

Problem 1 and 2

Matthew Leong

8/11/2020

Problem 1:

There were a number of things that sounded off about the “Excel guru”’s analysis but first, we wanted to check if there was an unacceptable amount of data lost.

```
## The original dataset originally had 7894 observations.
```

```
## The new cleaned dataset has 7679 observations.
```

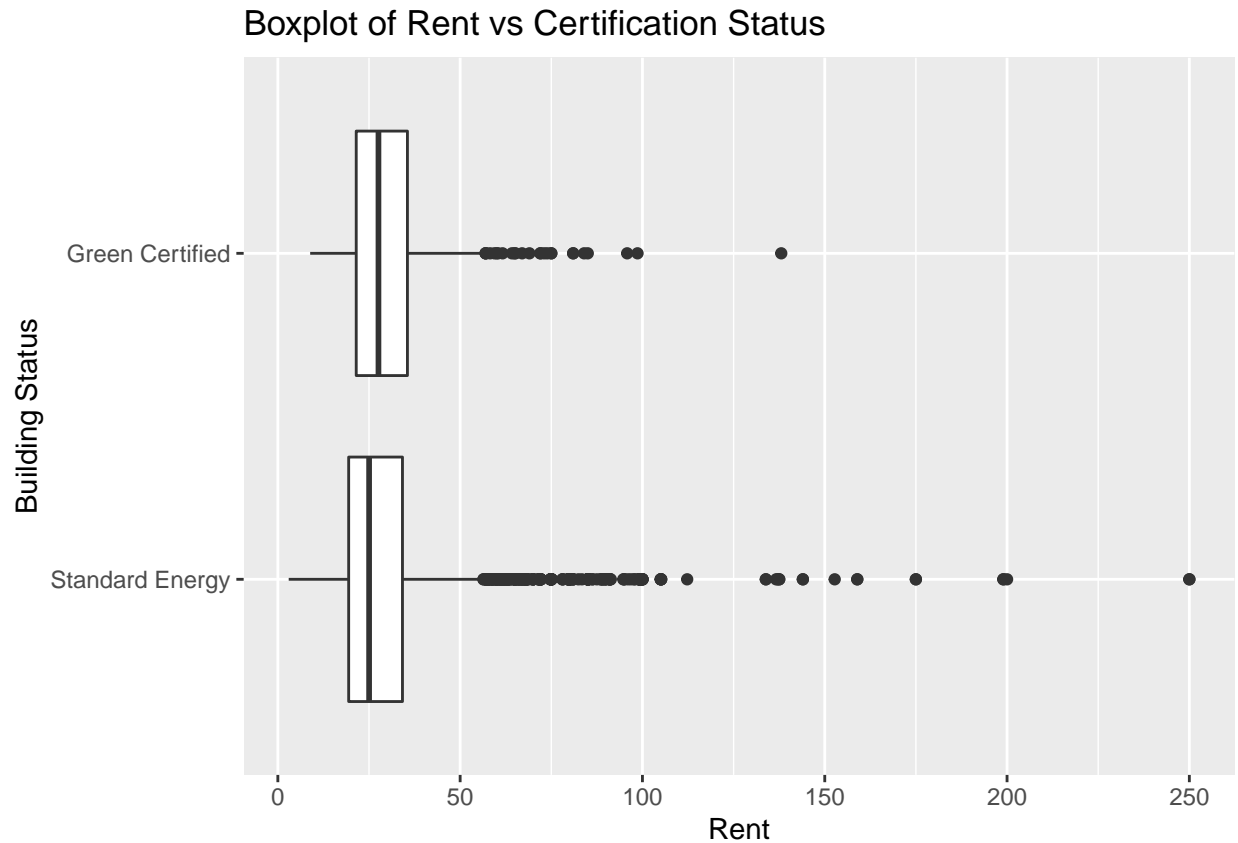
```
## They differ by 215 observations
```

```
## which is a percentage difference of 2.723588 percent.
```

