

# Code For Project 1 - Applied Statistics

Created by: Nachi Lieder 314399114

## Introduction

This is a report which its main purpose is to examine the variables affecting the BMI of the respondents. We will go through several stages in this report

- Data formatting - appending the data into one set that is manageable.
- Feature Exploration
- Dealing With Missing Data
- Response Variable Exploration
- Univariate regressions - initial screening
- Multivariate regressions - in depth screening & attempt to describe the BMI target using the given parameters.

We will explore the data , attempt to understand the behavior of different predictors, and try to use them to explain the behavior of the BMI target. We will be using univariate regressions to understand how each individual predictor behaves , whether it has a potential to contribute in the analysis ,and filter otherwise.

Later we will perform a Stepwise AIC to filter the recommended features that will aid us in the prediction. To enhance this we will add to this analysis the interactions between the different features and re-assess the process. We will find that this is proven to be helpful and indeed will create a good picture to the BMI target and its prediction , as well as the characterization of the Y variable.

In addition , I will supply some added depth to enrich the analysis by attempting to cluster and group the observants using the input features , and test whether its definable and applicable to characterize the BMI per group using the other features. (more on this will follow in the report.)

It should be noted that in this analysis I used the initial set of features , and tested of scope (within the code) the potential addition of the remaining feautres . I have found these features to be fruitful for our analysis and would recomend attempting to add tem in the next phase( spcifically the interactions between the survey data and the demographic data)

## Part 1

In this part I recieve as an input the 4 datasets , merge them , and combine 1 whole dataset. You can see in the result the dimensions of the resulted dataframe which we will explore. I merged each genders two datasets together (the demographic data and survey data) and appended it together. The final result here is a single dataframe with the dimension of 4036 x 57.

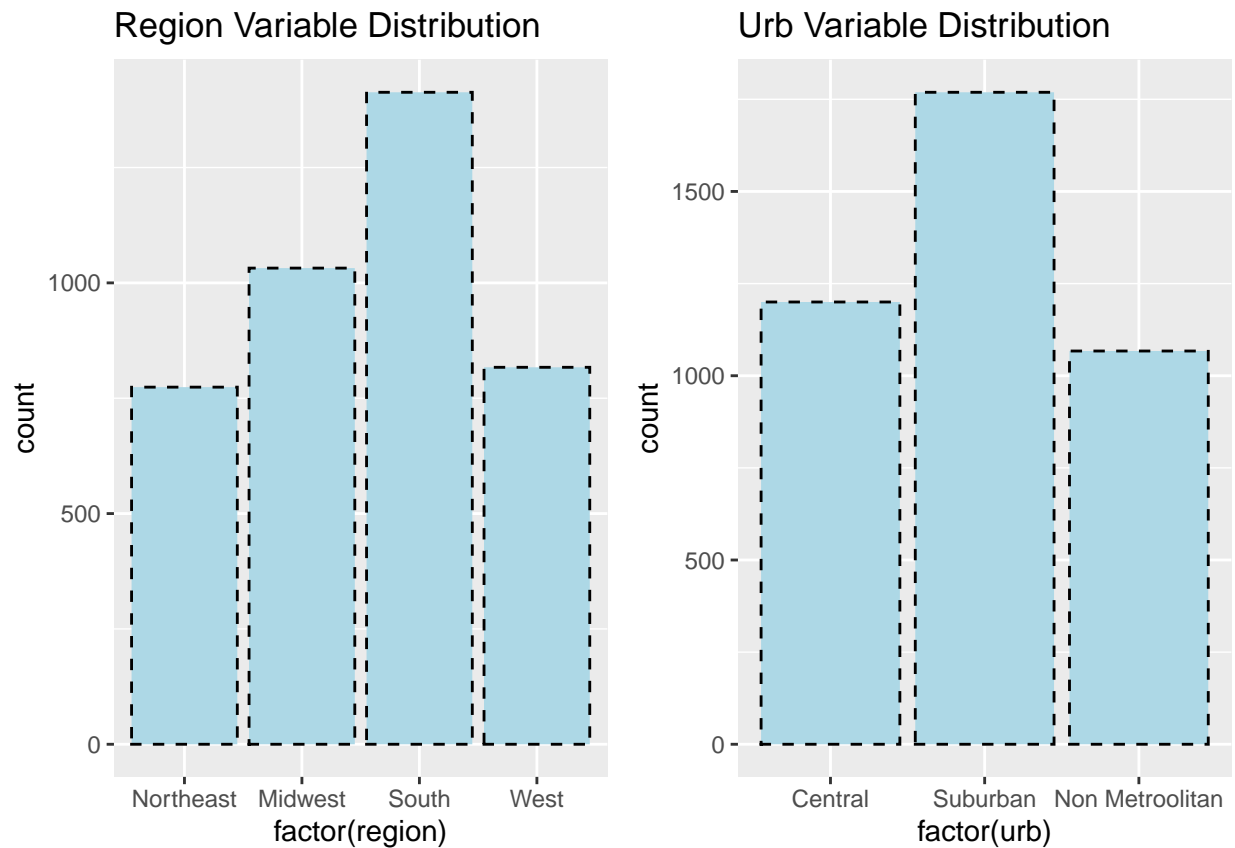
## Part 2 - Exploration

In this part I will perform an initial exploration on the datasets features (predictors) , and attempt to understand the characteristics of the predictors. Lets start off with looking at the key features:

