

- 01 Introduction
- 02 Data
- 03 Methods & Results
- 04 Discussion
- 05 Conclusion

- 01 Introduction
- 02 Data
- 03 Methods & Results
- 04 Discussion
- 05 Conclusion

Introduction

Research Questions and Hypothesis

Does a person's household income, race, or educational attainment have a relationship to the percentage of their income they pay toward rent each month?

My hypotheses:

People in minority racial groups (people who do not identify as white) have a higher GRPIP (gross rent as a percentage of household income.)

People with a lower level of education have a higher GRPIP (gross rent as a percentage of household income.)

Lower-income households have higher GRPIP (gross rent as a percentage of household income.)

- 01 Introduction
- 02 Data
- 03 Methods & Results
- 04 Discussion
- 05 Conclusion

Defining Sample Set

The unit of analysis is **person-level data** from the **2018 American Community Survey**.

The definition of the sample population is people **over in the age of 18 in rental households** in the state of Massachusetts.

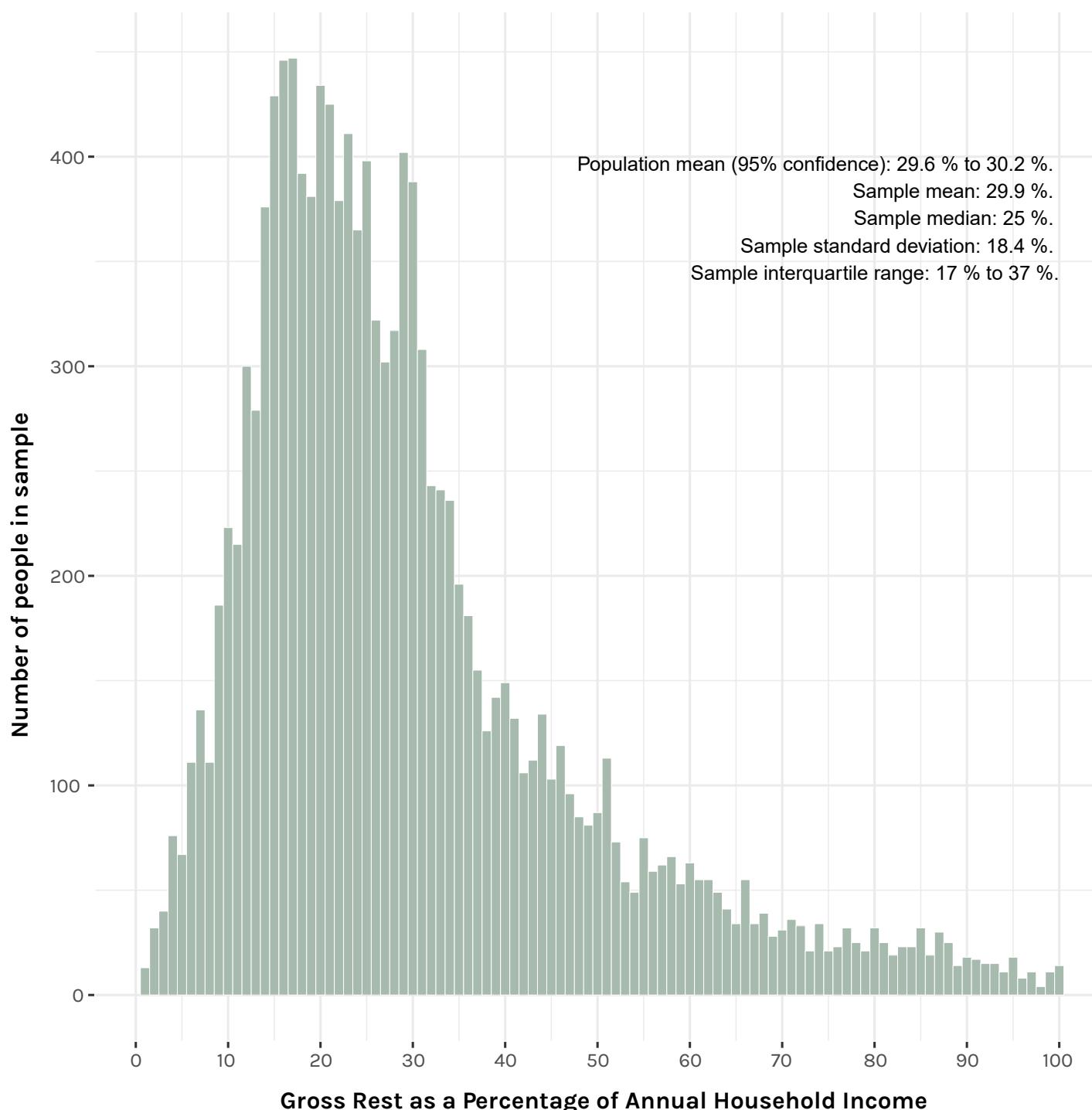
For the purposes of this exercise, I am defining a "rent-burdened" individual as someone who spends **30% or more of their household income on rent**. I used 30% because this is the percentage HUD uses to define cost-burdened families.

Data

Distributions and Proportions

Gross Rent as a Percentage of Household Income

Key takeaway: the population mean of about 30% tells us that, on average, renters in MA are right at the threshold of being rent burdened.

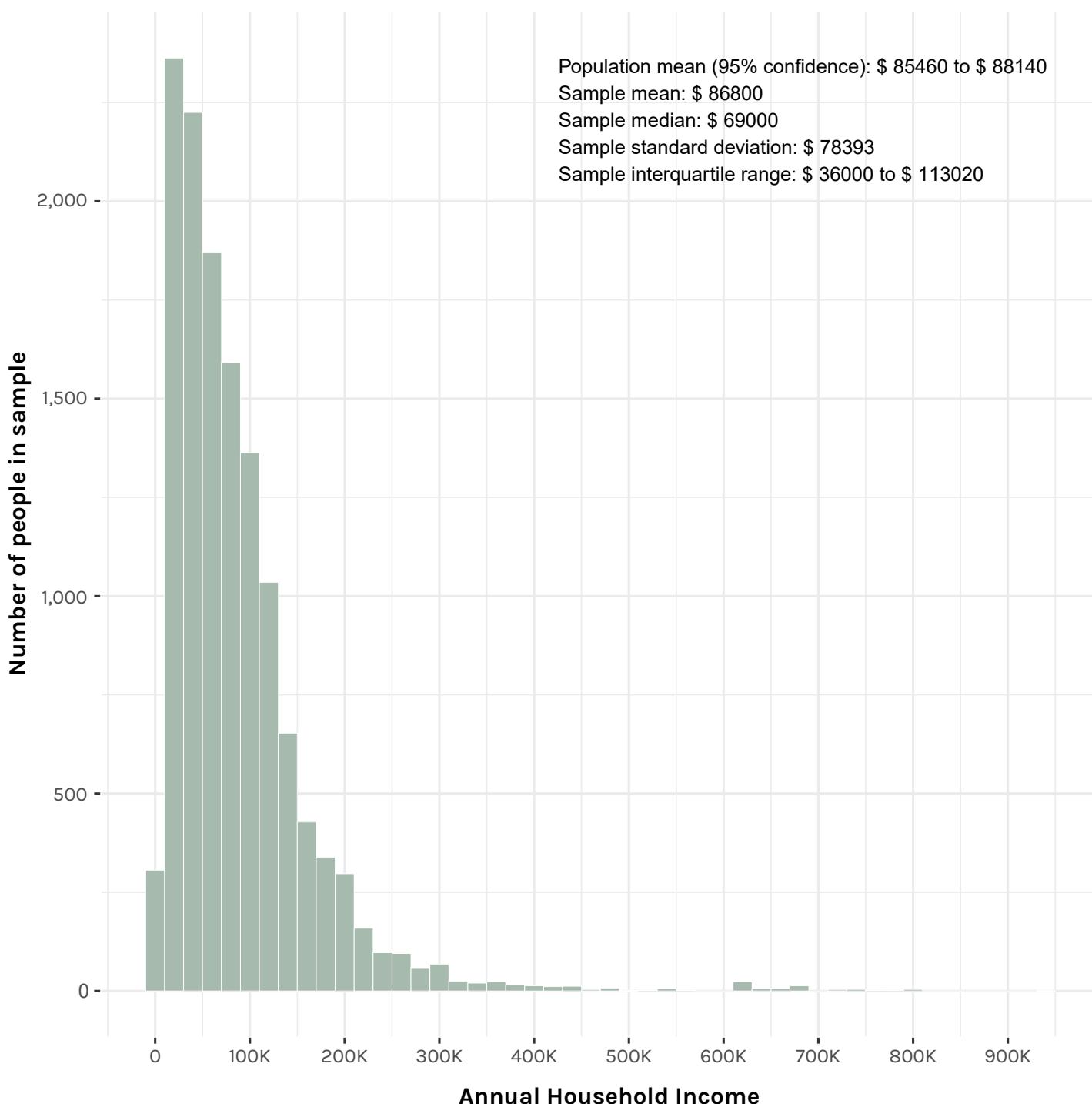


Data

Distributions and Proportions

Annual Household Income

Key takeaway: the population mean of \$85,460 - \$88,140 median of \$69,000 is high compared to the median household income in the U.S. of \$60,000.

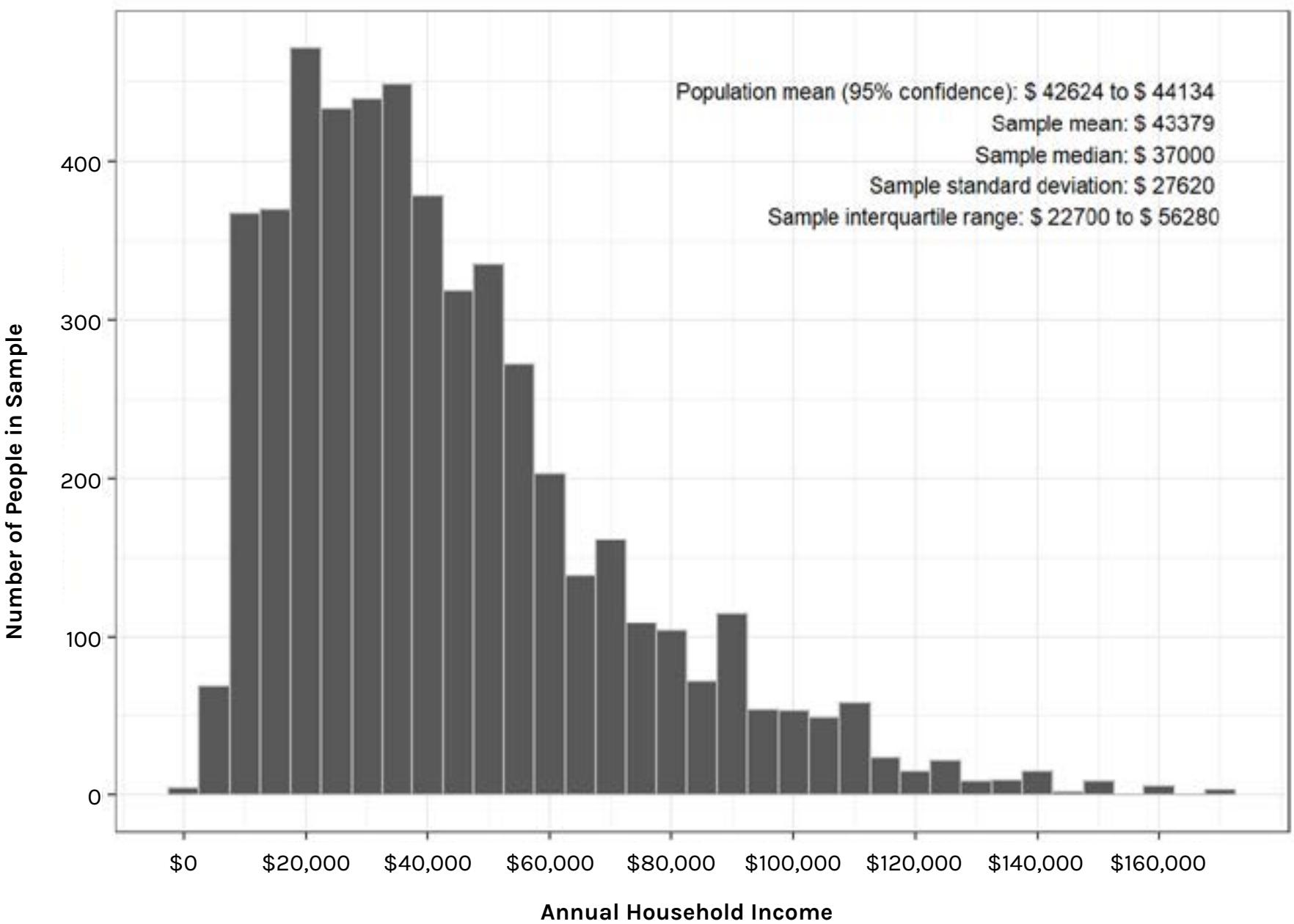


Data

Distributions and Proportions

Annual Household Income (of rent burdened people)

Key takeaway: when you look at the population of only rent burdened individuals ($GRPIP > 30\%$), the mean and median income decreases drastically.

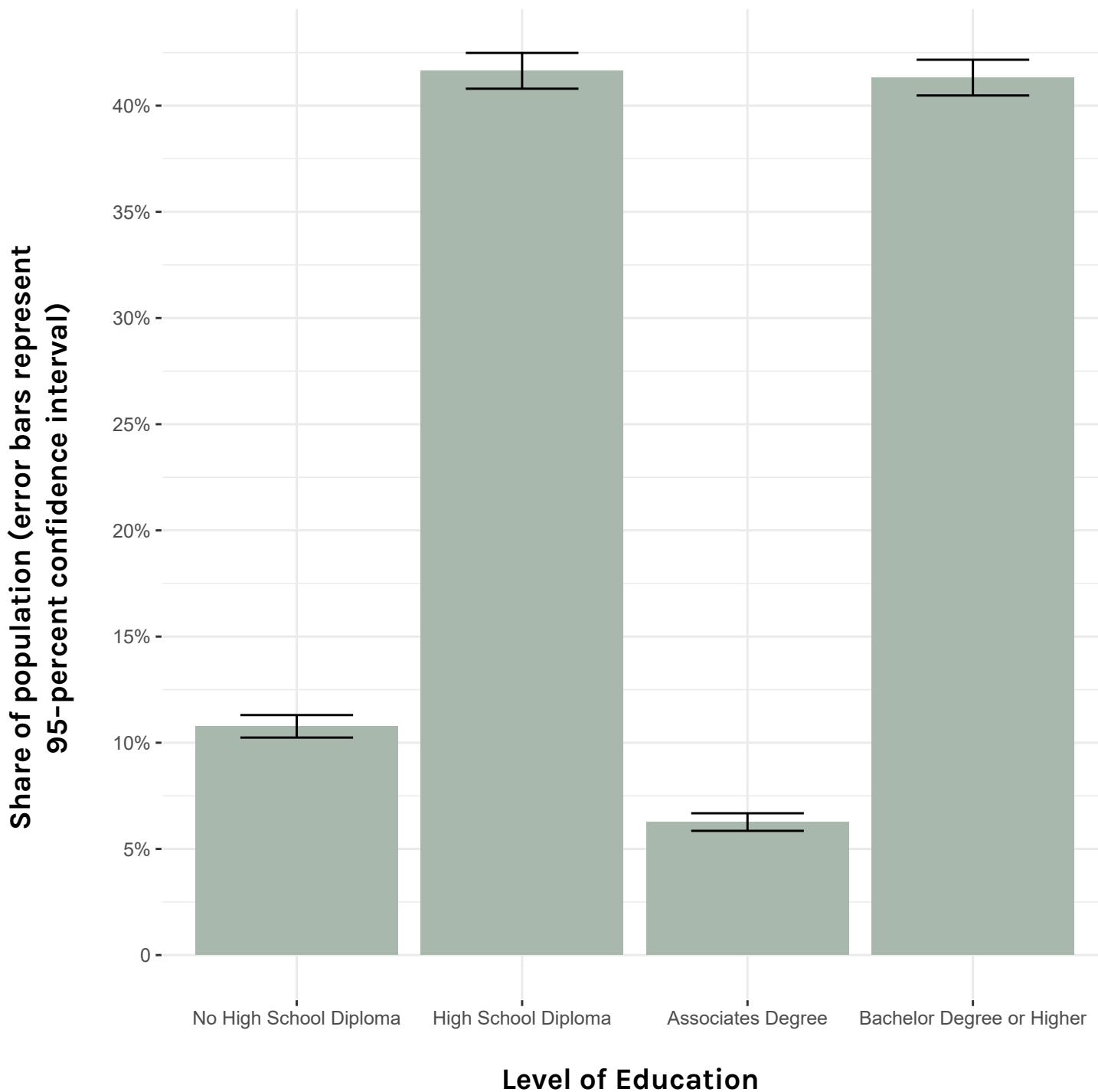


Data

Distributions and Proportions

Level of Education

Key takeaway: these are broad categories I designated by coding the original ACS variable. They are probably too broad and I would want to break them down more in the future. Over 40% of the population has a Bachelor's degree or higher, which is slightly higher than the national average of 33%.

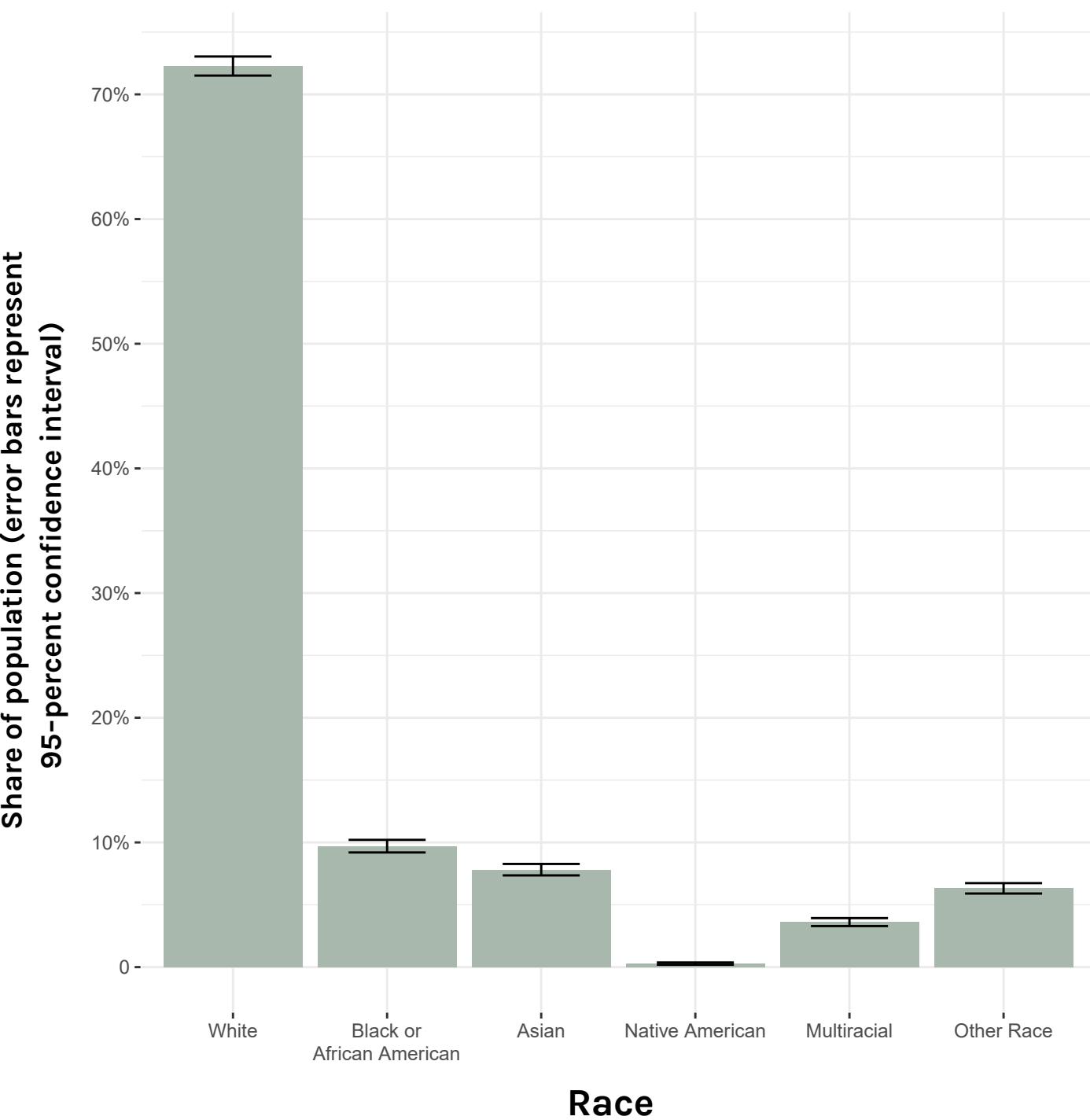


Data

Distributions and Proportions

Race

Key takeaway: these are broad categories I designated by coding the original ACS variable. Over 70% of the population is white alone. I didn't realize this, but this is fairly representative of Massachusetts as a whole. Census.gov reports 80% of Massachusetts as white alone.



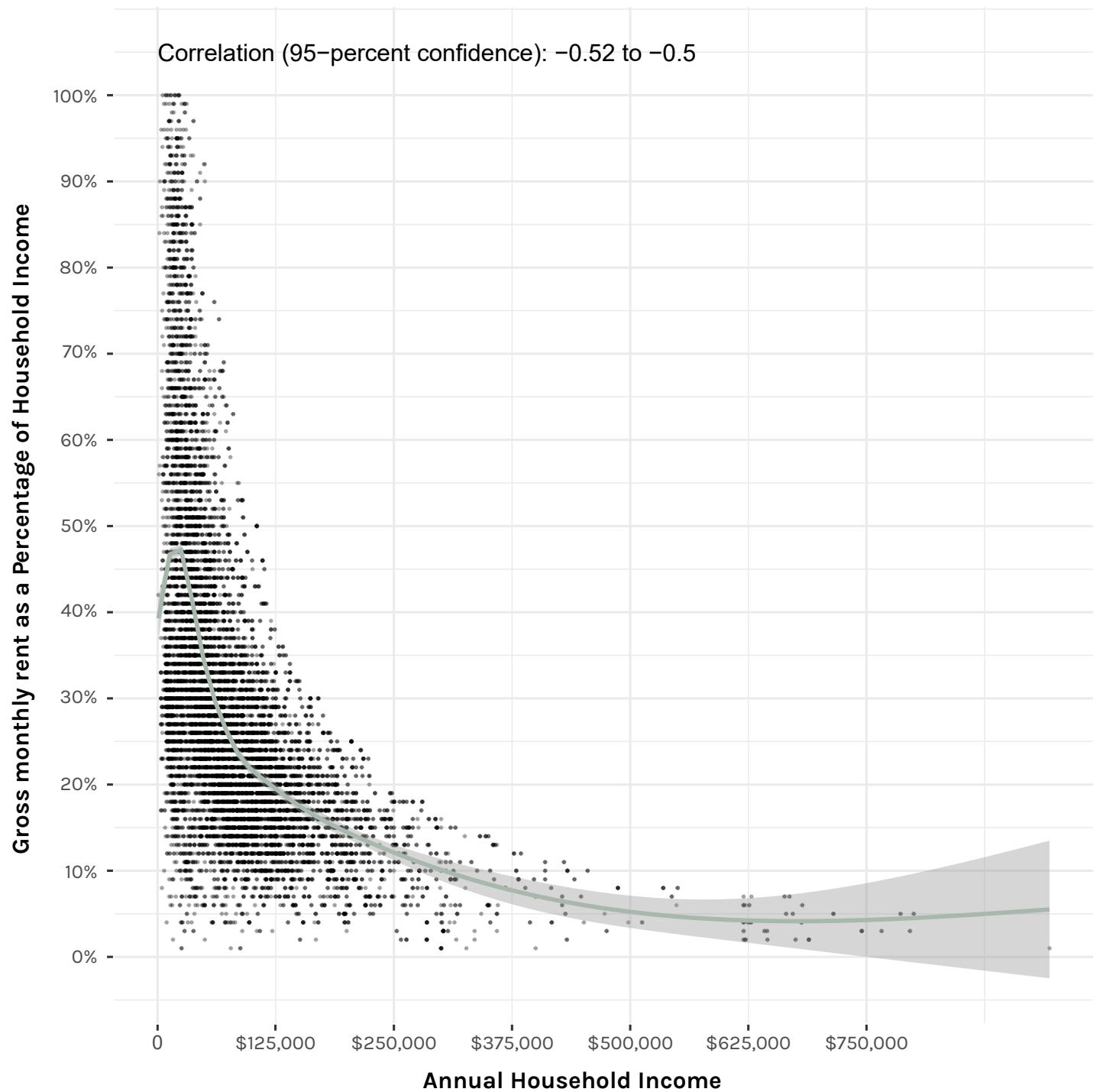
- 01 Introduction
- 02 Data
- 03 Methods & Results
- 04 Discussion
- 05 Conclusion

Methods & Results

Pearson's Test

Is there a correlation between Household Income and Gross Rent as a Percentage of Household Income?

This correlation is only moderate (around 0.5) but statistically significant. There is a negative correlation at a 95% confidence level.

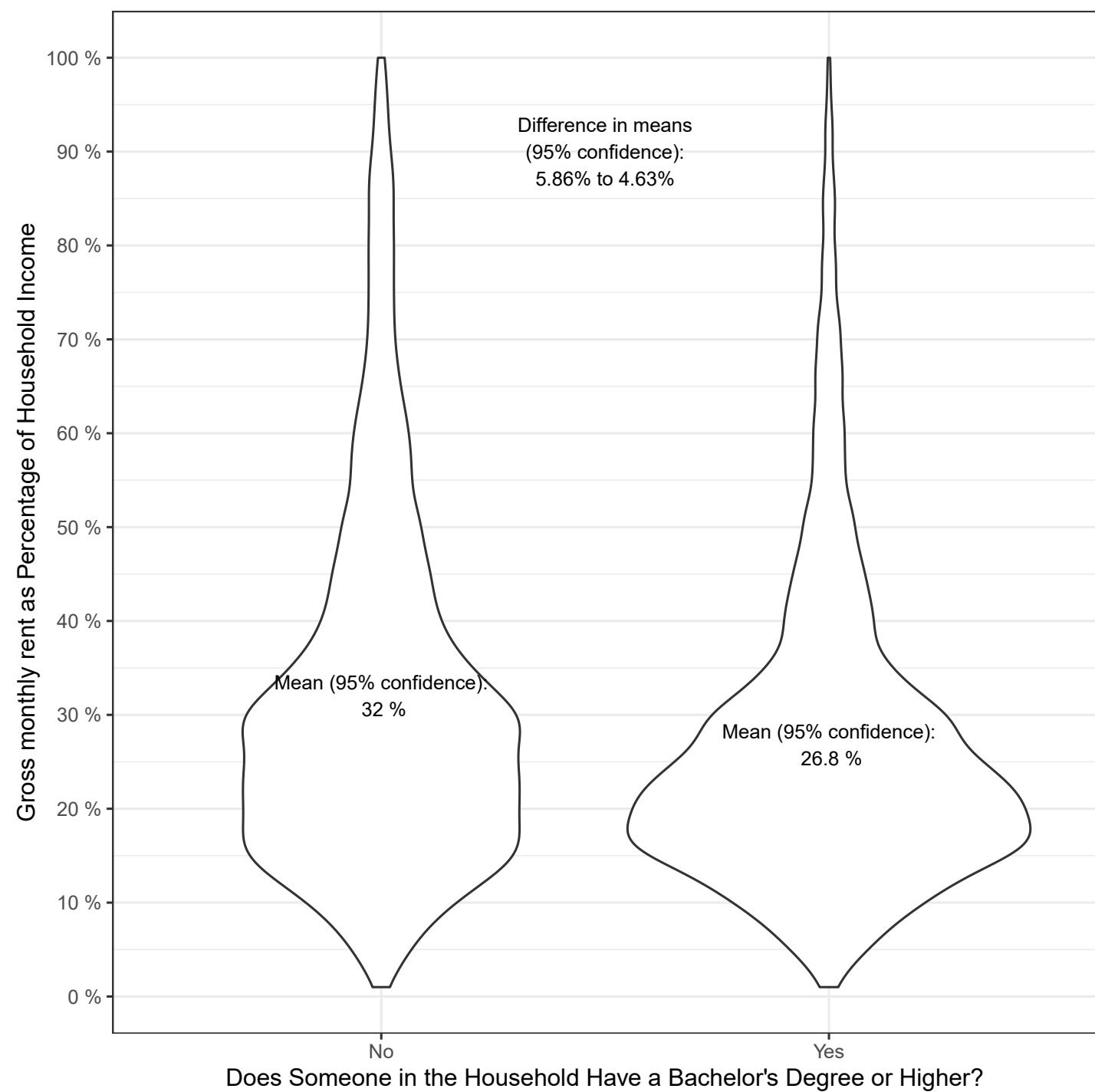


Methods & Results

Two-Sample t-tests

Are people with a Bachelor's Degree less likely to live in rent burdened households?

There is a significant relationship between these two variables. We can be 95% confident that people in this sample set with a Bachelor's degree or higher live in households that pay 5% - 6% less of their income on rent than those without a Bachelor's degree. It looks like many people with a college degree live in households that pay less than 20% of their income on rent.

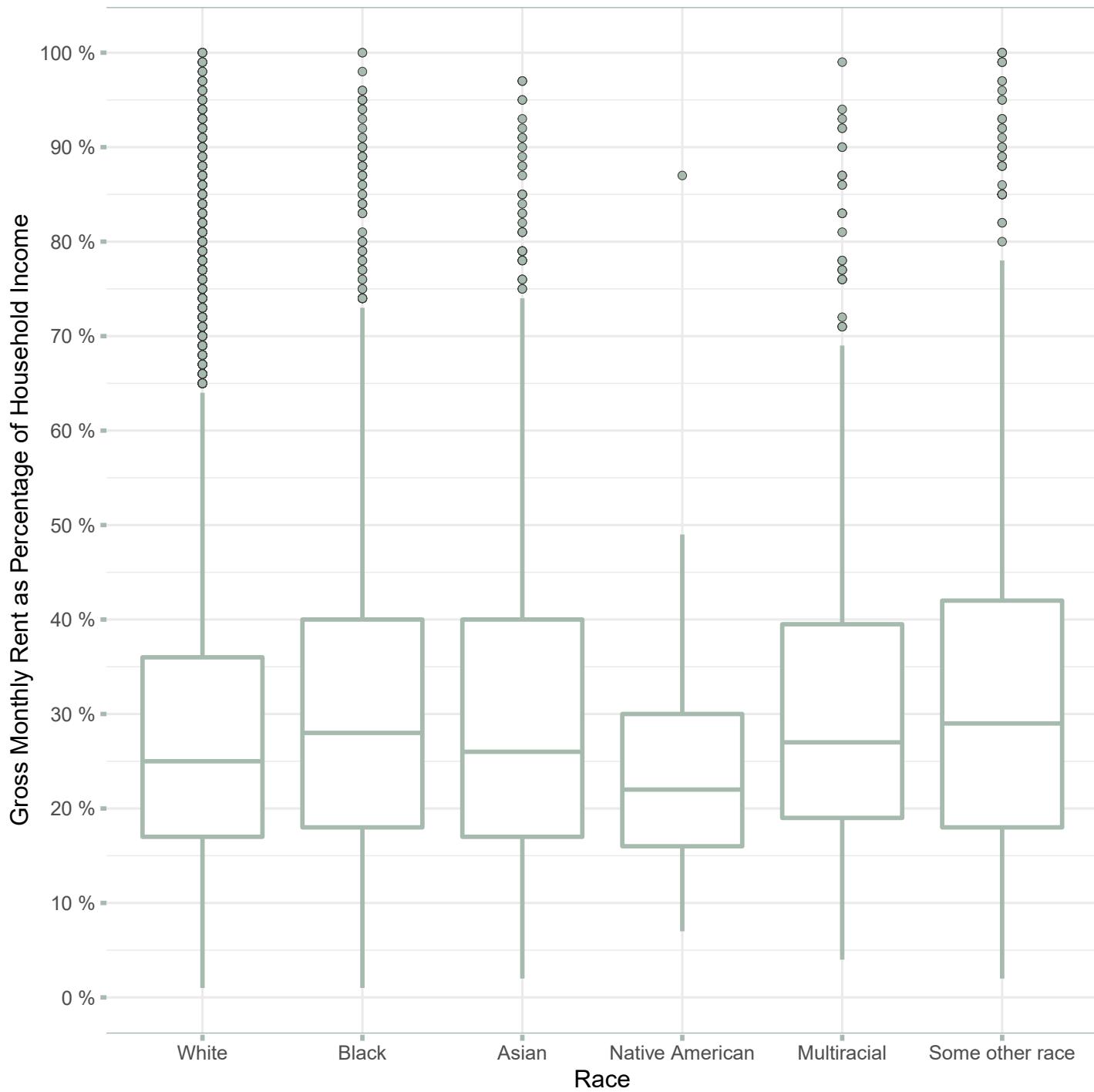


Methods & Results

Tukey / ANOVA

Is there a correlation between Gross Rent as a Percentage of Household Income and Race?

There is a honestly significant relationship between these variables. Here is their distribution shown in a box box. The upper and lower quartiles of each racial category are not dramatically different. The lowest median GRPIP is seen in the Native American or Pacific Islander population, followed by the White population.



Methods & Results

Linear Regression | First Model

My dependant variable is Gross Rent as a Percentage of Household Income.

My independant variables are Educational Attainment, Race, Age and Number of People in the Household.

The **multiple R-squared value for this model is 0.044**, which tells us that this model predicts about 4% of the variation in gross rent as a percentage of household income in this dataset.

```
##  
## Call:  
## lm(formula = GRPIP ~ EDU + RACE + AGEP + NP, data = person_data)  
##  
## Residuals:  
##      Min       1Q   Median       3Q      Max  
## -36.125 -11.969  -4.419   6.972  75.003  
##  
## Coefficients:  
##                               Estimate Std. Error t value Pr(>|t|)  
## (Intercept)            39.887140  0.868523 45.925 < 2e-16 ***  
## EDUAssoc. Degree     -5.229010  0.797713 -6.555 5.77e-11 ***  
## EDUBach Degree or Higher -8.695419  0.575215 -15.117 < 2e-16 ***  
## EDUHS Diploma        -3.146645  0.547474 -5.748 9.25e-09 ***  
## RACEAsian             3.121303  0.596580  5.232 1.70e-07 ***  
## RACEBlack              2.785634  0.542295  5.137 2.84e-07 ***  
## RACEMultiracial       2.393221  0.851052  2.812  0.00493 **  
## RACENative Am.        -4.890339  2.961429 -1.651  0.09869 .  
## RACEOther Race         2.610134  0.670054  3.895 9.85e-05 ***  
## AGEP                  -0.014633  0.009619 -1.521  0.12824  
## NP                    -1.884364  0.121329 -15.531 < 2e-16 ***  
## ---  
## Signif. codes:  0 '****' 0.001 '***' 0.01 '**' 0.05 '*' 0.1 '.' 1  
##  
## Residual standard error: 17.97 on 13142 degrees of freedom  
## Multiple R-squared:  0.04425,    Adjusted R-squared:  0.04352  
## F-statistic: 60.85 on 10 and 13142 DF,  p-value: < 2.2e-16
```

Methods & Results

Linear Regression | First Model

I am comparing **education levels** to the least amount of education: No High School Diploma. Compared to a household with no high school diploma:

A household with someone with a high school diploma pays, on average, 3.2 percentage points less on rent as a percentage of their household income.

A household with someone with an associate's degree pays, on average, 5.2 percentage points less on rent as a percentage of their household income.

A household with someone with a Bachelor's degree or higher pays, on average, 8.7 percentage points less on rent as a percentage of their household income.

These coefficients are all statistically significant.

```
##  
## Call:  
## lm(formula = GRPIP ~ EDU + RACE + AGEP + NP, data = person_data)  
##  
## Residuals:  
##      Min       1Q   Median       3Q      Max  
## -36.125 -11.969  -4.419   6.972  75.003  
##  
## Coefficients:  
##                               Estimate Std. Error t value Pr(>|t|)  
## (Intercept)            39.887140  0.868523 45.925 < 2e-16 ***  
## EDUAssoc. Degree     -5.229010  0.797713 -6.555 5.77e-11 ***  
## EDUBach Degree or Higher -8.695419  0.575215 -15.117 < 2e-16 ***  
## EDUHS Diploma        -3.146645  0.547474 -5.748 9.25e-09 ***  
## RACEAsian             3.121303  0.596580  5.232 1.70e-07 ***  
## RACEBlack              2.785634  0.542295  5.137 2.84e-07 ***  
## RACEMultiracial       2.393221  0.851052  2.812  0.00493 **  
## RACENative Am.         -4.890339  2.961429 -1.651  0.09869 .  
## RACEOther Race         2.610134  0.670054  3.895 9.85e-05 ***  
## AGEP                  -0.014633  0.009619 -1.521  0.12824  
## NP                   -1.884364  0.121329 -15.531 < 2e-16 ***  
## ---  
## Signif. codes:  0 '****' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1  
##  
## Residual standard error: 17.97 on 13142 degrees of freedom  
## Multiple R-squared:  0.04425,    Adjusted R-squared:  0.04352  
## F-statistic: 60.85 on 10 and 13142 DF,  p-value: < 2.2e-16
```

Methods & Results

Linear Regression | First Model

I am comparing race to the most common race in my sample set: White alone. Compared to a White household:

A household with at least one person who is **Asian** pays, on average, **3.1 percentage points more** on rent as a percentage of household income.

A household with at least one person who is **Black or African American** pays, on average, **2.8 percentage points more** on rent as a percentage of household income.

A household with at least one person who is **multiracial** pays, on average, **2.4 percentage points more** on rent as a percentage of household income.

A household with at least one person who is **some other race** pays, on average, **2.6 percentage points more** on rent as a percentage of household income.

These coefficients are statistically significant. The Native American / Pacific Islander variable was not.

```
##  
## Call:  
## lm(formula = GRPIP ~ EDU + RACE + AGEP + NP, data = person_data)  
##  
## Residuals:  
##    Min     1Q   Median     3Q    Max  
## -36.125 -11.969  -4.419   6.972  75.003  
##  
## Coefficients:  
##                                         Estimate Std. Error t value Pr(>|t|)  
## (Intercept)                   39.887140  0.868523 45.925 < 2e-16 ***  
## EDUAssoc. Degree          -5.229010  0.797713 -6.555 5.77e-11 ***  
## EDUBach Degree or Higher -8.695419  0.575215 -15.117 < 2e-16 ***  
## EDUHS Diploma            -3.146645  0.547474 -5.748 9.25e-09 ***  
## RACEAsian                  3.121303  0.596580  5.232 1.70e-07 ***  
## RACEBlack                  2.785634  0.542295  5.137 2.84e-07 ***  
## RACEMultiracial           2.393221  0.851052  2.812  0.00493 **  
## RACENative Am.             -4.890339  2.961429 -1.651  0.09869 .  
## RACEOther Race            2.610134  0.670054  3.895 9.85e-05 ***  
## AGEP                      -0.014633  0.009619 -1.521  0.12824  
## NP                        -1.884364  0.121329 -15.531 < 2e-16 ***  
## ---  
## Signif. codes:  0 '****' 0.001 '***' 0.01 '**' 0.05 '*' 0.1 '.' 1  
##  
## Residual standard error: 17.97 on 13142 degrees of freedom  
## Multiple R-squared:  0.04425,    Adjusted R-squared:  0.04352  
## F-statistic: 60.85 on 10 and 13142 DF,  p-value: < 2.2e-16
```

Methods & Results

Linear Regression | Best Fit

I was able to get a better model fit by doing three things:

1. Added Household Income (HINCP) to the regression.
2. Converted Age from a continuous variable to a categorical variable. In this case, I have two groups: 30 years old or younger and over 30 years old.
3. Added an interaction term between Race and Household Income.

I am hypothesizing that Race has an effect on Household Income when considering how much of their income one pays on rent.

The R-squared value increased from 0.04 to 0.27, which tells us that this model predicts about 27% of the variation in gross rent as a percentage of household income in this dataset.

```
##  
## Call:  
## lm(formula = GRPIP ~ HINCP + EDU + RACE + age_cat + NP + HINCP:RACE,  
##      data = person_data)  
##  
## Residuals:  
##      Min       1Q   Median       3Q      Max  
## -37.857  -9.988  -3.553   5.581  87.821  
##  
## Coefficients:  
##                               Estimate Std. Error t value Pr(>|t|)  
## (Intercept)                3.996e+01  5.277e-01  75.729 < 2e-16 ***  
## HINCP                  -1.217e-04  2.241e-06 -54.306 < 2e-16 ***  
## EDUAssoc. Degree        -1.357e+00  6.972e-01 -1.947 0.051615 .  
## EDUBach Degree or Higher 5.958e-01  5.117e-01  1.164 0.244333  
## EDUHS Diploma            -9.521e-01  4.762e-01 -1.999 0.045575 *  
## RACEAsian                 1.632e+00  7.672e-01  2.127 0.033441 *  
## RACEBlack                 4.566e-01  6.511e-01  0.701 0.483162  
## RACEMultiracial          4.553e+00  1.173e+00  3.882 0.000104 ***  
## RACENative Am.           -3.341e+00  4.463e+00 -0.749 0.454044  
## RACEOther Race            6.796e+00  9.019e-01  7.535 5.21e-14 ***  
## age_catyounger            1.077e+00  3.000e-01  3.591 0.000331 ***  
## NP                      -2.404e-02  1.050e-01 -0.229 0.818846  
## HINCP:RACEAsian           1.169e-05  5.541e-06  2.110 0.034904 *  
## HINCP:RACEBlack            8.482e-06  6.029e-06  1.407 0.159540  
## HINCP:RACEMultiracial     -4.175e-05  1.144e-05 -3.649 0.000264 ***  
## HINCP:RACENative Am.      -1.732e-05  4.448e-05 -0.389 0.696926  
## HINCP:RACEOther Race      -9.305e-05  1.066e-05 -8.728 < 2e-16 ***  
## ---  
## Signif. codes:  0 '****' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1  
##  
## Residual standard error: 15.7 on 13136 degrees of freedom  
## Multiple R-squared:  0.2713  Adjusted R-squared:  0.2704  
## F-statistic: 305.6 on 16 and 13136 DF,  p-value: < 2.2e-16
```