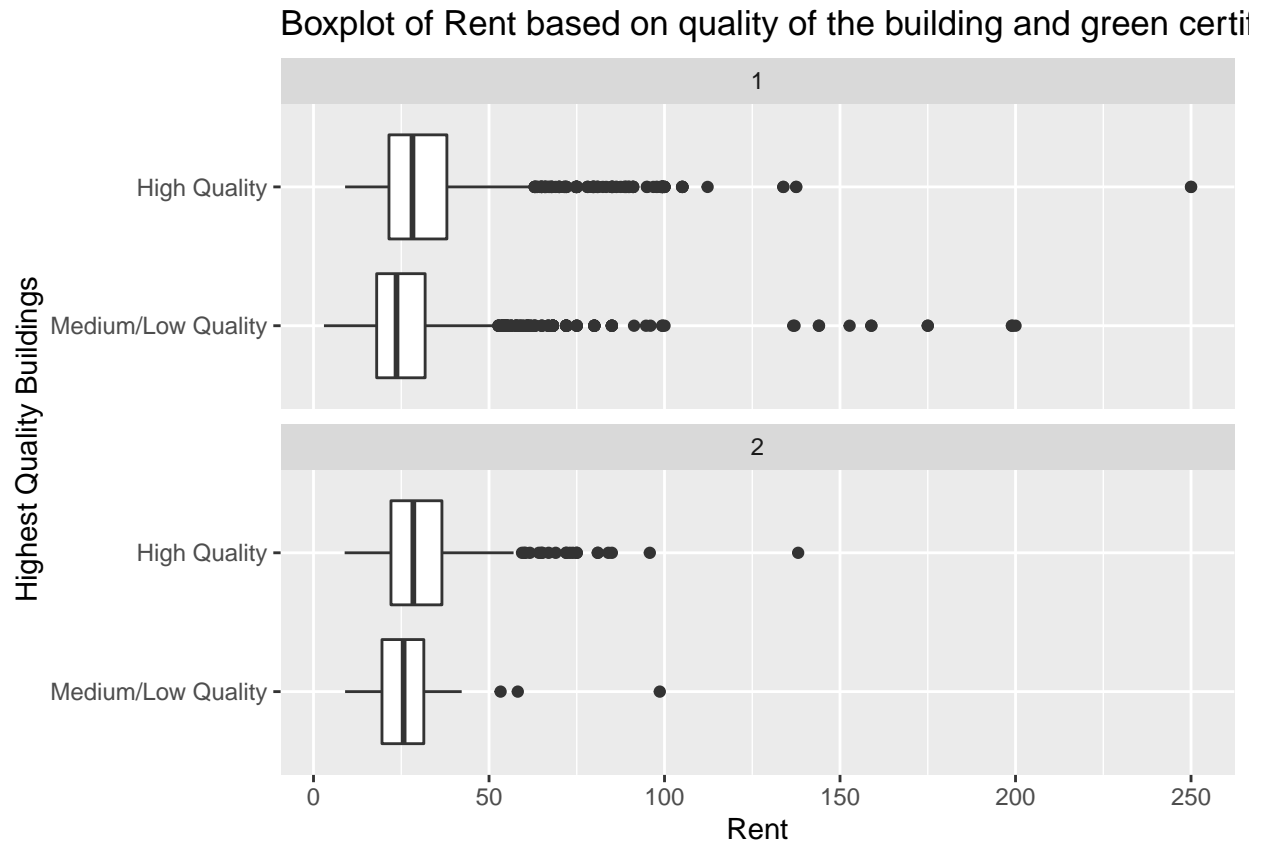
```
## There are 684 green certified buildings in the data set.
```

```
## There are 6995 non green certified buildings in the data set.
```

After running some code really quickly, we can see that the guru lost about 2.72% of the data which is honestly an acceptable amount. However, there is a key issue with this analysis in general. He assumes a causal relationship or strong positive correlation between rent and green rating without considering the other variables. To demonstrate this issue, let’s look at the boxplot obtained when the other variables are not considered.