1. Region
2. Urb
3. Income
4. Age
5. Gender
6. Grade
7. Exercise

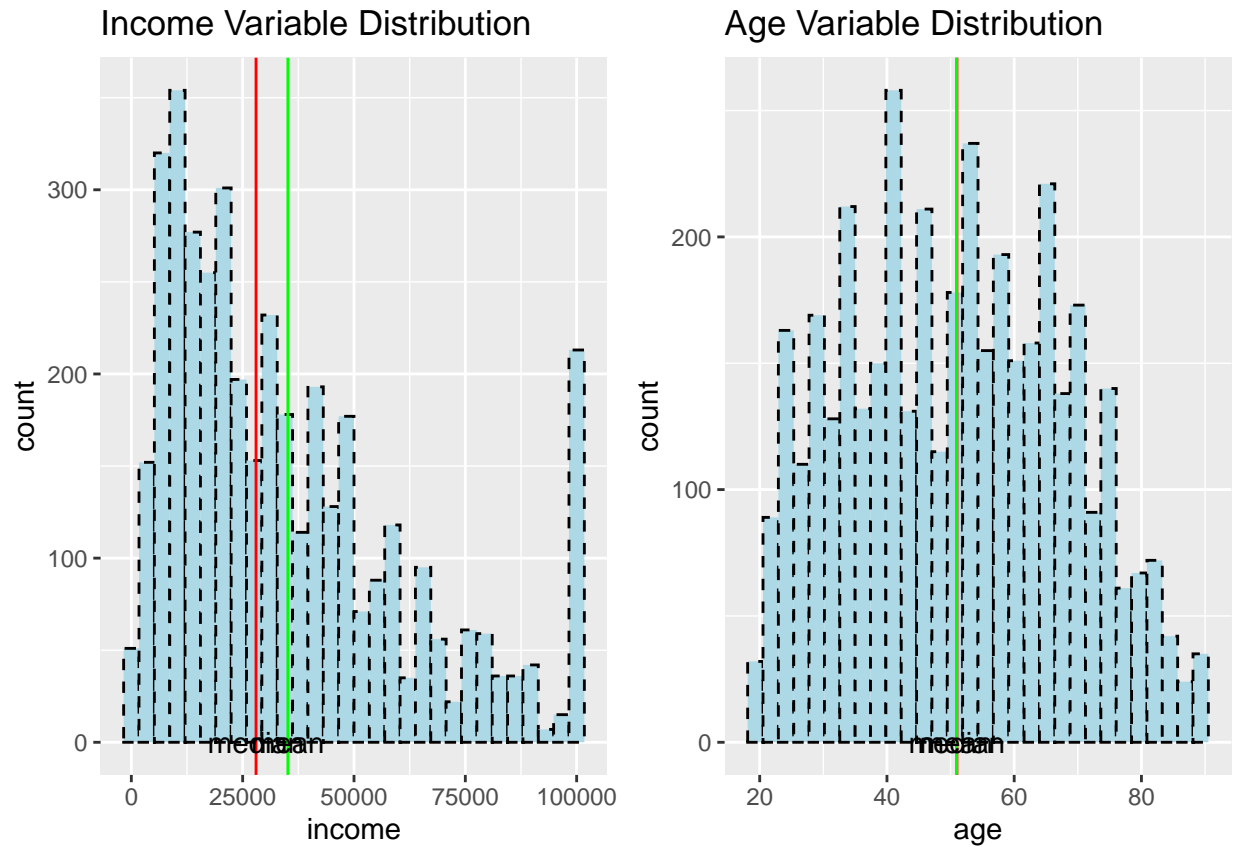
Lets observe the Region feature: We can see that overall there is a slight advantage to observations of patients coming from the South.

Next, lets observe the Urb feature: Here what we are observing is a slight advantage in the count for the Suburban originated patients.



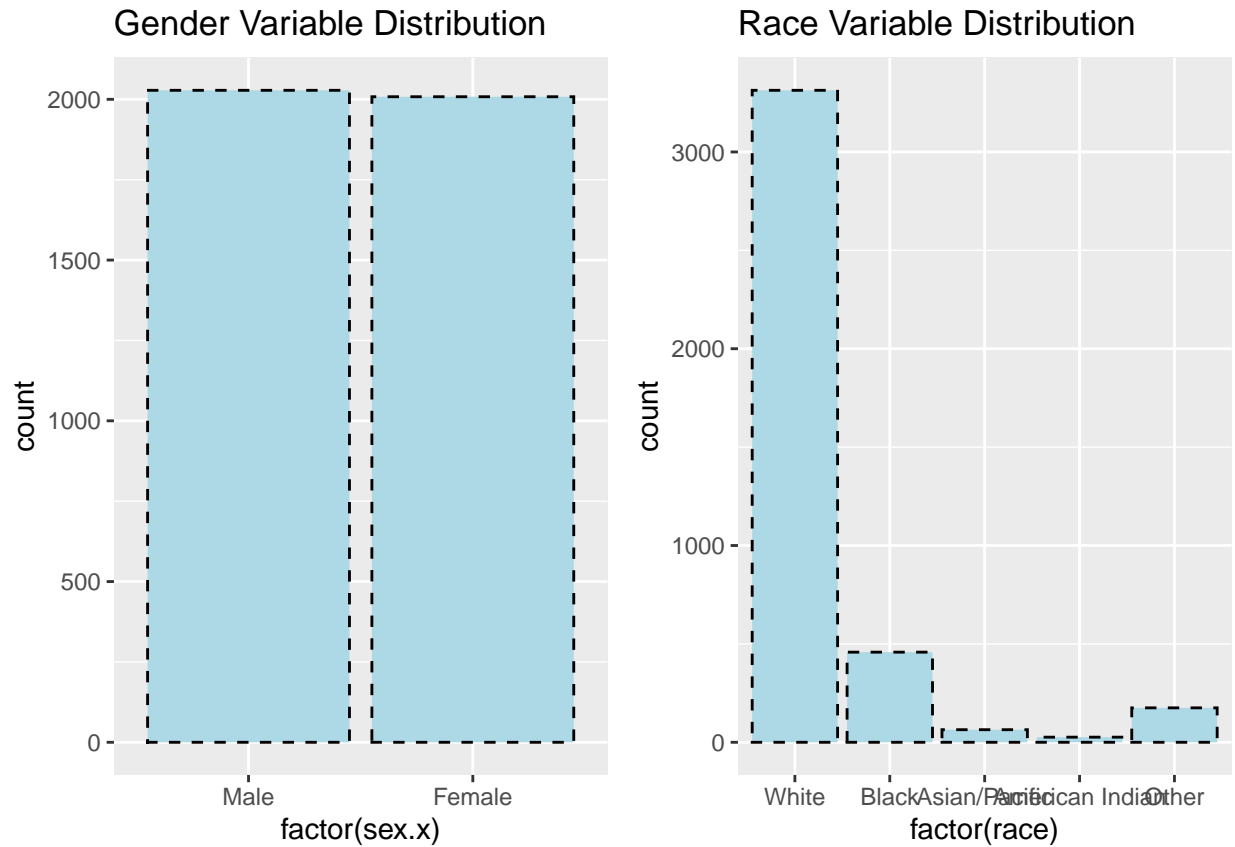
Next , lets observe the income variable. (red line - Median , Green line - mean) We can see that here there is a long right tail and what seems to be abnormal value around 100000. Due to the scaling this right tail spike can argueably make sense , therefor I have decided to leave it in the anlysis and not treat it as an outlier.

Lets review the next parameter - Age: Here the distribution surrounds the median and mean which is around 50.



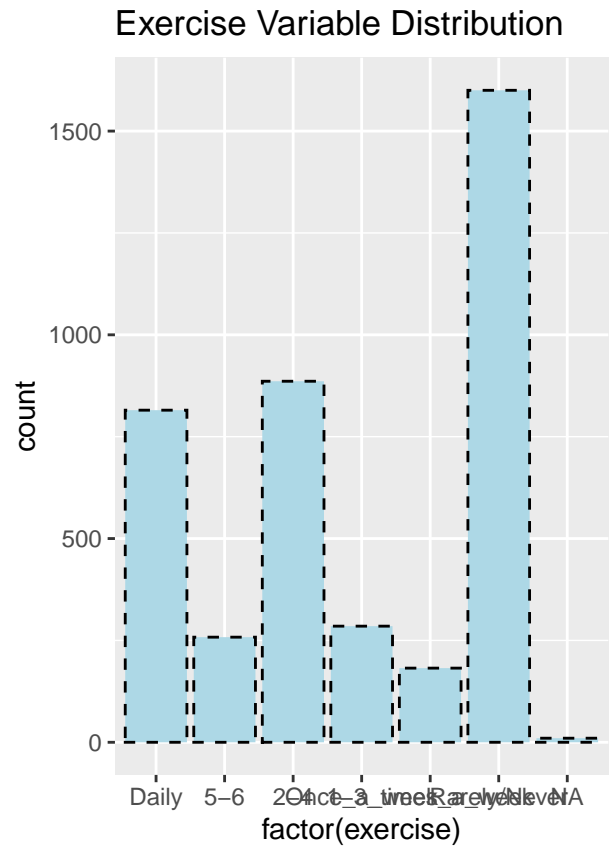
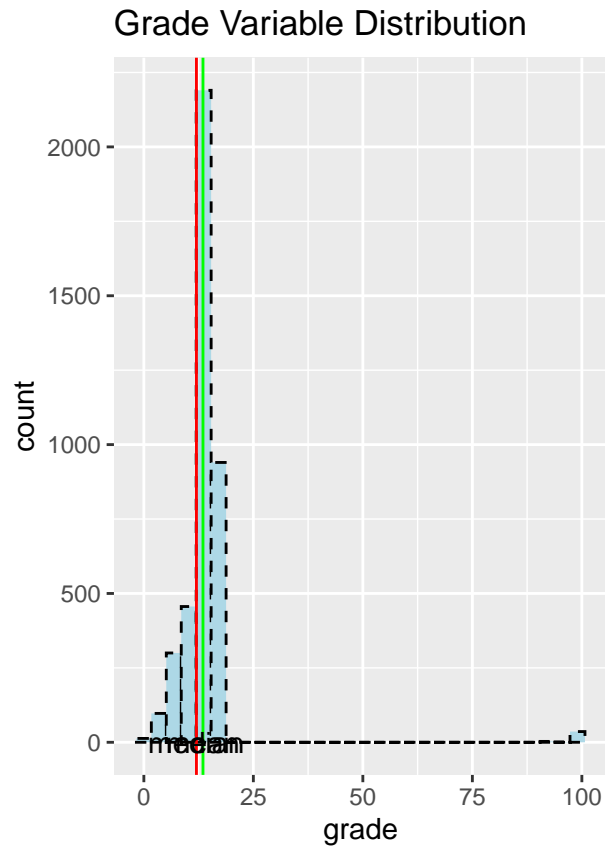
Lets observe the Gender distribution, where we can see an even distribution.

Lets look at the race feature: Here we can see an unbalanced set where there is an oversample of patients categorized under White. We will explore later on whether this feature can assist us in describing our target (BMI).



Lets view the Grade feature: Here we are able to catch the outliers where their grade is above 20. We will treat these outliers in the following parts. In addition we can see that as expected , the median and mean are near 12. which intuitively makes sense since this represents people who finished high school.

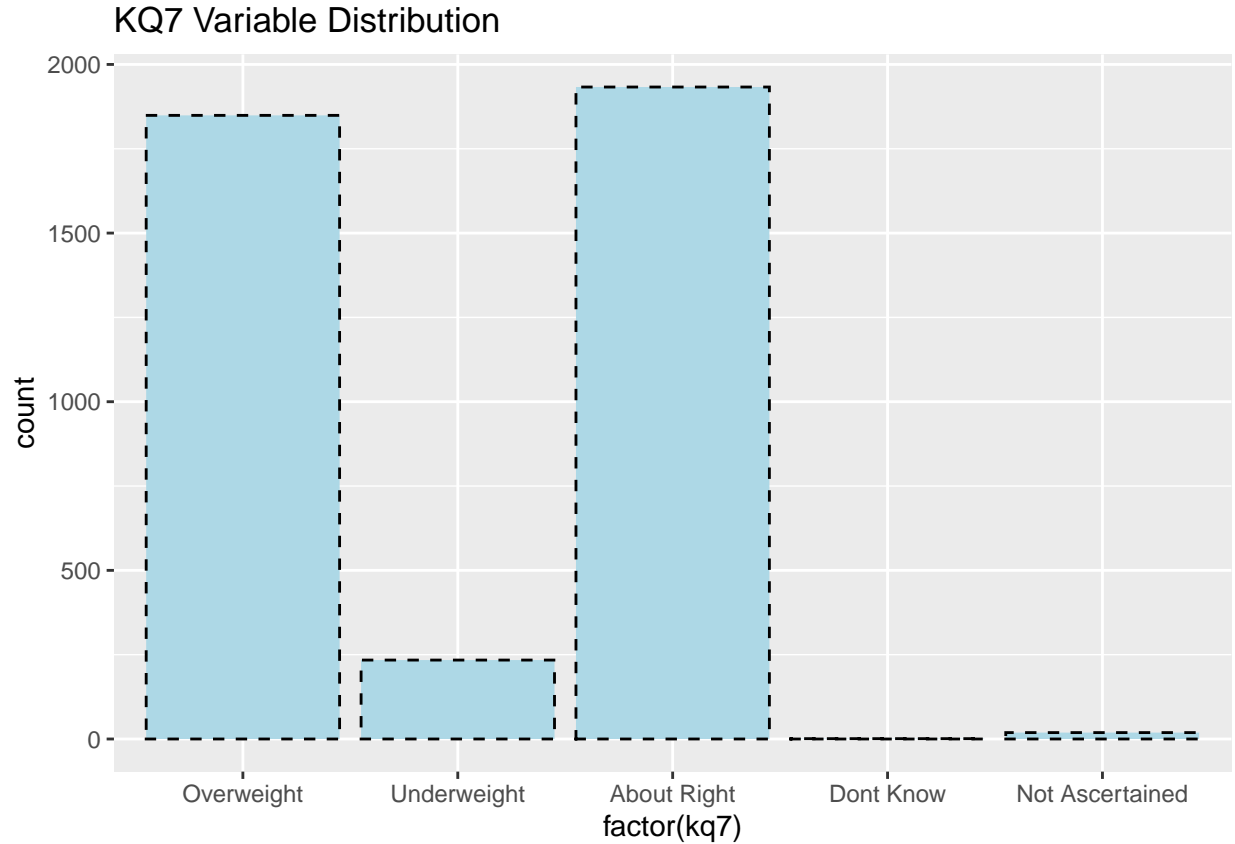
Lets explore the “Exercise” variable: here too we can see some NaN obsevation that will be treated in the following parts of the report. In terms of equally distributed data , the data here is not balanced.



Lets view the kq7 feature: Here we can see some null or insignificant values which will be treated later on.

