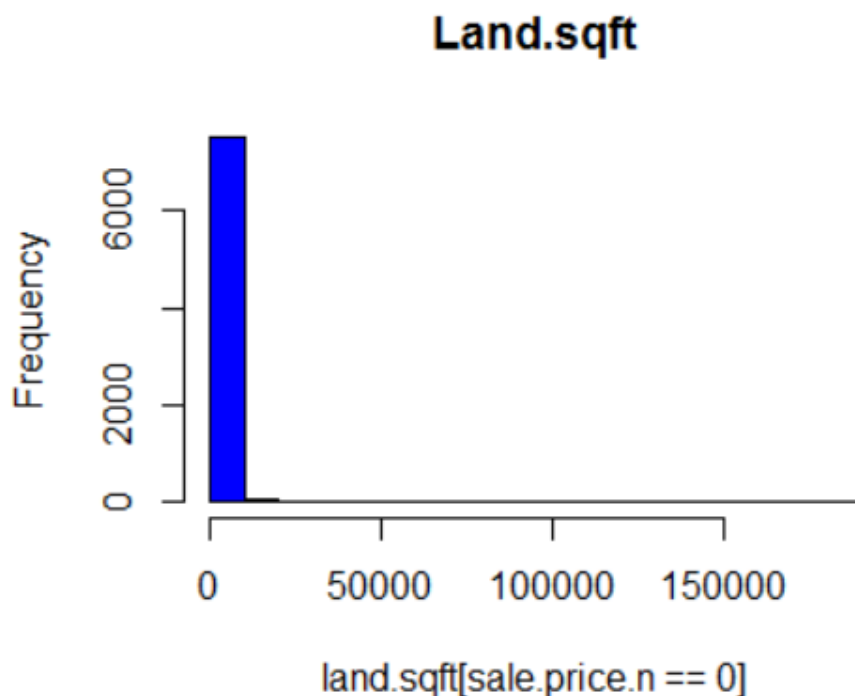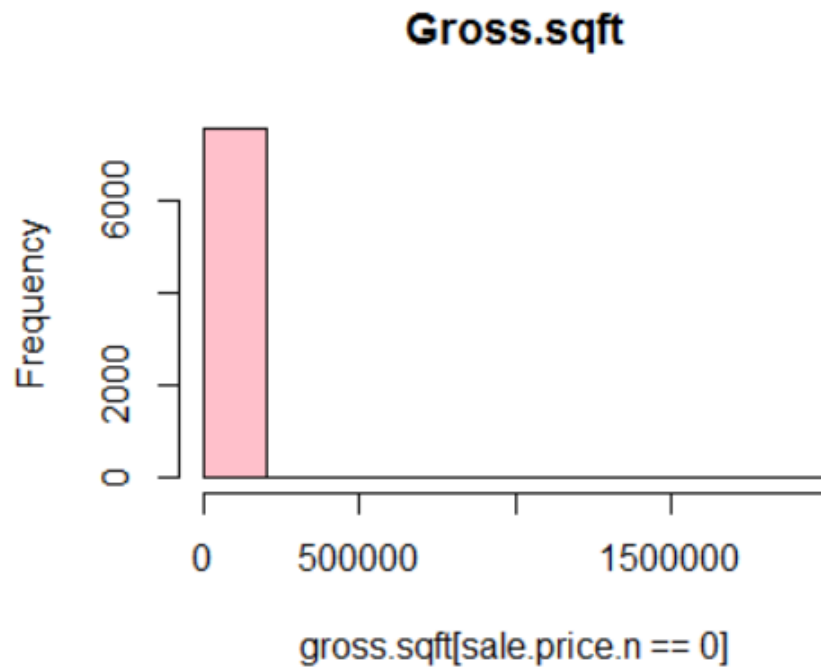# Problem 3 Case Study

This question is quite similar to Problem 2. The difference is in the given data and also a slight variation in the approach. In this question, the main task is to perform EDA on a simple data set. The first part of this question demands an EDA on the sales of one of the Bronx of New York state. To do that, I first ensured that the data provided is free from redundant data and then followed the Linear Regression process to obtain EDA. Also to mention, in order to fetch the .xls file, Windows machine required **Perl** to be installed along with the absolute path of perl.exe mentioned as one the arguments in read.xls() command.

