

# Project

## Part A(page 1-5)

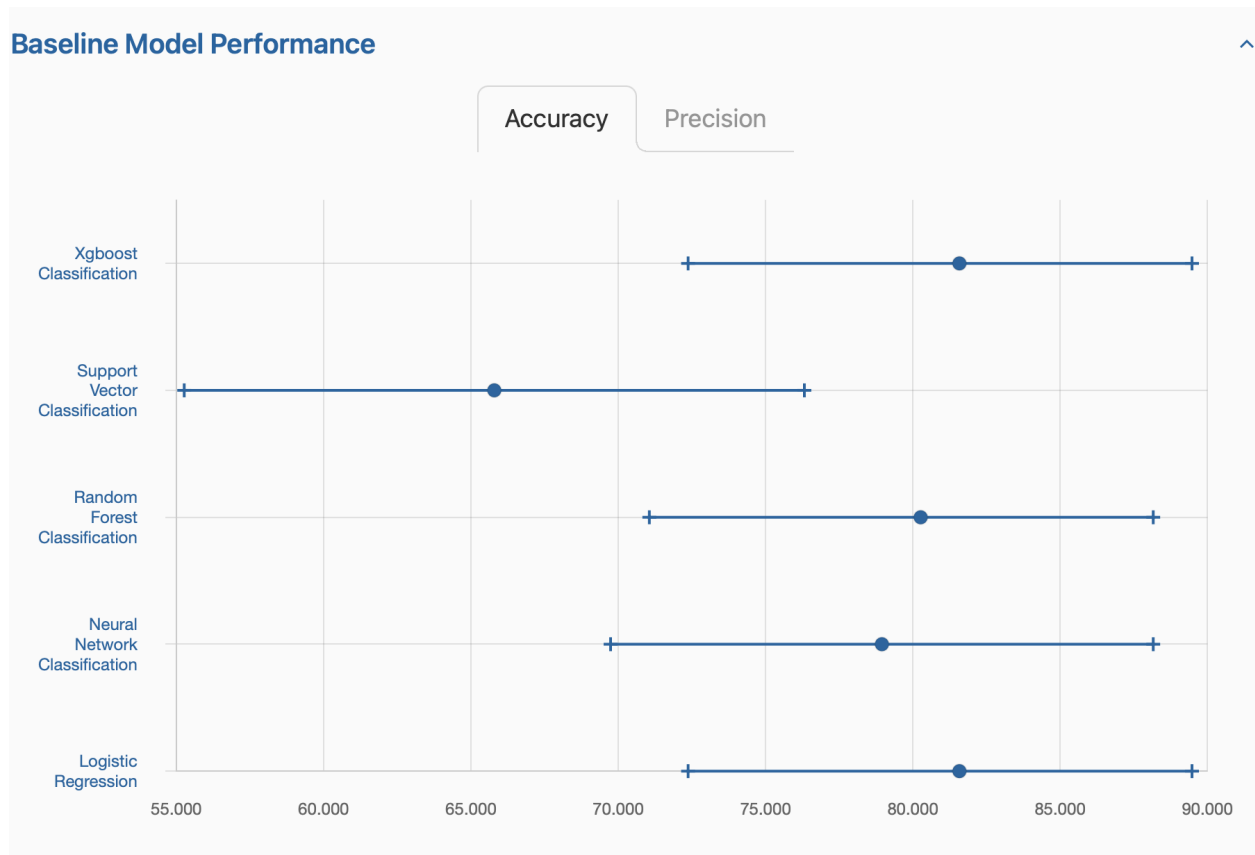


Figure 1: source:<https://archive.ics.uci.edu/dataset/45/heart+disease>

This graphic is from webpage of Heart Disease Datasets which is one of the most popular datasets in UC Irvine Machine Learning Repository. The graphic shows the accuracy of each machine learning method. The point in the middle shows the mean performance of each ML method while the 2 boundaries show the best and the worst performance. The graphic is brief and direct. We can see Logistic Regression and Xgboost Classification are likely to perform better than other methods. And in the case Support Vector Classification is likely to perform worst.

The advantages of the graphic are showing the average performance directly and showing random effect in testing machine learning. If the graphic could sort the method from top to bottom, the comparison may be more direct.

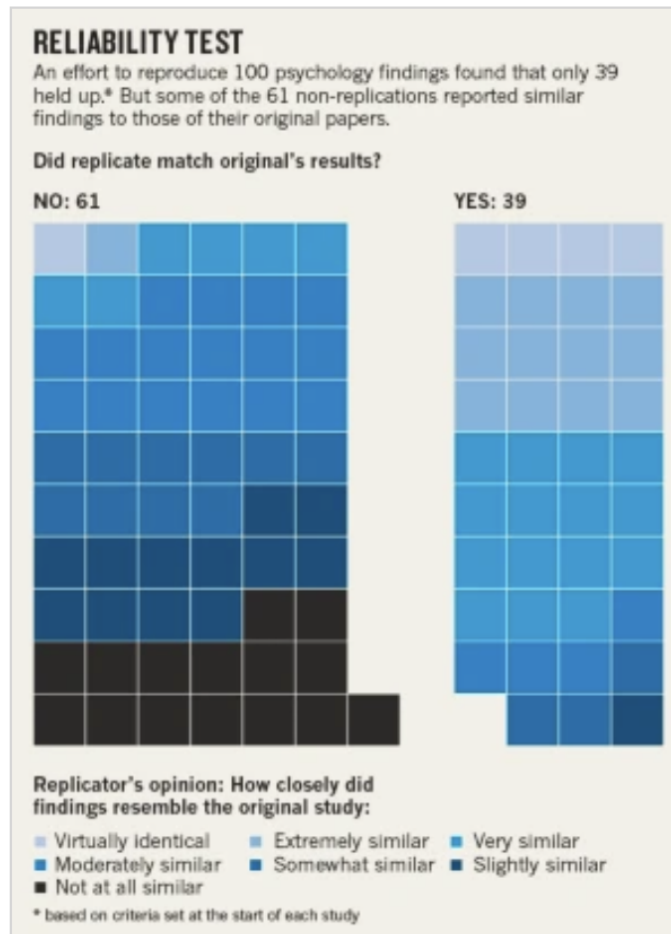


Figure 2: source:<https://www.nature.com/articles/nature.2015.18248>

The picture directly shows the ratio of successful replication results which strongly suggest that many psychology studies should be doubted. Besides, it also shows the similarity between original paper and reproducing results by color from light to dark (which always means negative), which is also very direct and brief.