Table 1: Frequency Matrix - Doctor Questions

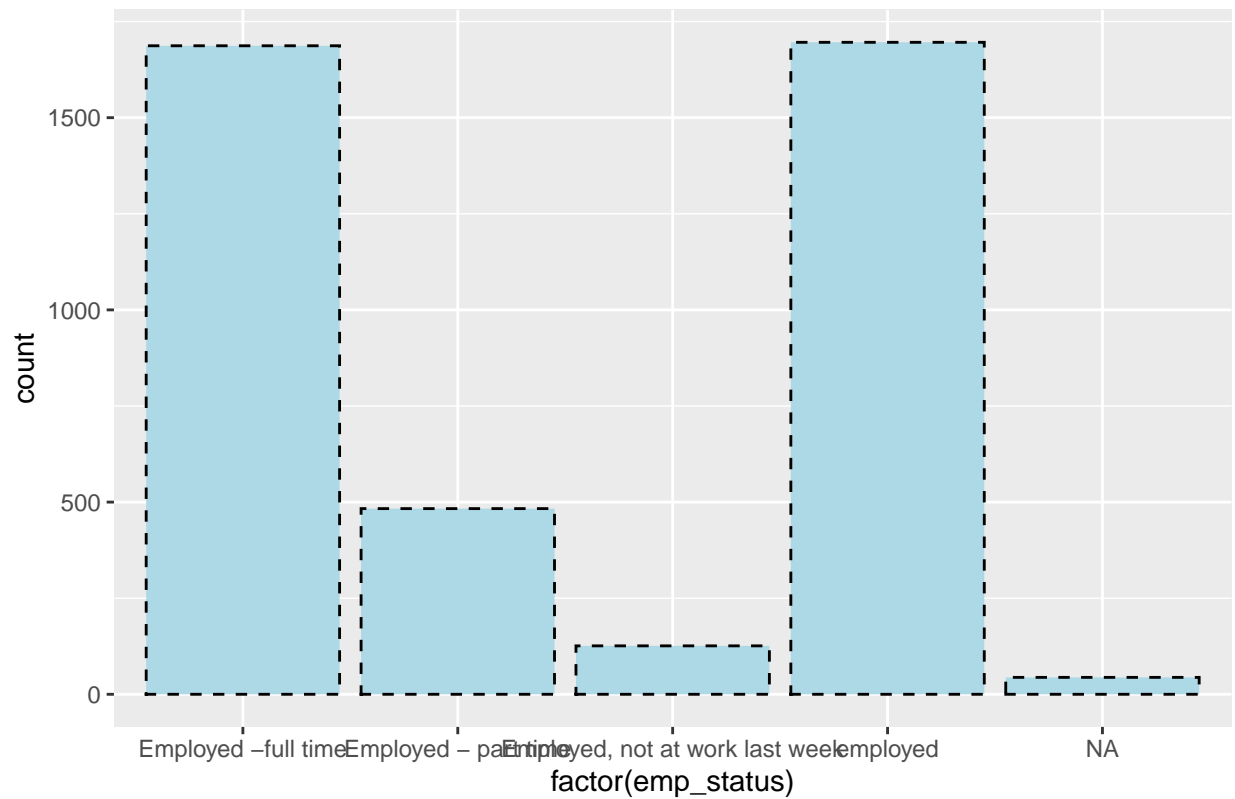
	dt01	dt02	dt03	dt06	dt07
1:	245	361	197	70	138
2:	3791	3675	3839	3966	3898



Lets view the binary questions known as dt01,02,03,06,07 In the following diagram we have the frequency matrix per question x category answer. These too are not balanced and mostly are answered with a one-sided answer, we will try to take this under consideration in the following steps.

Lets look at the Employee Status feature. There are some values to be treated , and it seems that the majority comes from two classes of the four.

### Employee Status Variable Distribution

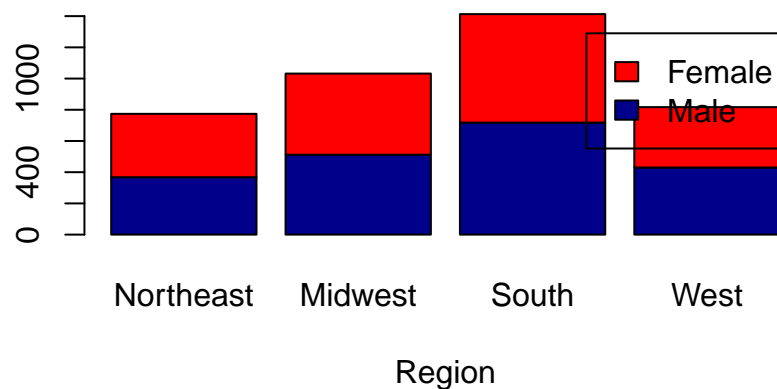


Lets explore some additional characterisites and relationships that may trigger some additional analysis.

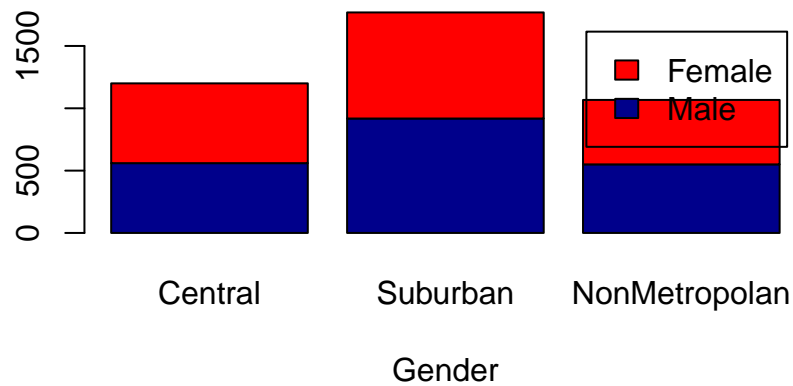
In this following plot we describe the distribution of gender & Region It seems that the porportions of the gender are equal across all regions .

Next lets view the split between the gender and the Urb feature. Here too there seems to be an equal split.

