

1 **MA 354: Data Analysis I – Fall 2019**

2 **Exam 2:**

3 **Instructions:**

- 4 • You have 45 minutes to complete the conceptual part of this exam.
- 5 • The data analysis is take home and due 12/06 by 11:59p.
- 6 • Take a deep breath. You're going to do well and the worst case is that it will be productive.

7 **R/L^AT_EX Sweave notes – this should be all that you need.**

- 8 • To run R and print the output.

```
<<>>=
      #Rcode goes here
      #Output is automatically printed in the .pdf
@
```

9 **Remark:** All R chunks must have no spaces preceding the <<>>= or @ syntax.

- 10 • Provide R code for plot and place the plot into our document.

```
<<plotName,eval=FALSE>>=
      #Rcode for plot
      #We will call this later so make sure it has a unique name
@
\begin{figure}[H]
  \centering
  <<fig=TRUE,echo=FALSE>>=
  library("graphics")
  <<plotName>>
  @
  \caption{Some information about our plot} \label{Fig:plot1}
\end{figure}
```

11 You can then reference a graph in latex using `\ref{Fig:plot1}`.

12 **Remark:** All R chunks must have no spaces preceding the <<>>= or @ syntax.

- 13 • If you wanted a one line equation that is centered like this,

$$\hat{y}_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \epsilon$$

14 you can use this L^AT_EX.

```
\[\widehat{y}_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \epsilon\]
```

- If you wanted a multiple line equation that is centered like this,

$$\begin{aligned} f_X(x) &= 90x^8(1-x) \\ &= 90x^8 - 90x^9 \end{aligned}$$

15 you can use this L^AT_EX.

```
\begin{align*}
f_X(x) &= 90 x^8(1-x)\\
&= 90x^8 - 90x^9\\
\end{align*}
```

Help: You can ask for information about any of the following functions that we've used by asking **R**. For example, if I wanted help with the `lm()` function I would run `?lm()` in the **R** console. Note that if you're asking a question about a function, its library must be loaded.

20	• Stock R functions	60	• ggplot2 Package Plotting	97	• boot Package
21	– which()	61	– ggplot()	98	– boot()
22	– subset()	62	– geom_bar()	99	– boot.ci()
23	– summary()	63	– coord_polar()	100	• BSDA Package
24	– names()	64	– geom_hline()	101	– SIGN.test()
25	– cumsum()	65	– geom_text()	102	• simpleboot Package
26	– apply()	66	– geom_histogram()	103	– two.boot()
27	– lapply()	67	– geom_density()	104	• RVAideMemoire Package
28	– sapply()	68	– geom_freqpoly()	105	– mood.medtest()
29	– tapply()	69	– geom_boxplot()	106	– cramer.test()
30	– table()	70	– geom_jitter()	107	• rcompanion Package
31	– prop.table()	71	– geom_violin()	108	– pairwiseMedianTest()
32	– pie()	72	– geom_point()	109	– cldList()
33	– barplot()	73	– geom_line()	110	– phi()
34	– hist()	74	– facet_grid()	111	– cramerV()
35	– density()	75	– coord_flip()	112	• multcomp Package
36	– boxplot()	76	– theme_bw()	113	– glht()
37	– lines()	77	– xlab()	114	– cld()
38	– points()	78	– ylab()	115	• FSA Package
39	– jitter()	79	– ggtitle()	116	– dunnTest()
40	– legend()	80	• Probability Distribution	117	• DescTools Package
41	– optim()	81	– dbinom()	118	– StuartTauC()
42	– prop.test()	82	– dhyper()		
43	– t.test()	83	– dnbinom()		
44	– var.test()	84	– dpois()		
45	– aov()	85	– dunif()		
46	– lm()	86	– dnorm()		
47	– anova()	87	– dlnorm()		
48	– tukeyHSD()	88	– dchisq()		
49	– p.adjust()	89	– dt()		
50	– fisher.test()	90	– df()		
51	– chisq.test()	91	• gridExtra Package		
52	– cor()	92	– grid.arrange()		
53	– cor.test()	93	• qqplotr Package		
54	• stringr Package	94	– stat_qq_band()		
55	– str_split()	95	– stat_qq_line()		
56	• extraDistr Package	96	– stat_qq_point()		
57	– dmnom()				
58	• nleqslv Package				
59	– nleqslv()				

- Bernoulli Distribution

$$f_X(x|p) = p^x(1-p)^{1-x} I(x \in \{0, 1\}) \quad \text{[PMF]}$$

$$E(X) = p \quad \text{[Expected Value]}$$

$$\text{var}(X) = p(1-p) \quad \text{[Variance]}$$

- Binomial Distribution

$$f_X(x|n, p) = \binom{n}{x} p^x (1-p)^{n-x} I(x \in \{0, 1, \dots, n\}) \quad \text{[PMF]}$$

$$E(X) = np \quad \text{[Expected Value]}$$

$$\text{var}(X) = np(1-p) \quad \text{[Variance]}$$

- Hypergeometric Distribution

$$f_X(x|N, n, m, k) = \frac{\binom{m}{x} \binom{n}{k-x}}{\binom{N}{k}} I(x \in \mathcal{X}) \quad \text{[PMF]}$$

$$E(X) = \frac{km}{m+n} \quad \text{[Expected Value]}$$

$$\text{var}(X) = \frac{km}{m+n} \frac{-n}{m+n} \frac{m+n-k}{m+n-1} \quad \text{[Variance]}$$

- Negative Binomial Distribution

$$f_X(x|n, p) = \binom{n+x-1}{x} p^n (1-p)^x I(x \in \{0, 1, \dots\}) \quad \text{[PMF]}$$

$$E(X) = \frac{n(1-p)}{p} \quad \text{[Expected Value]}$$

$$\text{var}(X) = \frac{n(1-p)}{p^2} \quad \text{[Variance]}$$

- Poisson Distribution

$$f_X(x|\lambda) = \frac{\lambda^x e^{-\lambda}}{x!} I(x \in \{0, 1, \dots\}) \quad \text{[PMF]}$$

$$E(X) = \lambda$$

$$\text{var}(X) = \lambda$$

- Uniform Distribution

$$f_X(x|a, b) = \frac{1}{b-a} I(x \in [a, b]) \quad \text{[PDF]}$$

$$E(X) = \frac{a+b}{2} \quad \text{[Expected Value]}$$

$$\text{var}(X) = \frac{(b-a)^2}{12} \quad \text{[Variance]}$$

- Gaussian Distribution

$$f_X(x|\mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{\frac{-(x-\mu)^2}{2\sigma^2}} I(x \in \mathbb{R}) \quad \text{[PDF]}$$

$$E(X) = \mu \quad \text{[Expected Value]}$$

$$\text{var}(X) = \sigma^2 \quad \text{[Variance]}$$

- Log-Normal Distribution

$$f_X(x|\mu, \sigma) = \frac{1}{x\sigma\sqrt{2\pi}} e^{\frac{(\ln(x)-\mu)^2}{2\sigma^2}} I(x \in (0, \infty)) \quad \text{[PDF]}$$

$$E(X) = e^{\mu+\sigma^2/2} \quad \text{[Expected Value]}$$

$$\text{var}(X) = e^{2\mu+\sigma^2} e^{\sigma^2-1} \quad \text{[Variance]}$$

- Chi-squared Distribution

$$f_X(x) = \frac{1}{\Gamma(\frac{v}{2}) 2^{v/2}} x^{\frac{v}{2}-1} e^{-\frac{x}{2}} \quad \text{[PDF]}$$

$$E(X) = v \quad \text{[Expected Value]}$$

$$\text{var}(X) = 2v \quad \text{[Variance]}$$

- Student T distribution

$$f_T(t) = \frac{\Gamma(\frac{v+1}{2})}{\sqrt{\pi} \Gamma(\frac{v}{2})} \left(1 + \frac{t^2}{2}\right)^{-(v+1)/2} \quad \text{[PDF]}$$

$$E(X) = 0 \quad \text{[Expected Value for } v > 1]$$

$$\text{var}(X) = \frac{v}{v-2} \quad \text{[Variance for } v > 2]$$

- F distribution

$$f_W(w) = \frac{\Gamma(\frac{u+v}{2})}{\Gamma(\frac{u}{2})\Gamma(\frac{v}{2})} \left(\frac{u}{v}\right)^{u/2} \frac{w^{\frac{u}{2}-1}}{[1 + (\frac{u}{v})w]^{(u+v)/2}} I(w > 0) \quad \text{[PDF]}$$

$$E(W) = \frac{v}{v-2} \quad \text{([Expected Value for } v > 2])}$$

$$\text{var}(W) = \left(\frac{u-2}{u}\right) \left(\frac{v}{v+2}\right) \quad \text{([Variance])}$$

1 In-exam Portion:

Part I (30 points)

In Part I, I'm simply evaluating your engagement with the material. If you've worked through the material, there should be clear distinctions to make. I have provided as much room as I think is necessary to answer these questions. Take a minute to think or do some scratch work – your answer should fit in the space provided, only keep the important distinctions. I do not expect you to recite the formulas but to explain the procedures, their hypotheses, conclusions and/or their differences.

Submit your exam by emailing the following to wcipolli@colgate.edu

1. A LASTNAME_FIRSTNAME.pdf file just containing your answers (pages 5-6)

2 Out-of-exam Portion:

Part II (70 points)

In Part II, you're completing a data analysis. In this analysis you should provide numerical and graphical summaries that provide information for the researcher related to their research question.

Submit your exam by emailing the following to wcipolli@colgate.edu

1. A LASTNAME_FIRSTNAME.pdf of your final draft data analysis
2. Your .Rnw file.

Part III (Optional with likely increased score)

Shortly after the exam, you will receive an email to anonymously review two exams. You should review their data analysis for completeness, correctness, and communication. You will type up **constructive** notes to make the response better. The idea is to provide guidance for what's needed for the full data analysis to be effectively communicated to where you can understand the logic and the conclusions made about the data analysis. The format is discussed below.

- Write a paragraph about the general pros and cons of the paper you're reviewing. There is something good about every paper – find it and discuss that part. Also provide, in broad strokes a **constructive** critique of the response.
- Provide a list of major issues.
- Provide a list of minor issue.
- Provide a list of typographical errors you've found while reading.
- Ensure to provide specific line item comments where applicable.

Part VI (Optional with likely increased score)

After you receive comments about your work, revisit your analysis from the exam. Write a final draft of your analysis and provide responses to reviewer comments.

- Write a revision of your original solution which incorporates comments made in the reviews you've received.
- Provide a list of responses to specific line item comments; e.g.,
 - On page 1, line 2, you appear to interpret the statistics incorrectly.
Response: This was actually done correctly because I was treating the predictor as categorical and not continuous. I've added a sentence to make this distinction clear when fitting the model.
 - On page 2, line 4, you're missing a period at the end of the sentence.
Response: Thank you for pointing this out; I've added the missing period to the end of the sentence.

Part II

Suburban areas play an integral role in the development of sustainable cities; however, developers often do not consider sustainability in the construction of subdivisions and the subsequent adoption of homeowner's association (HOA) covenants. While there are multiple actions homeowners can take to contribute to personal sustainability on the plot-by-plot level, these actions are not always adopted or supported by greater neighborhood norms.

The current literature provides assessments of individual sustainability indicators at the homeowner and neighborhood level as well as multi-indicator sustainability assessments of cities and larger metropolitan areas but lacks such multi-indicator analyses at the homeowner and neighborhood level. This study assesses the relationship among multiple sustainability indicators of homeowner behavior including recycling habits, lawn care, tree planting, and home gardening, and compares these behaviors between neighborhoods with HOAs and those without. Data metrics were collected from twelve neighborhoods in Greenville, South Carolina through on-site observation, analysis of Google Earth images, and qualitative assessment of HOA covenants.

Use the data collected by the researcher to extract information required for their study. The data consist of 1,616 observations of homes in Greenville, SC and the seven variables recorded for each. Basic descriptions of the variables and other important information can be found below.

- **Neighborhood:** This reports which of the 12 neighborhoods in Greenville, SC each observed home is located. There are no missing values for this variable.
- **Lot Number:** This reports the lot number of the homes. There are no missing values for this variable.
- **HOA:** This reports whether or not the homes are part of a homeowners' association (1 = yes; 2 = no). There are no missing values for this variable.
- **Recycle:** This reports the recycling status of the homes (1 = both recycling bin and trash bin present at the home; 2 = only a trash bin was present). There are 274 missing values for this variable. These missing values correspond to neighborhoods with no curb side pick up (i.e., Brownstone, Edgewood, Glastonbury, and Fox Springs).
- **Lawn Care:** This represents a likert variable on a scale of 1 to 4 (1 = excellent; 2 = good; 3 = poor; 4 = none). This measure maps onto how artificially managed the lawn is where 1 means the lawns were highly artificially managed with presumed regular chemical application and 4 means the lawns were naturally managed (i.e., not managed at all). There are no missing values for this variable.
- **Trees:** This represents the number of trees in the front yard. There are 14 missing values for this variable.
- **Garden:** This represents whether a garden was present in the front or back of the homes (1 = yes; 2 = no). There are no missing values for this variable.

The data can be accessed in R as follows.

```
> dat.HOA<-read.table("https://cipolli.com/students/data/Exam1Data.txt",  
+                      header=T,sep=",")
```

Copy and paste your analysis from Exam I below and complete a third draft. This will involve

1. Making changes related to comments on your revised Exam I.
2. Reading through your Exam I, noting where you can make improvements with what you've learned since then. You can now provide more than just visual "evidence" by using the appropriate tests to justify your visual insights.

I expect a highly polished final draft that is correct, communicated well and succinct.

Solution:

```
> dat.HOA<-read.table("https://cipolli.com/students/data/Exam1Data.txt",header=T,sep=",")
```

With this data we have multiple goals. The first is to help researchers compare sustainability behaviors between neighborhoods with HOAs and those without. The second is to perform a multi-indicator analyses at the homeowner and neighborhood level. After visualizing the data, we can perform hypothesis tests and find correlations for the data. However, before we can move forward with our analyses, we need to have a representative sample.

I want to first clean the data. It is first important to view the *Trees* variable. While it has missing values, we will not remove these values because these observations could have information for analyses for other variables. However, the code below would allow someone to remove the missing values if they chose to for their analysis.

```
> #creates a subset of the data
> #removes missing values from Trees
> dat.HOA.notree<-subset(dat.HOA,is.na(Trees)==FALSE,
+                         select=c(Neighborhood,Lot.Number,HOA,Recycle,LawnCare,Trees,Garden))
```

Recycle also has data that needs to be cleaned. Currently, neighborhoods without curbside pickup have the values *NA*. Since we may want to investigate the relationship between those homes with some curbside pickup and those with no curbside pickup at all, we can create a data set that includes these values as a characteristic we can analyze.

```
> for (i in 1:1616){ #for each obs in the dataset
+   if (is.na(dat.HOA$Recycle[i])==TRUE){ #if Recycling is a missing var
+     dat.HOA$Recycle[i]=0} #change missing var to zero
+ }
> #create second dataset with numeric values instead of factors
> dat.HOA1<-dat.HOA
```

After doing this data cleaning, we can confirm that the updated dataframe has recycling variables where observation that were once missing are now equal to zero, meaning that observation has no curbside pickup. Had we left these values as *NA*, these values would have been excluded from our analyses in R.

Now that we've cleaned our data, we can summarize it so that we can get a better understanding of the problems at hand.

```
> #summarize data
> summary(dat.HOA$HOA)
```

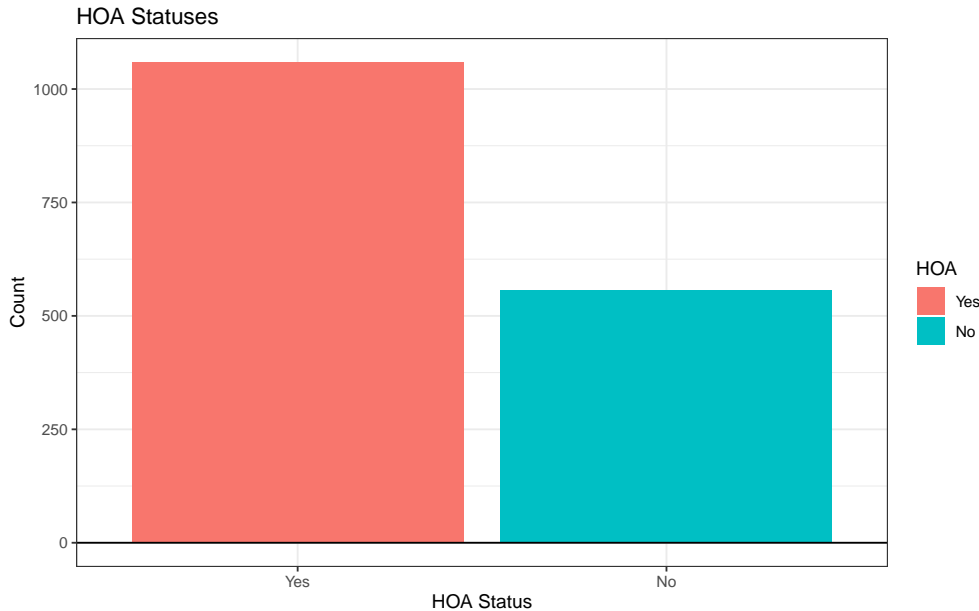
Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
1.000	1.000	1.000	1.345	2.000	2.000

The summary function is important because it allows us to have an understanding of how information is weighted in the data set before proceeding with the analyses. This is especially important for the HOA category, as we can see that because the median is 1 and the mean is 1.341, we can see that we have an unequal amount of HOA and non HOA homes.

Let's look at a preliminary plot comparing neighborhoods with HOAs and those without using ggplot2 Wickham (2016) and gridExtra Auguie (2017).

```
> library(ggplot2)
> library(gridExtra)
```

Summary of HOA Statuses



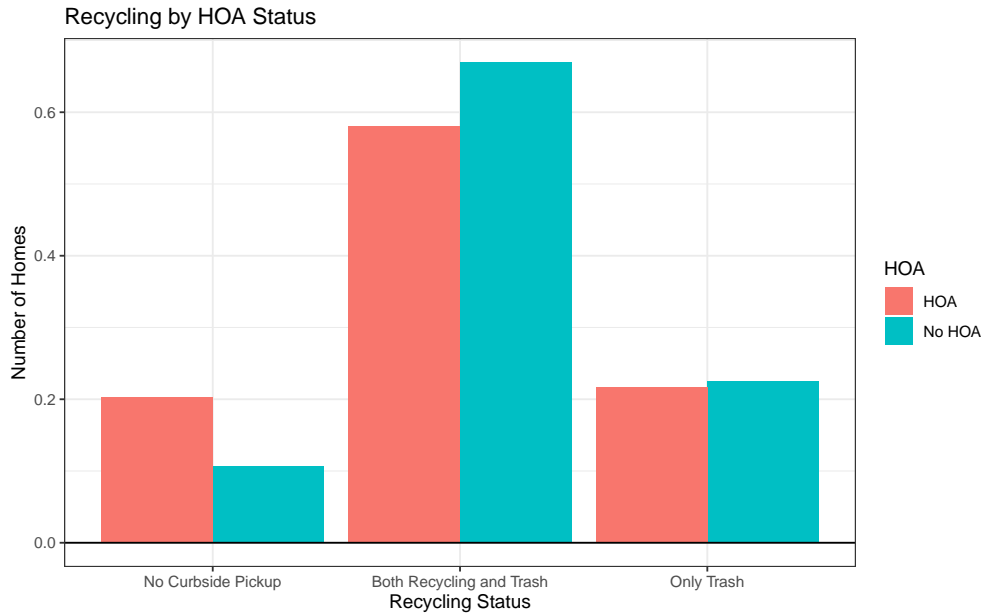
HOA Status	Frequency
Yes	0.659176
No	0.340824

There seems to be almost twice as many homes in our sample that are a part of an HOA association as homes that are not a part of an HOA association. To be precise, 65.91 percent of homes in our sample are HOA homes and 34.08 percent of homes in our sample are non HOA homes.