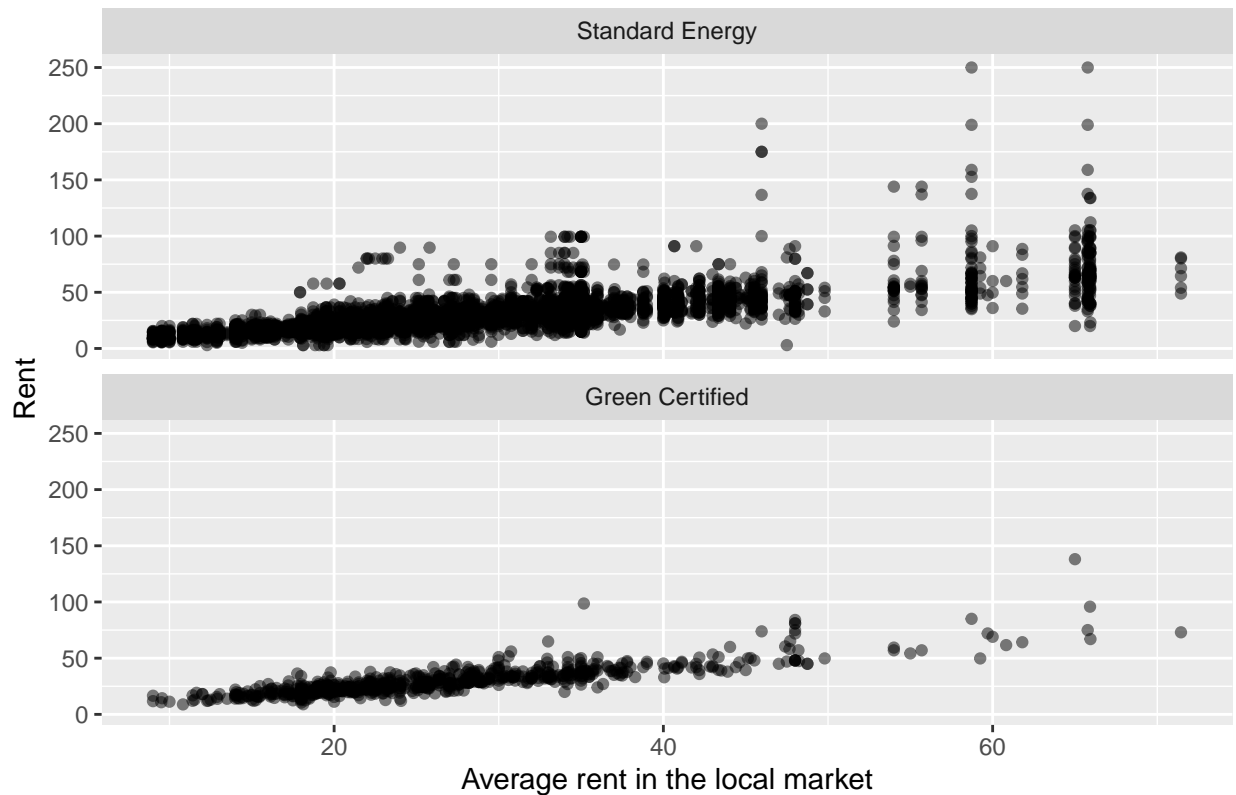


The box plot demonstrates the distribution of the split that the guru did. As we can see from the box plot, rent does not seem to differ that much between green certified and standard energy. The guru ignores the overall distribution though and bases his analysis on the median. The median here is misleading as while it is higher on this distribution, it will not necessarily hold true when we add in additional factors. For instance, let us see what would happen when we add in whether or not a building is high quality.