### Distribution by Gender and Region



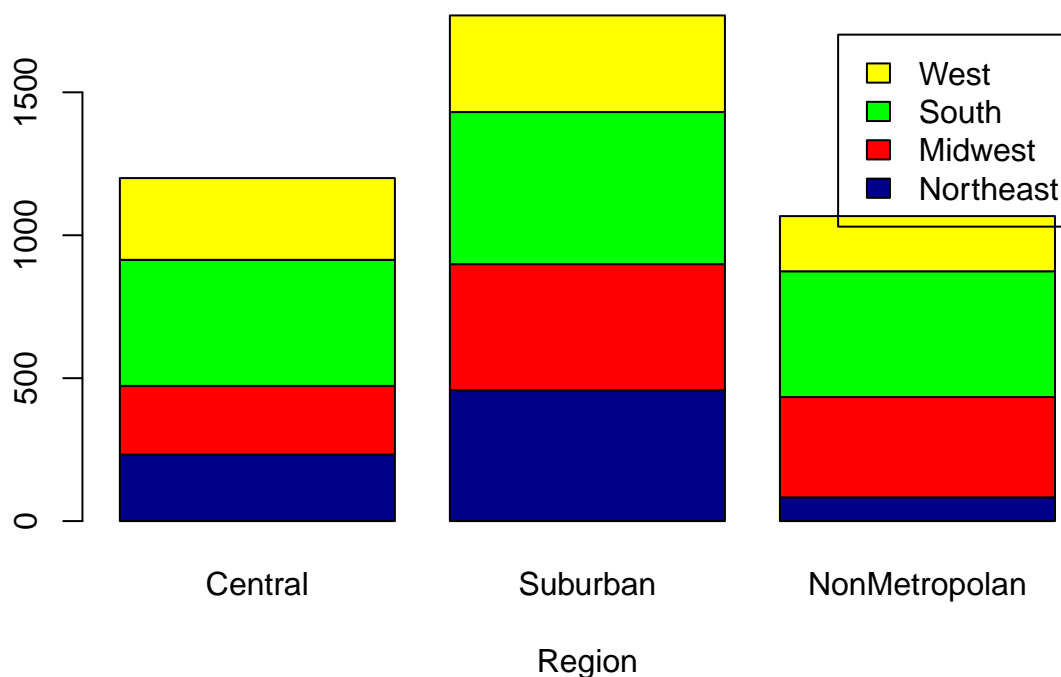
### Gender Distribution by Urb



Now lets observe the split between Urb and Region. This is to understand the different distributions between inner classes. In this example we can see that the Northeast under the Non Metropolitan has a lower frequency , than the other two Urb sections, and is porportionally less.

In this plot following the former one we observe the relationship and linear trend between the Age and the Income. The scatter is pretty well spread , though there is a slight negative trend/correlation overall between the two. This might be due to the fact that the elderly dont work and possibly have a lower income.(Thoughts to consider)

### Region Distribution by Urb







## Part 3 - Missing and Problematic Data

Given the insights from the prior section , we decide now what data to be changed to NA and follow by removing for regressional tests etc. In this part we deal with NaNs for exercise.

Deal with NaNs and problematic answers that are uninformative for the kq7 variable.

We deal with grades that are above 18 which will be considered outliers.

We deal with Null values for the 5 questions from the doctors.

We deal with Null and missing values for the employee status.

In total we remove about 100 observations from the initial dataset, given that this is a very small amount , we can afford to give up these observations and continue the analysis without them.

We can see that the dimension after the data cleaning is slightly reduced to 3700 x 57

## Part 4 - The Response Variable - Need For Transformation?

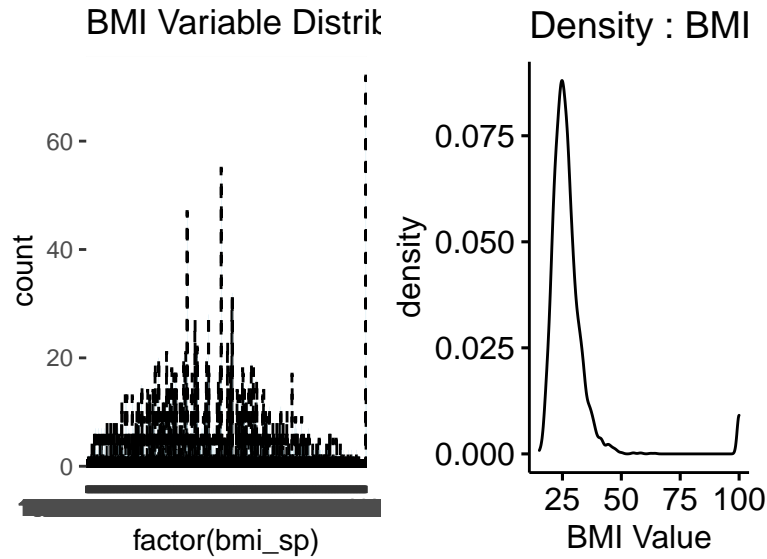
Lets look at the target variable. In this step , we will observe the BMI Measurement and evaluate the metrics , understand whether we will need to transform the value for further perdictive analysis and deal with outlier data.

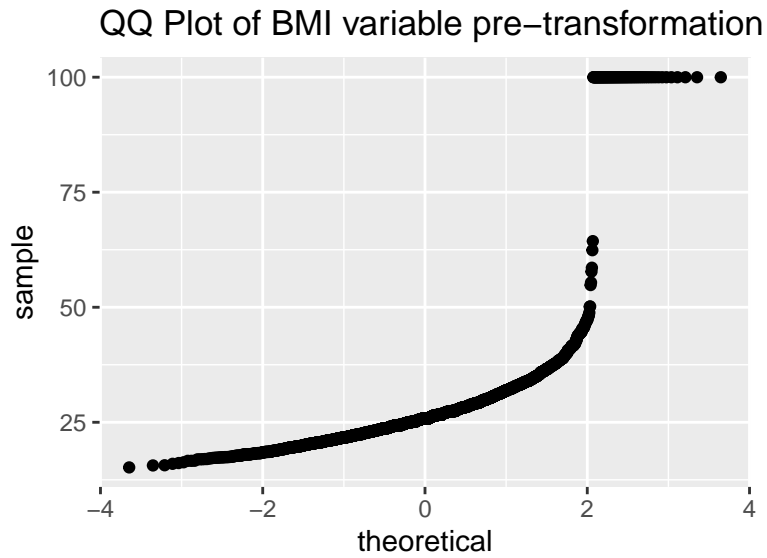
First off we know that values of BMI *tend* to range between 0-30 , where over 30 is considered obese. We can see here that there are certain outliers , where an entire quartile is above 30 ( and the top 2 percentile are above 75). We will consider the top two percentile as outliers and deal with them in the following sections.

