count | | |
|---|---|---|
| SnowfallDummy | N Obs | Sum |
| 0 | 8721 | 6165957.00 |
| 1 | 39 | 6357.00 |



**Exploratory Data Anlysis**

Distribution of Rented_Bike_Count by SnowfallDummy

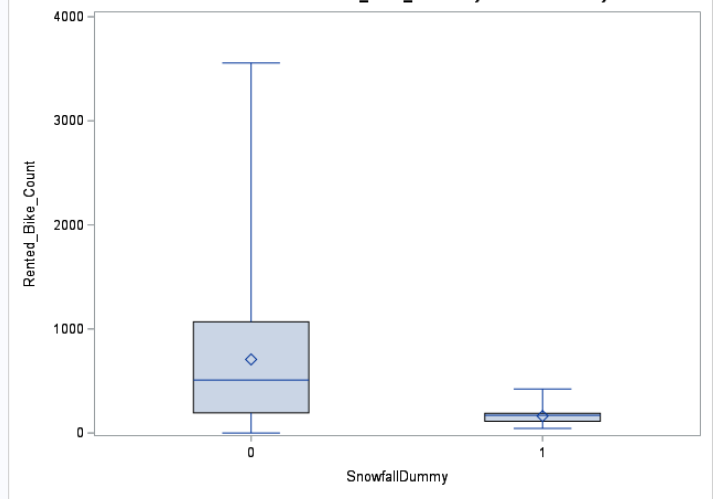| 3298 | 3309 | 3365 | 3380 | 3384 | 3404 | 3418 | 3556 | Total |
|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 295 |
| 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 3.37 |
| 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | |
| 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | |
| 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 8465 |
| 0.01 | 0.01 | 0.01 | 0.01 | 0.01 | 0.01 | 0.01 | 0.01 | 96.63 |
| 0.01 | 0.01 | 0.01 | 0.01 | 0.01 | 0.01 | 0.01 | 0.01 | |
| 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | |
| 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 8760 |
| 0.01 | 0.01 | 0.01 | 0.01 | 0.01 | 0.01 | 0.01 | 0.01 | 100.00 |

```
/* Sort the data */
PROC SORT data=SeoulBikeSharing;
BY IsFunctionalDay;
RUN;

/* Create a boxplot for Rented Bike Counts by Functioning Day */
PROC BOXPLOT data=SeoulBikeSharing;
PLOT Rented_Bike_Count*IsFunctionalDay;
RUN;

/* Create a frequency table for Rented Bike Counts by Functioning Day */
proc freq data=SeoulBikeSharing;
  tables IsFunctionalDay * Rented_Bike_Count;
run;
```

# Boxplot
Functioning vs Non-Functioning Day

**Exploratory Data Anlysis**

**Distribution of Rented_Bike_Count by IsFunctionalDay**

# Scatterplot

```
/* ScatterPlot */
proc sgscatter data=Test;
   matrix Rented_Bike_Count Hour Temperature Humidity Visibility Dew_Point_Temperature
   title "Scatterplot Matrix";
run;
```