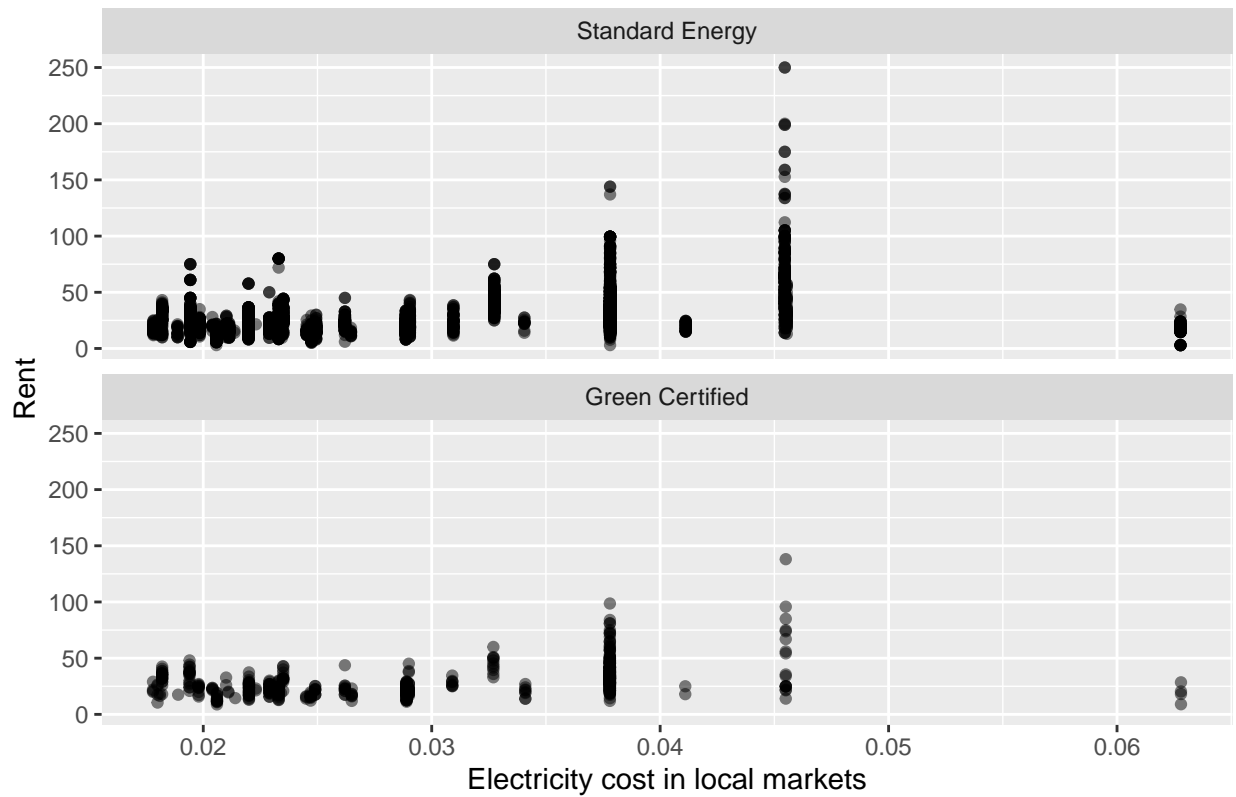
Naturally, as one would expect, high class buildings have a higher rent associated with them. Basing our analysis off of median, we can see in the box plot that while green certified medium/low quality buildings have a higher rent, standard energy (non certified green buildings) and green certified high class buildings have about the same median rent. This is essentially a confounding variable problem which entails how another different variable could affect rent and/or green certification status.