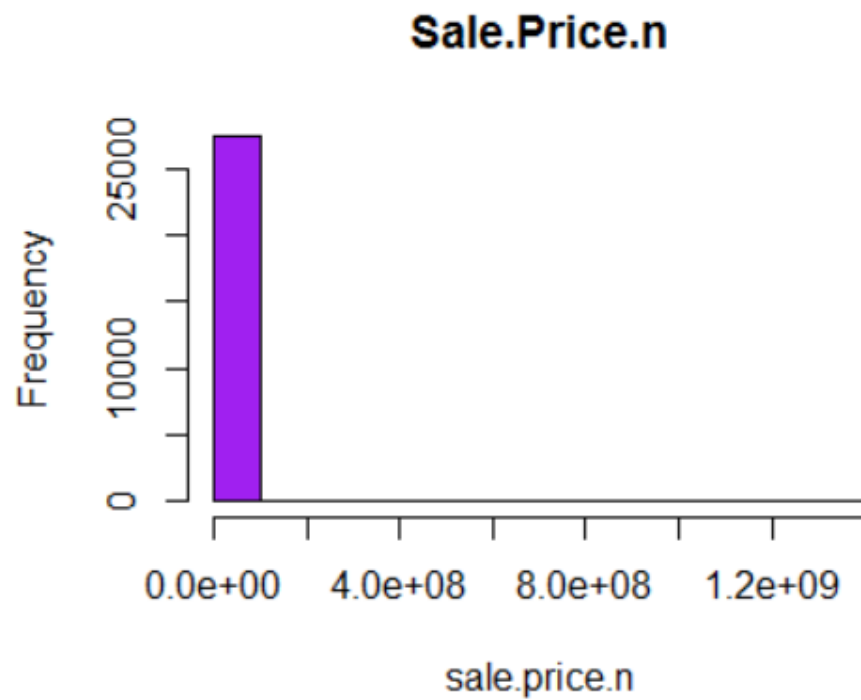Following were the plots obtained:
The first scenario represents Land.sqft vs Frequency when **sale.price.n==0 for one Borough of New York**

**Land.sqft**



land.sqft[sale.price.n == 0]

The next scenario represents Gross.sqft vs Frequency when sale.price.n==0 for one Borough of New York
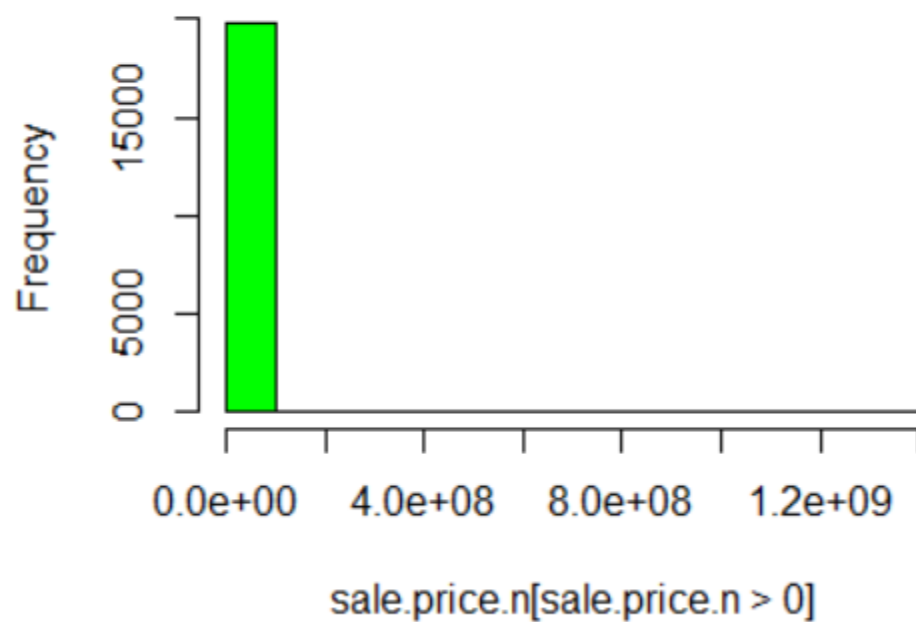
**Gross.sqft**