Methods & Results

Linear Regression | Best Fit

When looking at educational attainment, the coefficients have changed quite a bit.

One coefficient is statistically significant. In this model, a household with someone with a high school diploma pays, on average, 0.9 percentage points less on rent as a percentage of their household income than someone without a high school diploma.

The other coefficients are not statistically significant. We can assume one reason these results have changed is that we are now controlling for income. Previously, the coefficients could possibly be attributed to the fact that people with higher education have higher incomes.

Now, controlling for income, a household with someone with a Bachelor's Degree or higher pays, on average, 0.6 percentage points more on rent as a percentage of their household income than someone without a high school diploma. But again, the results are not statistically significant.

```
##  
## Call:  
## lm(formula = GRPIP ~ HINCP + EDU + RACE + age_cat + NP + HINCP:RACE,  
##      data = person_data)  
##  
## Residuals:  
##      Min       1Q   Median       3Q      Max  
## -37.857  -9.988  -3.553   5.581  87.821  
##  
## Coefficients:  
##                               Estimate Std. Error t value Pr(>|t|)  
## (Intercept)                3.996e+01  5.277e-01  75.729 < 2e-16 ***  
## HINCP                  -1.217e-04  2.241e-06 -54.306 < 2e-16 ***  
## EDUAssoc. Degree        -1.357e+00  6.972e-01 -1.947 0.051615 .  
## EDUBach Degree or Higher 5.958e-01  5.117e-01  1.164 0.244333  
## EDUHS Diploma           -9.521e-01  4.762e-01 -1.999 0.045575 *  
## RACEAsian                1.632e+00  7.672e-01  2.127 0.033441 *  
## RACEBlack                4.566e-01  6.511e-01  0.701 0.483162  
## RACEMultiracial          4.553e+00  1.173e+00  3.882 0.000104 ***  
## RACENative Am.           -3.341e+00  4.463e+00 -0.749 0.454044  
## RACEOther Race           6.796e+00  9.019e-01  7.535 5.21e-14 ***  
## age_catyounger           1.077e+00  3.000e-01  3.591 0.000331 ***  
## NP                      -2.404e-02  1.050e-01 -0.229 0.818846  
## HINCP:RACEAsian          1.169e-05  5.541e-06  2.110 0.034904 *  
## HINCP:RACEBlack           8.482e-06  6.029e-06  1.407 0.159540  
## HINCP:RACEMultiracial    -4.175e-05  1.144e-05 -3.649 0.000264 ***  
## HINCP:RACENative Am.     -1.732e-05  4.448e-05 -0.389 0.696926  
## HINCP:RACEOther Race     -9.305e-05  1.066e-05 -8.728 < 2e-16 ***  
## ---  
## Signif. codes:  0 '****' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1  
##  
## Residual standard error: 15.7 on 13136 degrees of freedom  
## Multiple R-squared:  0.2713, Adjusted R-squared:  0.2704  
## F-statistic: 305.6 on 16 and 13136 DF, p-value: < 2.2e-16
```

Methods & Results

Linear Regression | Best Fit

Let's take a look at Household Income and the interaction term between Household Income and Race.

The negative and statistically significant coefficient for Household Income (HINCP) tells us that as one's income increases, they pay fewer percentage points of their household income on rent. This is controlling for all other variables.

When looking at the interaction term, I am still comparing race to the most common race in my sample set: White alone. Not all coefficients are statistically significant, but let's take a look at a plot to better understand...

```
##  
## Call:  
## lm(formula = GRPIP ~ HINCP + EDU + RACE + age_cat + NP + HINCP:RACE,  
##      data = person_data)  
##  
## Residuals:  
##    Min      1Q  Median      3Q     Max  
## -37.857 -9.988 -3.553  5.581 87.821  
##  
## Coefficients:  
##                               Estimate Std. Error t value Pr(>|t|)  
## (Intercept)            3.996e+01 5.277e-01 75.729 < 2e-16 ***  
## HINCP                -1.217e-04 2.241e-06 -54.306 < 2e-16 ***  
## EDUAssoc. Degree     -1.357e+00 6.972e-01 -1.947 0.051615 .  
## EDUBach Degree or Higher 5.958e-01 5.117e-01  1.164 0.244333  
## EDUHS Diploma        -9.521e-01 4.762e-01 -1.999 0.045575 *  
## RACEAsian             1.632e+00 7.672e-01  2.127 0.033441 *  
## RACEBlack              4.566e-01 6.511e-01  0.701 0.483162  
## RACEMultiracial       4.553e+00 1.173e+00  3.882 0.000104 ***  
## RACENative Am.        -3.341e+00 4.463e+00 -0.749 0.454044  
## RACEOther Race         6.796e+00 9.019e-01  7.535 5.21e-14 ***  
## age_catyounger        1.077e+00 3.000e-01  3.591 0.000331 ***  
## NP                   -2.404e-02 1.050e-01 -0.229 0.818846  
## HINCP:RACEAsian        1.169e-05 5.541e-06  2.110 0.034904 *  
## HINCP:RACEBlack         8.482e-06 6.029e-06  1.407 0.159540  
## HINCP:RACEMultiracial -4.175e-05 1.144e-05 -3.649 0.000264 ***  
## HINCP:RACENative Am. -1.732e-05 4.448e-05 -0.389 0.696926  
## HINCP:RACEOther Race -9.305e-05 1.066e-05 -8.728 < 2e-16 ***  
## ---  
## Signif. codes:  0 '****' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1  
##  
## Residual standard error: 15.7 on 13136 degrees of freedom  
## Multiple R-squared:  0.2713, Adjusted R-squared:  0.2704  
## F-statistic: 305.6 on 16 and 13136 DF, p-value: < 2.2e-16
```

Methods & Results

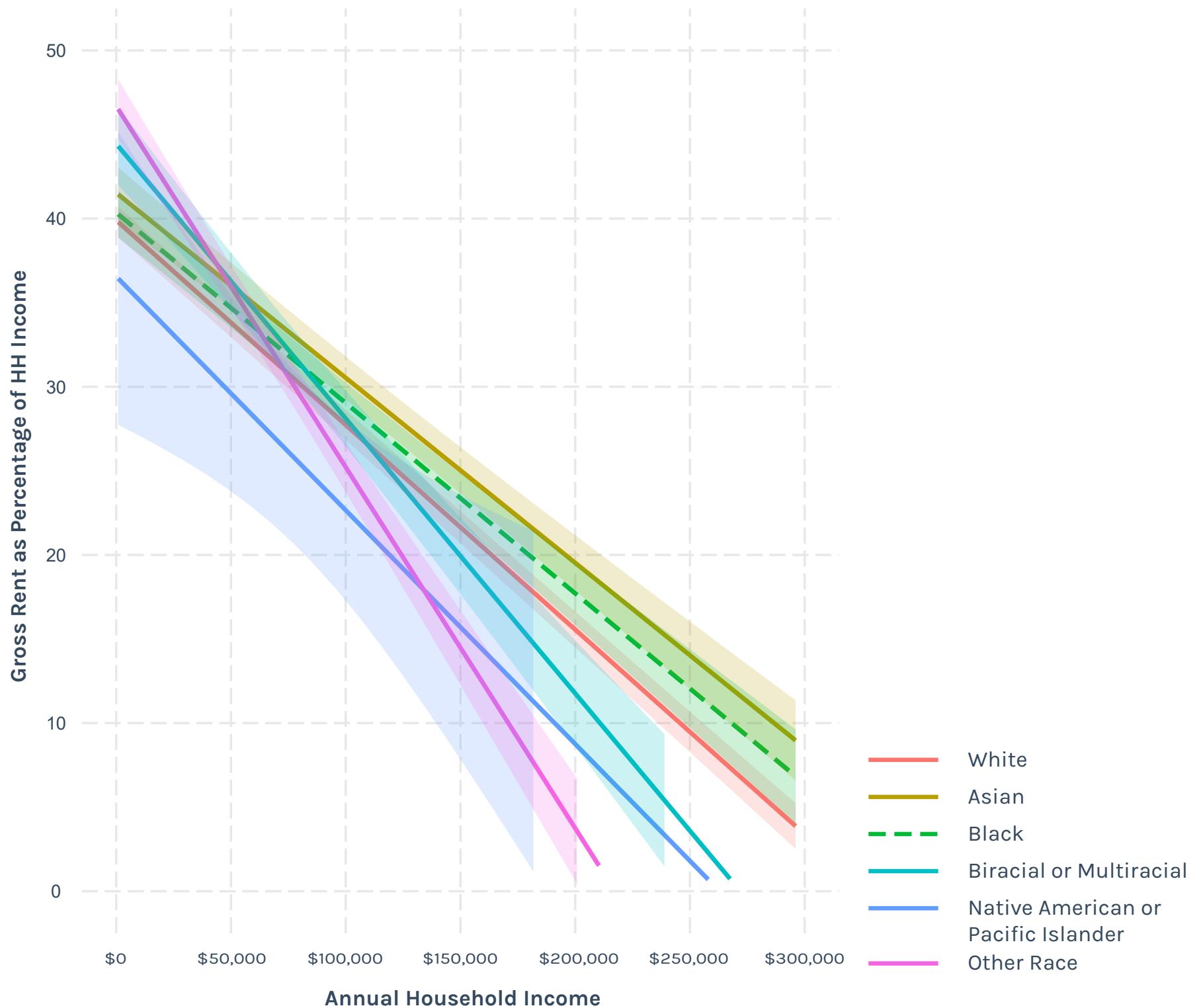
Linear Regression | Best Fit

When looking at middle-income households (renters who make around \$50,000), there is not much of a difference in GRPIP between the different racial categories. Everyone pays around 35%-38% of their income on rent. The real differences in the effect of race are seen at very low incomes and very high incomes.

Let's take Asian renters as an example. The positive and statistically significant coefficient for the interaction between this race and household income tells us that income has slightly more of an effect on GRPIP for Asian renters than it does for White renters. As we see here, the slopes representing White renters and Asian renters start off around the same point, but as we move up to a higher household income (\$300,000), there is a bigger difference between the GRPIP of White and Asian renters.

The same is true when comparing White renters and Black renters.

Opposely, there is a negative and statistically significant coefficient for the interaction between household income and multiracial as well as some other race.



Methods & Results

Linear Regression | Best Fit

Finally, we should quickly take a look at our new variable. If you are 30 years old or younger, you pay 1 percentage point more on rent than if you are over the age of 30. This coefficient is statistically significant.

```
##  
## Call:  
## lm(formula = GRPIP ~ HINCP + EDU + RACE + age_cat + NP + HINCP:RACE,  
##      data = person_data)  
##  
## Residuals:  
##    Min      1Q  Median      3Q     Max  
## -37.857 -9.988 -3.553  5.581 87.821  
##  
## Coefficients:  
##                               Estimate Std. Error t value Pr(>|t|)  
## (Intercept)            3.996e+01 5.277e-01 75.729 < 2e-16 ***  
## HINCP                 -1.217e-04 2.241e-06 -54.306 < 2e-16 ***  
## EDUAssoc. Degree      -1.357e+00 6.972e-01 -1.947 0.051615 .  
## EDUBach Degree or Higher 5.958e-01 5.117e-01  1.164 0.244333  
## EDUHS Diploma          -9.521e-01 4.762e-01 -1.999 0.045575 *  
## RACEAsian              1.632e+00 7.672e-01  2.127 0.033441 *  
## RACEBlack               4.566e-01 6.511e-01  0.701 0.483162  
## RACEMultiracial        4.553e+00 1.173e+00  3.882 0.000104 ***  
## RACENative Am.         -3.341e+00 4.463e+00 -0.749 0.454044  
## RACEOther Race          6.796e+00 9.019e-01  7.535 5.21e-14 ***  
## age_catyounger          1.077e+00 3.000e-01  3.591 0.000331 ***  
## NP                    -2.404e-02 1.050e-01 -0.229 0.818846  
## HINCP:RACEAsian         1.169e-05 5.541e-06  2.110 0.034904 *  
## HINCP:RACEBlack          8.482e-06 6.029e-06  1.407 0.159540  
## HINCP:RACEMultiracial   -4.175e-05 1.144e-05 -3.649 0.000264 ***  
## HINCP:RACENative Am.    -1.732e-05 4.448e-05 -0.389 0.696926  
## HINCP:RACEOther Race    -9.305e-05 1.066e-05 -8.728 < 2e-16 ***  
## ---  
## Signif. codes:  0 '****' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1  
##  
## Residual standard error: 15.7 on 13136 degrees of freedom  
## Multiple R-squared:  0.2713, Adjusted R-squared:  0.2704  
## F-statistic: 305.6 on 16 and 13136 DF, p-value: < 2.2e-16
```

- 01 Introduction
- 02 Data
- 03 Methods & Results
- 04 Discussion
- 05 Conclusion

- 01 Introduction
- 02 Data
- 03 Methods & Results
- 04 Discussion
- 05 Conclusion