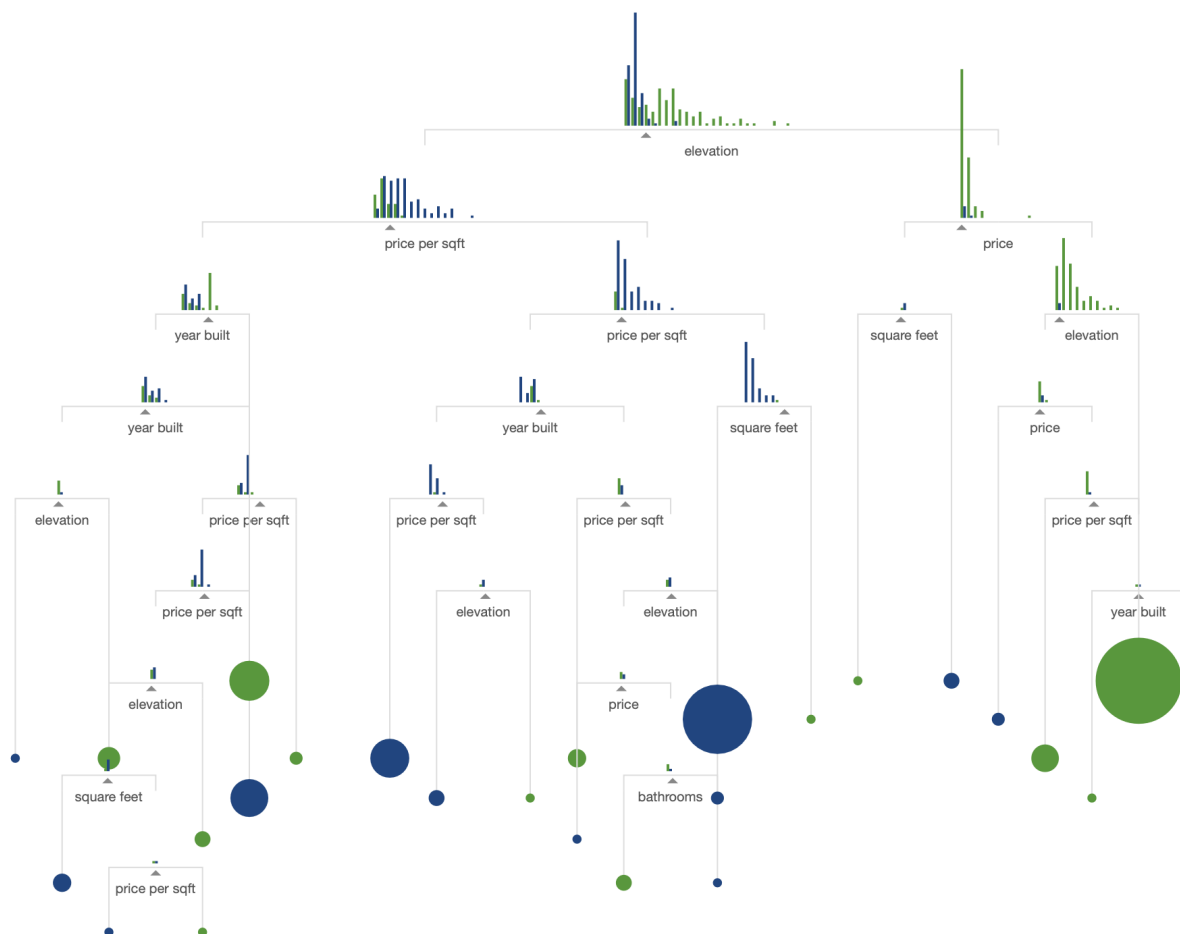


Figure 3: source:<http://www.r2d3.us/visual-intro-to-machine-learning-part-1/>

The graphic show how tree decision method which is one of the classical machine learning method works. The 2 colored histogram on each nodes clearly shows how it split a data subset into 2 partition depending on point on x( different covariates ) axis.

The original graphic is dynamic which really help ML starter learn the principle of tree decision quickly.

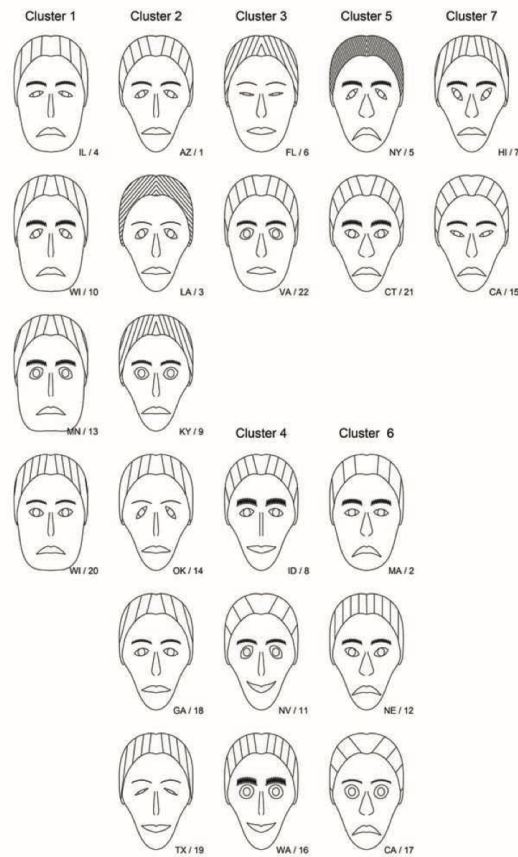


Figure 2. Chernoff faces for 22 public utilities

Figure 4: [source:https://journals.sagepub.com/doi/pdf/10.1177/1536867X0900900302](https://journals.sagepub.com/doi/pdf/10.1177/1536867X0900900302)

This is also an interesting graphic which use chernoff face to show different clusters. We can also find companies in the same cluster share what in similarity quickly because people are sensitive to face features .

It may cost a little time to describe these features with raw data..However, This graph may really be useful for people to see the hidden structure of data , how well cluster works and a good and quick approach to show high dimensional(more than 2D) data as well.