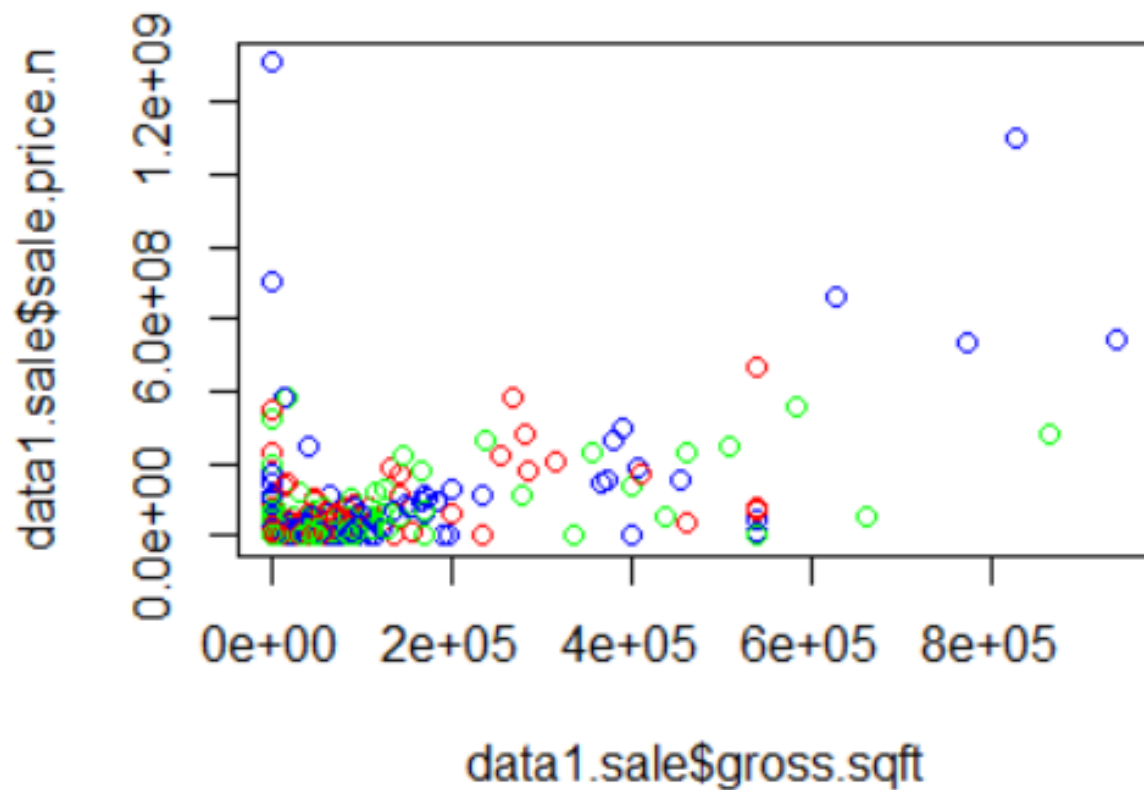The next scenario represents Sale.Price.n vs Frequency for one Borough of New York

## Sale.Price.n



The next scenario represents Sale.Price.n vs Frequency when **Sale.Price.n>0 for one Borough of New York**

# Sale.Price.n
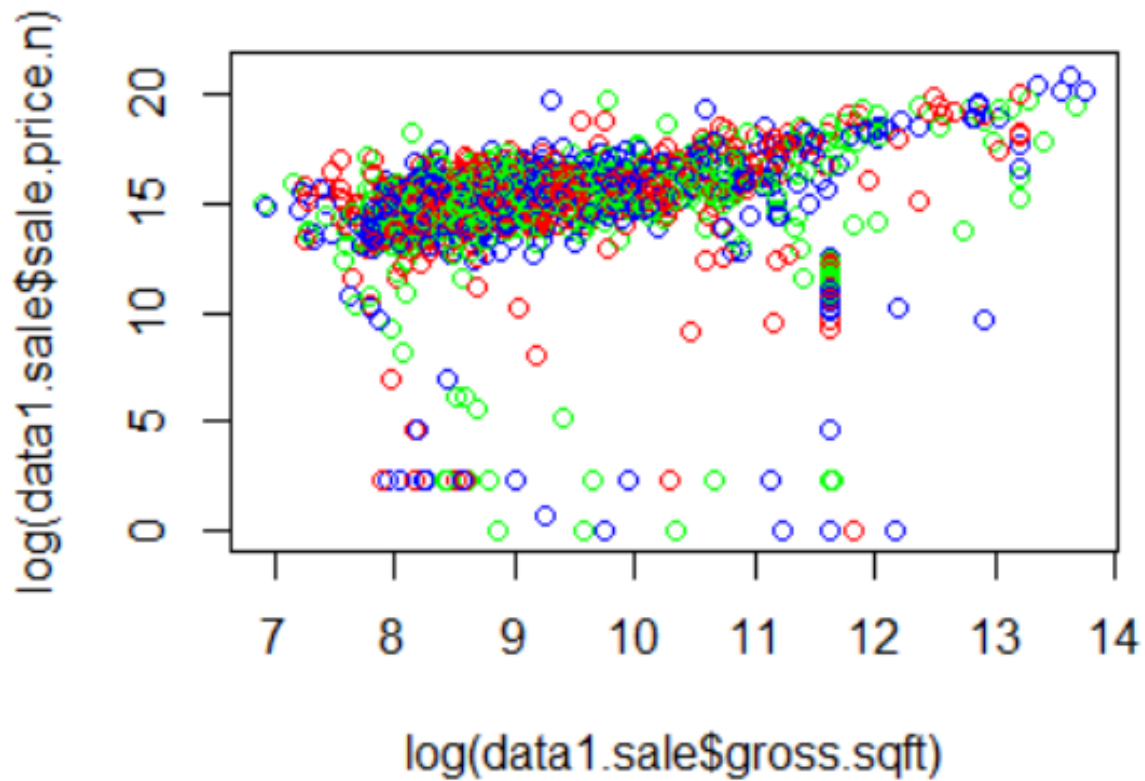


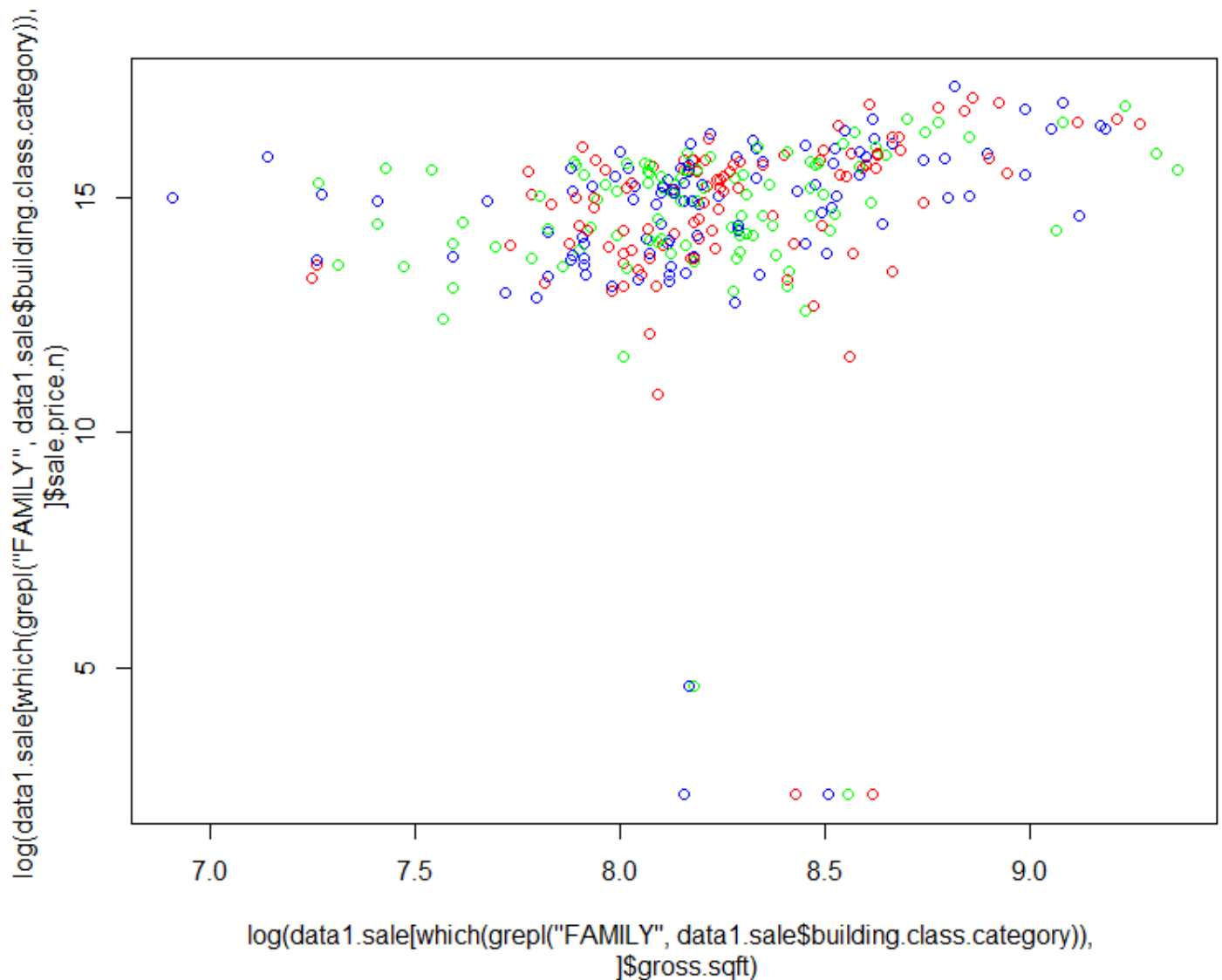Frequency vs sale.price.n[sale.price.n > 0]