Scatterplot Matrix

# Train and Test Data

- I choose to do 80/20 split for the dataset
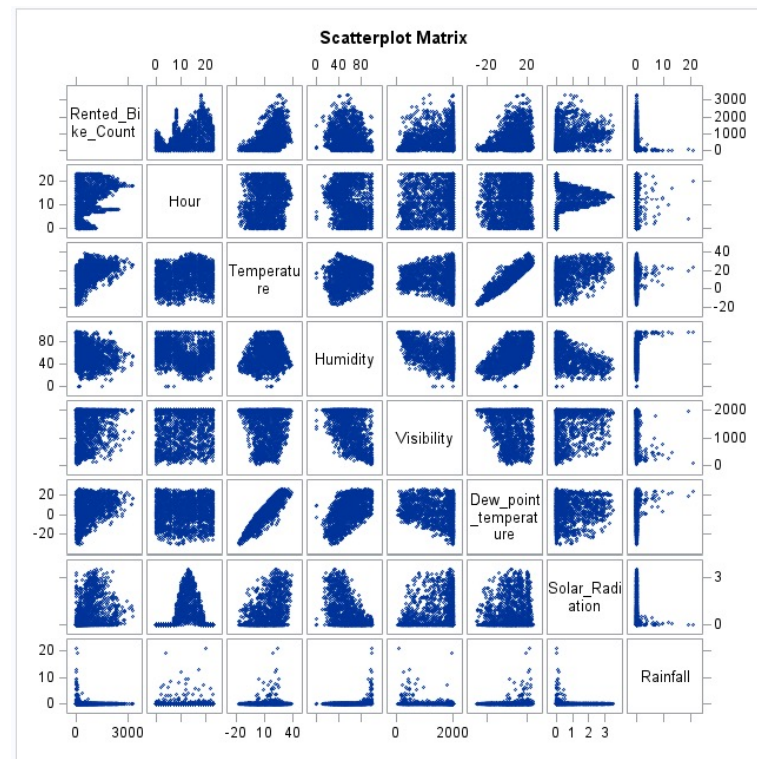- (samprate=0.8)
- Seed = 592587
- Training: Approx 1700 Rows of Data
- Testing: Approx 7000 Rows of Data

```
/* Create a training and test dataset split with an 80/20 ratio */
title 'Split';
proc surveyselect data=SeoulBikeSharing out=train outall seed=592587 samprate=0.8;
run;
```

**Split**

| Obs | Selected | Date | Rented_Bike_Count | Hour | Temperature | Humidity | Wind_speed | Visibility | Dew_point_temperature | Solar_Radiation | Rainfall | Dum_Winter | Dum_Spring | Dum_Summer | IsHol |
|-----|----------|---------|-------------------|------|-------------|----------|------------|------------|-----------------------|-----------------|----------|------------|------------|------------|-------|
| 1 | 0 | 01/12/20 | 204 | 1 | -5.5 | 38 | 0.8 | 2000 | -17.6 | 0.00 | 0.0 | 1 | 0 | 0 | |
| 2 | 0 | 01/12/20 | 100 | 5 | -6.4 | 37 | 1.5 | 2000 | -18.7 | 0.00 | 0.0 | 1 | 0 | 0 | |
| 3 | 0 | 01/12/20 | 449 | 12 | 1.7 | 23 | 1.4 | 2000 | -17.2 | 1.11 | 0.0 | 1 | 0 | 0 | |
| 4 | 0 | 01/12/20 | 463 | 15 | 2.1 | 36 | 3.2 | 2000 | -11.4 | 0.54 | 0.0 | 1 | 0 | 0 | |
| 5 | 0 | 02/12/20 | 219 | 8 | -4.2 | 79 | 2.1 | 1436 | -7.3 | 0.01 | 0.0 | 1 | 0 | 0 | |
| 6 | 0 | 02/12/20 | 479 | 12 | 4.3 | 41 | 1.3 | 1666 | -7.8 | 1.09 | 0.0 | 1 | 0 | 0 | |
| 7 | 0 | 02/12/20 | 385 | 19 | 5.0 | 52 | 2.3 | 1666 | -4.0 | 0.00 | 0.0 | 1 | 0 | 0 | |
| 8 | 0 | 02/12/20 | 359 | 20 | 4.6 | 51 | 1.2 | 1585 | -4.6 | 0.00 | 0.0 | 1 | 0 | 0 | |

**Split**

| Obs | Selected | Date | Rented_Bike_Count | Hour | Temperature | Humidity | Wind_speed | Visibility | Dew_point_temperature | Solar_Radiation | Rainfall | Dum_Winter | Dum_Spring | Dum_Summer | IsHol |
|-----|----------|---------|-------------------|------|-------------|----------|------------|------------|-----------------------|-----------------|----------|------------|------------|------------|-------|
| 1 | 1 | 01/12/20 | 254 | 0 | -5.2 | 37 | 2.2 | 2000 | -17.6 | 0.00 | 0.0 | 1 | 0 | 0 | |
| 2 | 0 | 01/12/20 | 204 | 1 | -5.5 | 38 | 0.8 | 2000 | -17.6 | 0.00 | 0.0 | 1 | 0 | 0 | |
| 3 | 1 | 01/12/20 | 173 | 2 | -6.0 | 39 | 1.0 | 2000 | -17.7 | 0.00 | 0.0 | 1 | 0 | 0 | |
| 4 | 1 | 01/12/20 | 107 | 3 | -6.2 | 40 | 0.9 | 2000 | -17.6 | 0.00 | 0.0 | 1 | 0 | 0 | |
| 5 | 1 | 01/12/20 | 78 | 4 | -6.0 | 36 | 2.3 | 2000 | -18.6 | 0.00 | 0.0 | 1 | 0 | 0 | |
| 6 | 0 | 01/12/20 | 100 | 5 | -6.4 | 37 | 1.5 | 2000 | -18.7 | 0.00 | 0.0 | 1 | 0 | 0 | |
| 7 | 1 | 01/12/20 | 181 | 6 | -6.6 | 35 | 1.3 | 2000 | -19.5 | 0.00 | 0.0 | 1 | 0 | 0 | |
| 8 | 1 | 01/12/20 | 460 | 7 | -7.4 | 38 | 0.9 | 2000 | -19.3 | 0.00 | 0.0 | 1 | 0 | 0 | |
| 9 | 1 | 01/12/20 | 930 | 8 | -7.6 | 37 | 1.1 | 2000 | -19.8 | 0.01 | 0.0 | 1 | 0 | 0 | |
| 10 | 1 | 01/12/20 | 490 | 9 | -6.5 | 27 | 0.5 | 1928 | -22.4 | 0.23 | 0.0 | 1 | 0 | 0 | |
| 11 | 1 | 01/12/20 | 339 | 10 | -3.5 | 24 | 1.2 | 1996 | -21.2 | 0.65 | 0.0 | 1 | 0 | 0 | |
| 12 | 1 | 01/12/20 | 360 | 11 | -0.5 | 21 | 1.3 | 1936 | -20.2 | 0.94 | 0.0 | 1 | 0 | 0 | |
| 13 | 0 | 01/12/20 | 449 | 12 | 1.7 | 23 | 1.4 | 2000 | -17.2 | 1.11 | 0.0 | 1 | 0 | 0 | |
| 14 | 1 | 01/12/20 | 451 | 13 | 2.4 | 25 | 1.6 | 2000 | -15.6 | 1.16 | 0.0 | 1 | 0 | 0 | |
| 15 | 1 | 01/12/20 | 447 | 14 | 3.0 | 26 | 2.0 | 2000 | -14.6 | 1.01 | 0.0 | 1 | 0 | 0 | |
| 16 | 0 | 01/12/20 | 463 | 15 | 2.1 | 36 | 3.2 | 2000 | -11.4 | 0.54 | 0.0 | 1 | 0 | 0 | |
| 17 | 1 | 01/12/20 | 484 | 16 | 1.2 | 54 | 4.2 | 793 | -7.0 | 0.24 | 0.0 | 1 | 0 | 0 | |
| 18 | 1 | 01/12/20 | 555 | 17 | 0.8 | 58 | 1.6 | 2000 | -6.5 | 0.08 | 0.0 | 1 | 0 | 0 | |

# Selection Process

```
title 'Selection Process';
/* Perform forward selection with the binary response variable */
proc reg data = train;
   model Rented_Bike_Count = Hour Temperature Humidity Wind_speed Visibility Dew_poin
run;

proc reg data = train;
   model Rented_Bike_Count = Hour Temperature Humidity Wind_speed Visibility Dew_poin
run;

proc reg data = train;
   model Rented_Bike_Count = Hour Temperature Humidity Wind_speed Visibility Dew_poin
run;
```

**1.Forward Selection:**
1. Start with an empty set of features.
2. Iteratively add features one at a time, selecting the one that provides the best improvement in model performance.

**2.Backward Elimination:**
1. Start with all features.
2. Iteratively remove features one at a time, eliminating the one that has the least impact on model performance.

**3.Stepwise Selection:**
1. This method combines forward selection and backward elimination.
2. It starts with an empty set of features and adds features in a forward selection fashion.
3. At each step, it also checks if removing any previously added feature (backward elimination) would improve model performance.

# Outliers , Influential Points, Significance values and VIFs>10

```
/* Fit a multiple regression model and compute VIF */
title 'Multiple Regression and VIF';
proc reg data = train;
    model Rented_Bike_Count = Hour Temperature Humidity Wind_speed Visibility Dew_poin
run;
```

```
proc reg data = train;
    model Rented_Bike_Count = Hour Temperature Humidity Wind_speed Dew_point_temperat
run;
```

```
proc reg data = train;
    model Rented_Bike_Count = Hour Temperature Humidity Wind_speed  Solar_Radiation F
run;
```

## Multiple Regression and VIF

### The REG Procedure
### Model: MODEL1
### Dependent Variable: Rented_Bike_Count

| Number of Observations Read | 8760 |
|---|---|
| Number of Observations Used | 8760 |

#### Analysis of Variance

| Source | DF | Sum of Squares | Mean Square | F Value | Pr > F |
|---|---|---|---|---|---|
| Model | 13 | 2004173836 | 154167218 | 822.28 | <.0001 |
| Error | 8746 | 1639760527 | 187487 | | |
| Corrected Total | 8759 | 3643934363 | | | |

| Root MSE | 432.99759 | R-Square | 0.5500 |
|---|---|---|---|
| Dependent Mean | 704.60205 | Adj R-Sq | 0.5493 |
| Coeff Var | 61.45279 | | |

#### Parameter Estimates

| Variable | DF | Parameter Estimate | Standard Error | t Value | Pr > |t| | Variance Inflation |
|---|---|---|---|---|---|---|
| Intercept | 1 | 10.08233 | 96.18448 | 0.10 | 0.9165 | 0 |
| Hour | 1 | 27.61298 | 0.73449 | 37.59 | <.0001 | 1.20777 |
| Temperature | 1 | 16.59975 | 3.65933 | 4.54 | <.0001 | 89.25760 |
| Humidity | 1 | -10.43131 | 1.02192 | -10.21 | <.0001 | 20.22905 |
| Wind_speed | 1 | 19.31135 | 5.09731 | 3.79 | 0.0002 | 1.30358 |
| Visibility | 1 | 0.00970 | 0.00988 | 0.98 | 0.3265 | 1.68843 |
| Dew_point_temperature | 1 | 10.31851 | 3.82564 | 2.70 | 0.0070 | 116.62782 |
| Solar_Radiation | 1 | -75.45111 | 7.56960 | -9.97 | <.0001 | 2.02029 |
| Rainfall | 1 | -58.88760 | 4.26991 | -13.79 | <.0001 | 1.08415 |
| Dum_Winter | 1 | -362.37915 | 19.67649 | -18.42 | <.0001 | 3.36061 |
| Dum_Spring | 1 | -138.24854 | 13.84919 | -9.98 | <.0001 | 1.68945 |
| Dum_Summer | 1 | -154.51586 | 17.21802 | -8.97 | <.0001 | 2.61134 |
| IsHoliday | 1 | -119.35433 | 21.60543 | -5.52 | <.0001 | 1.02253 |
| IsFunctionalDay | 1 | 933.37414 | 26.66037 | 35.01 | <.0001 | 1.08070 |

---

| Number of Observations Read | 8760 |
|---|---|
| Number of Observations Used | 8760 |

#### Analysis of Variance

| Source | DF | Sum of Squares | Mean Square | F Value | Pr > F |
|---|---|---|---|---|---|
| Model | 12 | 2003993317 | 166999443 | 890.73 | <.0001 |
| Error | 8747 | 1639941045 | 187486 | | |
| Corrected Total | 8759 | 3643934363 | | | |

| Root MSE | 432.99667 | R-Square | 0.5500 |
|---|---|---|---|
| Dependent Mean | 704.60205 | Adj R-Sq | 0.5493 |
| Coeff Var | 61.45265 | | |

#### Parameter Estimates

| Variable | DF | Parameter Estimate | Standard Error | t Value | Pr > |t| | Variance Inflation |
|---|---|---|---|---|---|---|
| Intercept | 1 | 43.19878 | 90.06877 | 0.48 | 0.6315 | 0 |
| Hour | 1 | 27.56318 | 0.73273 | 37.62 | <.0001 | 1.20201 |
| Temperature | 1 | 16.39517 | 3.65337 | 4.49 | <.0001 | 88.96787 |
| Humidity | 1 | -10.69002 | 0.98732 | -10.83 | <.0001 | 18.88251 |
| Wind_speed | 1 | 19.78163 | 5.07472 | 3.90 | <.0001 | 1.29205 |
| Dew_point_temperature | 1 | 10.54079 | 3.81892 | 2.76 | 0.0058 | 116.21890 |
| Solar_Radiation | 1 | -76.94661 | 7.41456 | -10.38 | <.0001 | 1.93839 |
| Rainfall | 1 | -59.02476 | 4.26761 | -13.83 | <.0001 | 1.08298 |
| Dum_Winter | 1 | -365.88923 | 19.34855 | -18.91 | <.0001 | 3.24953 |
| Dum_Spring | 1 | -141.26176 | 13.50441 | -10.46 | <.0001 | 1.60639 |
| Dum_Summer | 1 | -153.61635 | 17.19356 | -8.93 | <.0001 | 2.60394 |
| IsHoliday | 1 | -118.95628 | 21.60158 | -5.51 | <.0001 | 1.02217 |
| IsFunctionalDay | 1 | 932.98346 | 26.65734 | 35.00 | <.0001 | 1.08046 |

---

### Dependent Variable: Rented_Bike_Count

| Number of Observations Read | 8760 |
|---|---|
| Number of Observations Used | 8760 |

#### Analysis of Variance

| Source | DF | Sum of Squares | Mean Square | F Value | Pr > F |
|---|---|---|---|---|---|
| Model | 11 | 2002564965 | 182051360 | 970.28 | <.0001 |
| Error | 8748 | 1641369397 | 187628 | | |
| Corrected Total | 8759 | 3643934363 | | | |

| Root MSE | 433.16043 | R-Square | 0.5496 |
|---|---|---|---|
| Dependent Mean | 704.60205 | Adj R-Sq | 0.5490 |
| Coeff Var | 61.47590 | | |

#### Parameter Estimates

| Variable | DF | Parameter Estimate | Standard Error | t Value | Pr > |t| | Variance Inflation |
|---|---|---|---|---|---|---|
| Intercept | 1 | -184.69236 | 36.00604 | -5.13 | <.0001 | 0 |
| Hour | 1 | 27.48078 | 0.73240 | 37.52 | <.0001 | 1.20001 |
| Temperature | 1 | 26.19223 | 0.86547 | 30.26 | <.0001 | 4.98902 |
| Humidity | 1 | -8.09285 | 0.29913 | -27.05 | <.0001 | 1.73196 |
| Wind_speed | 1 | 19.31556 | 5.07383 | 3.81 | 0.0001 | 1.29062 |
| Solar_Radiation | 1 | -81.71191 | 7.21349 | -11.33 | <.0001 | 1.83330 |
| Rainfall | 1 | -60.38737 | 4.24057 | -14.24 | <.0001 | 1.06849 |
| Dum_Winter | 1 | -367.62035 | 19.34569 | -19.00 | <.0001 | 3.24612 |
| Dum_Spring | 1 | -142.81440 | 13.49780 | -10.58 | <.0001 | 1.60360 |
| Dum_Summer | 1 | -149.30870 | 17.12906 | -8.72 | <.0001 | 2.58249 |
| IsHoliday | 1 | -118.23412 | 21.60816 | -5.47 | <.0001 | 1.02202 |
| IsFunctionalDay | 1 | 930.66254 | 26.65415 | 34.92 | <.0001 | 1.07939 |

# Outliers , Influential Points, Significance values and VIFs>10

```
title 'Outliers & Influential Points';
proc reg data=train;
    model Rented_Bike_Count = Hour Temperature Humidity Wind_speed Solar_Radi
    plot student.*(Rented_Bike_Count = Hour Temperature Humidity Wind_speed S
    plot npp.*student.;
run;
```

**The REG Procedure**

## Outliers & Influential Points

Rented_Bike_Count = -184.69 +27.481 Hour +26.192 Temperature -8.0929 Humidity +19.316 Wind_speed
-81.712 Solar_Radiation -60.387 Rainfall -367.62 Dum_Winter -142.81 Dum_Spring
-149.31 Dum_Summer -118.23 IsHoliday +930.66 IsFunctionalDay

N 8760
Rsq 0.5496
AdjRsq 0.5490
RMSE 433.16

Normal Cumulative Distribution vs CDF of Studentized Residual

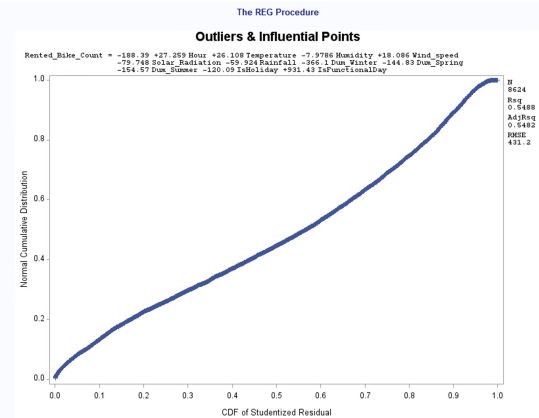# Outliers , Influential Points, Significance values and VIFs>10

```
title 'Outliers & Influential Points';
proc reg data = train;
model Rented_Bike_Count = Hour Temperature Wind_speed  Solar_Radiation Rainfall Dum
plot student.*(Rented_Bike_Count = Hour Temperature Wind_speed  Solar_Radiation Rai
plot npp.*student.;

data train_02;
SET train;
IF _N_ IN (3283, 3297, 3499, 3513, 3523, 3537, 3547, 3619, 3681, 3705, 3715, 3825,
4161, 4171, 4185, 4195, 4219, 4281, 4291, 4305, 4339, 4340, 4353, 4363, 4377, 4387,
4674, 4699, 4713, 4723, 4724, 4785, 4809, 4819, 4820, 4833, 4843, 4844, 4857, 4867,
5313, 5347, 5371, 5395, 5491, 6331, 6571, 6667, 6681, 6691, 6705, 6729, 6739, 6811,
7651, 7665, 7675, 7689, 7713, 7737, 7809, 7843, 7857, 7867, 7881, 8001, 8025, 8049,
RUN;

proc reg data = train_02;
model Rented_Bike_Count = Hour Temperature Wind_speed  Solar_Radiation Rainfall Dum
plot student.*(Rented_Bike_Count = Hour Temperature Wind_speed Solar_Radiation Rain
plot npp.*student.;
run;

data train_03;
SET train;
IF _N_ IN ( 2947, 2961, 2971, 2995, 3115, 3177, 3273, 3473, 3487, 3602, 3682, 3790,
4617, 4618, 4662, 4687, 4688, 4709, 4710, 4711, 4755, 4756, 4776, 4797, 4798, 4799,
6788, 7098, 7193, 7415, 7429, 7599, 7622, 7693, 7761, 7847, 8020, 8187, 8201, 8248,
RUN;

proc reg data = train_03;
model Rented_Bike_Count = Hour Temperature Wind_speed  Solar_Radiation Rainfall Dum
plot student.*(Rented_Bike_Count = Hour Temperature Wind_speed Solar_Radiation Rain
plot npp.*student.;
run;
```

# Challenges

- Since the dataset was really large with 8760 datasets it was a bit difficult to handle it. Even after 80-20 split, it was more >1000 datasets to handle.

- Even after removing All the Influential Points, Outliers Significance values and VIFs>10, It was still a lot of datapoints to find the outliers and influential points.

- Also due to the same reason the computational time was a lot compared to smaller datasets, especially on a virtual machine.

# Conclusion

- Rented Bike Count = -188.39 + 27.259 Hour +26.108 Temperature -7.9786 Humidity +18.085 Wind_Speed -79.748 Solar_Raditation -59.924 Rainfall -366.1 Dum_Winter +144.83 Dum_Spring +154.57 Dum_Summer -120.09 IsHoliday +931.43 IsFunctionalDay

- We observed that during the day, most demand was there during 8AM and 6PM, but the demand started to grow +- 1 hour before that.

- There is more demand for a bike during the Weekdays compared to Weekends

- The demand for bikes increased by 18 as windspeed increases.

- The demand decreases by 8 bikes when the humidity decreases.

- The demand decreased by 80 bikes when the solar radiation increases.

- When winter kicks in the demand decreases drastically by 366 bikes, during sprint it increases by 145 bikes, also during summer it increases by 155 bikes.

- For functioning days when it is a holdiday the demand decreases by 120 bikes but when it is a functional day it increases drastically by 931 bikes.