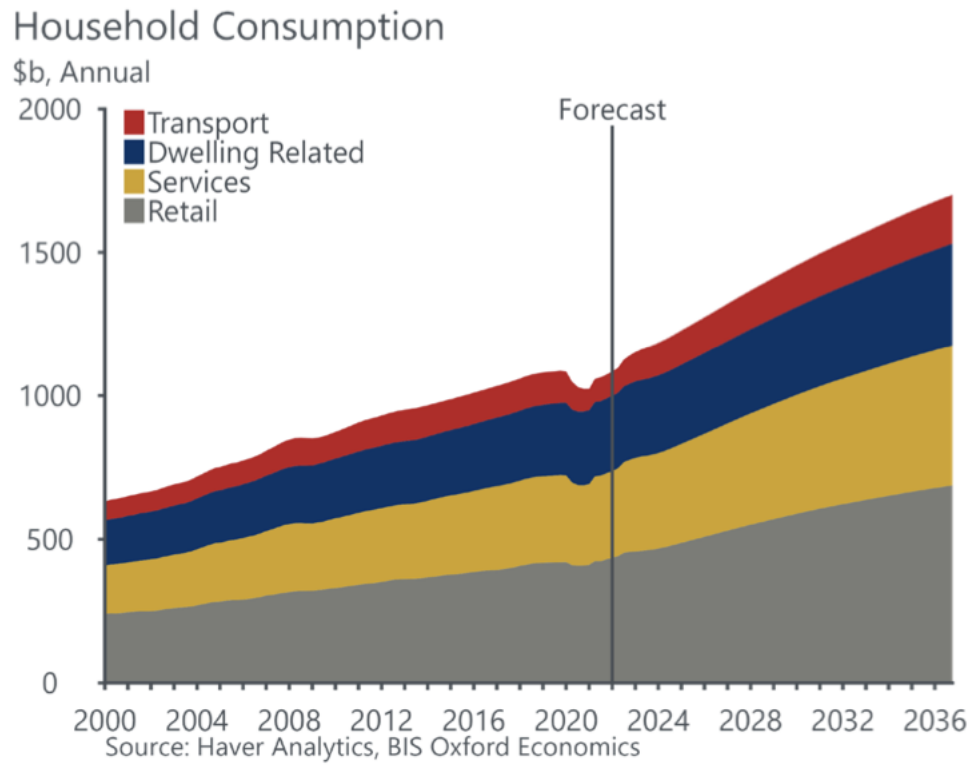


Figure 5: source:<https://www.portofmelbourne.com/wp-content/uploads/BISOE-Port-of-Melbourne-Container-Forecasts-and-Sensitivities-20221223.pdf>

It shows the household consumption in different part and total in one graphic. Besides, it also shows the change and the rate of change. The data on 2020 which showed an abnormal down because of the pandemic.

The author wanted to show the increasing trend of the household consumption therefore plot the forecast. It may be clear to see the total growth and retail growth but a little harder to get the growth of other part due to the graphic type. But since the retail accounts most, this may be reasonable. And we can compare the growth by slope.

## Part B (page 6-13)

### Introduction

This part focus on the analysis on Insurance availability in Chicago. The data has full description in GDA project description. The original variables include Zip, Race, Fire, Theft, Age, Volun, Invol and Income.

Zip which represents different neighborhood is the identity of each neighborhood. Race, Fire, Theft, Age and Income may be our important covariates espically Race and Age since the objectives of the study are to explore the extent to which racial composition and age of housing affect underwriting practices.

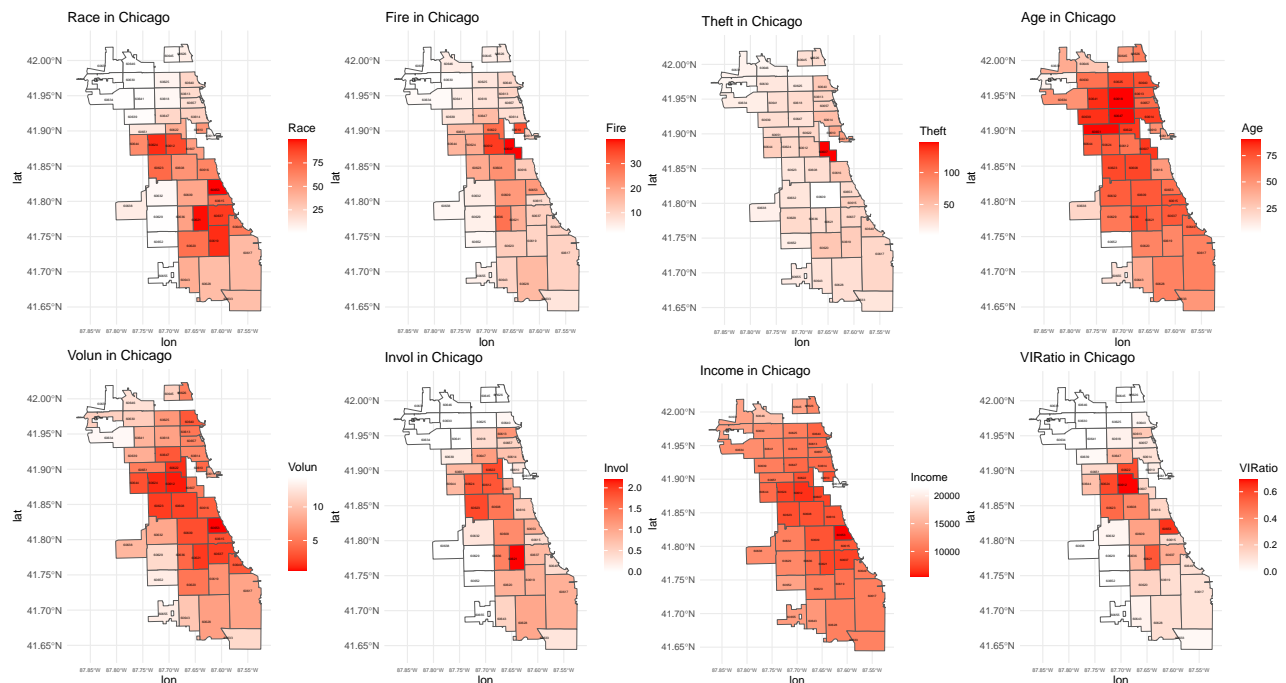
Volun represents the a rate of “accepted” requests to purchase insurance while Invol is likely to play the opposite role meaning rate of “rejected”. Since these 2 variables are rate, to some extent, they can show the rate of acceptance and rejection regardless of size of neighborhood. However, the ratio between  $Invol$  and  $Volun + Invol$  may be more reasonable to represent the extent of rejection. So here I produce a variable  $VIRatio$  which is  $\frac{Invol}{Volun + Invol}$ . Besides, insurance was probably not compulsory so we can see the  $max(\frac{Volun}{100})$  is around 0.1 while  $max(\frac{Invol}{100})$  is around 0.01.

Let's first have a general picture of the data.