Now in the below plot, it can be clearly seen that the density of data is very high near 0 square feet which is due to the presence of noise.
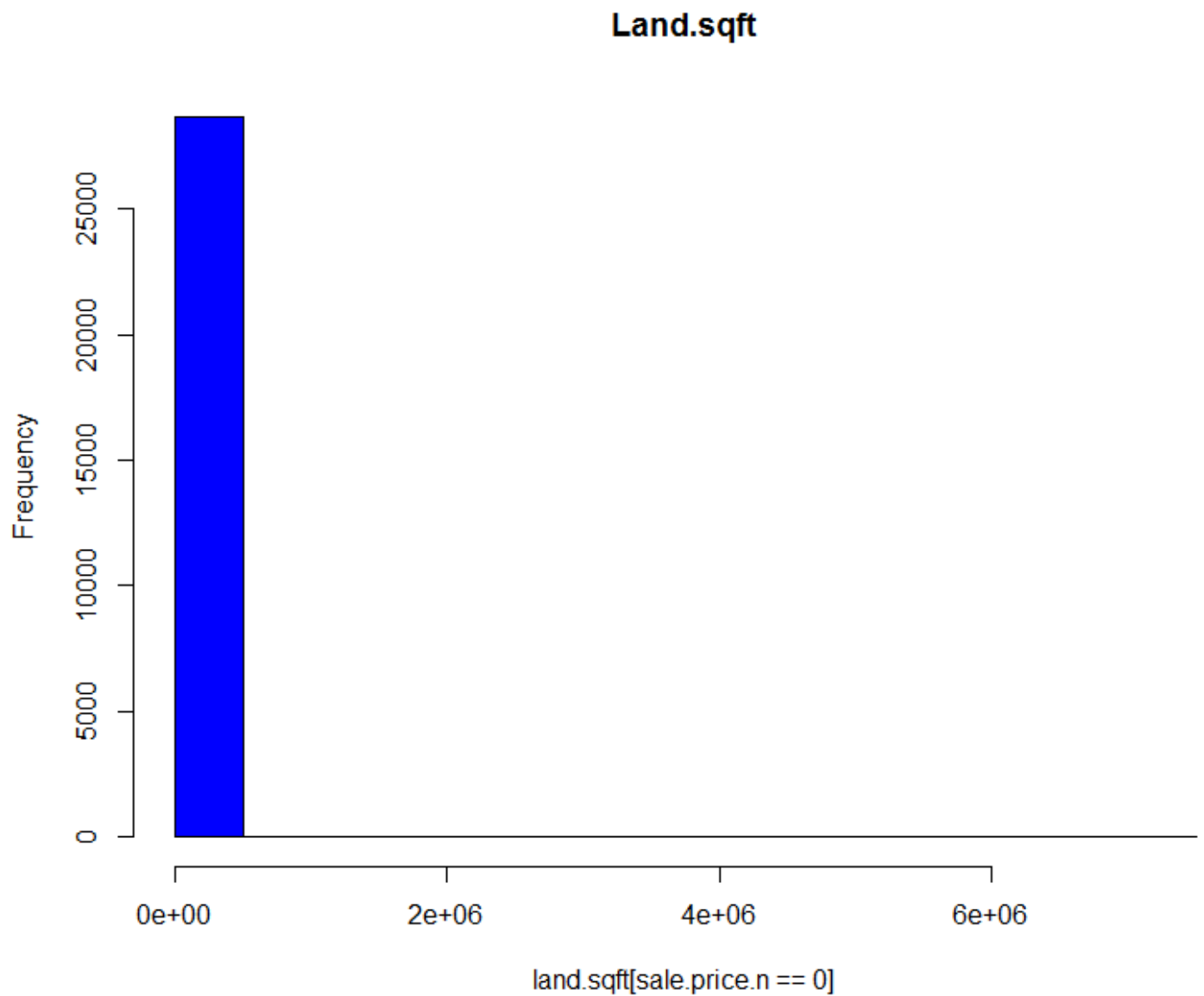
