

# ASA Data Challenge Expo

## Helping Communities During the COVID-19 Pandemic

### Entry Details

- ASA Data Challenge Expo (<https://community.amstat.org/dataexpo/home>)
- Name: Walter Yu
- Date: Fall 2020

### Executive Summary

This entry aims to help disadvantaged communities during the COVID-19 pandemic by answering the following questions:

1. Explore the relationship between socioeconomic features of the U.S. population and disadvantaged communities.
2. Identify disadvantaged communities based on their median household income and socioeconomic factors.
3. Identify which disadvantaged communities have been most impacted by the COVID-19 pandemic and in need of public services.
4. Provide recommendations on helping these communities based on data analysis results.

The intended audience are state/local governments, non-governmental organizations (NGOs) and volunteers which are able to provide aid and services to these communities.

### Scope

This entry focuses on California communities to control its scope since several questions are being considered, and data analysis of all U.S. communities would expand the scope and length of this report. This limited scope provides for more detail and attention to be paid to analysis, documentation and recommendations.

## Part 1: Overview

### Methodology

This entry is designed to be interpretable, so it clearly outlines data analysis steps into the following modules:

1. Overview: Outline approach, assumptions and data sources
2. Data Processing: Data preparation and manipulation for analysis
3. Data Analysis: Model fit, coefficient interpretation and diagnostics
4. Data Visualization: Communicate findings through data plots
5. Recommendations: Document key findings from data analysis

### Assumptions

This entry makes the following assumptions:

1. Although the scope is limited to California communities, the methodology can be applied to other states since it is based on data extracted from the U.S. Census for the state/county level and do not contain any characteristics specific to California.
2. State and federal guidelines ([https://www.hud.gov/topics/rental\\_assistance/phprog](https://www.hud.gov/topics/rental_assistance/phprog)) typically define disadvantaged communities as being low-income, so median household income was used to identify such communities.
3. Data analysis was documented to be clear and easily interpretable, so linear regression and the Law of Parsimony ([https://en.wikipedia.org/wiki/Occam%27s\\_razor](https://en.wikipedia.org/wiki/Occam%27s_razor)) were applied whenever possible

### Data Summary

This entry analyzes core and supplemental datasets from the data challenge problem statement (<https://opportunity.census.gov/assets/files/covid-19-top-asa-problem-statement.pdf>) as follows:

- Core Dataset: 2019 American Community Survey (ACS) Single-Year Estimates
- Supplemental Dataset: COVID-19 Data from the National Center for Health Statistics

Data was downloaded from portal websites as follows:

1. U.S. Census Website: Advanced search feature (<https://data.census.gov/cedsci/advanced>) was used to filter data in the following order: Surveys > Years > Geography > Topics.
2. U.S. Census COVID-19 Website: CA state data was downloaded from the categorical dataset search page (<https://covid19.census.gov/>).
3. National Center for Health Statistics (NCHS) Website: Death counts by county and race downloaded from their data portal (<https://www.cdc.gov/nchs/covid19/index.htm>).

## Core Datasets

Datasets of interest were identified from the U.S. Census data portal and extracted using the advanced search tool. Table ID numbers are listed for reference.

1. 2019 American Community Survey (ACS) Single-Year Estimates - Language Spoken
  - Description: PLACE OF BIRTH BY LANGUAGE SPOKEN AT HOME AND ABILITY TO SPEAK ENGLISH IN THE UNITED STATES
  - Survey/Program: American Community Survey
  - Years: 2019
  - Table: B06007 (<https://data.census.gov/cedsci/table?q=B06007&tid=ACSDT1Y2019.B06007>)
2. 2019 American Community Survey (ACS) Single-Year Estimates - Household Income
  - Description: HOUSEHOLD INCOME IN THE PAST 12 MONTHS (IN 2019 INFLATION-ADJUSTED DOLLARS)
  - Survey/Program: American Community Survey
  - Years: 2019
  - Table: B19001 (<https://data.census.gov/cedsci/table?text=B19001&tid=ACSDT1Y2019.B19001>)
3. 2019 American Community Survey (ACS) Single-Year Estimates - Median Household Income
  - Description: MEDIAN HOUSEHOLD INCOME IN THE PAST 12 MONTHS (IN 2019 INFLATION-ADJUSTED DOLLARS)
  - Survey/Program: American Community Survey
  - Years: 2019
  - Table: B19013 (<https://data.census.gov/cedsci/table?text=B19013&tid=ACSDT1Y2019.B19013>)
4. 2019 American Community Survey (ACS) Single-Year Estimates - Poverty Status
  - Description: POVERTY STATUS IN THE PAST 12 MONTHS BY SEX BY AGE
  - Survey/Program: American Community Survey
  - Years: 2019
  - Table: B17001 (<https://data.census.gov/cedsci/table?text=B17001&tid=ACSDT1Y2019.B17001>)
5. 2019 American Community Survey (ACS) Single-Year Estimates - Housing Cost
  - Description: MONTHLY HOUSING COSTS
  - Survey/Program: American Community Survey
  - Years: 2019
  - Table: B25104 (<https://data.census.gov/cedsci/table?text=B25104&tid=ACSDT1Y2019.B25104>)
6. 2019 American Community Survey (ACS) Single-Year Estimates - Education Attainment
  - Description: EDUCATIONAL ATTAINMENT FOR THE POPULATION 25 YEARS AND OVER
  - Survey/Program: American Community Survey
  - Years: 2019
  - Table: B15003 (<https://data.census.gov/cedsci/table?text=B15003&tid=ACSDT1Y2019.B15003>)
7. 2019 American Community Survey (ACS) Single-Year Estimates - Commute Mode
  - Description: MEANS OF TRANSPORTATION TO WORK BY AGE
  - Survey/Program: American Community Survey
  - Years: 2019
  - Table: B08101 (<https://data.census.gov/cedsci/table?text=B08101&tid=ACSDT1Y2019.B08101>)
8. 2019 American Community Survey (ACS) Single-Year Estimates - Race
  - Description: RACE
  - Survey/Program: American Community Survey
  - Years: 2019
  - Table: B02001 (<https://data.census.gov/cedsci/table?text=B02001&tid=ACSDT1Y2019.B02001>)

## Supplemental Datasets (U.S. Census)

Datasets of interest were identified from the U.S. Census COVID-19 data portal under the categorical dataset section.

1. U.S. Census - COVID-19 Demographic and Economic Resources
  - Dataset: California Counties DP02 Social (<https://covid19.census.gov/datasets/california-counties-dp02-social>)
2. U.S. Census - COVID-19 Demographic and Economic Resources
  - Dataset: California Counties DP03 Economic (<https://covid19.census.gov/datasets/california-counties-dp03-economic>)
3. U.S. Census - COVID-19 Demographic and Economic Resources
  - Dataset: California Counties DP04 Housing (<https://covid19.census.gov/datasets/california-counties-dp04-housing>)
4. U.S. Census - COVID-19 Demographic and Economic Resources
  - Dataset: California Counties DP05 Demographic (<https://covid19.census.gov/datasets/california-counties-dp05-demographic>)
5. U.S. Census - COVID-19 Demographic and Economic Resources
  - Dataset: Household Pulse Survey Public Use File (<https://www.census.gov/programs-surveys/household-pulse-survey/datasets.html>)
6. U.S. Census - COVID-19 Demographic and Economic Resources
  - Dataset: COVID-19 Case Surveillance Public Use Data (<https://data.cdc.gov/Case-Surveillance/COVID-19-Case-Surveillance-Public-Use-Data/vbim-akqf>)

## Supplemental Datasets (NCHS)

Datasets of interest were identified from the National Center for Health Statistics (NCHS) data portal.