We should also look at these numbers in reference to different factors and compare breakdowns with homes that are part of an HOA and those that are not.

The next graphs provide a side-by-side comparisons of households in HOAs and households not in HOAs, representing the relationships between these households' HOA statuses and multiple sustainability indicators. These graphs all utilize relative frequencies to compare proportionally between HOA households and non HOA households. As was illustrated by the graph above, over 60 percent of the households in the sample are in an HOA while about 30 percent of households in the sample are not in an HOA. This will be important as we begin to analyze the homes by different factors.

Recycling and HOA Statuses

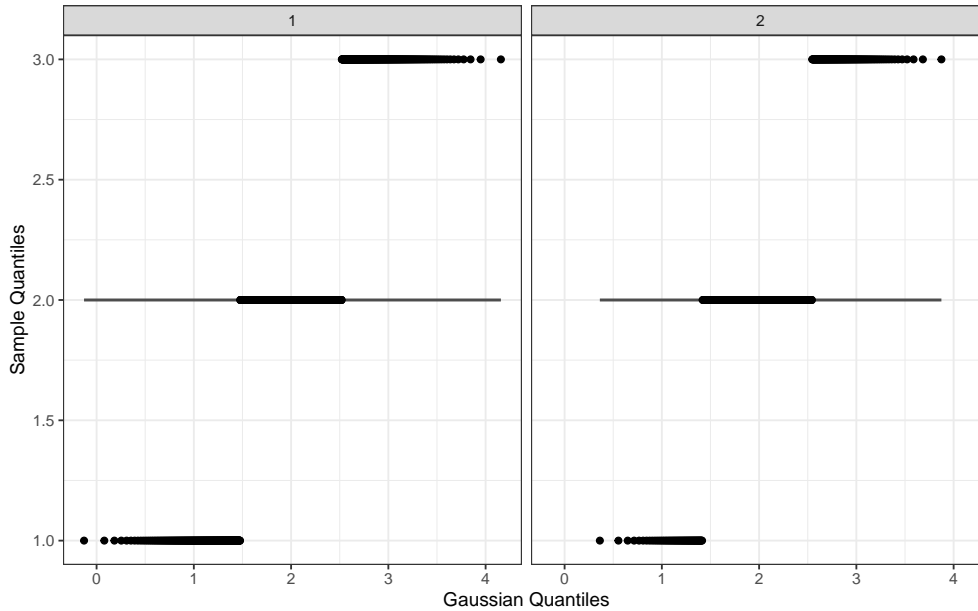


```
> #dat.HOA$HOA<-factor(dat.HOA$HOA, levels=c(1,2),labels = c("Yes", "No"))
> dat.HOA$Recycle<-factor(dat.HOA$Recycle, levels = c(0,1,2),
+                           labels = c("No Curbside Pickup",
+                                       "Both Recycling and Trash",
+                                       "Only Trash"))
> prop.table(table(dat.HOA$Recycle, dat.HOA$HOA),margin=2)*100
```

	Yes	No
No Curbside Pickup	20.30217	10.59246
Both Recycling and Trash	57.97923	66.96589
Only Trash	21.71860	22.44165

As we can see, regardless of HOA status homes the majority typically had both recycling and trash present. About 58 percent of HOA households and 67 percent of non HOA households had both trash and recycling pickup present. There were about the same proportion of homes that had only trash pickup for HOA and non HOA households (21.7 percent and 22.4 percent, respectively). About double the proportion of HOA households have no curbside pickup (20.3 percent) as non HOA households (10.6 percent). We see that there is a higher proportion of non HOA homes to HOA homes that have either trash pickup or both recycling and trash pick up rather than no curbside pickup. This suggests that non HOA homes may be better for the environment compared to HOA homes in regards to recycling and waste management.

We can perform a hypothesis test to determine whether or not there is a statistically significant difference between HOA homes and non HOA homes and across recycling statuses.



As we knew before and according to the above plot (Almeida et al., 2017), our data is not normally distributed, but rather is discretely distributed by recycling status. Therefore, for robustness, we perform a Mood's Median Test (Hervé, 2019). We want to generally identify if there is a difference across treatments, where the treatment is HOA versus non HOA. We look at these treatments specifically in the context of recycling status. We assume a representative sample. While these observations are not necessarily independent, since neighbors may base their recycling choices off of each other, for the sake of our study we will assume them to be so.

Our null hypothesis is that there is no significant difference between HOA and non HOA households with regards to recycling statuses; the population medians are all equal. Our alternative hypothesis is that there is a significant difference between HOA and non HOA households with regard to recycling statuses; at least one of the population medians is different. Since the plots of the data do not point towards normality, we will use a Mood's Median Test which is a nonparametric alternative to the ANOVA.

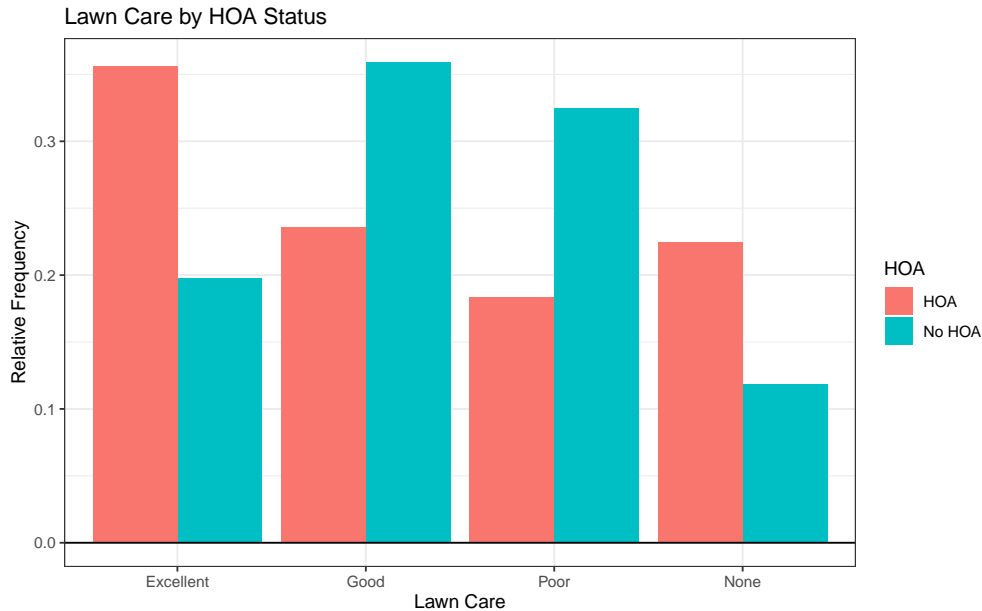
```
> #moods median test to determine if there are any significant differences
> #across treatments
> library(RVAideMemoire)
> mood.medtest(Recycle~HOA,data=dat.HOA1)
```

Mood's median test

```
data: Recycle by HOA
X-squared = 0.073138, df = 1, p-value = 0.7868
```

The Mood's median test is a chi-squared test that tests for differences across medians. Our chi-squared variable is 0.073138 and our p-value is 0.7868 which is greater than a significance level of 0.05. Since our p value is greater than 0.05, we fail to reject the null.

Lawn Care and HOA Statuses

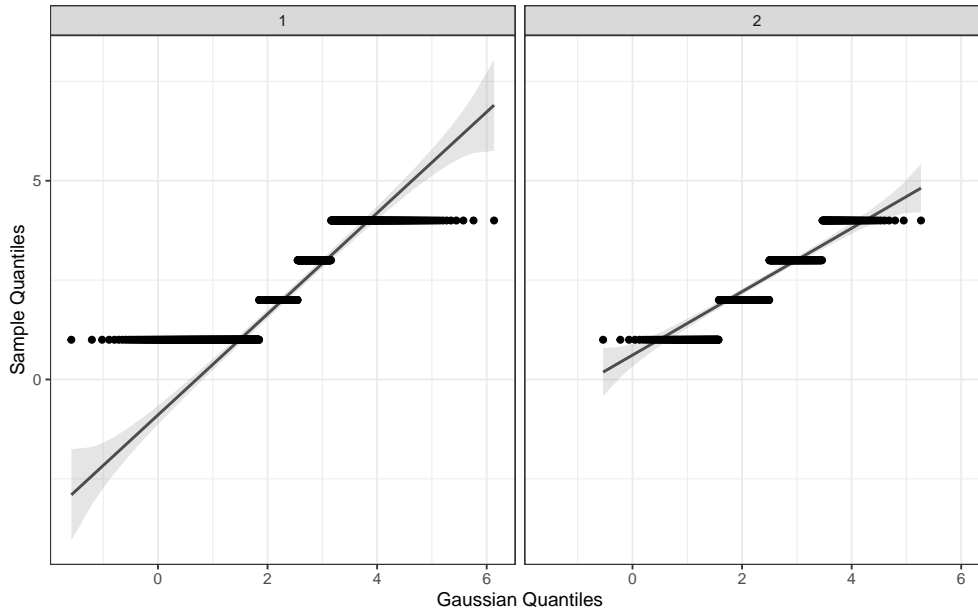


```
> dat.HOA$LawnCare<-factor(dat.HOA$LawnCare, levels = c(1,2,3,4),
+                           labels = c("Excellent",
+                                       "Good",
+                                       "Poor",
+                                       "None"))
> prop.table(table(dat.HOA$LawnCare, dat.HOA$HOA),margin=2)*100
```

	Yes	No
Excellent	35.59962	19.74865
Good	23.60718	35.90664
Poor	18.31917	32.49551
None	22.47403	11.84919

As you can see from the charts above, households in HOAs had a higher percent of excellent lawn care (35.6 percent) than households without HOAs (19.7 percent). From this, we can infer that these households with HOAs typically use more chemical application than those without HOAs, which is worse for the environment. We also see that the highest percentage of homes not in an HOA have good (35.9 percent) or poor (32.5 percent) lawn care relative to homes in HOAs (23.7 percent and 18.4 percent, respectively). How much more damaging non HOA homes are compared to HOA homes in regards to lawn care is hard to say from the graph alone, but we can see from the table that more HOA households had excellent or good lawn care (about 59 percent) compared to non HOA households (about 56 percent). Again though, lawn care status alone may not be enough to determine if HOA households or non HOA households are better or worse for the environment.

Again, now that we have a better picture of the data we can perform statistical analyses to determine whether or not Lawn Care differs significantly by HOA status.



As we knew before and according to the above plot (Almeida et al., 2017), our data is not normally distributed, but rather is discretely distributed by recycling status. Therefore, for robustness, we perform a Mood's Median Test. We want to generally identify if there is a difference across treatments, where the treatment is HOA versus non HOA. We look at these treatments specifically in the context of lawn care. We assume a representative sample. While these observations are not necessarily independent, since neighbors may base their lawn care choices off of each other and off of social pressure, for the sake of our study we will assume them to be so.

Our null hypothesis is that there is no significant difference between HOA and non HOA households with regards to lawn care; the population medians are all equal. Our alternative hypothesis is that there is a significant difference between HOA and non HOA households with regard to lawn care; at least one of the population medians is different. Since the plots of the data do not point towards normality, we will use a Mood's Median Test which is a nonparametric alternative to the ANOVA.

```
> #moods median test to determine if there are any significant differences
> #across treatments
> mood.medtest(LawnCare~HOA,data=dat.HOA1)
```

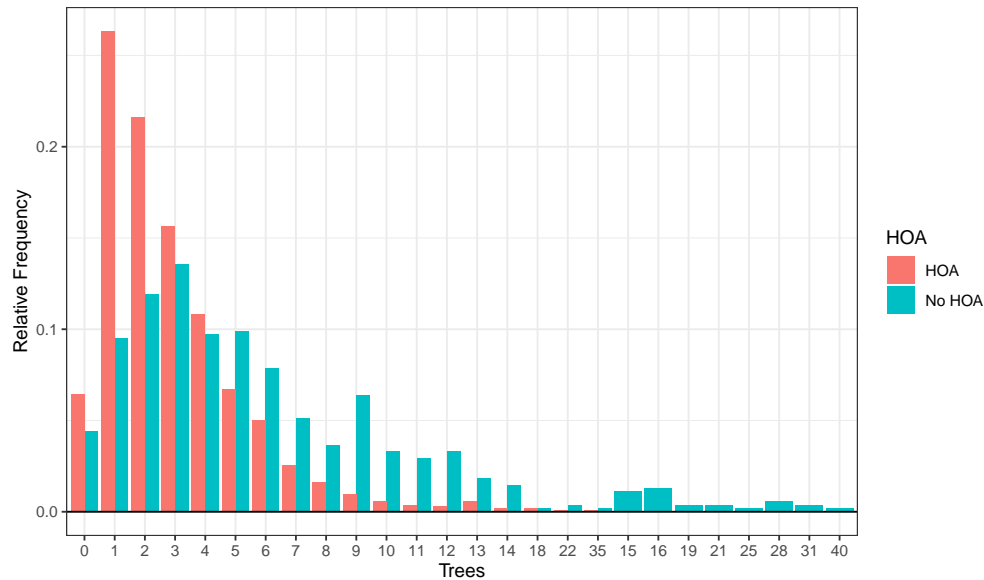
Mood's median test

```
data: LawnCare by HOA
X-squared = 1.7468, df = 1, p-value = 0.1863
```

The Mood's median test (Hervé, 2019) is a chi-squared test that tests for differences across medians. Our chi-squared variable is 1.7468 and our p-value is 0.1863 which is greater than a significance level of 0.05. Since our p value is greater than 0.05, we fail to reject the null.

Trees in Yard and HOA Statuses

Number of Trees by HOA Status



```
> prop.table(table(dat.HOA$Trees, dat.HOA$HOA),margin=2)*100
```

	Yes	No
0	6.43939394	4.39560440
1	26.32575758	9.52380952
2	21.59090909	11.90476190
3	15.62500000	13.55311355
4	10.79545455	9.70695971
5	6.72348485	9.89010989
6	5.01893939	7.87545788
7	2.55681818	5.12820513
8	1.60984848	3.66300366
9	0.94696970	6.41025641
10	0.56818182	3.29670330
11	0.37878788	2.93040293
12	0.28409091	3.29670330
13	0.56818182	1.83150183
14	0.18939394	1.46520147
15	0.00000000	1.09890110
16	0.00000000	1.28205128
18	0.18939394	0.18315018
19	0.00000000	0.36630037
21	0.00000000	0.36630037
22	0.09469697	0.36630037
25	0.00000000	0.18315018
28	0.00000000	0.54945055
31	0.00000000	0.36630037
35	0.09469697	0.18315018
40	0.00000000	0.18315018

```
> summary(dat.HOA$Trees)
```

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	NA's
0.000	1.000	3.000	3.952	5.000	40.000	14

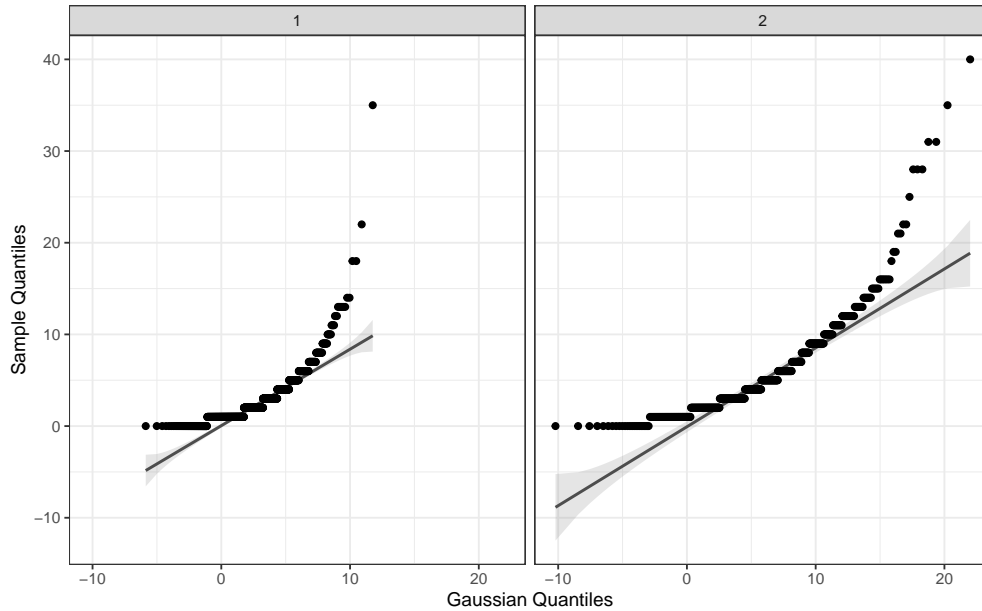
This figure illustrates the comparison between HOA homes and non HOA homes and the number of trees they have in their yard. It is important to note though that while most households in HOAs typically tend to have between 0 and 10 trees, it seems that households not in HOAs typically tend to have between 0 and 15 trees. In fact, about 62 percent of HOA homes have between 1 and 3 trees, where about half the proportion of non HOA homes (about 33 percent) have between 1 and 3 trees. The distribution of non HOA homes is less right skewed than the distribution of HOA homes, meaning a larger proportion of non HOA homes have more trees than HOA homes. From this, we can hypothesize that non HOA homes are more sustainable than HOA homes in the sense that non HOA homes have more trees planted than HOA homes, which is better for the environment.

We can try to determine the population mean and variances of trees in HOA and non HOA households by utilizing MOM and MLE estimation. MOM, or method of moments, simultaneously solves equations of moments where the expected value(s) is set equal to the sample value(s). The MLE, or Maximum Likelihood Estimator, assumes a distribution and then find the parameter value(s) that maximize the likelihood function of that distribution's PDF given the sample data we observed.

When looking at the plot of *Trees*, we see that the data follows an approximately poisson distribution. We therefore can utilize this knowledge to calculate the MOM and MLE estimates. This is where we utilize a dataset with the missing values of *Trees* removed, since these cannot pass through the *nleqslv* (Hasselman, 2018) package.

```
> #method of moments
> g<-function(x.data,theta) {#data, theta
+   beta = theta[1]
+   nu = theta[2]
+   EX = beta
+   varX = nu
+   m1 = EX - mean(x.data)
+   m2 = varX - var(x.data)
+   return(c(m1,m2))
+ }
> library(nleqslv) #load
> nleqslv(x=c(0,40), #best guess
+       fn=g, #function
+       x.data=dat.HOA.notree)
```

We can now perform a hypothesis test to determine numerically if the number of trees in a yard varies across treatments of HOA and non HOA homes. First, we must test for normality.



As we can see from the plots, the data is not normally distributed, but follows a parabolic shape. Therefore, for robustness, we perform a Mood's Median Test. We want to generally identify if there is a difference across treatments, where the treatment is HOA versus non HOA. We look at these treatments specifically in the context of number of trees. We assume a representative sample. While these observations are not necessarily independent, since neighbors may base their recycling choices off of each other, for the sake of our study we will assume them to be so.

Our null hypothesis is that there is no significant difference between HOA and non HOA households with regards to recycling statuses; the population medians are all equal. Our alternative hypothesis is that there is a significant difference between HOA and non HOA households with regard to recycling statuses; at least one of the population medians is different. Since the plots of the data do not point towards normality, we will use a Mood's Median Test which is a nonparametric alternative to the ANOVA.

```
> #moods median test to determine if there are any significant differences
> #across treatments
> mood.medtest(Trees~HOA,data=dat.HOA1)
```

Mood's median test

```
data: Trees by HOA
X-squared = 138.67, df = 1, p-value < 2.2e-16
```

The Mood's median test is a chi-squared test that tests for differences across medians. Our chi-squared variable is 138.67 and our p-value is essentially zero. Since our p value is less than 0.05, we reject the null.

We can perform a post-hoc test to adjust the p-value according to the number of groups. For Mood's Median Tests, we perform a Pairwise Median Test.

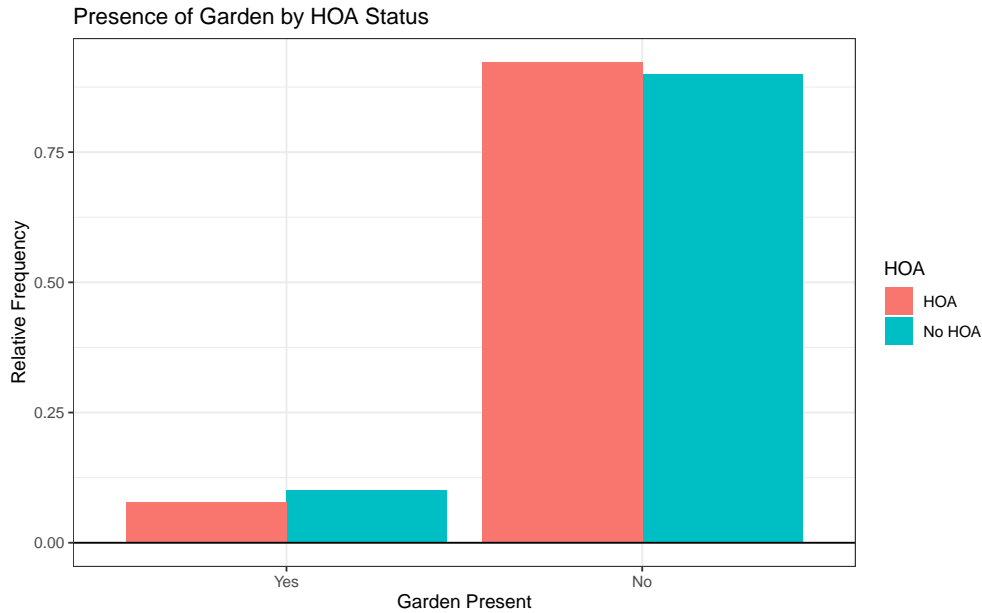
```
> #pairwise to find particular difference
> library(rcompanion)
> PTBH<-pairwiseMedianTest(Trees~HOA,
+                           data = dat.HOA1,
+                           method = "BH")

> cldList(p.adjust ~ Comparison,
+         data = PTBH,
+         threshold = 0.05)
```

	Group	Letter	MonoLetter
1	1	a	a
2	2	b	b

342 I utilize the Benjamini Hochberg approach (Mangiafico, 2019) to adjusting p-values because having a
 343 Type I error is not catatstrophic in this case. We see that with the p-value adjustemnt, HOA and non HOA
 344 households are significantly different from each other in terms of number of trees in a yard.

Garden and HOA Statuses

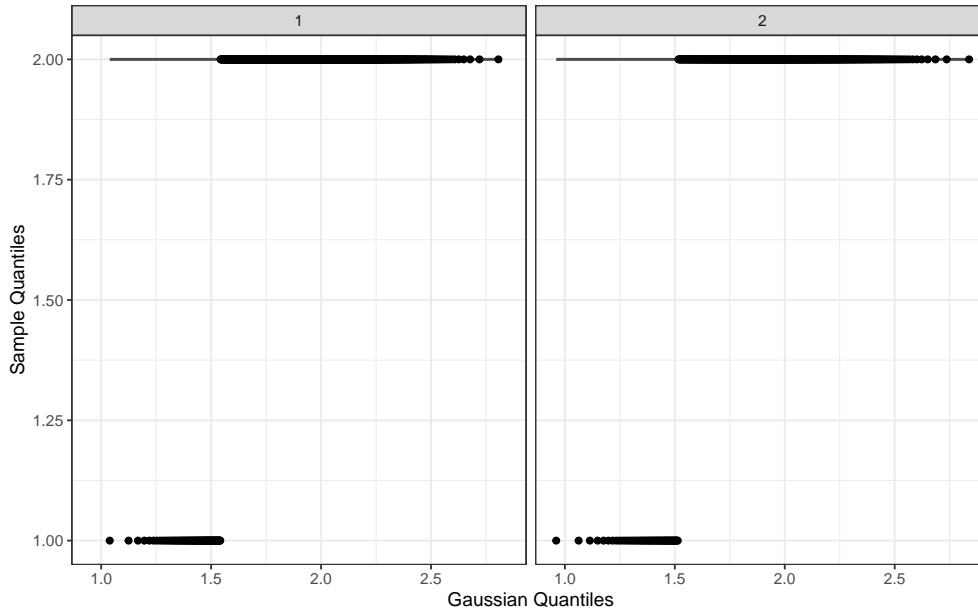


```
> dat.HOA$Garden<-factor(dat.HOA$Garden, levels = c(1,2),
+                          labels = c("Yes", "No"))
> prop.table(table(dat.HOA$Garden, dat.HOA$HOA),margin=2)*100
```

	Yes	No
Yes	7.743154	10.053860
No	92.256846	89.946140

Households in and not in HOAs tend to follow a similar pattern when it comes to gardens. Both tend to not have gardens present in the front or back of their homes. 92.2 percent of HOA homes and 89.9 percent of non HOA homes did not have gardens in the yards. Therefore, in terms of garden status, HOA homes and non HOA homes have a similar, positive effect on the environment by not having gardens. This is under the assumption that it is possible to overplant, which is bad for the environment. However, it is important to note that this is a judgement call as some gardens may benefit the environment, but we have no way to know from just viewing pictures.

We can perform a hypothesis test to determine whether or not there is a statistically significant difference between HOA homes and non HOA homes and across Garden status.



As we knew before and according to the above plot, our data is not normally distributed, but rather is discretely distributed by recycling status. Therefore, for robustness, we perform a Mood's Median Test. We want to generally identify if there is a difference across treatments, where the treatment is HOA versus non HOA. We look at these treatments specifically in the context of Garden status. We assume a representative sample. While these observations are not necessarily independent, since neighbors may base their recycling choices off of each other, for the sake of our study we will assume them to be so.

```
> #moods median test to determine if there are any significant differences
> #across treatments
> library(RVAideMemoire)
> mood.medtest(Garden~HOA,data=dat.HOA1)
```

Mood's median test

```
data: Garden by HOA
X-squared = 155.94, df = 1, p-value < 2.2e-16
```

Our null hypothesis is that there is no significant difference between HOA and non HOA households with regards to Garden statuses; the population medians are equal. Our alternative hypothesis is that there is a significant difference between HOA and non HOA households with regard to recycling statuses; at least one of the population medians is different. Since the plots of the data do not point towards normality, we will use a Mood's Median Test which is a nonparametric alternative to the ANOVA.

The Mood's median test is a chi-squared test that tests for differences across medians. Our chi-squared variable is 155.94 and our p-value is essentially zero, which is less than a significance level of 0.05. Since our p value is less than 0.05, we reject the null.

Since we are only comparing across 2 factors with 2 categories per factor, we did not need to perform a post-hoc test.

After examining sustainability factors at the household level, we can see that according to the number of trees in a home's yard and recycling status that non HOA homes are more sustainable than HOA homes. As a reminder, non HOA homes had more trees in their yard and better access to recycling and trash pickup compared to non HOA homes. Garden status and lawn care status gave results that did not clearly identify if HOA or non HOA homes were more sustainable. This is due to the fact that HOA and non HOA homes followed the same pattern of garden status. This indeterminate result also occurs because while HOA homes typically had more excellent lawn care than non HOA homes, they also had more lawns with no care than

non HOA homes. Without knowing how many/much chemicals were used for the different lawn statuses of "excellent", "good", and "poor", it is hard to determine the environmental effects of homes in each category.

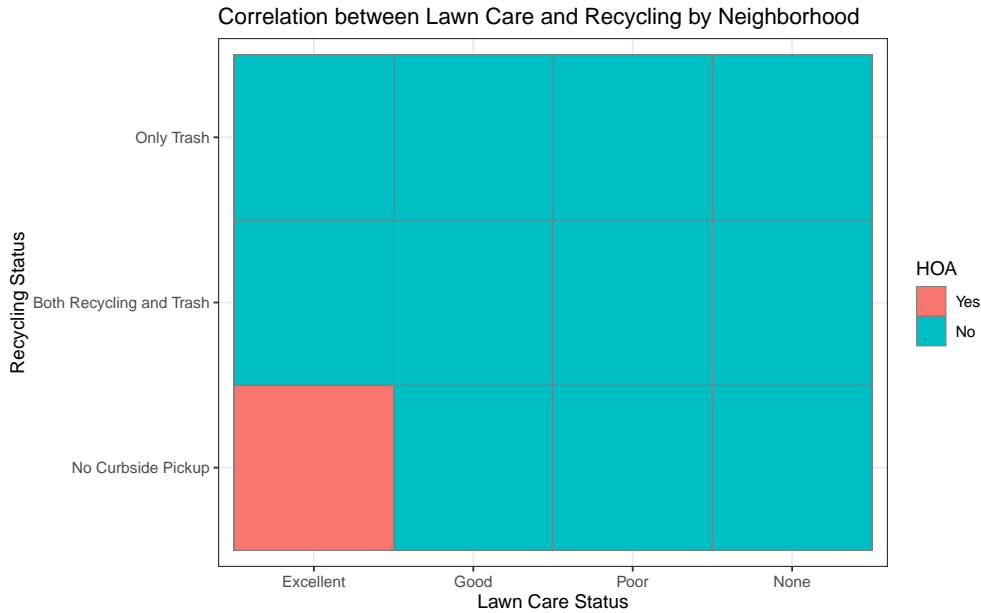
Now that we've seen the general count of observations with and without HOA, we can begin to examine multi-indicator sustainability analyses at the neighborhood level. We can do this by recreating the graphs that we previously created for each neighborhood.

After looking at the data, we can split up neighborhoods by their HOA status.

Neighborhoods with HOAs are: Brownstone Crossing, Edgewood at Paris, Glastonbury Village, Half Mile Lake, Northcliff, Patridge Ridge

Neighborhoods without HOAs are: Buxton, Croftstone Acres, Fox Springs, Liberty Park, Timberlake, Windermere

Relationship between Lawn Care and Recycling



As we can see from the tile plot above, HOA homes are associated with excellent lawn care and no curbside pickup for recycling, while non HOA homes are associated with all other categories. This points us to believe that HOA homes are less sustainable than non HOA homes because they are associated with unsustainable practices while non HOA homes are not. We can look at tables of the data get more insight.

	1	2	3	4	Sum
0	83	53	33	46	215
1	238	149	106	121	614
2	56	48	55	71	230
Sum	377	250	194	238	1059

	1	2	3	4
0	0.07837583	0.05004721	0.03116147	0.04343720
1	0.22474032	0.14069877	0.10009443	0.11425873
2	0.05288008	0.04532578	0.05193579	0.06704438

	1	2	3	4	Sum
0	12	19	18	10	59
1	82	129	119	43	373
2	16	52	44	13	125
Sum	110	200	181	66	557

	1	2	3	4
0	0.02154399	0.03411131	0.03231598	0.01795332
1	0.14721724	0.23159785	0.21364452	0.07719928
2	0.02872531	0.09335727	0.07899461	0.02333932

Recycling Status by Lawn Care Status for HOA

	None	Poor	Good	Excellent
No Trash Pickup	0.07859848	0.05018939	0.03125000	0.04356061
Trash and Recycling Pickup	0.22537879	0.14109848	0.10037879	0.11268939
Trash Pickup Only	0.05208333	0.04545455	0.05208333	0.06723485

Recycling Status by Lawn Care Status for non HOA

	None	Poor	Good	Excellent
No Trash Pickup	0.02197802	0.03479853	0.03296703	0.01831502
Trash and Recycling Pickup	0.14835165	0.23443223	0.21428571	0.06776557
Trash Pickup Only	0.02747253	0.09523810	0.08058608	0.02380952

As we can see from the tables, about 60 percent of non HOA households have both Trash and Recycling Pickup and none, poor, or good lawn care. About 46 percent of HOA households have both Trash and Recycling Pickup and none, poor, or good lawn care. This is important because it shows that non HOA neighborhoods have a larger proportion of their homes following good sustainability trends compared to HOA neighborhoods.

We can perform an association test to see if recycling and lawn care are associated by HOA. To do so, we can perform a chi-square independence test, a test that assesses the relationship between two categorical variables, which is also a nonparametric alternative to the Fisher test.

We can test whether observed dependence of recycling status and lawn care by HOA is due to random chance or not.

The null hypothesis is that the two categorical variables, recycling status and lawn care, are independent. The alternative hypothesis is that the two categorical variables, recycling status and lawn care, are dependent.

There are several assumptions we need to hold for HOA and non HOA households.

The first is that the two variables are categorical, which is true.

The second is that the observations are independent. While this isn't necessarily true, as neighbors can affect each other's behavior, for our purposes we will assume this to be true. Each household is only reported once.

The third is that the sample size is at least the number of cells in the table multiplied by 5. If there are 3 recycling statuses and 4 lawn care statuses, so 12 cells total, 1616 is more than 5 times larger than 12.

We need to check our expected counts, or the number of observations that we would expect to see in a cell if the observations were truly independent (if the null hypothesis is true). We need to ensure 80% of expected counts are greater than 5. This is calculated by multiplying the total of row i by the total of column j and dividing this quantity by n. We can do this for HOA and non HOA households.

$$(215 \times 377)/1616=50.16$$

$$(614 \times 377)/1616=143.24$$

$$(230 \times 377)/1616=53.66$$

$$(215 \times 250)/1616=33.26$$

$$(614 \times 250)/1616=94.99$$

$$(230 \times 250)/1616=35.58$$

$$(215 \times 194)/1616=25.81$$

$$(614 \times 194)/1616=73.71$$

$$(230 \times 194)/1616=27.61$$

$$(215 \times 238)/1616=31.66$$

$$(614 \times 238)/1616=90.42$$

$$(230 \times 238)/1616=33.87$$

$$(59 \times 110)/1616=4.02$$

$$(373 \times 110)/1616=25.39$$

$$(125 \times 110)/1616=8.51$$

$$(59 \times 200)/1616=7.30$$

$$(373 \times 200)/1616=46.16$$

$$(125 \times 200)/1616=15.47$$

$$(59 \times 181)/1616=6.61$$

$$(373 \times 181)/1616=41.78$$

$$(125 \times 181)/1616=14.00$$

```

451 (59 X 66)/1616=2.41
452 (373 X 66)/1616=15.23
453 (125 X 66)/1616=5.11 As we can see, all expected counts are greater than 5 for the HOA sample and 10 out
454 of 12 cells for the non HOA sample or about 83% of cells are greater than 5.
455     Lastly, we need to confirm that none of the expected counts are less than one, which is true in these
456 cases.
457     We can now calculate test statistics for the chi-square tests.

> #HOA homes
> chisq.test(x=dat.HOA.y$Recycle,y=dat.HOA.y$LawnCare)

Pearson's Chi-squared test

data:  dat.HOA.y$Recycle and dat.HOA.y$LawnCare
X-squared = 26.147, df = 6, p-value = 0.000209

> #nonHOA homes
> chisq.test(x=dat.HOA.n$Recycle,y=dat.HOA.n$LawnCare)

Pearson's Chi-squared test

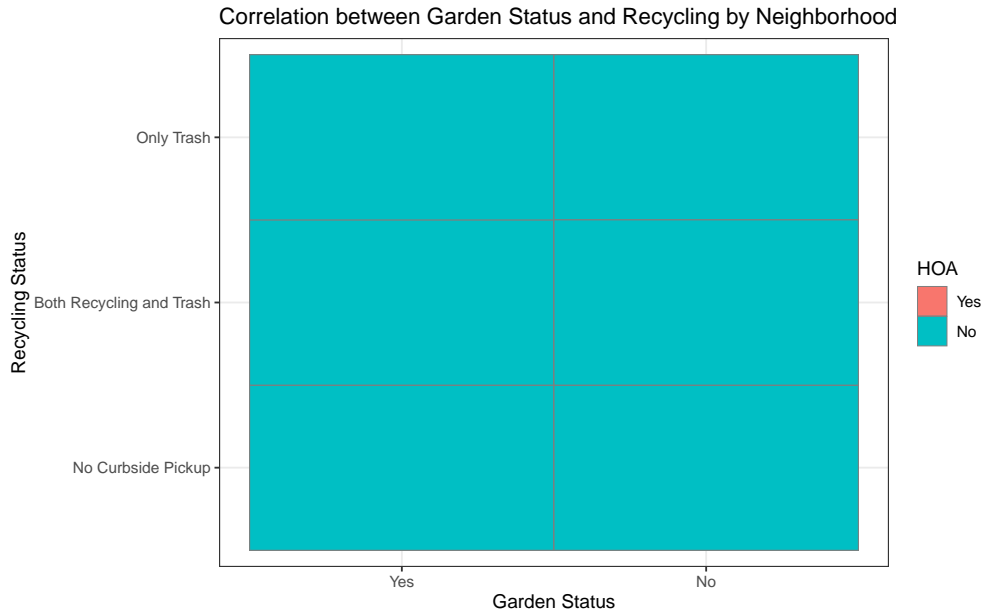
data:  dat.HOA.n$Recycle and dat.HOA.n$LawnCare
X-squared = 7.488, df = 6, p-value = 0.2781

```

458 As we can see for HOA households, the test statistic is 26.147 and the p-value is 0.0002, which is close to
459 zero. We therefore can reject the null hypothesis and conclude that for HOA households, there is sufficient
460 evidence to suggest that there is a relationship between recycling status and lawn care.

461 For non HOA households, the test statistic is 7.488 and the p-value is 0.2781, which is greater than
462 0.05. We therefore fail to reject the null hypothesis and conclude that for non HOA households, there is not
463 sufficient evidence to suggest that there is a relationship between recycling status and lawn care.

Relationship between Garden Status and Recycling



As we can see from the plot above, there is no correlation for garden status and recycling status for HOA households. However, all combinations of garden status and recycling status are correlated for non HOA households. We can further explore these relationships by looking at a table.

```
> rgy<-table(dat.HOA.y$Recycle,dat.HOA.y$Garden)
> rgy1<-addmargins(table(dat.HOA.y$Recycle,dat.HOA.y$Garden))
> rgy1
> prop.table(rgy)
> rgn<-table(dat.HOA.n$Recycle,dat.HOA.n$Garden)
> rgn1<-addmargins(table(dat.HOA.n$Recycle,dat.HOA.n$Garden))
> rgn1
> prop.table(rgn)
```

Recycling Status by Garden Status for HOA

	Yes	No
No Trash Pickup	0.025568182	0.178030303
Trash and Recycling Pickup	0.045454545	0.534090909
Trash Pickup Only	0.006628788	0.210227273

Recycling Status by Garden Status for non HOA

	Yes	No
No Trash Pickup	0.01282051	0.09523810
Trash and Recycling Pickup	0.07509158	0.58974359
Trash Pickup Only	0.01465201	0.21245421

Examining the tables, we see that about 59 percent of non HOA homes had trash and recycling pickup and did not have a garden compared to 53 percent of HOA homes. It is important to note that the effects of gardens for non HOA homes compared to HOA homes is minimal because both types of neighborhoods typically tend to not have gardens. Especially since gardens can be organic or not, contain different kinds

of plants that are better or worse for soil and can be watered at different frequencies, their net effects are probably hard to determine from visuals alone.

We can now perform association tests to see if garden and recycling status are correlated by HOA.

We can perform an association test to see if recycling and garden statuses are associated by HOA. To do so, we can perform a chi-square independence test, a test that assesses the relationship between two categorical variables, which is also a nonparametric alternative to the Fisher test.

We can test whether observed dependence of recycling status and garden status by HOA is due to random chance or not.

The null hypothesis is that the two categorical variables, recycling status and garden status, are independent. The alternative hypothesis is that the two categorical variables, recycling status and garden status, are dependent.

There are several assumptions we need to hold for HOA and non HOA households.

The first is that the two variables are categorical, which is true.

The second is that the observations are independent. While this isn't necessarily true, as neighbors can affect each other's behavior, for our purposes we will assume this to be true. Each household is only reported once.

The third is that the sample size is at least the number of cells in the table multiplied by 5. If there are 3 recycling statuses and 2 garden statuses, so 6 cells total, 1616 is more than 5 times larger than 6.

We need to check our expected counts, or the number of observations that we would expect to see in a cell if the observations were truly independent (if the null hypothesis is true). We need to ensure 80% of expected counts are greater than 5. This is calculated by multiplying the total of row *i* by the total of column *j* and dividing this quantity by *n*. We can do this for HOA and non HOA households.

HOA Households

$(215 \times 82)/1616=10.91$

$(614 \times 82)/1616=31.16$

$(230 \times 82)/1616=11.67$

$(215 \times 977)/1616=129.98$

$(614 \times 977)/1616=257.23$

$(230 \times 977)/1616=139.05$

Non HOA Households

$(59 \times 56)/1616=2.04$

$(373 \times 56)/1616=12.93$

$(125 \times 56)/1616=4.33$

$(59 \times 501)/1616=18.29$

$(373 \times 501)/1616=115.64$

$(125 \times 501)/1616=38.75$ As we can see, all expected counts are greater than 5 for the HOA sample and 4 out of 6 cells for the non HOA sample or about two-thirds of cells are greater than 5. Our non HOA household therefore does not meet the assumptions required for the chi-squared test.

Lastly, we need to confirm that none of the expected counts are less than one, which is true in these cases.

We can now calculate test statistics for the chi-square tests.

```
> #HOA homes
```

```
> chisq.test(x=dat.HOA.y$Recycle,y=dat.HOA.y$Garden)
```

Pearson's Chi-squared test

data: dat.HOA.y\$Recycle and dat.HOA.y\$Garden

X-squared = 14.094, df = 2, p-value = 0.0008701

```
> #nonHOA homes
```

```
> chisq.test(x=dat.HOA.n$Recycle,y=dat.HOA.n$Garden)
```

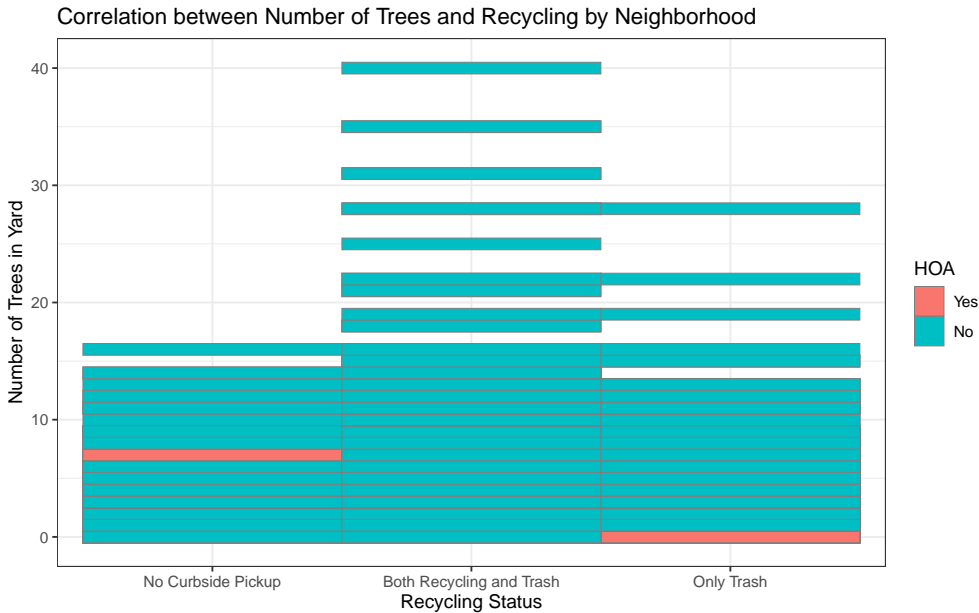
Pearson's Chi-squared test

```
data: dat.HOA.n$Recycle and dat.HOA.n$Garden  
X-squared = 2.4223, df = 2, p-value = 0.2979
```

519 As we can see for HOA households, the test statistic is 14.094 and the p-value is 0.00087, which is close to
520 zero. We therefore can reject the null hypothesis and conclude that for HOA households, there is sufficient
521 evidence to suggest that there is a relationship between recycling status and garden status.

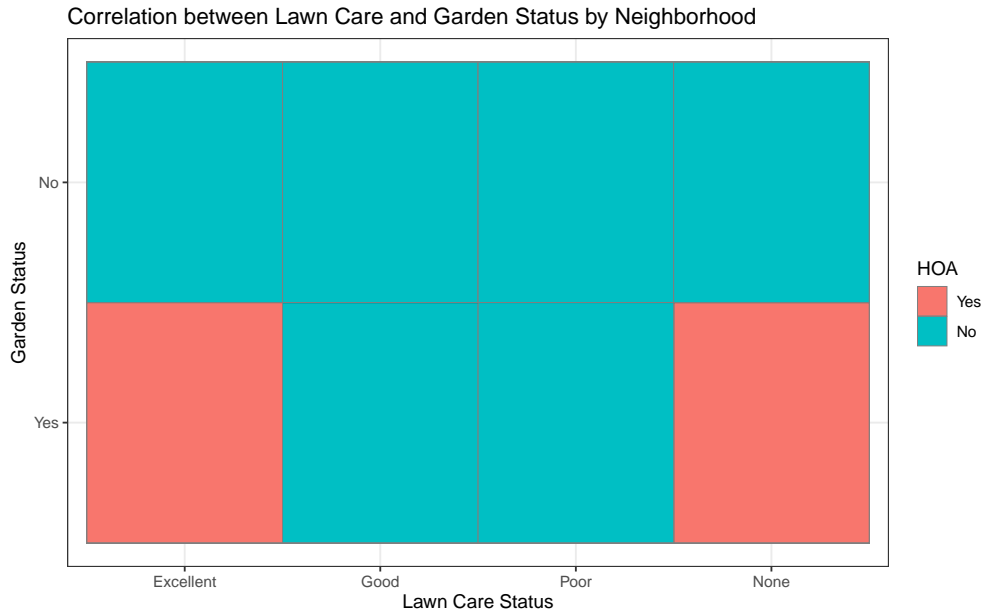
522 For non HOA households, our sample did not meet the assumptions. However, if we run the test re-
523 gardless, the test statistic is 2.4223 and the p-value is 0.2979, which is greater than 0.05. We therefore fail
524 to reject the null hypothesis and conclude that for non HOA households, there is not sufficient evidence to
525 suggest that there is a relationship between recycling status and lawn care.

Relationship between Trees and Recycling



This graph illustrates that recycling and trash pickup is correlated with more trees for both HOA and non HOA neighborhoods, although non HOA neighborhoods have more trees generally. Not only that but there is a stronger relationship between trash pickup only and number of trees for non HOA households than HOA households. From this graph, we can posit that non HOA neighborhoods may be more sustainable than HOA neighborhoods because they have greater numbers of trees and instances of recycling and trash pickup or trash pickup only. HOA neighborhoods seem to only have associations with only trash pickup and no trees, or no curbside pickup and 7 trees. We can perform an association test to see if there is correlation between trees and recycling status by HOA.

Relationship between Lawn Care and Garden Status



We can see from the plot above that HOA households are associated with excellent lawn care and having gardens, as well as having gardens and no lawn care. Non HOA households are associated with all other combinations. The correlation between having excellent lawn care and garden status is important because it points us to believe that HOA status homes could be worse for the environment than non HOA homes. We can further explore this relationship with tables.

```
> lcyg<-table(dat.HOA.y$LawnCare,dat.HOA.y$Garden)
> lcyg1<-addmargins(table(dat.HOA.y$LawnCare,dat.HOA.y$Garden))
> lcyg1
> prop.table(lcyg)
> lcn1<-table(dat.HOA.n$LawnCare,dat.HOA.n$Garden)
> lcn1<-addmargins(table(dat.HOA.n$LawnCare,dat.HOA.n$Garden))
> lcn1
> prop.table(lcn1)
```

Lawn Care Status by Garden Status for HOA

	Yes	No
None	0.02462121	0.33143939
Poor	0.01988636	0.21685606
Good	0.01325758	0.17045455
Excellent	0.01988636	0.20359848

Lawn Care Status by Garden Status for non HOA

	Yes	No
None	0.014652015	0.183150183
Poor	0.042124542	0.322344322
Good	0.042124542	0.285714286
Excellent	0.003663004	0.106227106

The tables show us that 55 percent of non HOA neighborhoods had no or poor lawn care while 58 percent of HOA neighborhoods had no or poor lawn care. These statistics alone may not be enough to show that one neighborhood is more sustainable than the other for this relationship. We see that about 10.6 percent of non HOA neighborhoods had excellent lawn care while about 22 percent of HOA neighborhoods had excellent lawn care. This indicates more that HOA neighborhoods may be less sustainable than non HOA neighborhoods because they use more chemicals for their lawn care. Again, since both HOA and non HOA neighborhoods have gardens around the same proportion, it is hard to compare these neighborhoods looking at both gardens status and lawn care status.

We can perform an association test to determine if there is correlation between garden status and lawn care status by HOA status.