Table 2: Quantile Matrix - BMI results

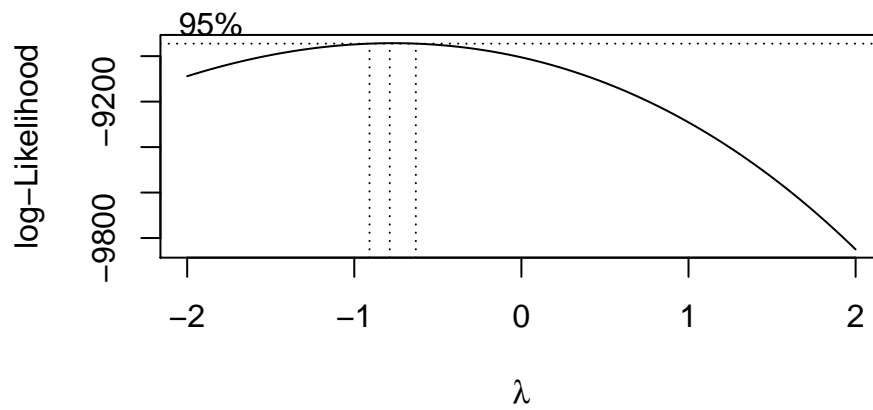
	x
0%	15.1900
25%	22.9600
50%	25.8200
75%	29.3700
90%	34.0140
98%	57.4656
100%	99.9900

In addition ,from the following figure, we are dealing with data that has a gaussian-like curve , though not specifically a normal distribution. We would like to verify that it is actually normal as a pre condition for future analysis such as log likelihood etc. For this we performed the following QQ plot where we can see the huge right tail of our outliers affecting this analysis. For the sake of the checkup , I removed these outliers and re-checked the QQ plot finding a somewhat smoother plot , though not a normal distribution.

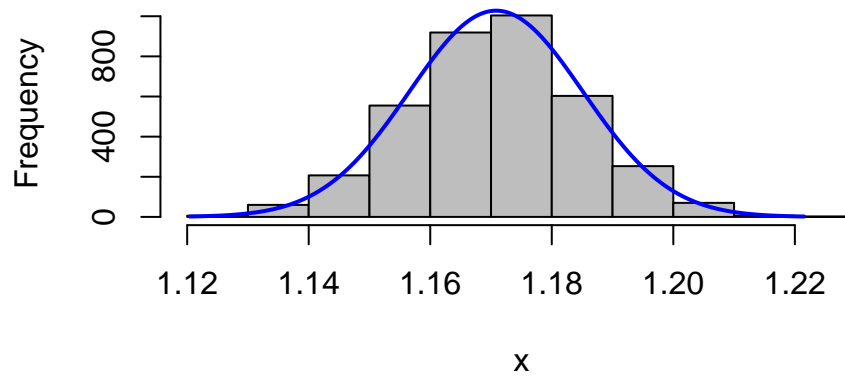




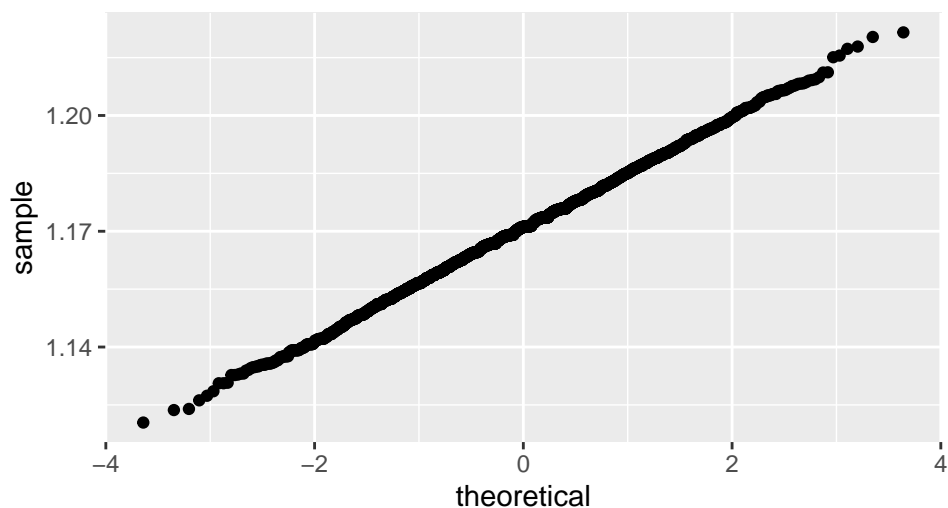
Given the results we saw regarding the BMI distribution , we will perform a Box-Cox transformation , and re-evaluate the distribution. In the following plot we perform the box-cox , and plot the updated histogram. We can see now that the distribution plot look a lot better and are able to take under normality assumptions.