### General picture

The 8 graphics below show Race, Fire, Theft, Age, Volun, Invol, Income and  $VIRatio$  on the map of Chicago. We may find that the  $VIRatio$  in the center part of the Chicago is more liekly to be lower. Besides, it seems that the graphics of Race, Fire, Invol and  $VIRatio$  share some pattern (deep red and light red seem to have similar areas) which suggests high positive correlation between variables.

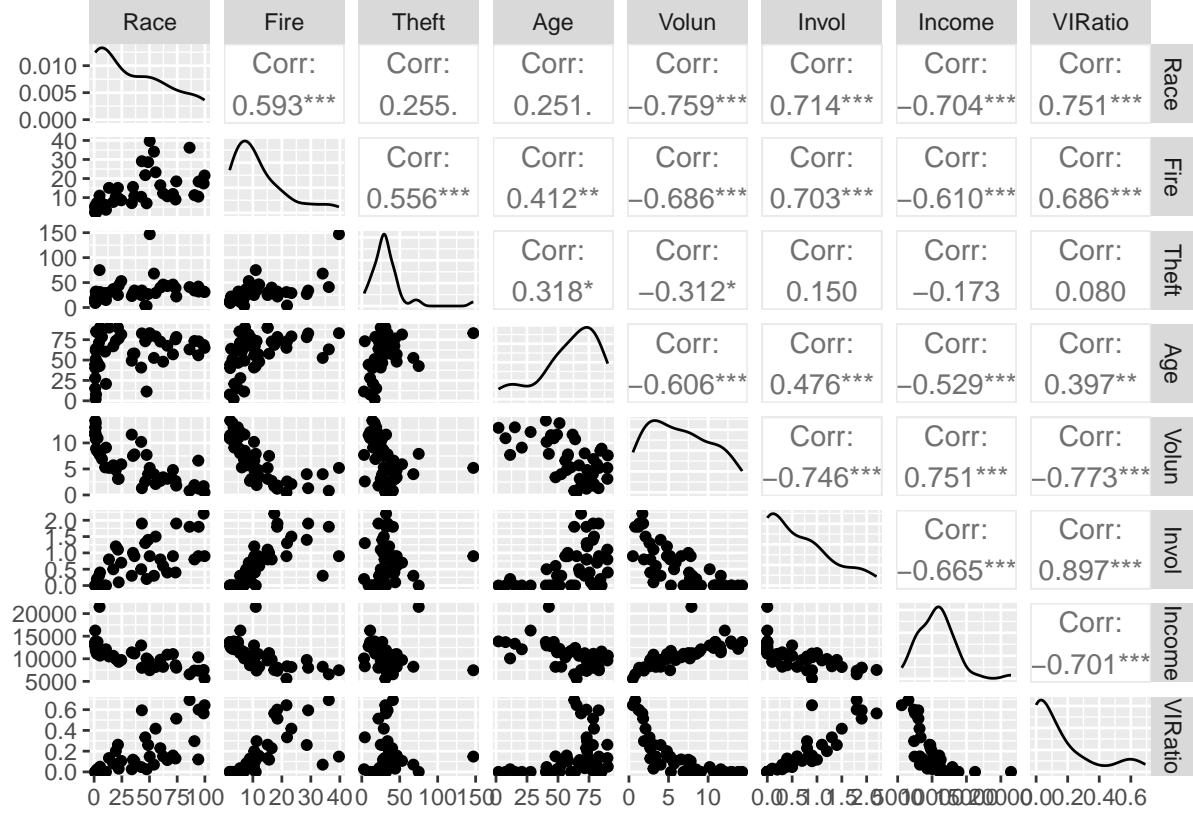
Finally, we also need to pay attention to the areas where  $Invol$  and  $VIRatio$  are zero. These areas are all in almost white color meaning very low in Race (almost 0) but show similar Fire, Theft and Age compared to areas around or near them. Meanwhile, these 0-Invol areas also show different house age which may suggests the relationship between “reject” and house age is weak. These features may suggest the race discrimination may happen in making the desicison of “accept” or “reject”.



## Matrix of correlation and scatter plot

Take a further look on the correlation between each variables. From the matrix of plots below, it seems that Race really show high correlation with VIRatio which is corresponding to our finding in general picture's pattern. We may also find Race show a strong linear relationship with Volun and Invol which may also suggest affect the underwriting practice to a large extent. We may found fire and Income show high correlation with IVRatio while Theft show low correlation with IVRatio.

Besides, we can also easily learn the correlation between covariates such as Race ,Fire and income. The high correlation and strong linear pattern of the scatter plot between Race and Income is also corresponding to the pattern of deeper red closer to center if we review the general picture.



## Regression model

Now, let us make a regression model to have a deeper view on the relationship among several variables. Our response variable may be one of Volun , Invol and VIRatio. VIRatio is a product combining Volun and Invo and reflect the close ratio between rejection and all insurance requests. So we may first try

$$VIRatio_i = \beta_0 + \beta_1 Race_i + \beta_2 Fire_i + \beta_3 Theft_i + \beta_4 Age_i + \beta_5 Income_i + \epsilon_i$$

and

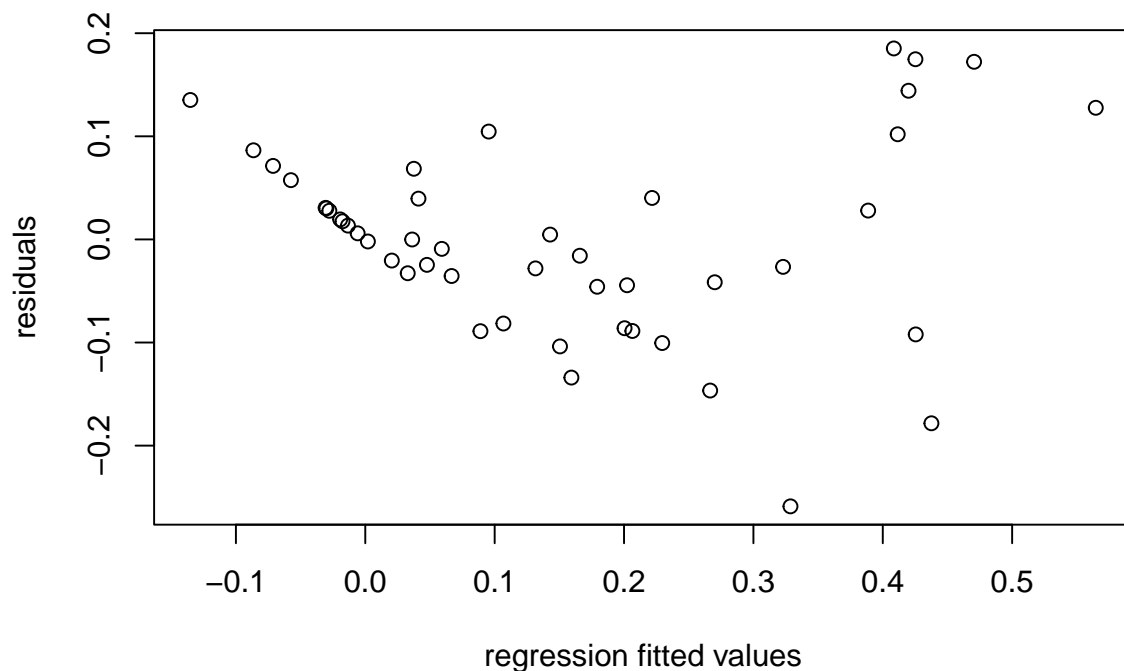
$$\epsilon_i \sim iid Norm(0, \sigma^2)$$

From the output we get

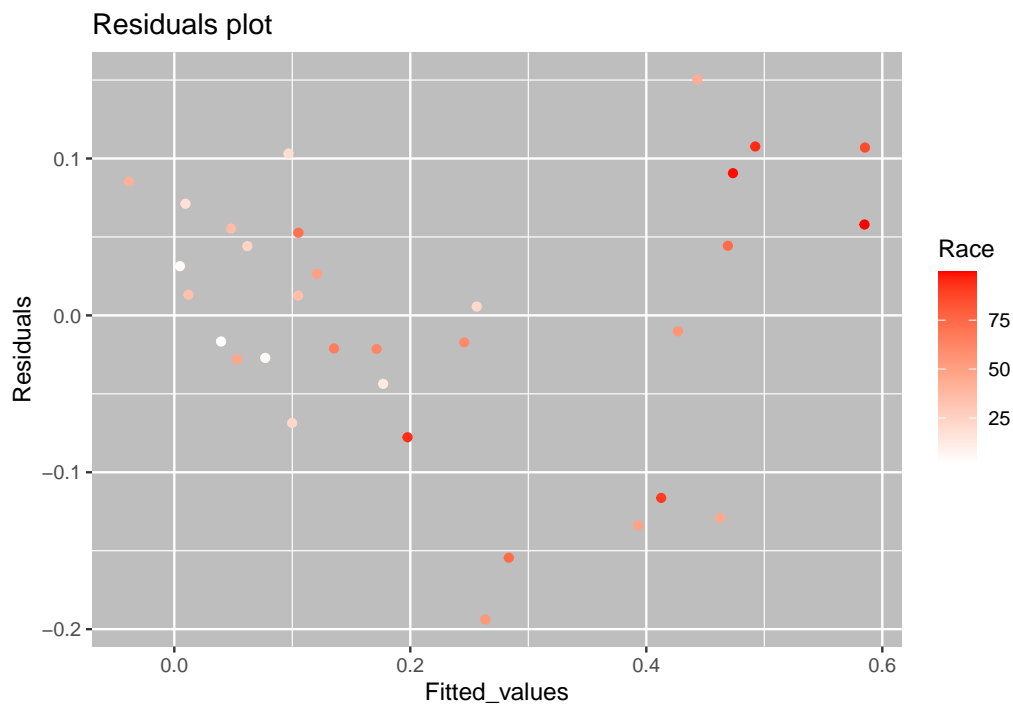
$$VIRatio_i = \frac{-698.9749 + 29.64232Race_i + 119.0461Fire_i - 36.59007Theft_i + 15.70569Age_i - 0.001532509Income_i}{10^4}$$

The output of linear regression show that the coefficient  $\beta_1, \beta_2, \beta_3$  are significant and indicating that in the model the coefficient of Race, Fire and Theft are important and should not be zero. However, when we plot the residuals ,we may find that the spread increases when the fitted value increases which may suggest transform

the  $VIRatio$  to  $\log(VIRaion)$ . Besides, we also found a interesting pattern where low fitted values have a strong linear relationship with residuals. The strange points mainly come from the true value of  $VIRatio_i$  which is 0. So the strange pattern of residuals can be explained because all these points have the same  $VIRatio$  but different covariates values and the regression model is a linear predictor.



**Refitting on data with  $VIRatio \neq 0$**



Let's first remove these  $VIRatio_i = 0$  points, fit the rest date with the same model again. The figure above shows the spread of residuals increase again when fitted values increase. So we apply a log or square-root

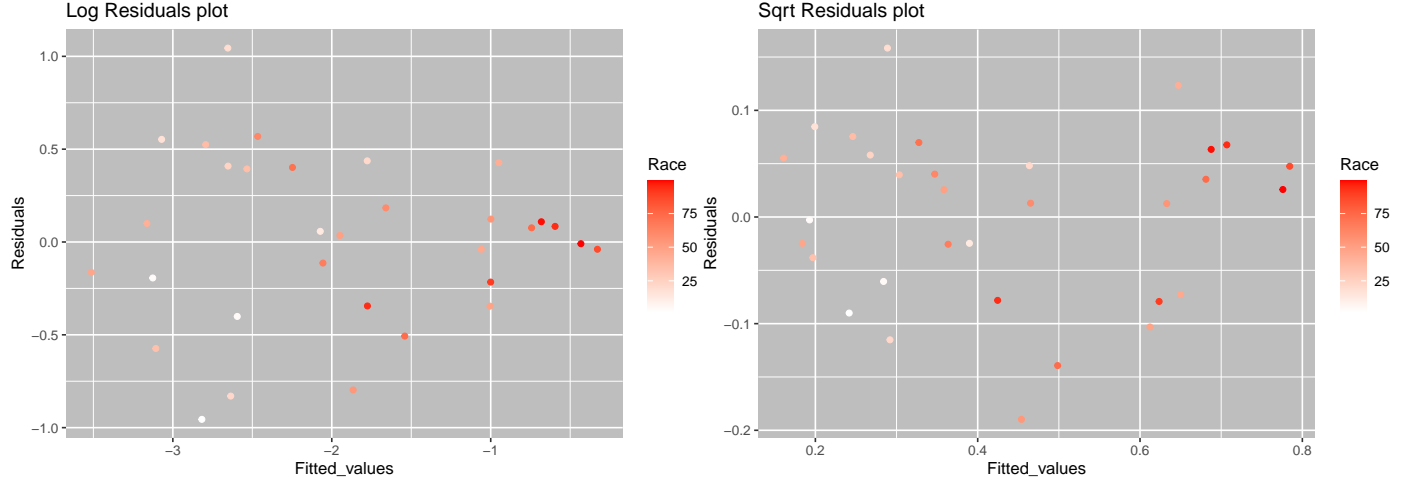


transformation to  $VIRatio_i$ . Then fit the following model:

$$\log(VIRatio_i) = \beta_0 + \beta_1 Race_i + \beta_2 Fire_i + \beta_3 Theft_i + \beta_4 Age_i + \beta_5 Income_i + \epsilon_i \quad (1)$$

$$\sqrt{VIRatio_i} = \beta_0 + \beta_1 Race_i + \beta_2 Fire_i + \beta_3 Theft_i + \beta_4 Age_i + \beta_5 Income_i + \epsilon_i \quad (2)$$

After refitting, from the plot below we may find the residuals pattern of fitted model (2) perform better than fitted model (1). So we choose the model (2).



Then we may want to see the coefficients' significance to know whether we can reduce numbers of covariates to make model easier.

```
##
## Call:
## lm(formula = sqrt(yc) ~ xc)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.18975 -0.06353  0.01926  0.05597  0.15839
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  2.350e-01  2.342e-01   1.003  0.324874
## xcRace       3.383e-03  7.872e-04   4.297  0.000215 ***
## xcFire       9.002e-03  2.441e-03   3.688  0.001049 **
## xcTheft     -3.772e-03  7.675e-04  -4.915  4.2e-05 ***
## xcAge        4.635e-03  1.128e-03   4.109  0.000352 ***
## xcIncome    -3.135e-05  1.563e-05  -2.006  0.055348 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.08678 on 26 degrees of freedom
## Multiple R-squared:  0.856, Adjusted R-squared:  0.8283
## F-statistic: 30.91 on 5 and 26 DF, p-value: 3.736e-10
```

The outputs suggest that Intercept can not be rejected to be zero. So we may have a try to get a new reduced model:

$$\sqrt{VIRatio_i} = \beta_1 Race_i + \beta_2 Fire_i + \beta_3 Theft_i + \beta_4 Age_i + \beta_5 Income_i + \epsilon_i \quad (3)$$

```
##              Estimate Std. Error t value Pr(>|t|)
## xcRace       3.967136e-03 5.294928e-04  7.492333 4.644002e-08
## xcFire       1.008156e-02 2.191624e-03  4.600044 8.921816e-05
```

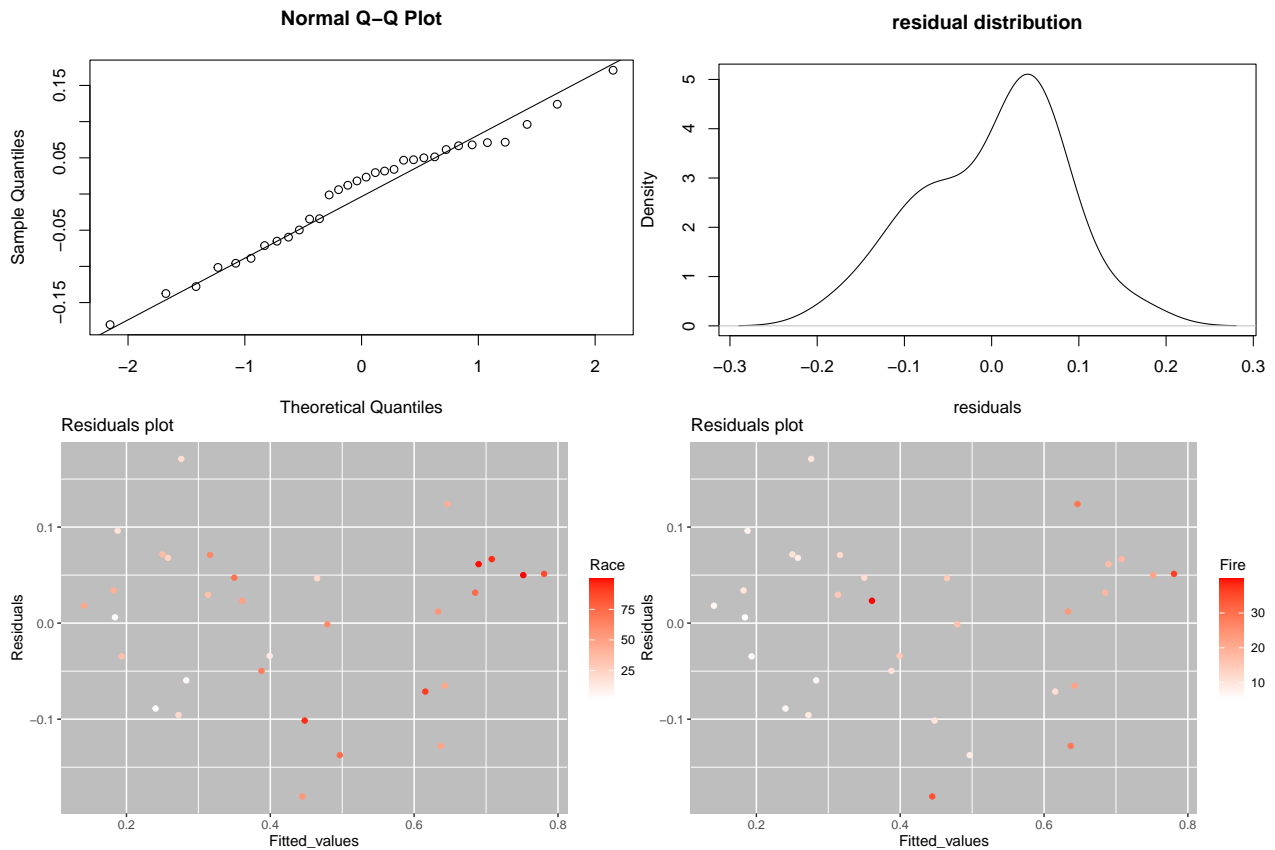
```
## xcTheft -3.816870e-03 7.663213e-04 -4.980769 3.207096e-05
## xcAge 5.378401e-03 8.503346e-04 6.325040 9.025759e-07
## xcIncome -1.661485e-05 5.345794e-06 -3.108023 4.401691e-03
```

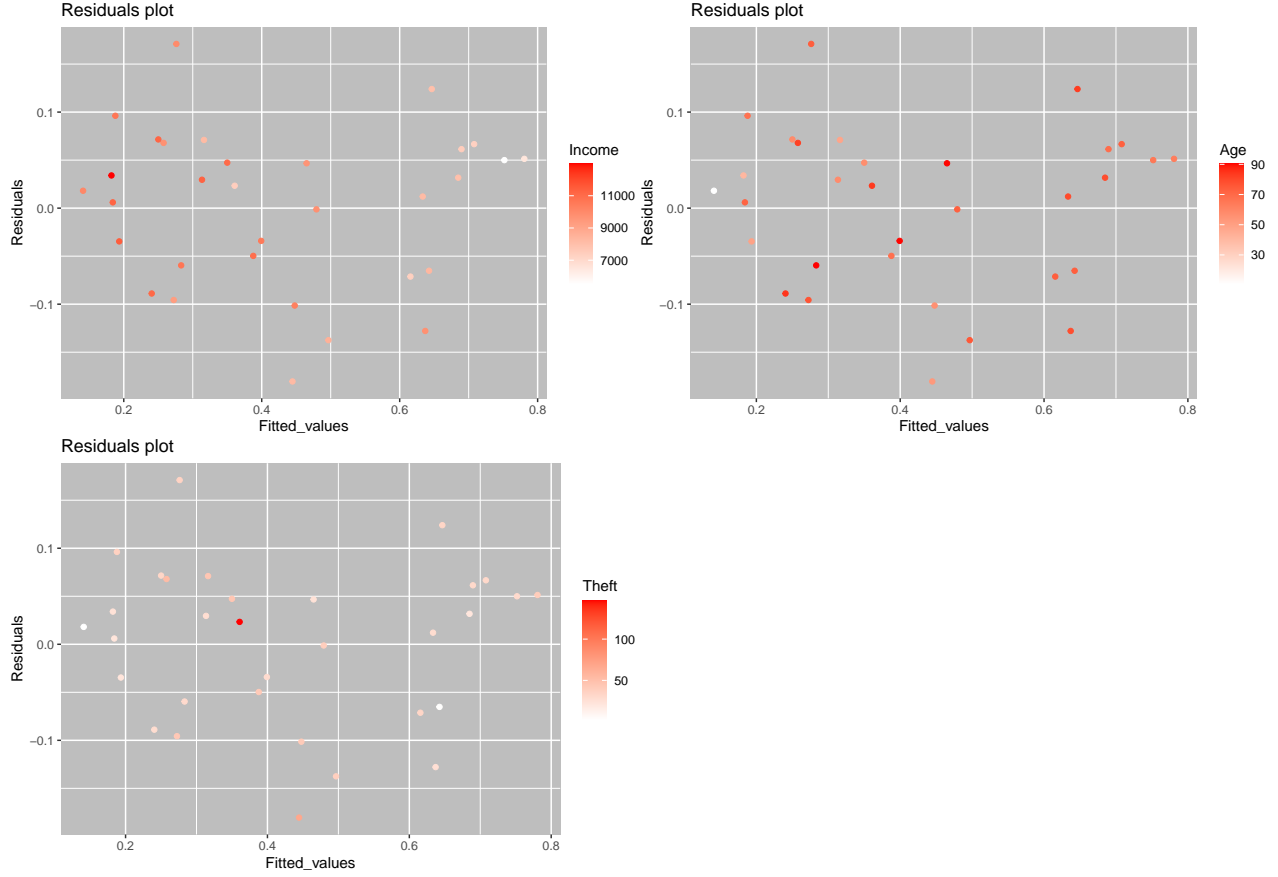
The output seems show that all coefficients are significant in this reduced model 3 so we can't reject each coefficients equals to zero. To see more residuals plot to see whether the residuals have a good pattern and how well they fits iid normal distribution assumption, we usually apply a QQ plot. The qqnorm seems to perform neither good nor bad because the point is not so near the qqline near zero while the overall performance seems not bad. We may also take a look at the residuals density plot which also explain the QQ plot results. The reasons for the not such good results of residuals could be small volume of data.

$$\sqrt{\hat{VIRatio}_i} = 0.003967Race_i + 0.010082Fire_i - 0.003817Theft_i + 0.005378Age_i + -0.000017Income_i \quad (4)$$

However ,the performance of this model 4 may be good enough to explain the data and residuals and show how Race and Age affect VIRatio. Meanwhile, from the colored residuals plot we ma also find the Race , Fire and Income play more important roles because the dark red and light red data points seem to be classified by fitted values. But the different colored points of Age and Theft are much more widely distributed on Fitted values which may suggest Age and Theft are weak factors.

Finally, the coffeicient of Theft was negative which may be contradict to the reality. The reasons may be various such as reduced data ,not good enough model and so on.





## Modeling on Volun as response variable

The above model we pay much attention on “reject” and let us view the “accept” which may not completely showed by “reject” because the problems with data eg.  $Invol = 0$  or much more complex fact in reality. Similarly, our model is

$$\sqrt{Volun} = \beta_0 + \beta_1 Race_i + \beta_2 Fire_i + \beta_3 Theft_i + \beta_4 Age_i + \beta_5 Income_i + \epsilon_i \quad (5)$$

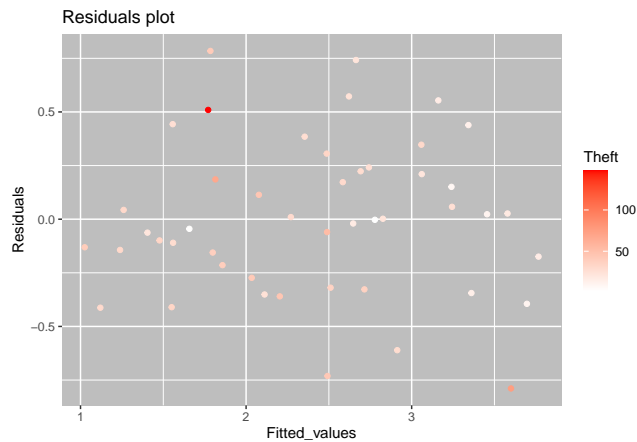
```
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  3.3869149122 5.615180e-01  6.031712 3.899719e-07
## xRace        -0.0130384752 2.625647e-03 -4.965814 1.251931e-05
## xFire        -0.0257127244 9.564085e-03 -2.688467 1.032835e-02
## xTheft        0.0053555598 3.234536e-03  1.655743 1.054093e-01
## xAge         -0.0114763086 3.153662e-03 -3.639042 7.586059e-04
## xIncome       0.0000301161 3.593704e-05  0.838024 4.068762e-01
```

The fitted results of model 5 seems good

The fitted model is given by:

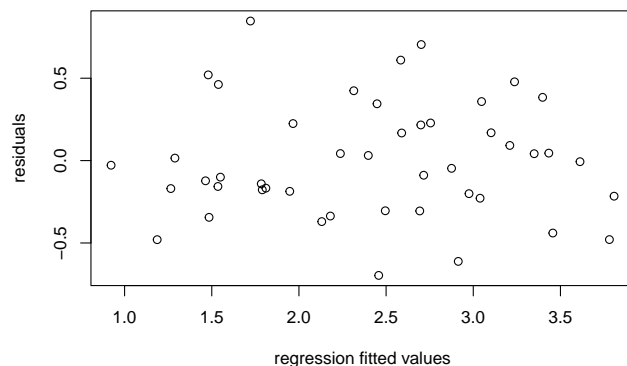
$$\sqrt{Volun} = 3.386915 - 0.013038 Race_i - 0.025713 Fire_i + 0.005356 Theft_i - 0.011476 Age_i + 0.000030 Income_i \quad (6)$$

Similarly, we find a strange positive value as coefficient of Theft. The reason may be the extreme data point (deep red) as shown in the below figure cause overfitting.