1. NCHS - COVID-19 Data from the National Center for Health Statistics
  - Dataset: Provisional COVID-19 Death Counts by County and Race (<https://data.cdc.gov/NCHS/Provisional-COVID-19-Death-Counts-by-County-and-Race/k8wy-p9cg>)

## Geospatial Datasets

Datasets of interest were identified from the U.S. Census COVID-19 data portal under the categorical dataset section.

1. U.S. Census - COVID-19 Demographic and Economic Resources
  - Description: American Community Survey (ACS) about household income ranges and cutoffs and Poverty Status.
  - These are 5-year estimates shown by state and county boundaries.
  - Link: Dataset (<https://uscensus.maps.arcgis.com/home/item.html?id=b2ba19b4cce04a9796d9cdeecaba2f18>)
2. U.S. Census - COVID-19 Demographic and Economic Resources
  - Description: American Community Survey (ACS) about household income ranges and cutoffs.
  - These are 5-year estimates shown by county, and state boundaries.
  - Link: Dataset (<https://uscensus.maps.arcgis.com/home/item.html?id=0fbb1571e5b6458f941580d1d64a6693>)
3. U.S. Census - COVID-19 Demographic and Economic Resources
  - Description: American Community Survey (ACS) about total population count by sex and age group.
  - These are 5-year estimates shown by state and county boundaries.
  - Link: Dataset (<https://uscensus.maps.arcgis.com/home/item.html?id=eab0f44ba5184c609175caa7ae317f0c>)

## Part 2: Data Processing

### Methodology

Core and supplemental datasets of interest were processed and joined as listed below prior to model fit and data visualization.

This module completes the tasks listed below; however, all text output, warnings and messages are silenced to minimize report length so please check the code repository for details about implementation.

1. Import csv files as dataframes
2. Remove first record from each dataframe (header data)

3. Import selected columns from full dataset
4. Relabel selected columns from full dataset
5. Join tables together into single dataframe

## Part 3: Data Analysis

### Methodology

Processed data was used to fit a linear regression model to identify key attributes associated with disadvantaged communities and those communities most impacted by the pandemic.

Data analysis was conducted as follows:

1. Fit 2019 ACS data features into separate linear regression models to identify which had the best fit and association with median household income.
2. Conduct coefficient interpretation and model diagnostics to refine models.
3. Identify communities with lowest median income and highest COVID-19 death counts.

### 2019 ACS - Commute Mode

Commute mode as a whole has a good relationship with median household income as verified with the low p-value. Features were split into two groups for analysis: car-based modes and transit-based modes.

Features within each group were fit individually into their own model. Carpool had a good fit in the car-based modes. Transit/remote work had good fits with transit-based modes. This finding implies that certain commute modes (i.e. carpool, transit and remote work) have a better relationship to income than other ones (i.e. walking).

```
# features:
# B08101_009E: Car, truck, or van - drove alone
# B08101_017E: Car, truck, or van - carpooled
# B08101_025E: Public transportation (excluding taxicab)
# B08101_033E: Walked
# B08101_041E: Taxicab, motorcycle, bicycle, or other means
# B08101_049E: Worked from home
```

```
# model fit for transit-based commute modes
# linear regression - transit
fit_acs19_commute_transit <- lm(
  hh_median ~ as.numeric(commute_transit),
  data=acs19_nofactor
)
round(summary(fit_acs19_commute_transit)$coeff, 6)
```

```
##               Estimate Std. Error   t value Pr(>|t|)
## (Intercept)    91637.6726   7286.8220  12.575808 0.000000
## as.numeric(commute_transit) -698.1878    320.7253  -2.176903 0.035448
```

```
# linear regression - transit + walk
fit_acs19_commute_transit_walk <- lm(
  hh_median ~ as.numeric(commute_transit) +
    as.numeric(commute_walk),
  data=acs19_nofactor
)
round(summary(fit_acs19_commute_transit_walk)$coeff, 6)
```

```
##               Estimate Std. Error  t value Pr(>|t|)
## (Intercept)      87881.6979   7634.8470  11.510604 0.000000
## as.numeric(commute_transit) -1176.2170   455.7074  -2.581080 0.013725
## as.numeric(commute_walk)     664.0563   455.7074   1.457199 0.153067
```

```
# Linear regression - transit + walk + remote
fit_acs19_commute_transit_remote <- lm(
  hh_median ~ as.numeric(commute_transit) +
  as.numeric(commute_walk) +
  as.numeric(commute_remote),
  data=acs19_nofactor
)
round(summary(fit_acs19_commute_transit_remote)$coeff, 6)
```

```
##               Estimate Std. Error  t value Pr(>|t|)
## (Intercept)      99606.0378   8501.6970  11.716018 0.000000
## as.numeric(commute_transit) -1064.6047   428.9094  -2.482120 0.017596
## as.numeric(commute_walk)     747.7170   427.9278   1.747297 0.088664
## as.numeric(commute_remote)   -775.9597   304.5944  -2.547518 0.015018
```

```
# variance analysis
anova(
  fit_acs19_commute_transit,
  fit_acs19_commute_transit_walk,
  fit_acs19_commute_transit_remote
)
```

```
## Analysis of Variance Table
##
## Model 1: hh_median ~ as.numeric(commute_transit)
## Model 2: hh_median ~ as.numeric(commute_transit) + as.numeric(commute_walk)
## Model 3: hh_median ~ as.numeric(commute_transit) + as.numeric(commute_walk) +
##   as.numeric(commute_remote)
##   Res.Df      RSS Df Sum of Sq    F Pr(>F)
## 1      40 1.8756e+10
## 2      39 1.7788e+10  1  968488352 2.4223 0.12791
## 3      38 1.5193e+10  1 2594746246 6.4898 0.01502 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
# model fit for car-based commute modes
# Linear regression - car
fit_acs19_commute_car <- lm(
  hh_median ~ as.numeric(commute_car),
  data=acs19_nofactor
)
round(summary(fit_acs19_commute_car)$coeff, 6)
```

```
##               Estimate Std. Error  t value Pr(>|t|)
## (Intercept)      91837.2973   7274.4999  12.62455 0.000000
## as.numeric(commute_car)  -708.0749   320.1829  -2.21147 0.032781
```

```
# Linear regression - car + carpool
fit_acs19_commute_car_carpool <- lm(
  hh_median ~ as.numeric(commute_car) +
  as.numeric(commute_carpool),
  data=acs19_nofactor
)
round(summary(fit_acs19_commute_car_carpool)$coeff, 6)
```

```
##
##           Estimate Std. Error  t value Pr(>|t|)
## (Intercept)      100655.1880   7708.8667 13.057067 0.000000
## as.numeric(commute_car)      -284.0894    346.1756 -0.820651 0.416833
## as.numeric(commute_carpool)   -860.7206    346.1756 -2.486370 0.017295
```

```
# Linear regression - car + carpool + other
fit_acs19_commute_car_other <- lm(
  hh_median ~ as.numeric(commute_car) +
  as.numeric(commute_carpool) +
  as.numeric(commute_other),
  data=acs19_nofactor
)
round(summary(fit_acs19_commute_car_other)$coeff, 6)
```

```
##
##           Estimate Std. Error  t value Pr(>|t|)
## (Intercept)      105124.7550   8272.4716 12.707781 0.000000
## as.numeric(commute_car)     -121.3169    361.7098 -0.335398 0.739169
## as.numeric(commute_carpool) -779.2708    347.1480 -2.244780 0.030683
## as.numeric(commute_other)   -465.5924    335.6560 -1.387112 0.173490
```

```
# variance analysis
anova(
  fit_acs19_commute_car,
  fit_acs19_commute_car_carpool,
  fit_acs19_commute_car_other
)
```

```
## Analysis of Variance Table
##
## Model 1: hh_median ~ as.numeric(commute_car)
## Model 2: hh_median ~ as.numeric(commute_car) + as.numeric(commute_carpool)
## Model 3: hh_median ~ as.numeric(commute_car) + as.numeric(commute_carpool) +
##           as.numeric(commute_other)
##   Res.Df    RSS Df Sum of Sq    F Pr(>F)
## 1      40 1.8693e+10
## 2      39 1.6135e+10  1 2557652979 6.3285 0.01623 *
## 3      38 1.5358e+10  1 777611785 1.9241 0.17349
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

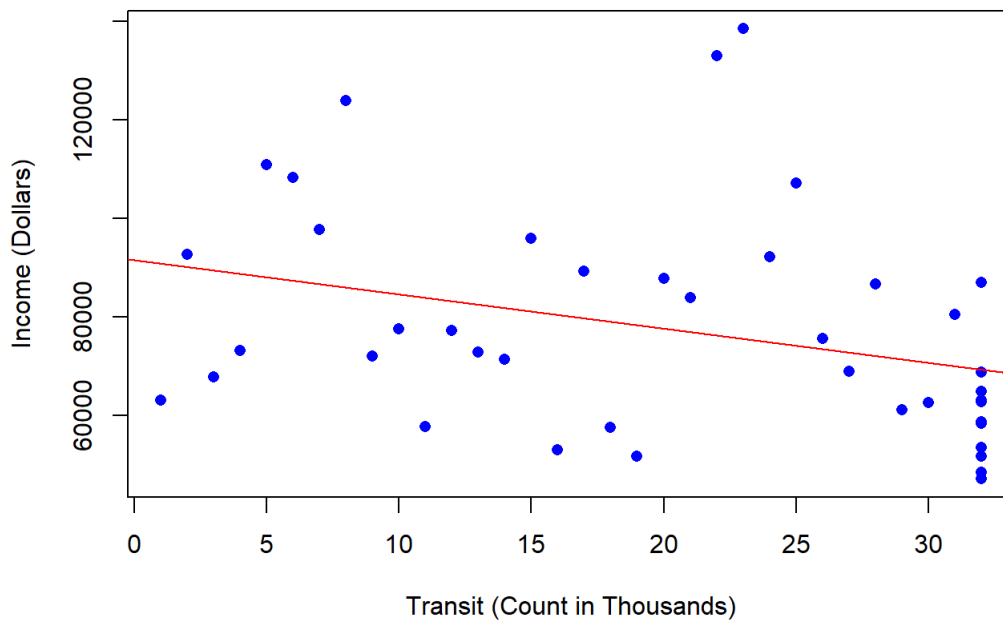
```
# visualize income and transit variables
# ref: https://www.theanalysisfactor.com/linear-models-r-plotting-regression-lines/
plot(hh_median ~ as.numeric(commute_transit),
  main="Median Household Income and Transit Variables",
  xlab="Transit (Count in Thousands)",
  ylab="Income (Dollars)",
  pch=16,
  col="blue",
  data=acs19_nofactor
)

# Linear regression - all
fit_acs19_commute_all <- lm(
  hh_median ~ as.numeric(commute_transit) +
  as.numeric(commute_walk) +
  as.numeric(commute_remote) +
  as.numeric(commute_car) +
  as.numeric(commute_carpool) +
  as.numeric(commute_other),
  data=acs19_nofactor
)
summary(fit_acs19_commute_all)
```

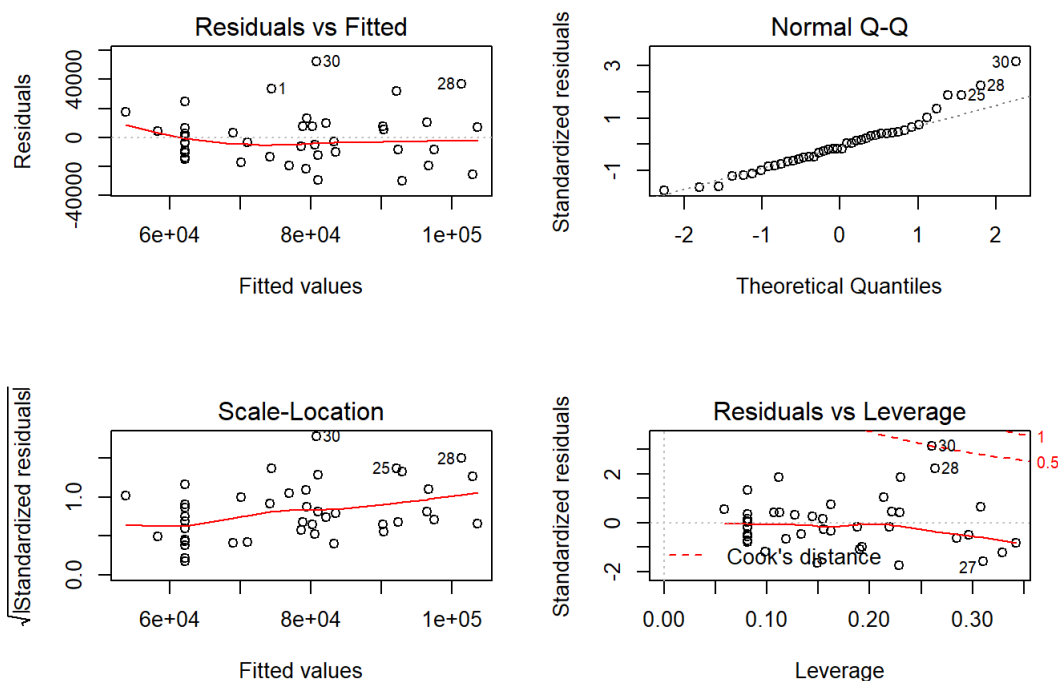
```
##
## Call:
## lm(formula = hh_median ~ as.numeric(commute_transit) + as.numeric(commute_walk) +
##   as.numeric(commute_remote) + as.numeric(commute_car) + as.numeric(commute_carpool) +
##   as.numeric(commute_other), data = acs19_nofactor)
##
## Residuals:
##   Min     1Q  Median     3Q      Max
## -29868 -11591  -2974   7657  52359
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    103765.67    8491.27   12.220 3.49e-14 ***
## as.numeric(commute_transit)    -563.87    468.69   -1.203  0.2370
## as.numeric(commute_walk)      982.97    430.39    2.284  0.0286 *
## as.numeric(commute_remote)   -331.85    406.48   -0.816  0.4198
## as.numeric(commute_car)      -42.06    379.73   -0.111  0.9124
## as.numeric(commute_carpool)  -649.67    406.68   -1.597  0.1191
## as.numeric(commute_other)   -694.39    416.06   -1.669  0.1040
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 19360 on 35 degrees of freedom
## Multiple R-squared:  0.3749, Adjusted R-squared:  0.2677
## F-statistic: 3.498 on 6 and 35 DF, p-value: 0.008148
```

```
# trend line plot
abline(
  fit_acs19_commute_transit,
  col="red"
)
```

## Median Household Income and Transit Variables



```
# diagnostic plot
par(mfrow=c(2,2))
plot(fit_acs19_commute_all)
```



## 2019 ACS - Housing Cost

Housing cost as a whole has a good relationship with median household income as verified with the low p-value. In particular, higher monthly housing costs were more closely associated with median household income.

This finding implies that certain higher housing cost is associated with median household income. However, lower housing cost features were not analyzed due to their large number, values with little or no rent and not being adjusted for cost of living (i.e. some counties may have significantly lower rents). As a result, only the higher housing cost features were analyzed.

```
# features:
# B25104_002E: Less than $100
# B25104_003E: 100 to $199
# B25104_004E: 200 to $299
# B25104_005E: 300 to $399
# B25104_006E: 400 to $499
# B25104_007E: 500 to $599
# B25104_008E: 600 to $699
# B25104_009E: 700 to $799
# B25104_010E: 800 to $899
# B25104_011E: 900 to $999
# B25104_012E: 1,000 to $1,499
# B25104_013E: 1,500 to $1,999
# B25104_014E: 2,000 to $2,499
# B25104_015E: 2,500 to $2,999
# B25104_016E: 3,000 or more
# B25104_017E: No cash rent

# linear regression model fit - housing cost
fit_acs19_housing_3000 <- lm(hh_median ~ housing_more_3000,
  data=acs19_nofactor
)
round(summary(fit_acs19_housing_3000)$coeff, 6)
```



```
##               Estimate Std. Error  t value Pr(>|t|)
## (Intercept)    76593.836875  3.6601e+03  20.926706  0.000000
## housing_more_3000    0.009326  1.0552e-02   0.883789  0.382092
```

```
fit_acs19_housing_2500 <- lm(hh_median ~ housing_more_3000 +
  housing_2500,
  data=acs19_nofactor
)
round(summary(fit_acs19_housing_2500)$coeff, 6)
```

```
##               Estimate Std. Error  t value Pr(>|t|)
## (Intercept)    76735.542976  2773.177922  27.670617   0e+00
## housing_more_3000    0.740007    0.132152   5.599655  2e-06
## housing_2500       -1.311729    0.236807  -5.539232  2e-06
```

```
fit_acs19_housing_2000 <- lm(hh_median ~ housing_more_3000 +
  housing_2500 +
  housing_2000,
  data=acs19_nofactor
)
round(summary(fit_acs19_housing_2000)$coeff, 6)
```

```
##               Estimate Std. Error  t value Pr(>|t|)
## (Intercept)    76846.250921  2696.486566  28.498659  0.000000
## housing_more_3000    0.378466    0.237616   1.592764  0.119499
## housing_2500        0.663281    1.115978   0.594349  0.555801
## housing_2000       -0.903725    0.499667  -1.808654  0.078423
```

```
# model fit summary
summary(fit_acs19_housing_2000)
```

```
##
## Call:
## lm(formula = hh_median ~ housing_more_3000 + housing_2500 + housing_2000,
##     data = acs19_nofactor)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -28127 -11690  -1139   13088   33075
##
## Coefficients:
##               Estimate Std. Error t value Pr(>|t|)
## (Intercept)    76846.2509  2696.4866   28.499  <2e-16 ***
## housing_more_3000    0.3785    0.2376    1.593   0.1195
## housing_2500        0.6633    1.1160    0.594   0.5558
## housing_2000       -0.9037    0.4997   -1.809   0.0784 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 16700 on 38 degrees of freedom
## Multiple R-squared:  0.4946, Adjusted R-squared:  0.4547
## F-statistic: 12.39 on 3 and 38 DF, p-value: 8.47e-06
```

```
# anova analysis to identify change
anova(
  fit_acs19_housing_3000,
  fit_acs19_housing_2500,
  fit_acs19_housing_2000
)
```

```
## Analysis of Variance Table
##
## Model 1: hh_median ~ housing_more_3000
## Model 2: hh_median ~ housing_more_3000 + housing_2500
## Model 3: hh_median ~ housing_more_3000 + housing_2500 + housing_2000
##   Res.Df    RSS Df Sum of Sq    F    Pr(>F)
## 1      40 2.0577e+10
## 2      39 1.1516e+10  1 9060336538 32.4700 1.477e-06 ***
## 3      38 1.0603e+10  1 912795180  3.2712  0.07842 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

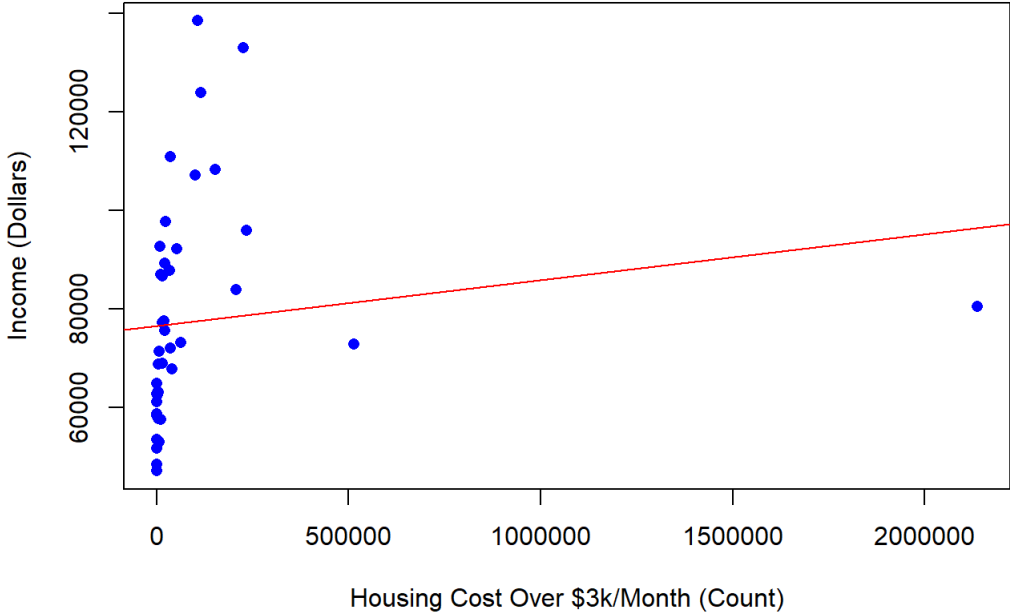
```
# calculate correlation
col_housing = c(
  "housing_more_3000",
  "housing_2500",
  "housing_2000"
)
cor(acs19_nofactor$hh_median, acs19_nofactor[col_housing])
```

```
##      housing_more_3000 housing_2500 housing_2000
## [1,]      0.1383946    0.09838385    0.07478075
```

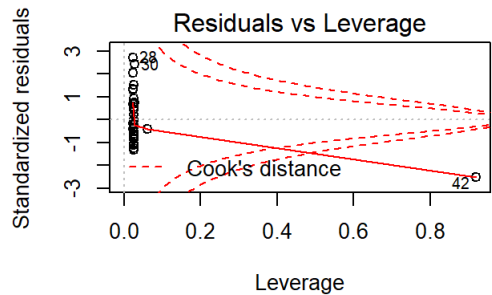
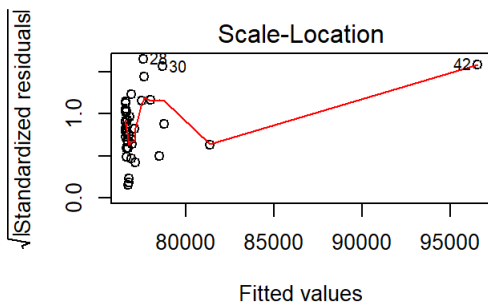
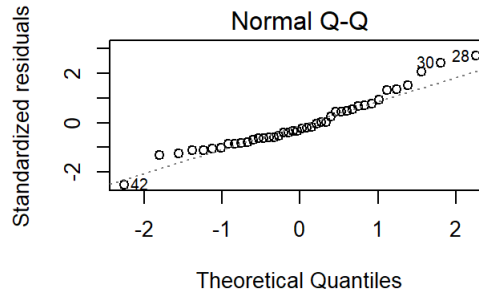
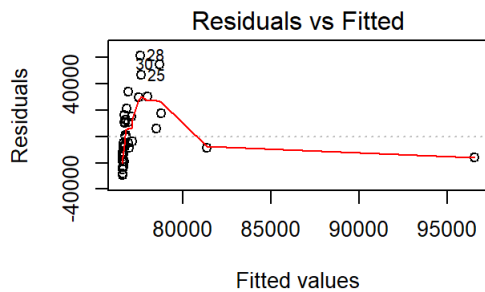
```
# visualize income and transit variables
# ref: https://www.theanalysisfactor.com/linear-models-r-plotting-regression-lines/
plot(hh_median ~ as.numeric(housing_more_3000),
     main="Median Household Income and Housing Cost",
     xlab="Housing Cost Over $3k/Month (Count)",
     ylab="Income (Dollars)",
     pch=16,
     col="blue",
     data=acs19_nofactor
)

# trend line plot
abline(
  fit_acs19_housing_3000,
  col="red"
)
```

## Median Household Income and Housing Cost



```
# diagnostic plot
par(mfrow=c(2,2))
plot(fit_acs19_housing_3000)
```



## 2019 ACS - Education Attainment

Education attainment as a whole has a good relationship with median household income as verified with the low p-value. In particular, education attainment including and beyond a bachelors degree was more closely associated with median household income.

This finding implies that certain higher education attainment is associated with median household income. When variables were added the model, most maintained relatively lower p-values.

```
# features:
# B15003_002E: No schooling completed
# B15003_003E: Nursery school
# B15003_004E: Kindergarten
# B15003_005E: 1st grade
# B15003_006E: 2nd grade
# B15003_007E: 3rd grade
# B15003_008E: 4th grade
# B15003_009E: 5th grade
# B15003_010E: 6th grade
# B15003_011E: 7th grade
# B15003_012E: 8th grade
# B15003_013E: 9th grade
# B15003_014E: 10th grade
# B15003_015E: 11th grade
# B15003_016E: 12th grade, no diploma
# B15003_017E: Regular high school diploma
# B15003_018E: GED or alternative credential
# B15003_019E: Some college, less than 1 year
# B15003_020E: Some college, 1 or more years, no degree
# B15003_021E: Associate's degree
# B15003_022E: Bachelor's degree
# B15003_023E: Master's degree
# B15003_024E: Professional school degree
# B15003_025E: Doctorate degree

# linear regression model fit - education attainment
fit_acs19_edu_doctorate <- lm(hh_median ~ edu_doctorate, data=acs19_nofactor)
round(summary(fit_acs19_edu_doctorate)$coeff, 6)
```

```
##              Estimate Std. Error  t value Pr(>|t|)
## (Intercept)  7.657e+04 3660.053958  20.920456 0.000000
## edu_doctorate 4.390e-02   0.048655   0.902281 0.372309
```

```
fit_acs19_edu_professional <- lm(hh_median ~ edu_doctorate +
  edu_professional,
  data=acs19_nofactor
)
round(summary(fit_acs19_edu_professional)$coeff, 6)
```

```
##              Estimate Std. Error  t value Pr(>|t|)
## (Intercept)  76612.646777 3405.744215 22.495126 0.000000
## edu_doctorate    1.221055    0.441099   2.768209 0.008579
## edu_professional -0.817221    0.304609  -2.682853 0.010652
```

```
fit_acs19_edu_master <- lm(hh_median ~ edu_doctorate +
  edu_professional +
  edu_master,
  data=acs19_nofactor
)
round(summary(fit_acs19_edu_master)$coeff, 6)
```

```
##              Estimate Std. Error  t value Pr(>|t|)
## (Intercept)  76716.058097 3302.287047 23.231190 0.000000
## edu_doctorate    3.492346    1.288222   2.710982 0.010015
## edu_professional  0.390498    0.710432   0.549662 0.585768
## edu_master      -0.779620    0.417107  -1.869113 0.069329
```

```
fit_acs19_edu_bachelor <- lm(hh_median ~ edu_doctorate +
  edu_professional +
  edu_master +
  edu_bachelor,
  data=acs19_nofactor
)
round(summary(fit_acs19_edu_bachelor)$coeff, 6)
```

```
##              Estimate Std. Error  t value Pr(>|t|)
## (Intercept)  76431.200647 3135.960658 24.372500 0.000000
## edu_doctorate    2.116071   1.363125  1.552367 0.129087
## edu_professional  2.377478   1.101358  2.158678 0.037432
## edu_master      -0.258002   0.457080 -0.564457 0.575850
## edu_bachelor    -0.329807   0.144564 -2.281387 0.028376
```

```
# model fit summary
summary(fit_acs19_edu_bachelor)
```

```
##
## Call:
## lm(formula = hh_median ~ edu_doctorate + edu_professional + edu_master +
##     edu_bachelor, data = acs19_nofactor)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -28933 -15473  -4100  16759  36781
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  76431.2006  3135.9607  24.372  <2e-16 ***
## edu_doctorate    2.1161    1.3631   1.552   0.1291
## edu_professional  2.3775    1.1014   2.159   0.0374 *
## edu_master      -0.2580    0.4571  -0.564   0.5758
## edu_bachelor    -0.3298    0.1446  -2.281   0.0284 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 19410 on 37 degrees of freedom
## Multiple R-squared:  0.3357, Adjusted R-squared:  0.2639
## F-statistic: 4.675 on 4 and 37 DF, p-value: 0.003726
```

```
# anova analysis to identify change
anova(
  fit_acs19_edu_doctorate,
  fit_acs19_edu_professional,
  fit_acs19_edu_master,
  fit_acs19_edu_bachelor
)
```

```
## Analysis of Variance Table
##
## Model 1: hh_median ~ edu_doctorate
## Model 2: hh_median ~ edu_doctorate + edu_professional
## Model 3: hh_median ~ edu_doctorate + edu_professional + edu_master
## Model 4: hh_median ~ edu_doctorate + edu_professional + edu_master + edu_bachelor
##   Res.Df    RSS Df Sum of Sq    F Pr(>F)
## 1      40 2.0560e+10
## 2      39 1.7357e+10  1 3203277308 8.5053 0.005985 **
## 3      38 1.5895e+10  1 1461354034 3.8802 0.056380 .
## 4      37 1.3935e+10  1 1960219018 5.2047 0.028376 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

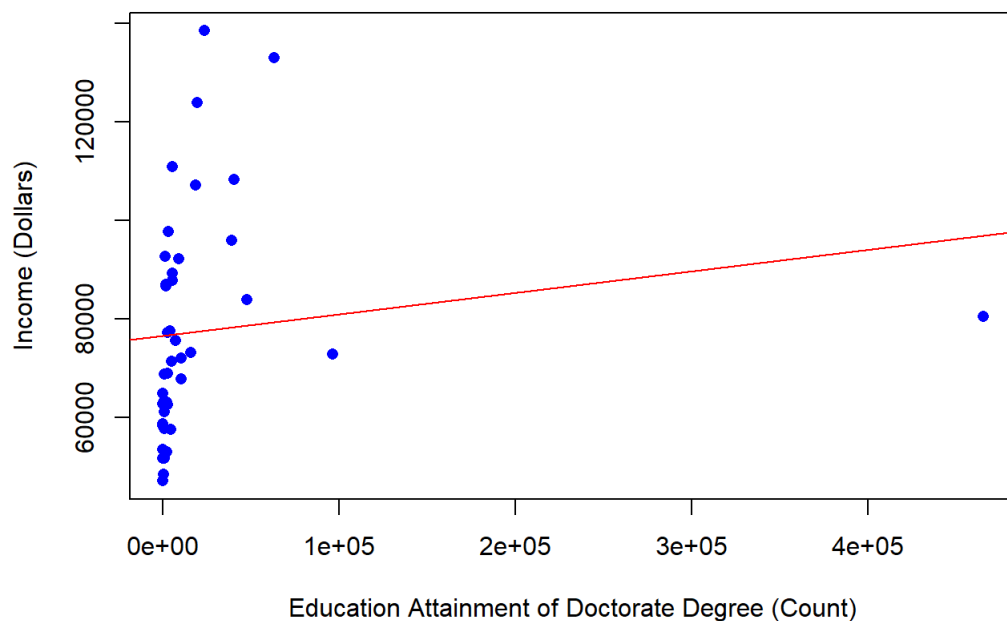
```
# calculate correlation
col_edu = c(
  "edu_doctorate",
  "edu_professional",
  "edu_master",
  "edu_bachelor"
)
cor(acs19_nofactor$hh_median, acs19_nofactor[col_edu])
```

```
##      edu_doctorate edu_professional edu_master edu_bachelor
## [1,]      0.1412332       0.1003801  0.1182039   0.09137493
```

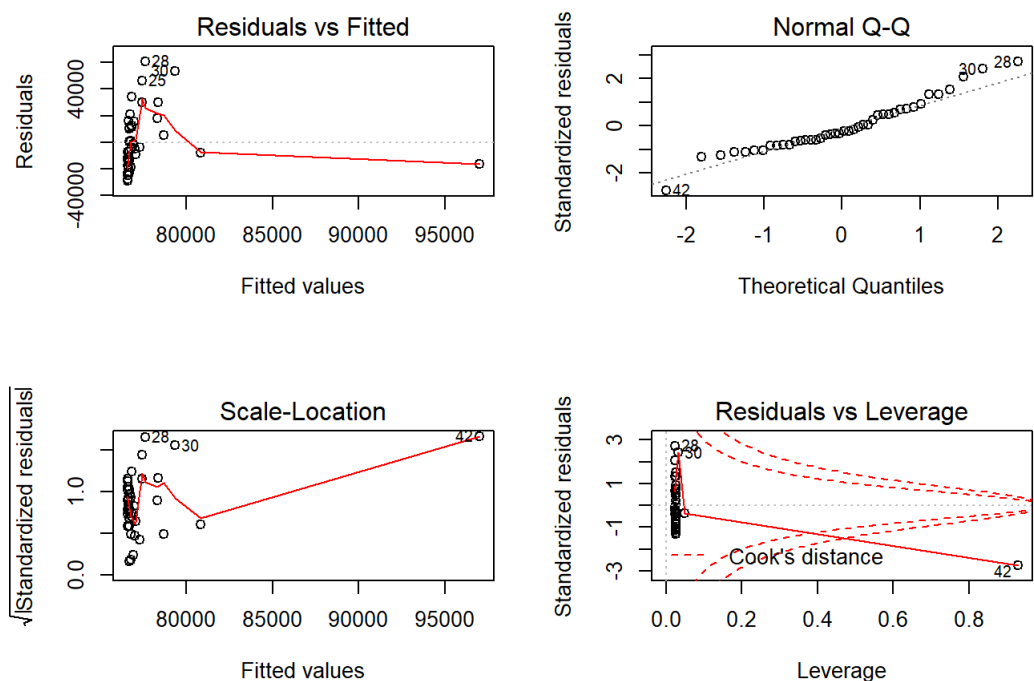
```
# visualize income and transit variables
# ref: https://www.theanalysisfactor.com/linear-models-r-plotting-regression-lines/
plot(hh_median ~ as.numeric(edu_doctorate),
     main="Median Household Income and Education Attainment",
     xlab="Education Attainment of Doctorate Degree (Count)",
     ylab="Income (Dollars)",
     pch=16,
     col="blue",
     data=acs19_nofactor
)

# trend line plot
abline(
  fit_acs19_edu_doctorate,
  col="red"
)
```

### Median Household Income and Education Attainment



```
# diagnostic plot
par(mfrow=c(2,2))
plot(fit_acs19_edu_doctorate)
```



## 2019 ACS - Race

Race as a whole has a good relationship with median household income as verified with the low p-value. In particular, “American Indian”, “Alaska Native alone” and “Asian alone” features were more closely associated with median household income.

```
# features:
# B02001_002E: White alone
# B02001_003E: Black or African American alone
# B02001_004E: American Indian and Alaska Native alone
# B02001_005E: Asian alone
# B02001_006E: Native Hawaiian and Other Pacific Islander alone
# B02001_007E: Some other race alone
# B02001_008E: Two or more races
```

```
# linear regression
fit_acs19_race1 <- lm(
  hh_median ~ as.numeric(B02001_002E),
  data=acs19_nofactor
)
round(summary(fit_acs19_race1)$coeff, 6)
```

	Estimate	Std. Error	t value	Pr(> t )
## (Intercept)	77215.039514	3693.175587	20.907492	0.00000
## as.numeric(B02001_002E)	0.000294	0.000979	0.299947	0.76577

```
fit_acs19_race2 <- lm(
  hh_median ~ as.numeric(B02001_002E) +
  as.numeric(B02001_003E),
  data=acs19_nofactor
)
round(summary(fit_acs19_race2)$coeff, 6)
```

##		Estimate	Std. Error	t value	Pr(> t )
##	(Intercept)	77185.280719	3739.116581	20.642652	0.000000
##	as.numeric(B02001_002E)	0.001819	0.006177	0.294561	0.769892
##	as.numeric(B02001_003E)	-0.015334	0.061278	-0.250245	0.803711

```
fit_acs19_race3 <- lm(
  hh_median ~ as.numeric(B02001_002E) +
  as.numeric(B02001_003E) +
  as.numeric(B02001_004E),
  data=acs19_nofactor
)
round(summary(fit_acs19_race3)$coeff, 6)
```

##		Estimate	Std. Error	t value	Pr(> t )
##	(Intercept)	76505.505408	3225.948817	23.715660	0.000000
##	as.numeric(B02001_002E)	0.045989	0.012742	3.609315	0.000883
##	as.numeric(B02001_003E)	0.162037	0.070342	2.303568	0.026811
##	as.numeric(B02001_004E)	-4.514352	1.183282	-3.815111	0.000487