Scatterplot of average local market rent vs rent sorted by green certification



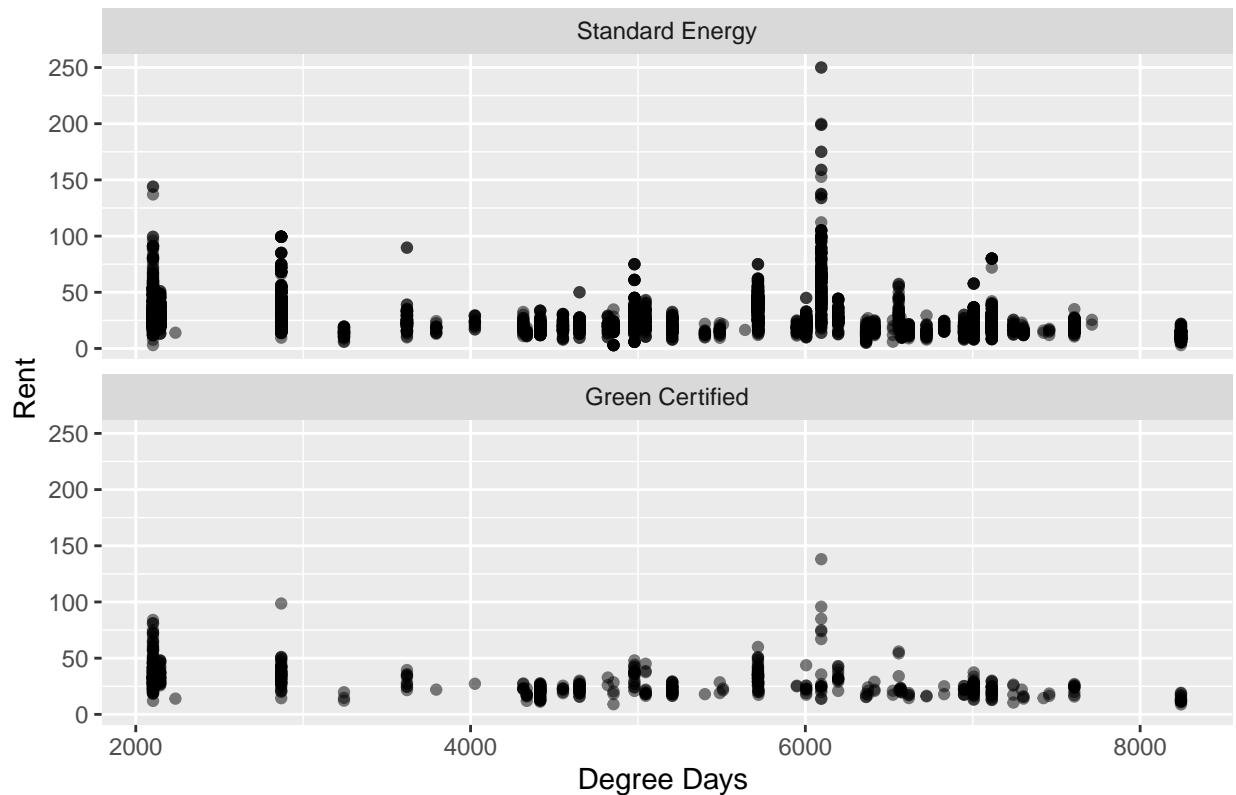
Another such confounding variable left out by the guru is the average rent in the local market. This variable has a key relationship with rent charged in that they essentially deal with the same thing. If green certified buildings netted a higher rent, this would be reflected in the graph when controlling for average rent in the local market. However according to the graph, it seems that green certified buildings perform worse but are more consistent overall which can be explained by the fact that there are a lot fewer green certified buildings than non green certified buildings. Additionally, the distributions are fairly similar implying that green certification does not actually affect the rent price. Let's look at a couple more potential confounding variables and see if the distributions differ.

Rent vs Electricity cost based on green certification



Again, when controlling for electricity costs, we can see that the distributions do not differ significantly further supporting the case that green certification does not necessarily mean higher rent. In fact when looking at the distribution, the rent seems lower in the green certified buildings. One possible explanation for this phenomenon is that green buildings have lower fixed costs than standard energy buildings which allow them to charge less rent based on energy usage.

Rent vs Degree Days based on green certification



Degree days in this data set are essentially measures of demand for energy where higher values mean greater demand. As this is another energy variable, we would expect to see and do actually see a similar distribution in that green certified buildings have lower rent than standard energy buildings.

In conclusion, the confounding variable problem makes the guru's analysis invalid. When controlling for such variables, we can see evidence that points towards green certified buildings have equal or even lower rent than non green certified ones. We can conclude that going green does not necessarily mean higher rent. TO further verify this conclusion, we recommend a linear regression or some other similar predictive modeling algorithm to control for all the variables and see which ones would actually lead to higher rent. Additionally, we would recommend running such a model on the green rating as well to determine other potential factors that are actually related to green buildings. As it stands now though, we do not recommend building a green building as it makes no difference being able to charge a higher rent to recoup the costs of certification.

Problem 2: