Briana Flynn
SES 5215 Quantitative Methods for Urban Planning
Exploring Predictors of Rent Burden in Suffolk County

# Introduction

This report will analyze gross percentage of household income spent on rent and its relationship to other selected household characteristics of renter households in Suffolk County, Massachusetts. I was interested in exploring the relationship between rent burden, household income, and other household characteristics in Suffolk County in order to better understand the prevalence of rent burden, and the characteristics of households are more or less rent-burdened. My broad question was whether there were certain characteristics that made renter households more likely to experience higher or lower rent burden.

## Hypotheses

1. Higher income is correlated with lower rent burden;
2. Higher rent is correlated with higher rent burden;
3. English proficiency is correlated with lower rent burden;
4. Internet access is correlated with lower rent burden; and
5. Presence of any number of children is correlated with higher rent burden as compared with no children.

While the combination of higher income and lower rent clearly results in lower rent burden, and the combination of lower income and higher rent results in higher burden, I was also interested in understanding rent burden's prevalence across income levels and rent levels to understand whether the experience was limited to certain populations or rental unit types.

# Data

To answer my research question, I relied on 1-year 2018 American Community Survey ("ACS") estimates via the Public Use Microdata Sample ("PUMS") data. The resulting dataset included results with 17,922 observations of the six variables of interest. In other words, the dataset includes a sample of 17,922 renter-occupied households in Suffolk County.

*Figure 1: Variables in Dataset*

| Variable | Data Label | Variable Type |
|---|---|---|
| Household Income | HINCP | Continuous |
| Gross Monthly Rent Payment | GRNTP | Continuous |
| Gross Monthly Rent as a Percentage of Household Income | GRPIP | Continuous |
| Access to the Internet | ACCESS_label | Categorical |
| Limited English-Speaking Household[1] | LNGI_label | Categorical |
| Household Presence and Age of Children | HUPAC_label | Categorical |

# Methods

## Distribution

*Figure 2: Testing the Distribution Statistics of Rent as a Percentage of Household Income ("GRPIP")*

| Test | Description |
|---|---|
| Sample GRPIP Mean | Sum of the GRPIP values divided by the number of values |
| Sample GRPIP Median | Value separating the higher half from the lower half of GRPIP values |
| Sample GRPIP Standard Deviation | Measure of the amount of variation or dispersion of a GRPIP values |
| 95 % GRPIP Confidence Interval | A range of GRPIP values that one can be 95% certain contains the true mean |
| Interquartile Range for GRPIP | A measure of variability, based on dividing GRPIP values into quartiles |
| GRPIP Histogram | An approximate graphical representation of the distribution of values |

---

[1] This binary categorical variable is transformed into a true or false variable for "English Proficient Household" with the label "Engl_prof" later in report.

This statistical test was run as an introduction to the variable of interest. The distribution statistics illustrate the overall prevalence of rent burden, what share of households are paying more than 30 and 50 percent of their income on rent, and what percentage of household income the "average" Suffolk County household is spending on rent, as well as the bottom quartile and top quartile.

## Correlation

Next I tested the relationship between the rent as a percentage of household income and each of my other variables, which required four types of relationship tests: Pearson test for the strength and significance of the relationship between two continuous variables, Two-sample T-test for the strength and significance of the relationship between a continuous variable and a binary categorical variable, ANOVA test for the significance of the relationship between a continuous variable and a categorical variable with three of more levels, and Tukey Test for the confidence intervals for the differences in continuous variable means between each pair of categories in a categorical variable.

*Figure 3: Testing Correlations and Differences*

| Test | Description |
|---|---|
| HINCP + GRPIP Pearson Test | Test the strength and significance of the relation between household income and rent as a percentage of household income |
| GRNTP + GRPIP Pearson Test | Test the strength and significance of the relationship between rent and rent as a percentage of household income |
| GRPIP + LNGI Two-Sample T-Test | Test the strength and the significance of the relationship between rent as a percentage of household income and English proficiency |
| GRPIP + ACCESS ANOVA Test | Test the significance of the relationship between rent as a percentage of household income and internet access categories |
| GRPIP + ACCESS Tukey Test | Test the confidence intervals for the differences in rent as a percentage of household income between each pair of categories. |
| GRPIP + HUPAC ANOVA Test | Test the significance of the relationship between rent as a percentage of household income and presence and age of children categories |
| GRPIP + HUPAC Tukey Test | Test the confidence intervals for the differences in rent as a percentage of household income between each pair of categories. |

These tests will allow me understand whether there is likely a statistically significant relationship between rent as a percentage of household income and other variables, and for certain relationships, the magnitude or strength of the relationship.

Pearson's correlation will indicate the relationship between two continuous variables. It can range from -1 to 1. A value of zero means there's no relationship between the two variables. Values closer to 1 indicate a stronger positive relationship, and values closer to -1 indicate a stronger negative relationship. A correlation will also have a p-value associated with it. The p-value is the likelihood that you would have calculated that correlation for your sample if the correlation for the full population was actually zero. When the p-value is less than 0.05, the 95-percent confidence interval will not include zero, and the correlation is significant.

The two-sample t-test will provide the 95-percent confidence interval for the difference in averages between two categorical variables, as well as the p-value (significance) of that difference. In other words, this will be the average difference between GRPIP points in the categories. If the 95-percent confidence interval for the two-sample t-test does not include zero (either all positive or all negative numbers), then there is a significant difference.

An Analysis of Variance (ANOVA) test will provide the significance of the relationship between a categorical variable and a continuous variable. In other words, it will give you a p-value representing the likelihood that there is being in any particular category has relationship with what the value of the continuous variable will be.

Tukey's Honestly Significant Difference (HSD) test will provide confidence intervals for the differences in averages (means) between each possible pair of categories, which indicates the magnitude of a relationship to supplement the ANOVA test's indication of significance.

## Regression

Finally, I used two initial regression models to attempt to predict percentage of household income based on my other variables. Because rent as a percentage of household income can be derived from my two other continuous variables, I ran one set of regression models that excluded gross monthly rent from my independent variable.

$$GRPIP = \beta_0 + \beta_1(HINCP) + \beta_2(ACCESS\_label) + \beta_3(HUPAC\_label) + \beta_4(engl\_prof)$$

Following this initial regression model, I modified the model with log-transformation of the continuous independent variable (household income and gross monthly rent, respectively), and added interactions between certain variables in order to improve model fit. I then ran a second set of regression models that excluded household income from my independent variables, but including gross monthly rent.

$$GRPIP = \beta_0 + \beta_1(GRNTP) + \beta_2(ACCESS\_label) + \beta_3(HUPAC\_label) + \beta_4(engl\_prof)$$

These regression models indicate the most likely coefficient values, how likely it is the coefficients are zero, how much of the variation of rent as a percentage of household income can be predicted from the regression tests, and how likely it is that the set of variables, taken together, have a linear relationship.

*Figure 4: Regression Tests and Selected Modifications*

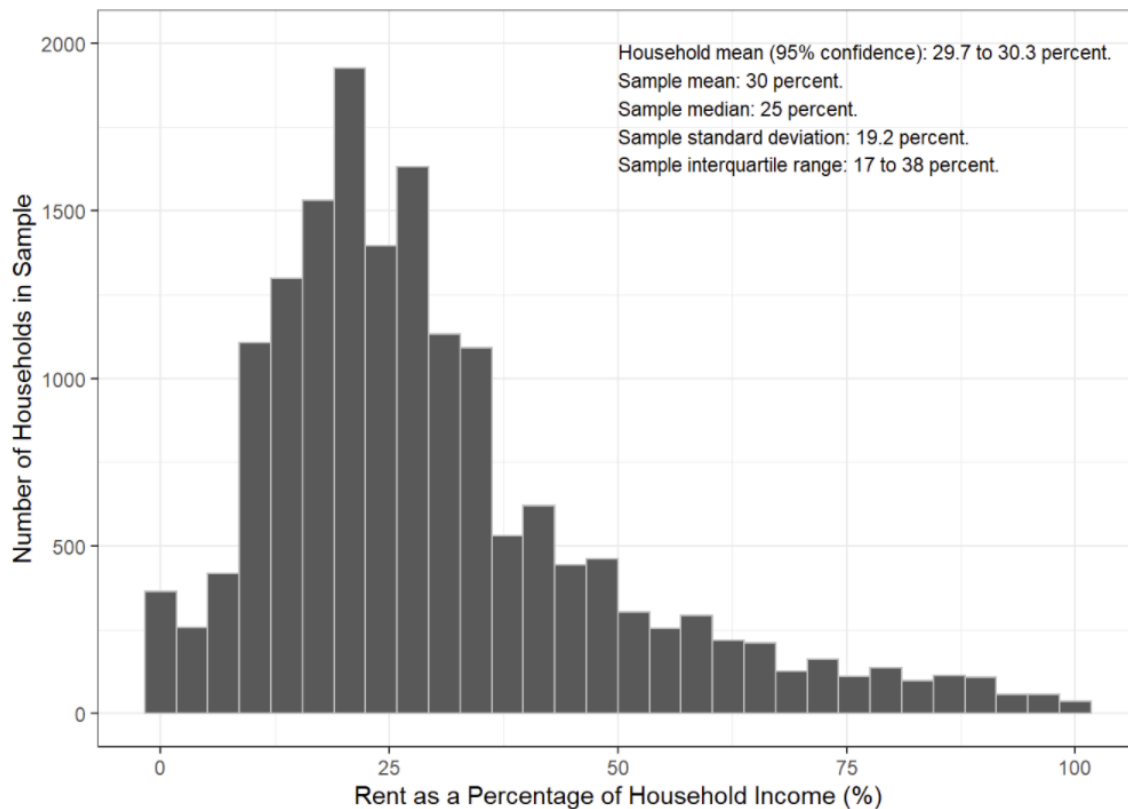| Name | Test | Description |
|---|---|---|
| Regression Model 1 | Initial Regression Test: with HINCP | A regression model to predict GRPIP based on the values of household income, internet access, presence and age of children, and English proficiency |
| Regression Model 1B | Regression Test with Log-transformed HINCP | A regression model to predict GRPIP based on the values of log-transformed household income, internet access, presence and age of children, and English proficiency |
| Regression Model 1C | Regression Test with Interactions between log-transformed HINCP | A regression model to predict GRPIP based on the values of log-transformed household income, internet access, presence and age of children, and English proficiency, with interactions between log-transformed household income and internet access and log-transformed household income and English proficiency |
| Regression Model 2 | Initial Regression Test: with GRNTP | A regression model to predict GRPIP based on the values of gross monthly rent, internet access, presence and age of children, and English proficiency |

# Results

## Distribution

*Distribution of Rent as a Percentage of Household Income*

*Figure 5: Sample Mean, Median, Standard Deviation, Interquartile Range, and Confidence Interval for Rent as a Percentage of Household Income*

| Statistic | Rent as a Percentage of Household Income |
|---|---|
| Sample Mean | 30% |
| Median | 25% |
| Standard Deviation | 19.2% |
| Interquartile range | 17% to 38% |
| Population mean (95% confidence) | 29.7% to 30.3% |

*Figure 6: Sample Mean, Standard Deviation, and Confidence Interval*



The distribution of monthly rent as a percentage of household income is right skewed. The average for the dataset is 30% of income spent on rent, while the median household spends 25% of household income on rent. A quarter of households spend more than 38% of their income on rent and a quarter spent less than 17% of their income on rent. The 95% confidence interval indicates that the mean is 95% likely to fall within 29.7 to 30.3 percent.
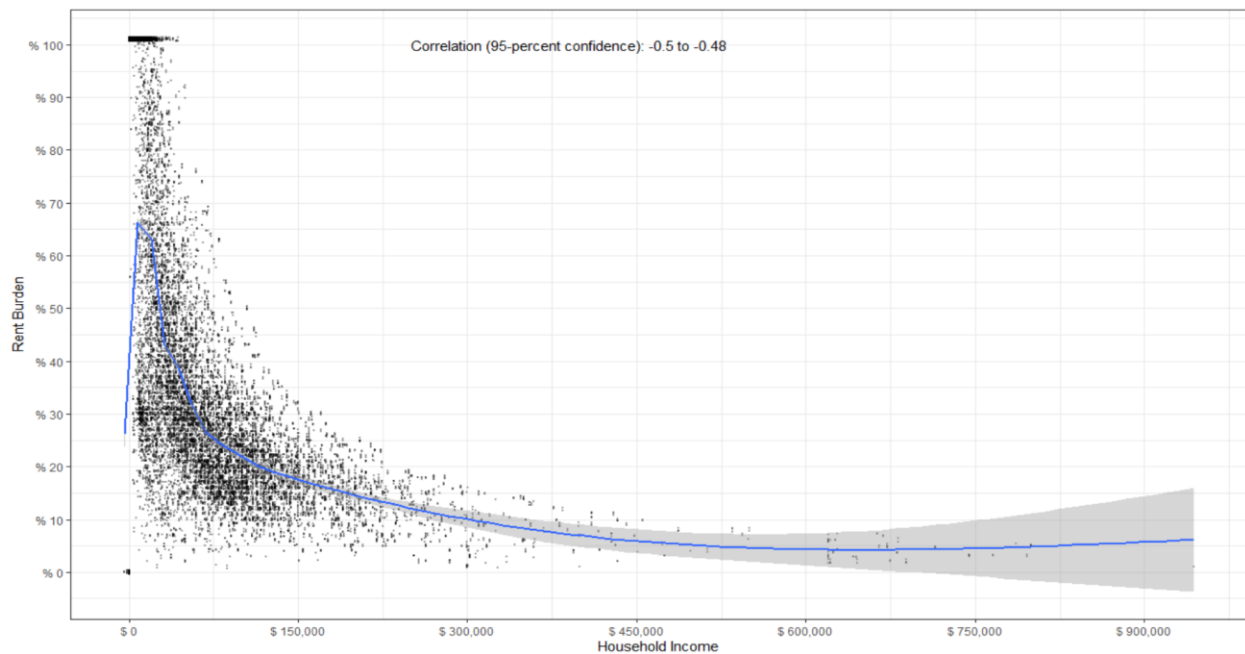
## Correlations, Differences, and Associations

*Pearson Tests Between Continuous Variables*

*Figure 7: Pearson Test Between Rent as a Percentage of Household Income and Household Income*

```
data:  households$HINCP and households$GRPIP
t = -74.547, df = 17920, p-value < 0.00000000000000022
alternative hypothesis: true correlation is not equal to 0
95 percent confidence interval:
 -0.4976206 -0.4752692
sample estimates:
       cor
-0.4865245
```

*Figure 8: Visualizing Correlation between Rent as a Percentage of Household Income and Household Income*
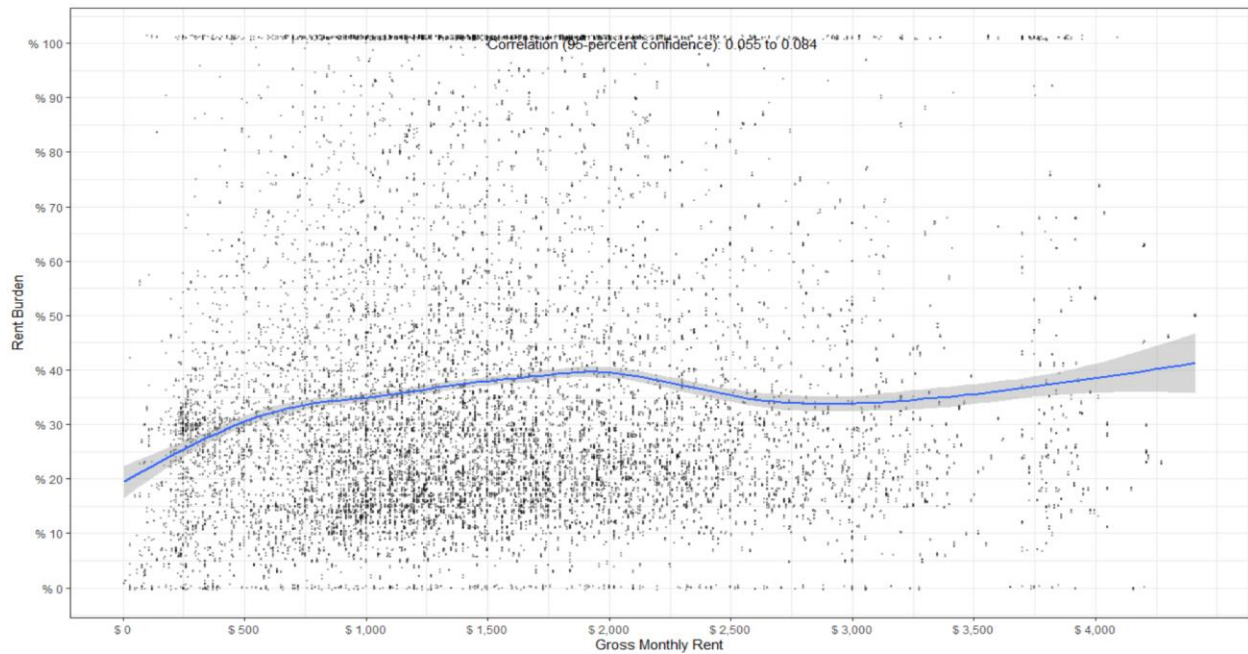


The results of a Pearson correlation test between household income and rent as a percentage of household income indicate that there is a statistically significant relationship between household income and rent as a percentage of household income. The correlation is -0.45, with a 95% confidence interval of -0.49 to -0.49. In general, you can see in Figure 8 that as household income increases, the percentage of household income spent on rent decreases.

*Figure 9: Pearson Test between Rent as a Percentage of Household Income and Gross Monthly Rent*

```
data:  households$GRNTP and households$GRPIP
t = 9.3251, df = 17920, p-value < 0.00000000000000022
alternative hypothesis: true correlation is not equal to 0
95 percent confidence interval:
 0.05490688 0.08404680
sample estimates:
       cor
0.06949166
```

*Figure 10: Visualizing Correlation between Rent as a Percentage of Household Income and Gross Monthly Rent*
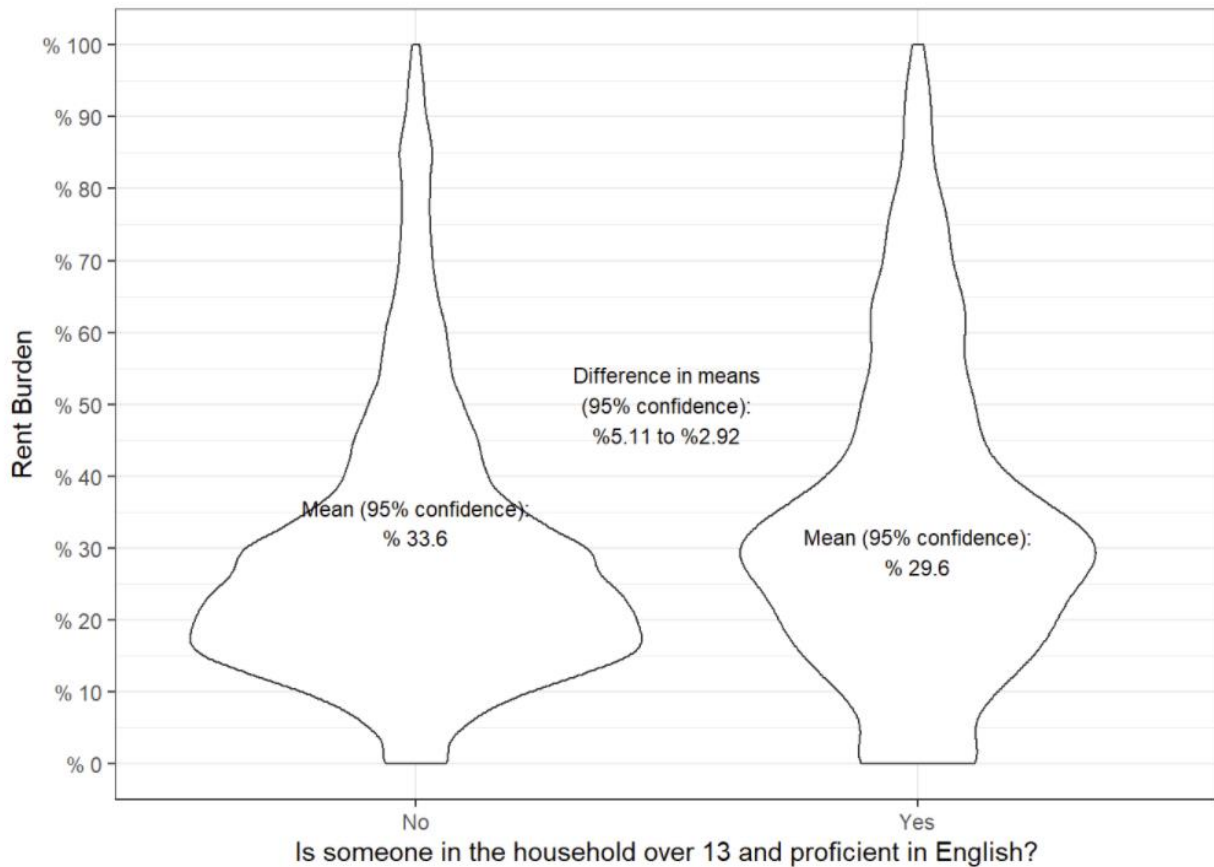


The results of a Pearson correlation test between gross monthly rent and rent as a percentage of household income indicate that there is also statistically significant relationship between household income and rent as a percentage of household income. However, the correlation appears to be very weak, with a correlation value of around 0.05, and a 95% confidence interval of 0.04 to 0.08.

*Two-sample T-Tests Between a Continuous and a Binary Categorical Variables*

*Figure 11: 2-Sample T-Test: Rent as a Percentage of Household Income and English Proficiency*

```
## t = 7.1972, df = 1923.7, p-value = 8.773e-13
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
##   2.921134 5.109425
## sample estimates:
## mean in group FALSE   mean in group TRUE
##            33.59207             29.57679
```

*Figure 12: Visualizing Difference in Rent as a Percentage of Household Income based on English Proficiency*



The results of the Two-Sample T-test suggest that there is a significant difference in rent as a percentage of household income means between English proficient and non-English proficient households, with a P-value near zero. However, the difference in means appears to be relatively small at 2.0%, with a 95% confidence interval of 2.92% to 5.11%. The results suggest that households that are English proficient spend 2.0% more of their income on rent than households that are not English proficient.
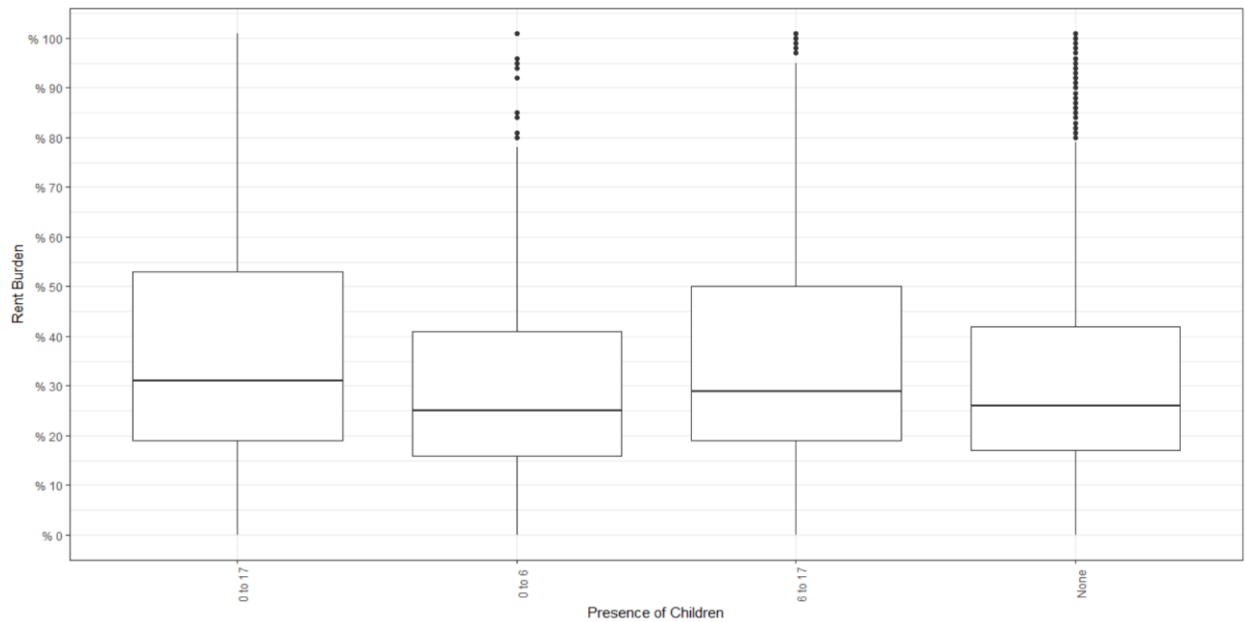
*ANOVA Tests Between a Continuous and a 3+ Result Categorical Variables*
*Figure 13: ANOVA Test: Rent as a Percentage of Household Income and Presence/Age of Children*

```
              Df    Sum Sq Mean Sq F value            Pr(>F)
HUPAC_label    3     95671   31890   45.26 <0.0000000000000002 ***
Residuals  17918 12625395     705
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

*Figure 14: Tukey Test to Visualize Difference in Rent as a Percentage of Household Income based on Presence of Children*



There appears to be a significant relationship between rent as a percentage of household income and presence/age of children in the household, with a P-value near 0. However, the differences appear to be very minor, with the differences in means not varying more than 5% between any of the categories.
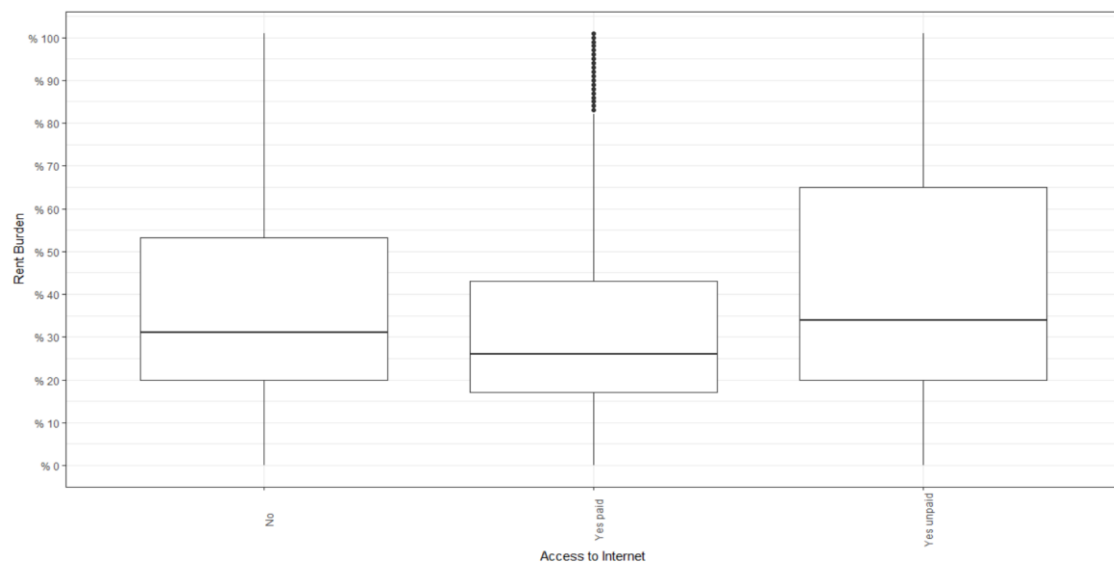
*Figure 15: ANOVA Test: Rent as a Percentage of Household Income and Internet Access*

```
                 Df    Sum Sq Mean Sq F value           Pr(>F)
ACCESS_label      2     95113   47556   67.49 <0.0000000000000002 ***
Residuals     17919 12625953     705
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

*Figure 16: Visualizing Difference in Rent Burden based on Access to Internet*

There appears to be a significant relationship between rent as a percentage of household income and internet access, with a P-value near 0 and differences up to 10% between those with paid and unpaid access to the internet.

## Regression

I first tested the results of a regression model in which rent as a percentage of household income is the dependent variable, and household income, English proficiency, presence and age of children, and internet access are the dependent variables.

*Figure 17: Model 1: Using a Regression Model with Household Income to Predict Rent as a Percentage of Household Income*

```
Residuals:
    Min      1Q  Median      3Q     Max
-48.914 -14.912  -6.386   7.270 122.405

Coefficients:
                          Estimate   Std. Error t value              Pr(>|t|)
(Intercept)            51.267005565  0.878995002  58.325 < 0.0000000000000002 ***
HINCP                  -0.000179961  0.000002254 -79.840 < 0.0000000000000002 ***
engl_profTRUE          -0.110522529  0.583561886  -0.189             0.849787
HUPAC_label0 to 6      -2.568027066  0.764662658  -3.358             0.000786 ***
HUPAC_label6 to 17     -0.830846984  0.658274075  -1.262             0.206908
HUPAC_labelNone        -2.750713400  0.588379444  -4.675          0.000002960415 ***
ACCESS_labelYes paid    1.386849786  0.606135767   2.288             0.022148 *
ACCESS_labelYes unpaid  7.296797442  1.168039841   6.247          0.000000000428 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 22.4 on 17565 degrees of freedom
Multiple R-squared:  0.2818,    Adjusted R-squared:  0.2815
F-statistic: 984.7 on 7 and 17565 DF,  p-value: < 0.00000000000000022
```

The multiple r-squared results of 0.2818 indicate that these variables likely predict about 28% of the variation in rent as a percentage of household income across the dataset.

*Figure 18: Model 1: Regression Using Log-transformed Household Income*

```
Residuals:
    Min      1Q  Median      3Q     Max
-88.566 -10.335  -1.309   8.341  66.575

Coefficients:
                       Estimate Std. Error  t value             Pr(>|t|)
(Intercept)            237.8383     1.4866  159.983 < 0.0000000000000002 ***
log(HINCP)             -19.4723     0.1366 -142.580 < 0.0000000000000002 ***
engl_profTRUE            3.8806     0.4654    8.338 < 0.0000000000000002 ***
HUPAC_label0 to 6       -1.2161     0.6075   -2.002             0.0453 *
HUPAC_label6 to 17       0.1851     0.5233    0.354             0.7236
HUPAC_labelNone         -2.5389     0.4673   -5.433          0.0000000561 ***
ACCESS_labelYes paid     8.7197     0.4872   17.898 < 0.0000000000000002 ***
ACCESS_labelYes unpaid   7.9227     0.9284    8.534 < 0.0000000000000002 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 17.8 on 17565 degrees of freedom
Multiple R-squared:  0.5463,    Adjusted R-squared:  0.5461
F-statistic:  3021 on 7 and 17565 DF,  p-value: < 0.00000000000000022
```

I wanted to improve the model fit, so I log-transformed household income in Figure 19 above. Log-transforming household income resulted in an increase the R-squared value from 0.28 to 0.55.
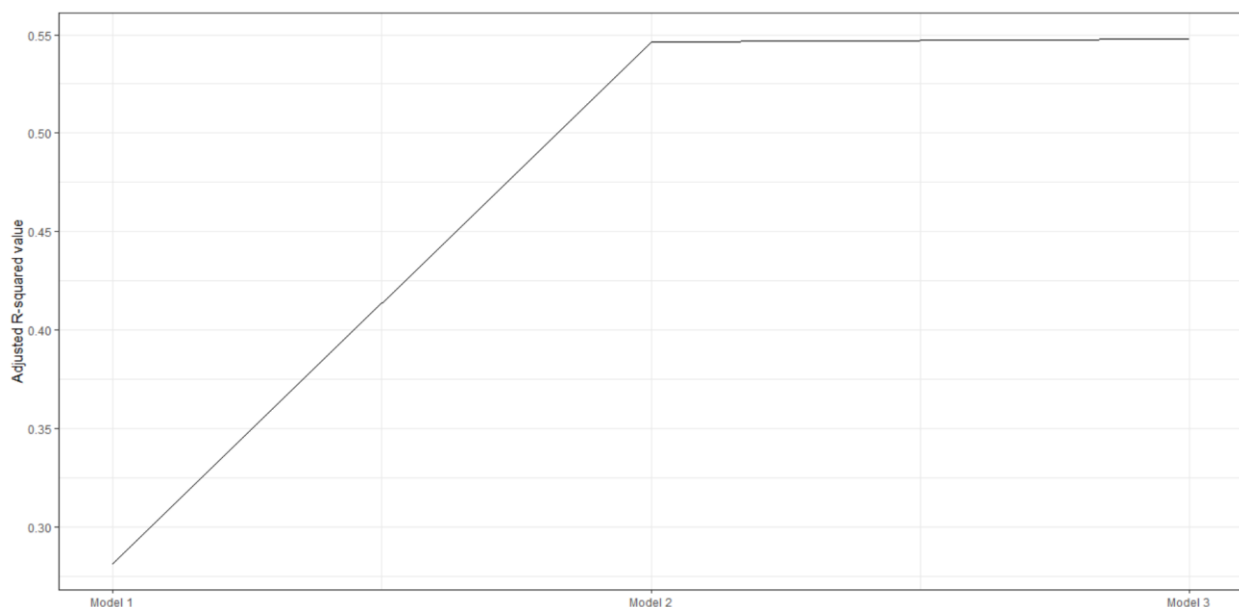
I then tested a third model, creating interaction terms between log-transformed household income and English proficiency, and log-transformed household income and internet access, further improved model fit, but only slightly. Because the interaction terms did not increase the R-squared value, I selected Model 1B as the best fit model. The results of all three model tests are shown below in Figure 20, and the relative r-squared values are illustrated in Figure 21.

*Figure 20: Comparing Model 1, 1B, and 1C Fit*

| | Model 1 | Model 1B | Model 1C |
|---|---|---|---|
| Household Income | -0.00 *** (p = 0.00) | | |
| Household Income (log-transformed) | | -19.47 *** (p = 0.00) | -15.08 *** (p = 0.00) |
| English proficiency (vs. no English proficiency) | -0.11 (p = 0.85) | 3.88 *** (p = 0.00) | 35.41 *** (p = 0.00) |
| Paid Internet access | 1.39 * (p = 0.02) | 8.72 *** (p = 0.00) | 27.36 *** (p = 0.00) |
| Unpaid Internet access | 7.30 *** (p = 0.00) | 7.92 *** (p = 0.00) | 25.07 ** (p = 0.00) |
| Presence of children age 0 to 6 | -2.57 *** (p = 0.00) | -1.22 * (p = 0.05) | -1.14 (p = 0.06) |
| Presence of children age 6 to 17 | -0.83 (p = 0.21) | 0.19 (p = 0.72) | 0.21 (p = 0.69) |
| No children | -2.75 *** (p = 0.00) | -2.54 *** (p = 0.00) | -2.31 *** (p = 0.00) |
| Interaction: log-transformed income and paid internet access | | | -1.82 *** (p = 0.00) |
| Interaction: log-transformed income and unpaid internet access | | | -1.68 * (p = 0.04) |
| Interaction: log-transformed income and English proficiency | | | -3.04 *** (p = 0.00) |
| N | 17573 | 17573 | 17573 |
| R2 | 0.28 | 0.55 | 0.55 |

*** p < 0.001; ** p < 0.01; * p < 0.05.

*Figure 21: Visualizing Model 1, 1B, and 1C Fit*

The results for the preferred regression model shown in detail in Figure 19 above revealed the following findings:

- The coefficient estimate for log-transformed household income is -19.47. This means that doubling household income is associated with an 19.47% percent decrease in rent as a percentage of household income. The P-value for log-transformed household income is extremely close to zero, suggesting that this relationship is statistically significant.
- The coefficient for English language proficiency (engl_profTRUE) is 3.88, meaning that households in which at least one person over the age of 14 speaks English very well spends 3.88 percent more of their household income on rent on average than households without English proficiency do, controlling for everything else in the model. The P-value is near zero, indicating statistical significance.
- Not all of the relationships for the categorical variable regarding presence of children were significant, but there was a significant finding for households with no children, who may spend 2.53% less of their income on rent than households with multiple children aged 0 to 17. There was also possible statistical significance for households with children 0 to 6, which appear to spend 1.22% less than households with multiple children aged 0 to 17.
- The internet access variable also appeared to be a significant predictor of rent burden. The coefficients for internet access range from 7.9 to 8.7, suggesting that households with internet access spend between about 8 or 9 percent more of their income on rent than those without internet access.

I then ran a regression model in which rent as a percentage of household income is the dependent variable, and gross monthly rent, English proficiency, presence and age of children, and internet access are the dependent variables, because I wanted to understand the extent to which rent might predict rent as a percentage of household income when controlling for other factors. The results of this regression are shown in Figure 22 below.