```
fit_acs19_race_all <- lm(
  hh_median ~ as.numeric(B02001_002E) +
  as.numeric(B02001_003E) +
  as.numeric(B02001_004E) +
  as.numeric(B02001_005E) +
  as.numeric(B02001_006E) +
  as.numeric(B02001_007E) +
  as.numeric(B02001_008E),
  data=acs19_nofactor
)
round(summary(fit_acs19_race_all)$coeff, 6)
```

##		Estimate	Std. Error	t value	Pr(> t )
##	(Intercept)	76528.547054	2767.810361	27.649491	0.000000
##	as.numeric(B02001_002E)	0.024350	0.020430	1.191882	0.241561
##	as.numeric(B02001_003E)	0.120353	0.103428	1.163641	0.252671
##	as.numeric(B02001_004E)	-3.378673	1.070627	-3.155791	0.003344
##	as.numeric(B02001_005E)	0.096149	0.027864	3.450717	0.001512
##	as.numeric(B02001_006E)	0.620357	1.435570	0.432133	0.668373
##	as.numeric(B02001_007E)	-0.021297	0.044566	-0.477881	0.635793
##	as.numeric(B02001_008E)	-0.154053	0.309610	-0.497570	0.621990

```
# model fit summary
summary(fit_acs19_race_all)
```

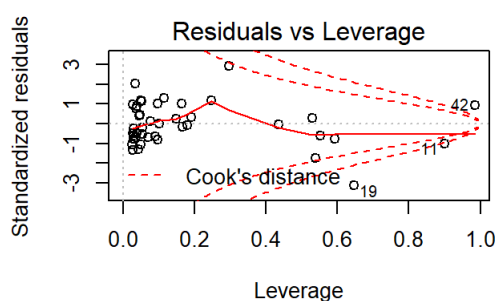
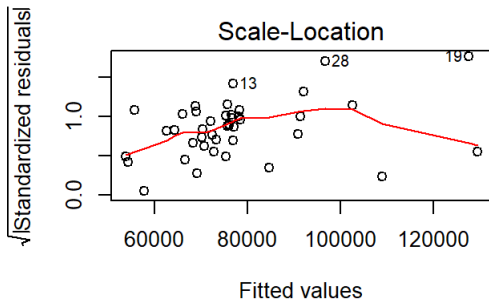
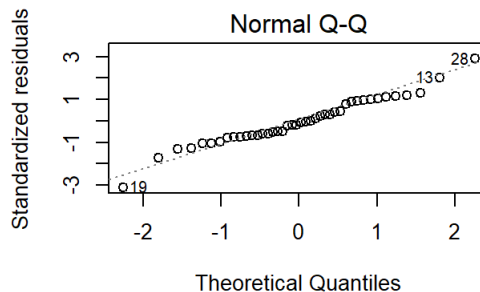
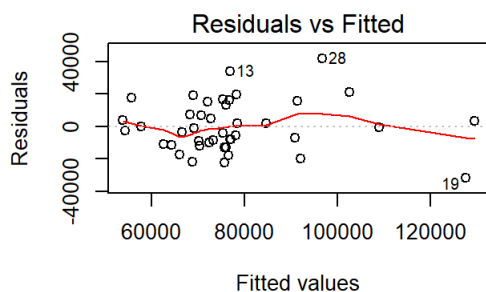


```
##
## Call:
## lm(formula = hh_median ~ as.numeric(B02001_002E) + as.numeric(B02001_003E) +
##   as.numeric(B02001_004E) + as.numeric(B02001_005E) + as.numeric(B02001_006E) +
##   as.numeric(B02001_007E) + as.numeric(B02001_008E), data = acs19_nofactor)
##
## Residuals:
##   Min       1Q   Median       3Q      Max
## -31736 -10621  -1961   11730   41733
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    7.653e+04  2.768e+03  27.649 < 2e-16 ***
## as.numeric(B02001_002E)  2.435e-02  2.043e-02   1.192  0.24156
## as.numeric(B02001_003E)  1.203e-01  1.034e-01   1.164  0.25267
## as.numeric(B02001_004E) -3.379e+00  1.071e+00  -3.156  0.00334 **
## as.numeric(B02001_005E)  9.615e-02  2.786e-02   3.451  0.00151 **
## as.numeric(B02001_006E)  6.204e-01  1.436e+00   0.432  0.66837
## as.numeric(B02001_007E) -2.130e-02  4.457e-02  -0.478  0.63579
## as.numeric(B02001_008E) -1.540e-01  3.096e-01  -0.498  0.62199
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 17090 on 34 degrees of freedom
## Multiple R-squared:  0.5264, Adjusted R-squared:  0.4289
## F-statistic: 5.399 on 7 and 34 DF, p-value: 0.0003163
```

```
# diagnostic plot
par(mfrow=c(2,2))
plot(fit_acs19_race_all)
```

```
## Warning in sqrt(crit * p * (1 - hh)/hh): NaNs produced
```

```
## Warning in sqrt(crit * p * (1 - hh)/hh): NaNs produced
```



Poverty level as a whole has a good relationship with median household income as verified with the low p-value. However when either of the two features are isolated, then their fit is poor so poverty is not particularly useful for the analysis.

The poor fit is likely due to the large size of each group (above or below poverty income level) which basically divides the dataset into two halves and results in poor fit due to large variance of each subset. Poverty level is intuitively a good indicator of median household income; however, the feature is not granular enough to use for additional analysis.

```
# features:
# B17001_002E: Income in the past 12 months below poverty level
# B17001_031E: Income in the past 12 months at or above poverty level

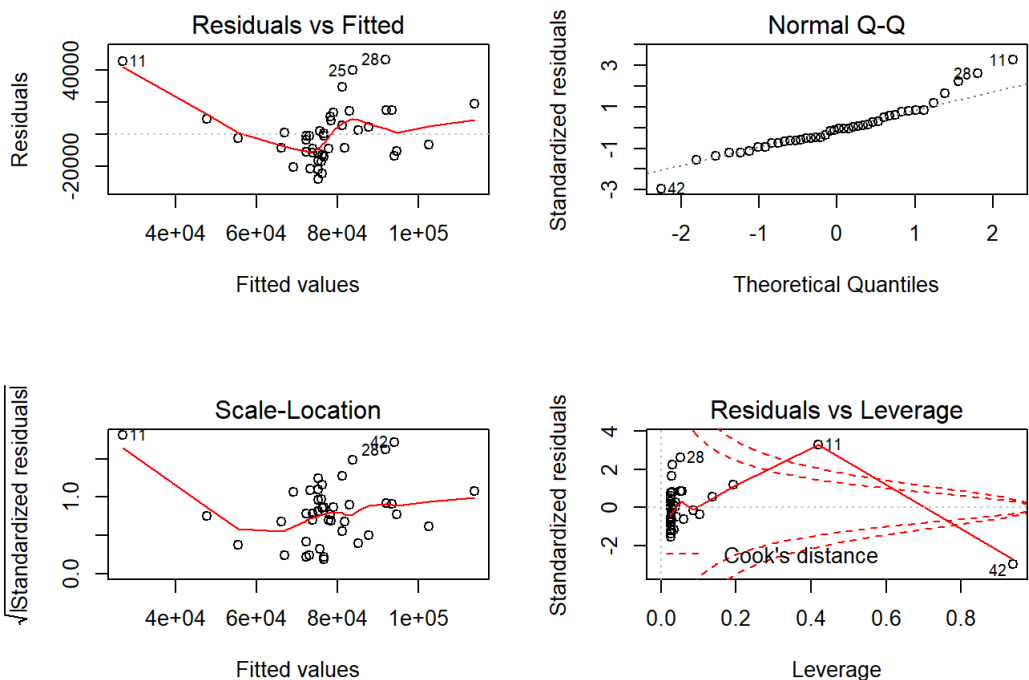
# Linear regression model fit - poverty below
fit_acs19_poverty_above <- lm(
  hh_median ~ poverty_above,
  data=acs19_nofactor
)
summary(fit_acs19_poverty_above)
```