Let's remove the point and refit model. After several practices on model fitting and selection.

```
##
## Call:
## lm(formula = sqrt(y) ~ x)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.69681 -0.21222 -0.03754  0.22263  0.84710
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  3.958106   0.154645  25.595 < 2e-16 ***
## xRace        -0.013179   0.002109  -6.250 1.73e-07 ***
## xFire        -0.031842   0.008641  -3.685 0.000649 ***
## xAge         -0.011837   0.002597  -4.558 4.41e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.362 on 42 degrees of freedom
## Multiple R-squared:  0.8283, Adjusted R-squared:  0.816
## F-statistic: 67.52 on 3 and 42 DF,  p-value: 4.115e-16
```



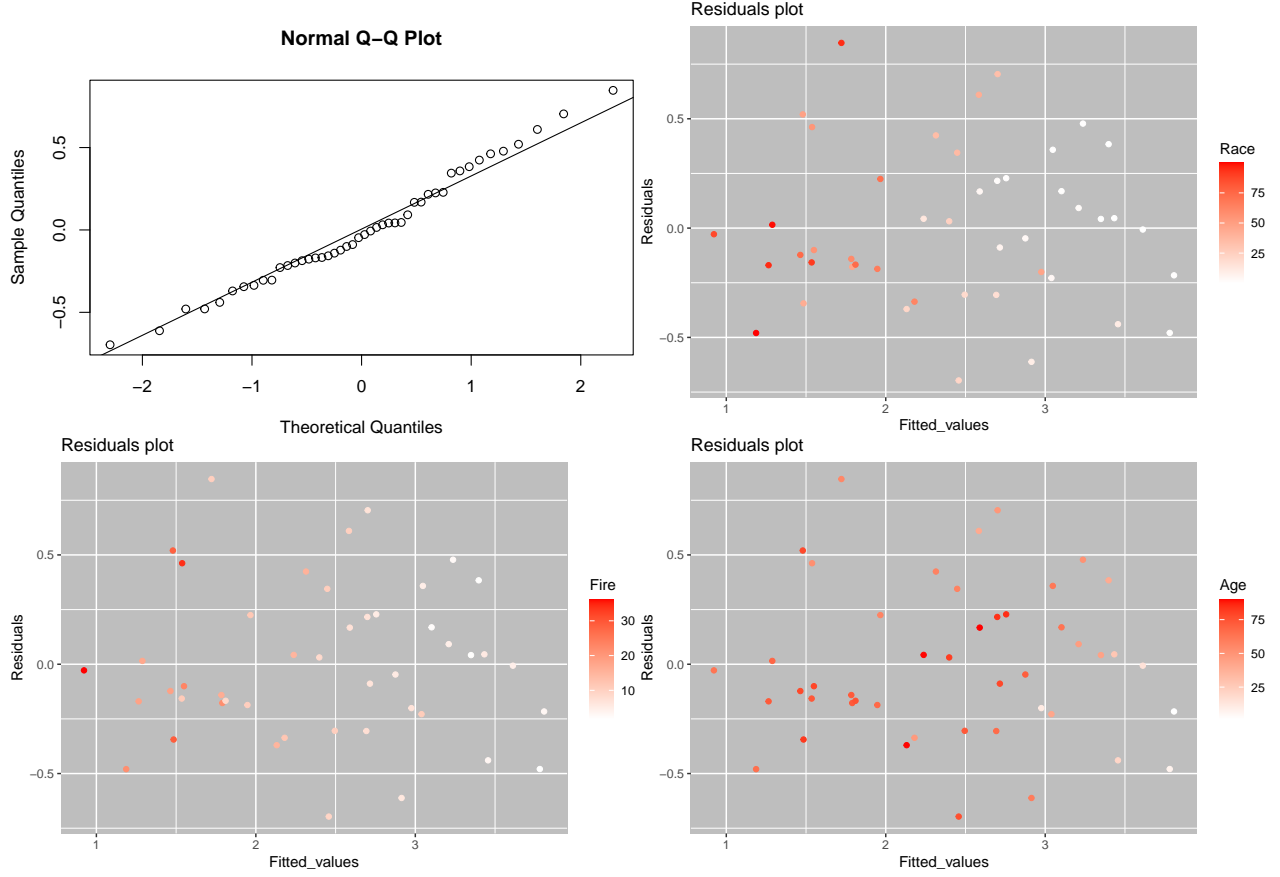
The final model is given by:

$$\sqrt{Volun} = \beta_0 + \beta_1 Race_i + \beta_2 Fire_i + \beta_3 Age_i + \epsilon_i \quad (7)$$

The fitted model is:

$$\sqrt{\hat{Volun}} = 3.95811 - 0.01318 Race_i - 0.03184 Fire_i - 0.01184 Age_i \quad (8)$$

The result of model (7) show that all three factors have negative impact on  $\sqrt{Volun}$ . The residuals show good patterns shown by QQ plot. We can see the general color change while fitted values increases clearly from the plot below which also suggests a good fitted results. So to conclude, Race, Fire and Age affect the acceptance to large extent.



## Conclusion

From the overall analysis, we may find racial composition affect the underwriting practice to a large extent while the age of house do not show such strong effect on “rejection”. However, when we model on Volun as response variable, we find both racial composition and age of housing affect the underwriting practices to an obvious large extent after controlling for factors like fire and others according to each coefficients meaning in regression model is showing how response variable change with the covariate increase when keeping the other covariates unchanging. The racial composition plays an important role may imply the race discrimination exists in the insurance company while the age of house also may be taken into account of decision of insurance company. To summarize, both of both racial composition and age of housing may affect the “accept” and “reject” significantly according to the 2 model (4 and 8) when the other condition keep the same. But model (4) has a strange interpretation on coefficient of Theft which may need to be improved.