*Figure 22: Model 2 Regression Model with Gross Monthly Rent to Predict Rent as a Percentage of Household Income*

```
Call:
lm(formula = GRPIP ~ GRNTP + engl_prof + HUPAC_label + ACCESS_label,
    data = hh_data)

Residuals:
    Min      1Q  Median      3Q     Max
-43.867 -17.679  -8.530   8.728  73.382

Coefficients:
                          Estimate Std. Error t value             Pr(>|t|)
(Intercept)             48.4202078  1.0330413  46.872 < 0.0000000000000002 ***
GRNTP                    0.0035508  0.0002486  14.281 < 0.0000000000000002 ***
engl_profTRUE           -6.6365058  0.6773636  -9.798 < 0.0000000000000002 ***
HUPAC_label0 to 6       -7.4030176  0.8855006  -8.360 < 0.0000000000000002 ***
HUPAC_label6 to 17      -1.4016725  0.7640379  -1.835             0.066588 .
HUPAC_labelNone         -5.6846320  0.6818269  -8.337 < 0.0000000000000002 ***
ACCESS_labelYes paid    -8.2897589  0.7140081 -11.610 < 0.0000000000000002 ***
ACCESS_labelYes unpaid   4.7415569  1.3579152   3.492             0.000481 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 26 on 17565 degrees of freedom
Multiple R-squared:  0.03243,    Adjusted R-squared:  0.03204
F-statistic:  84.1 on 7 and 17565 DF,  p-value: < 0.00000000000000022
```

The r-squared value was very low for this model, at 0.03, meaning that these variables can predict only 3% of the variation in rent as a percentage of household income, and that rent itself is a poor predictor of rent as a percentage of household income compared with household income. I went through the refinement process for this regression model, log-transforming gross monthly rent and testing interaction terms between gross monthly rent and English proficiency, and gross monthly rent and internet access. The R-squared value for the preferred model only reached 0.05, so I did not explore the results further, and focused on analyzing the results of the base model as supplemental to the results of the regression models that included household income. The model suggests that a $1,000 increase in gross monthly rent is associated with spending 3.6 percent more of household income on rent.

## Discussion

The results of these statistical tests begin to shed light on how Suffolk County renter households vary in terms of percentage of income spent on rent. The results were varied in terms of support for my hypothesis.

### Hypotheses Revisited
1.  Higher income is correlated with lower rent burden: TRUE
2.  Higher rent is correlated with higher rent burden: TRUE
3.  English proficiency is correlated with lower rent burden: FALSE
4.  Internet access is correlated with lower rent burden: FALSE
5.  Presence of any number of children is correlated with higher rent burden as compared with no children: TRUE AND FALSE

### Implications
In analyzing the key takeaways my results, first, and perhaps obviously, it is important to note that lower-income households experience higher rent burden. Second, the weak correlation between rent and rent burden demonstrate that renters across the rental market are cost-burdened: not just those living in the County's least expensive or most expensive units. Third, the lack of another amenity such as internet access may not predict higher rent burden. Fourth, being a native English speaking household is associated with higher rent burden, so non-English speaking households do not here appear to be an at-risk group for rent burden when controlling for other factors. Finally, having multiple children ages 0 to 17 predicts higher rent burden when compared to households with no children of only children ages 0 to 6.

More generally, it is notable that with a distribution test (Figure 6 and Figure 7) we learned that the average value of rent as a percentage of income spent on rent is 30%, which is the threshold for rent burden as defined by HUD.[2] By reviewing the interquartile range, we know that 25% of households are spending more than 38% of their income on rent, and the right skew of the histogram illustrates that many these top quartile households are in fact spending more than 50% of their income on rent. The distribution, in combination with HUD's definitions of rent burden and severe rent burden, raises the flag there is a high prevalence of both rent burden and severe rent burden in Suffolk County.

It is important to note that the results of this analysis are limited. The sample size for this study was only 17,922 renter households, whereas Suffolk County had 196,676 renter households according to the 2018 ACS 5-Year Estimate. With a sample size representing less than 10 percent of all study area households,

---

[2] A cost-burdened household is defined by the Department of Housing and Urban Development as households "who pay more than 30 percent of their income for monthly housing costs" and "may have difficulty affording necessities such as food, clothing, transportation, and medical care."

this sample may not be representative. It is also limited in that it does differentiate between neighborhoods or subareas within the County, which may have trends that do not reflect the County as a whole.

For future studies, I would be interested in investigating in more detail the larger, lower-income families spending more of their income on rent, as the results of my study raise questions around the affordability and availability of larger family-sized rental units in Suffolk County. I would be interested in looking at both bedroom size, structure type (single-family homes vs. attached and multifamily homes), and the year structures were built to study the relationship between housing inventory and rent burden.

## Conclusion

The study of rent burden is central to understanding and mitigating displacement risk for renter households amidst a nationwide shortage of affordable rental housing. Since 2000, the share of renters who are rent-burdened has risen dramatically nationally (from 39 to 50 percent overall). At least half of renters are currently rent burdened in 23 of the 100 largest cities.

While rent burden is an important measure of housing affordability and an area's economic vitality, it can often be difficult to characterize rent-burdened households: for example, a young professional living alone with no children earning $60,000 per year and spending $2,500 on monthly rent for a luxury one-bedroom apartment is spending the same percentage of income on rent (50%) as a single parent of three children earning $30,000 spending $1,250 on monthly rent for a three-bedroom. In order to better serve rent-burdened populations, it is important to understand certain other household characteristics that make a household more vulnerable to rent burden. I hope that this study highlights some of the nuanced elements of rent burden Suffolk County.