```
##
## Call:
## lm(formula = hh_median ~ poverty_above, data = acs19_nofactor)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -29986 -15651  -5750  11721  61199
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  7.711e+04  3.689e+03  20.903  <2e-16 ***
## poverty_above 2.660e-04  6.684e-04   0.398   0.693
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 22860 on 40 degrees of freedom
## Multiple R-squared:  0.003942, Adjusted R-squared: -0.02096
## F-statistic: 0.1583 on 1 and 40 DF, p-value: 0.6928
```

```
# Linear regression model fit - poverty below + above
fit_acs19_poverty_all <- lm(
  hh_median ~ poverty_below +
  poverty_above,
  data=acs19_nofactor
)
summary(fit_acs19_poverty_all)
```

```
##
## Call:
## lm(formula = hh_median ~ poverty_below + poverty_above, data = acs19_nofactor)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -28063 -11906  -1760    9444  46747
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  7.670e+04  2.963e+03  25.890 < 2e-16 ***
## poverty_below -2.958e-01  6.159e-02  -4.803 2.33e-05 ***
## poverty_above  3.991e-02  8.270e-03   4.825 2.17e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 18350 on 39 degrees of freedom
## Multiple R-squared:  0.3742, Adjusted R-squared:  0.3421
## F-statistic: 11.66 on 2 and 39 DF, p-value: 0.0001074
```

```
# diagnostic plot
par(mfrow=c(2,2))
plot(fit_acs19_poverty_all)
```



## 2019 ACS - Model Refinement

Improve linear model with best features and analyze results.

Model fit was improved by including the best attributes from the data analysis which improved the model as verified with a lower p-value and higher r-squared value as compared to the other models with fewer features.

```

# fit best features to evaluate their importance
fit_acs19_best_features = lm(hh_median ~ as.numeric(housing_more_3000) +
  as.numeric(commute_transit) +
  as.numeric(commute_remote) +
  as.numeric(commute_carpool) +
  as.numeric(edu_doctorate) +
  as.numeric(edu_professional) +
  as.numeric(edu_master) +
  as.numeric(edu_bachelor),
  data=acs19_nofactor
)
summary(fit_acs19_best_features)

```

```

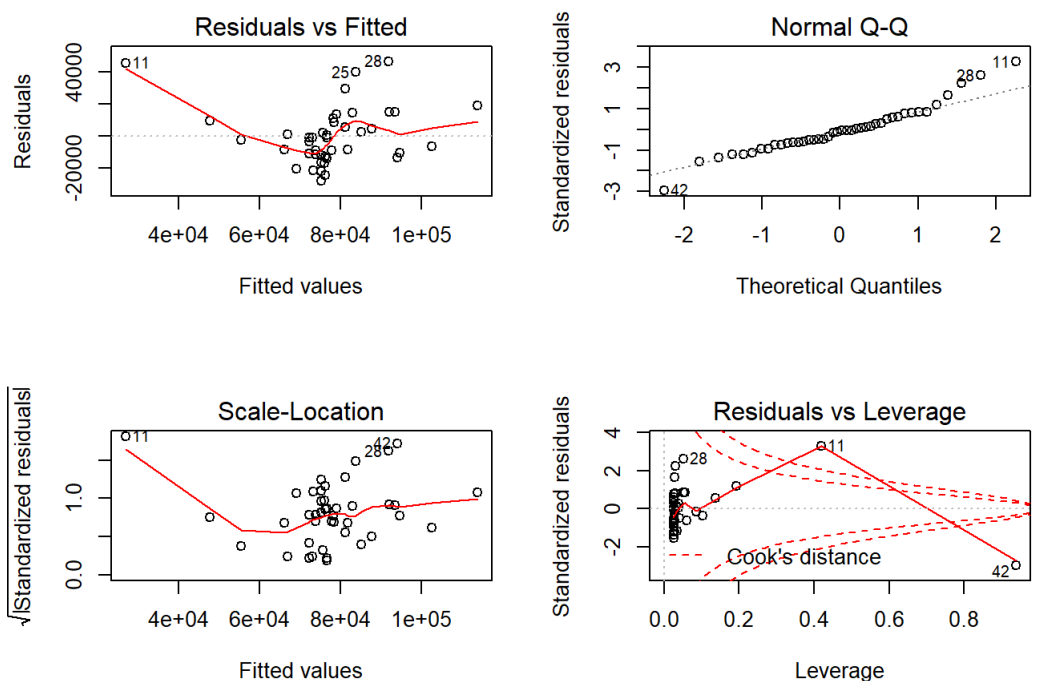
##
## Call:
## lm(formula = hh_median ~ as.numeric(housing_more_3000) + as.numeric(commute_transit) +
##   as.numeric(commute_remote) + as.numeric(commute_carpool) +
##   as.numeric(edu_doctorate) + as.numeric(edu_professional) +
##   as.numeric(edu_master) + as.numeric(edu_bachelor), data = acs19_nofactor)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -22574.4  -6485.2   -88.7   7790.4  24734.7
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      1.006e+05  5.571e+03  18.049  < 2e-16 ***
## as.numeric(housing_more_3000)  9.312e-01  1.924e-01   4.841  2.95e-05 ***
## as.numeric(commute_transit) -8.982e+02  2.258e+02  -3.979  0.000357 ***
## as.numeric(commute_remote)  -4.748e+02  2.656e+02  -1.788  0.082994 .
## as.numeric(commute_carpool)  1.997e+02  2.876e+02   0.694  0.492340
## as.numeric(edu_doctorate)    1.368e-01  1.068e+00   0.128  0.898860
## as.numeric(edu_professional)  2.929e-01  8.491e-01   0.345  0.732370
## as.numeric(edu_master)      -2.818e-01  3.039e-01  -0.927  0.360653
## as.numeric(edu_bachelor)    -2.654e-01  9.752e-02  -2.721  0.010296 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 12630 on 33 degrees of freedom
## Multiple R-squared:  0.7491, Adjusted R-squared:  0.6883
## F-statistic: 12.32 on 8 and 33 DF,  p-value: 5.775e-08

```

```

# diagnostic plot
par(mfrow=c(2,2))
plot(fit_acs19_poverty_all)

```



## NCHS COVID-19 Data

Import CDC COVID-19 data and compare with ACS data.

## Part 4: Data Visualization

### Summary Plots

Provide summary plots based on data analysis results.

## Part 5: Recommendations

### Key Findings

Provide recommendations based on data analysis results.

Some key findings are as follows:

1. Education attainment, commute mode and housing cost had good relationship with median household income.
2. Some features have a wide variance and outlier/high leverage points.
3. Combining best features resulted in better model fit (i.e. lower p-value and higher adjusted r-squared value).