Let us try to plot the data after removing the induced noise. The below graph validates the removal of noise.
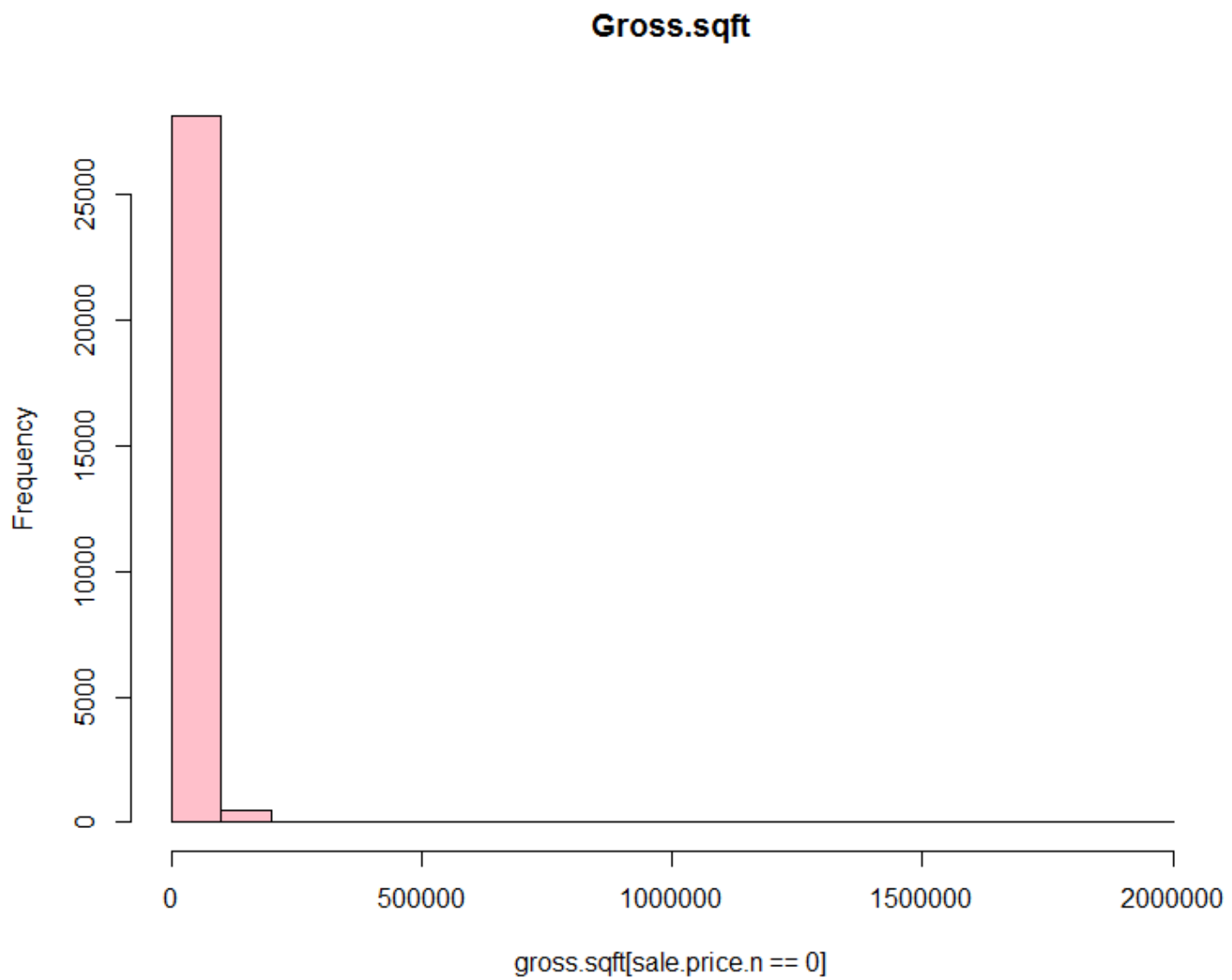
Now extending the above analysis to all the six boroughs of New York to have a combined analysis.