We can perform an association test to see if lawn care and garden statuses are associated by HOA. To do so, we can perform a chi-square independence test, a test that assesses the relationship between two categorical variables, which is also a nonparametric alternative to the Fisher test.

We can test whether observed dependence of lawn care status and garden status by HOA is due to random chance or not.

The null hypothesis is that the two categorical variables, lawn care status and garden status, are independent. The alternative hypothesis is that the two categorical variables, lawn care status and garden status, are dependent.

There are several assumptions we need to hold for HOA and non HOA households.

The first is that the two variables are categorical, which is true.

The second is that the observations are independent. While this isn't necessarily true, as neighbors can affect each other's behavior, for our purposes we will assume this to be true. Each household is only reported once.

The third is that the sample size is at least the number of cells in the table multiplied by 5. If there are 4 lawn care statuses and 2 garden statuses, so 8 cells total, 1616 is more than 5 times larger than 8.

We need to check our expected counts, or the number of observations that we would expect to see in a cell if the observations were truly independent (if the null hypothesis is true). We need to ensure 80% of expected counts are greater than 5. This is calculated by multiplying the total of row i by the total of column j and dividing this quantity by n. We can do this for HOA and non HOA households.

HOA Households

$(377 \times 82)/1616=19.13$

$(250 \times 82)/1616=12.68$

$(194 \times 82)/1616=9.84$

$(238 \times 82)/1616=12.08$

$(377 \times 977)/1616=227.93$

$(250 \times 977)/1616=151.14$

$(194 \times 977)/1616=117.29$

$(238 \times 977)/1616=143.89$

Non HOA Households

$(110 \times 56)/1616=3.81$

$(200 \times 56)/1616=6.93$

$(181 \times 56)/1616=6.27$

$(66 \times 56)/1616=2.29$

$(110 \times 501)/1616=34.10$

$(200 \times 501)/1616=62.00$

$(181 \times 501)/1616=56.11$

$(66 \times 501)/1616=20.46$ As we can see, all expected counts are greater than 5 for the HOA sample and 6 out of 8 cells for the non HOA sample or 75% of cells are greater than 5. Our non HOA household therefore does not meet the assumptions required for the chi-squared test.

Lastly, we need to confirm that none of the expected counts are less than one, which is true in these cases.

We can now calculate test statistics for the chi-square tests.

```
> #HOA homes
> chisq.test(x=dat.HOA.y$LawnCare,y=dat.HOA.y$Garden)
```

Pearson's Chi-squared test

```
data: dat.HOA.y$LawnCare and dat.HOA.y$Garden
X-squared = 0.99345, df = 3, p-value = 0.8028
```

```
> #nonHOA homes
> chisq.test(x=dat.HOA.n$LawnCare,y=dat.HOA.n$Garden)
```

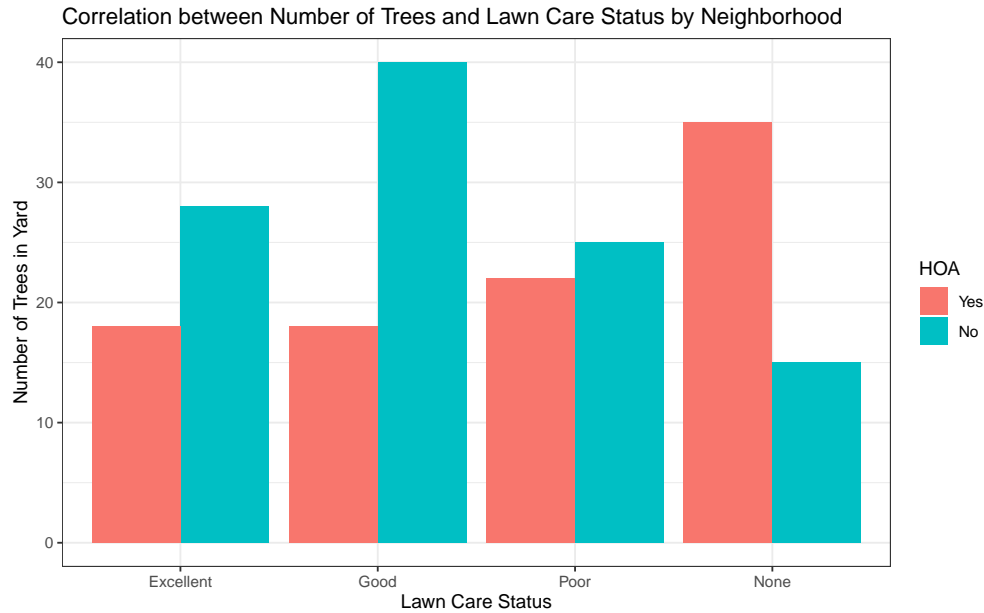
Pearson's Chi-squared test

```
data: dat.HOA.n$LawnCare and dat.HOA.n$Garden
X-squared = 6.4128, df = 3, p-value = 0.09317
```

601 As we can see for HOA households, the test statistic is 0.99345 and the p-value is 0.8028, which is greater
 602 than 0.05. We therefore fail to reject the null hypothesis and conclude that for HOA households, there is
 603 not sufficient evidence to suggest that there is a relationship between lawn care status and garden status.

604 For non HOA households, our sample did not meet the assumptions. However, if we run the test regard-
 605 less, the test statistic is 6.4128 and the p-value is 0.09317, which is greater than 0.05. We would therefore
 606 fail to reject the null hypothesis and conclude that for non HOA households, there is not sufficient evidence
 607 to suggest that there is a relationship between recycling status and lawn care.

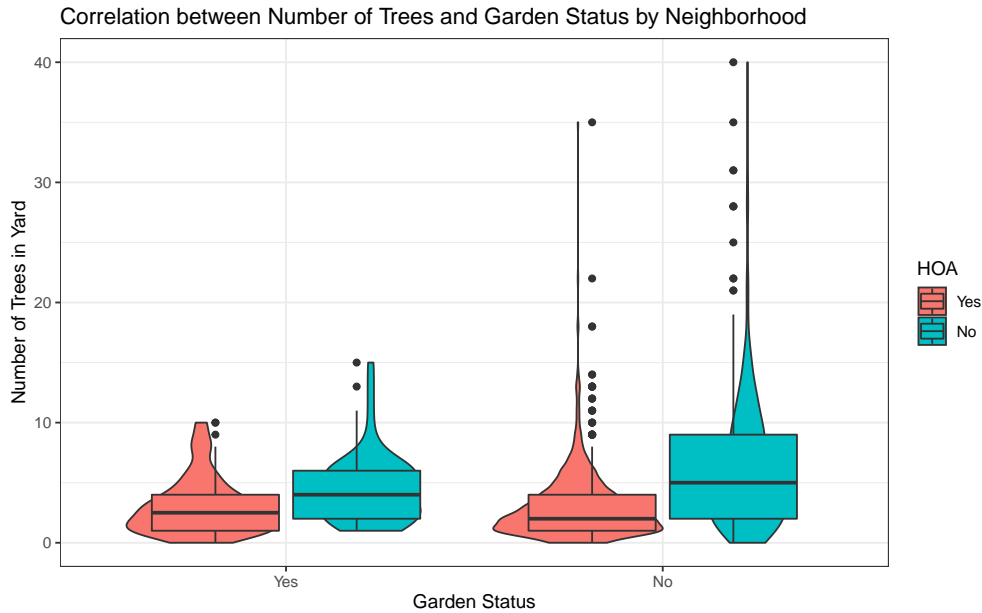
Relationship between Trees and Lawn Care



The graph above illustrates that non HOA neighborhoods could be more sustainable than HOA neighborhoods because non HOA homes are more likely to have more trees and less lawn care compared to HOA neighborhoods. As we can see, the two kinds of neighborhoods follow different patterns. Non HOA neighborhoods follow a right skew pattern, where most observations have a lot of trees and little lawn care and then decrease the amount of trees they have as they increase their lawn care. HOA neighborhoods follow an increasing exponential pattern where none and poor lawn care statuses have fewer trees and as lawn care improves, the amount of trees seen on the property increases.

We can examine this relationship further by performing an association test.

Relationship between Trees and Garden Status



The graph above illustrates that regardless of garden status, non HOA neighborhoods have greater medians of trees than HOA neighborhoods. If non HOA neighborhoods have more trees regardless of garden status, then they can be seen as more sustainable than HOA neighborhoods.

Overall, we can see that non HOA neighborhoods are more sustainable than HOA neighborhoods. While some factors such as gardens may be ambiguous in terms of their sustainability, our analysis shows that non HOA neighborhoods have more trees, less lawn care, and more trash and recycling pickup options than HOA homes. Our sample had more HOA observations than non HOA observations, so we compared these variables using relative frequency for comparing HOA statuses to the four variables of lawn care, recycling, garden status, and trees, as well as illustrating relative frequency through tables for discrete vs. discrete variables and graphs for discrete vs. continuous variables. This paper encourages future researchers to look into the sustainability of gardens and what specific factors can be used to tell what gardens may or may not be sustainable. This research is important for potential homebuyers looking to purchase a sustainable home as well as those who are environmentally conscious and may or may not live in an HOA neighborhood.

We can perform an association test to test our belief of correlation between number of trees and garden status.

References

- Almeida, A., Loy, A., and Hofmann, H. (2017). *qqplotr: Quantile-Quantile Plot Extensions for 'ggplot2'*. R package version 0.0.3 initially funded by Google Summer of Code 2017.
- Auguie, B. (2017). *gridExtra: Miscellaneous Functions for "Grid" Graphics*. R package version 2.3.
- Hasselman, B. (2018). *nleqslv: Solve Systems of Nonlinear Equations*. R package version 3.3.2.
- Hervé, M. (2019). *RVAideMemoire: Testing and Plotting Procedures for Biostatistics*. R package version 0.9-73.
- Mangiafico, S. (2019). *rcompanion: Functions to Support Extension Education Program Evaluation*. R package version 2.3.7.
- Wickham, H. (2016). *ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York.