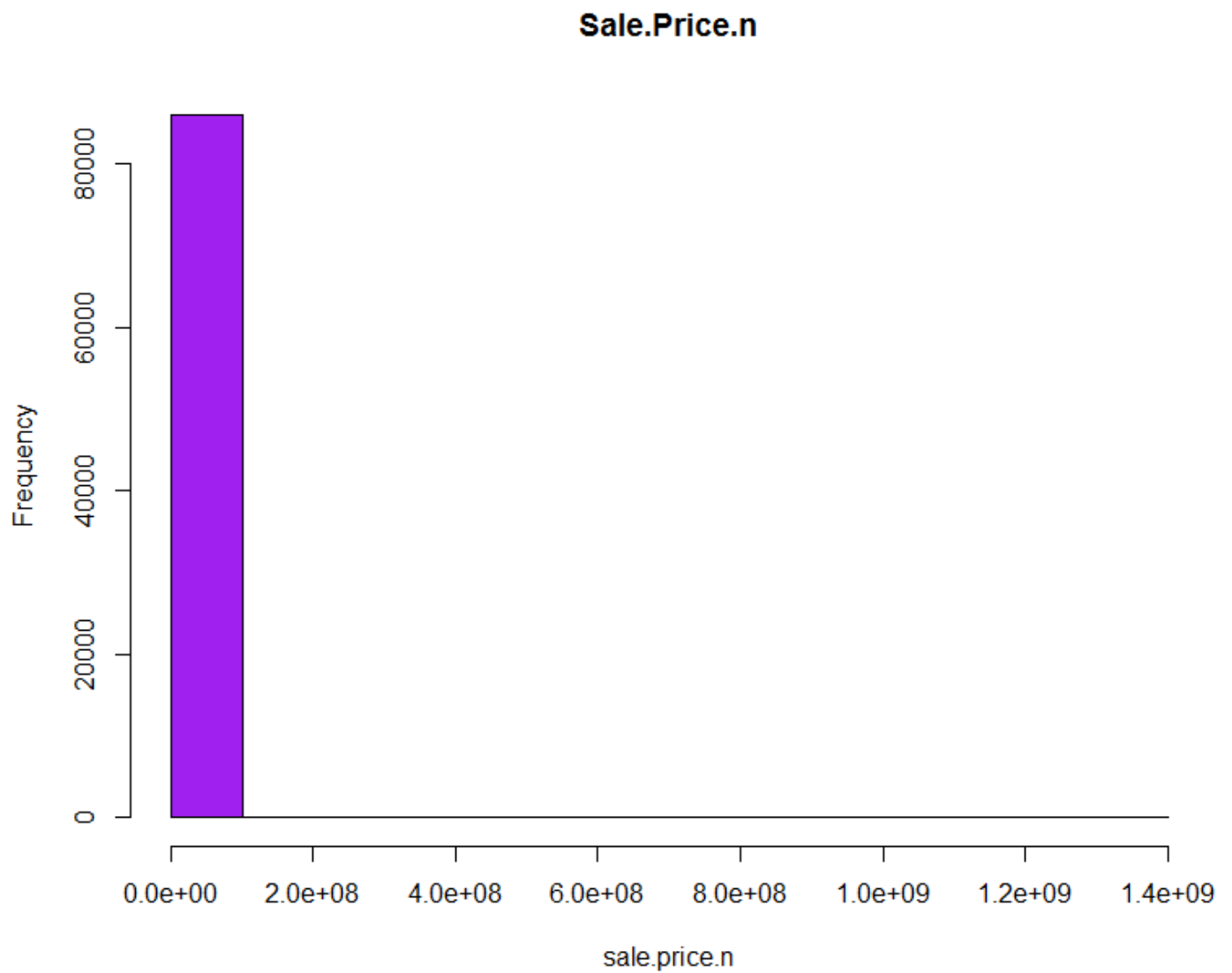
The first scenario represents Land.sqft vs Frequency when
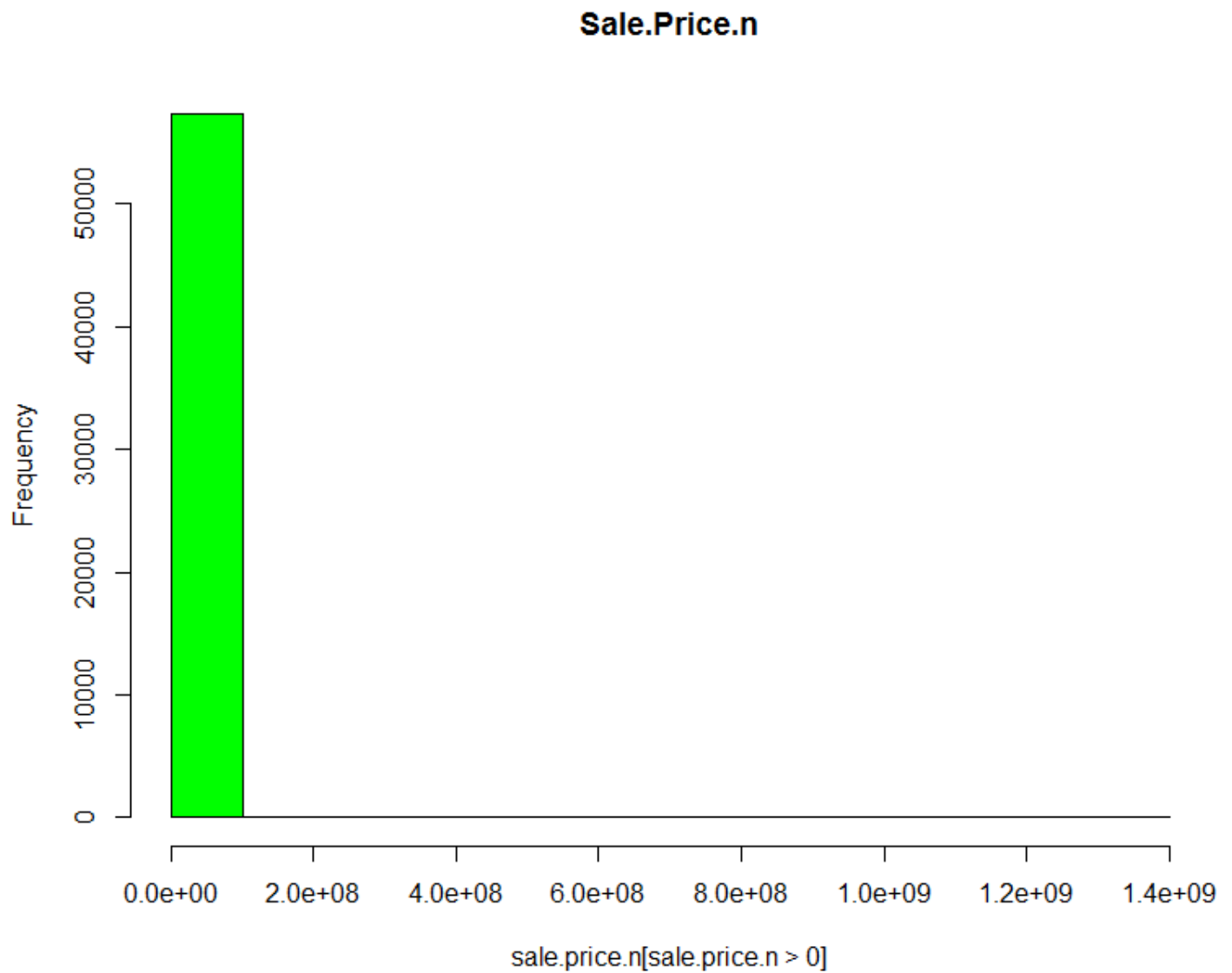**sale.price.n==0 for all 6 Boroughs of New York**

**Land.sqft**



The next scenario represents Gross.sqft vs Frequency when
sale.price.n==0 for all 6 Boroughs of New York
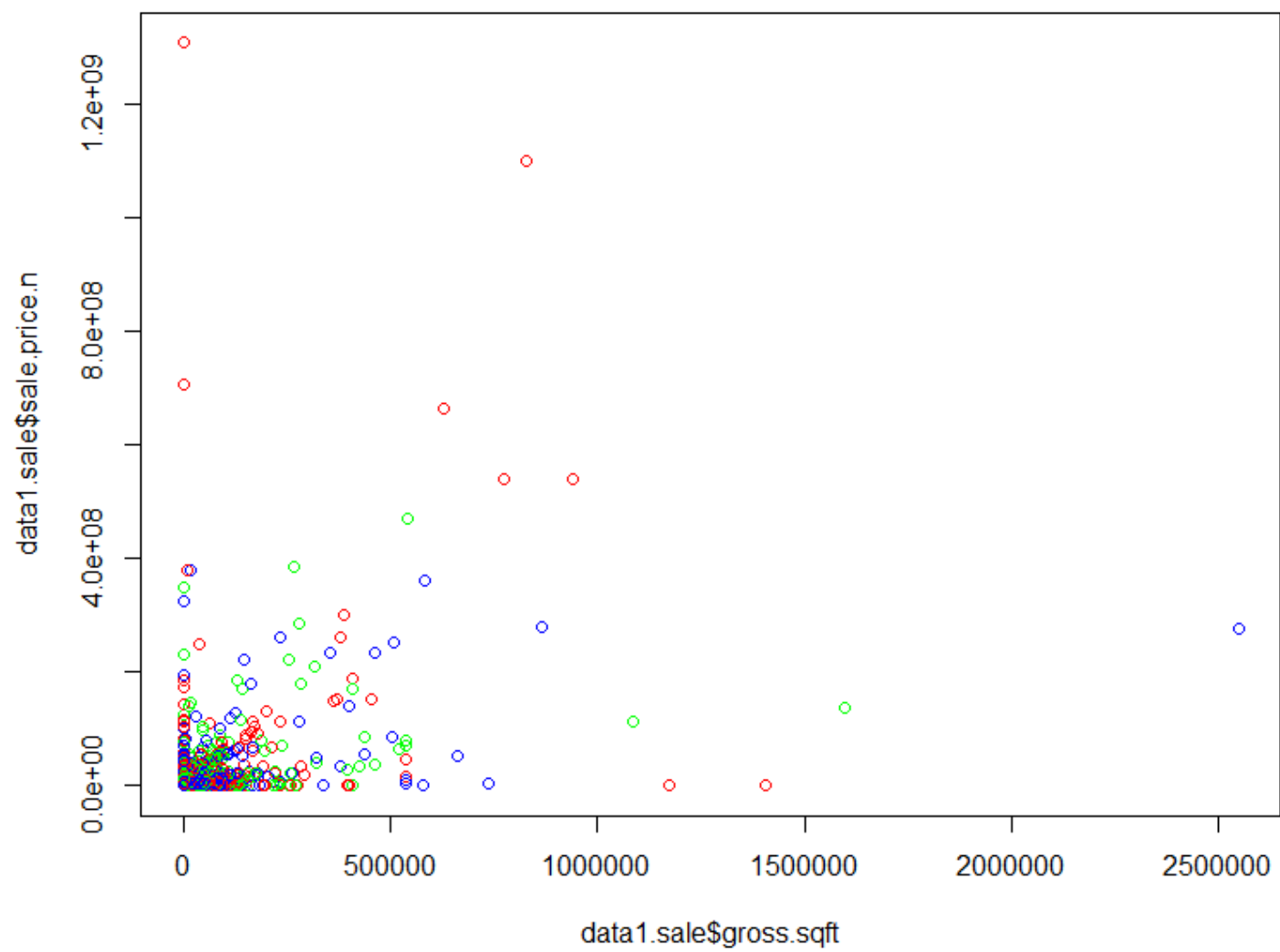
**Gross.sqft**



The next scenario represents Sale.Price.n vs Frequency for all 6 Boroughs
of New York

## Sale.Price.n



The next scenario represents Sale.Price.n vs Frequency when
Sale.Price.n>0 for all 6 Boroughs of New York

# Sale.Price.n



sale.price.n[sale.price.n > 0]

Now in the below plot for all the boroughs of New York, it can be clearly seen that the density of data is very high near 0 square feet which is due to the presence of noise.

After removing the erroneous noise the plot improves as shown below