## Normality Histogram – Transformed BMI Variable



## QQ Plot of BMI variable post-transformation



## Part 5 - Relation of BMI with other values

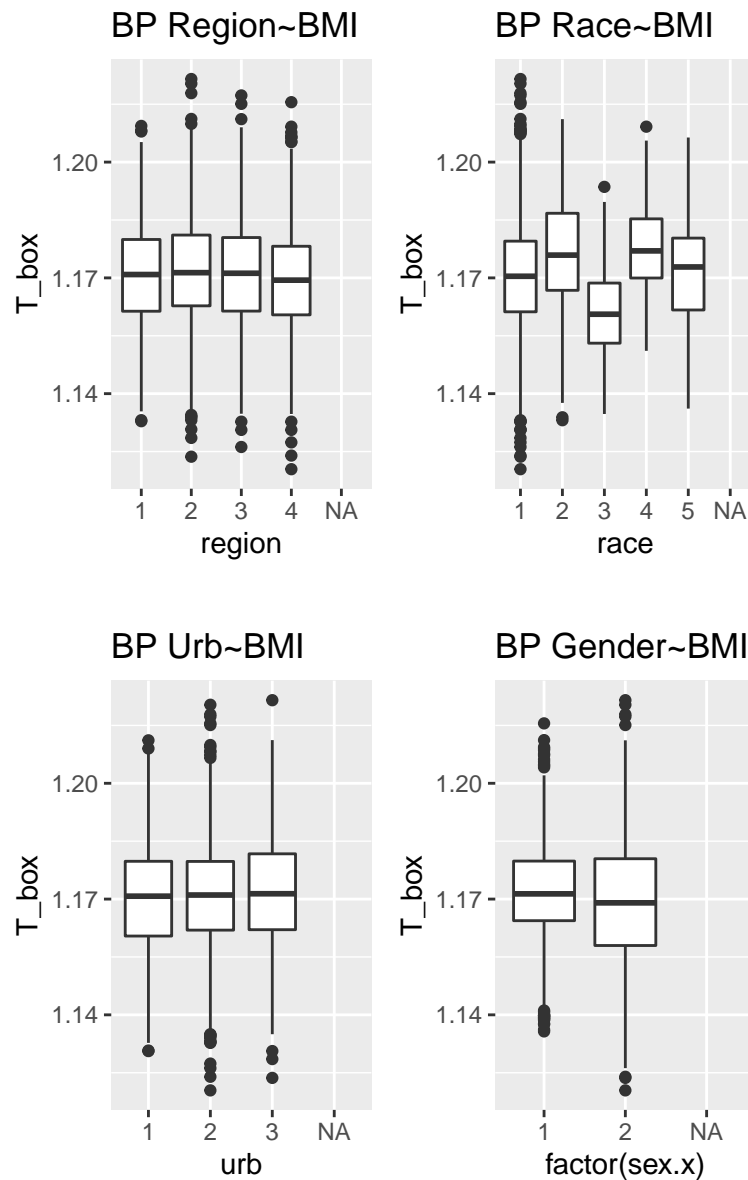
Lets review the relationship of the BMI with other features , identify some correlations and potential descriptive analysis.

Looking at the BMI grouped by regions , we can see that Northeast is with the smallest range , and least amount of outliers. We can also see that the medians are similar to all groups.

While plotting the box plots of the BMI , vs the Race feature ,we see different ranges through the group , having “white” as expected with the largest range . This is expected since it holds the majority of the observations. What is interesting is that most (if not all) the outliers that are potential from this set , are under the “white” category.

Looking at the relationship between BMI and Urb , we find equal medians with Class 2 (Suburban) containing some additional long tails , though we are still under standard acceptable ranges , so this is considered fine.

Next we can see the box plots per gender. The body of Group 2 is larger under almost equally sized groups, which is interesting to understand full relation and correlation between the two.

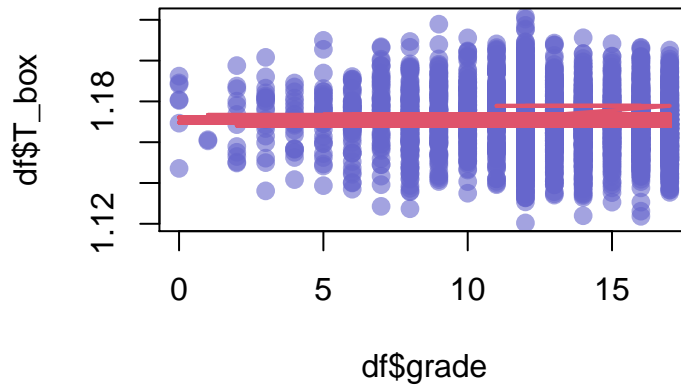


Next we will look at the negative correlation between the BMI transformed and the grade. The ranges are much larger for grades 12 +, this makes sense, since the age range is probably larger and more diverse, vs lower grades which probably contains a more unified set of patients. What is interesting here is using a polynomial regression line we can see a small peak around 11-12, and minimums around 3-6 and around 15.

In the following figure we demonstrate a similar analysis to the previous though with age vs BMI.

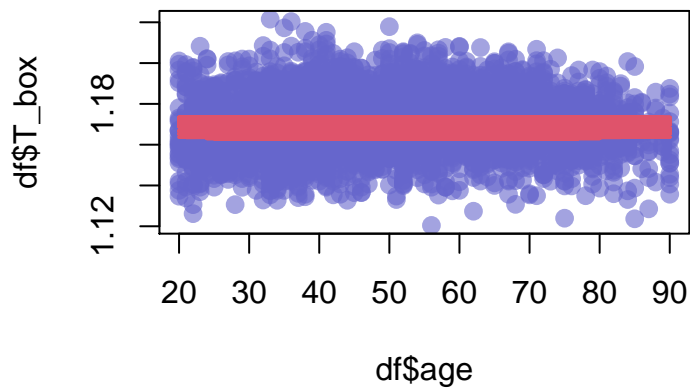
Here we can see that the highest BMI range around the ages of 50-70. Leveraging this, we can consider using a Transformation of this variable for the prediction of the BMI.

### Scatter plot Grade ~ BMI Transformed



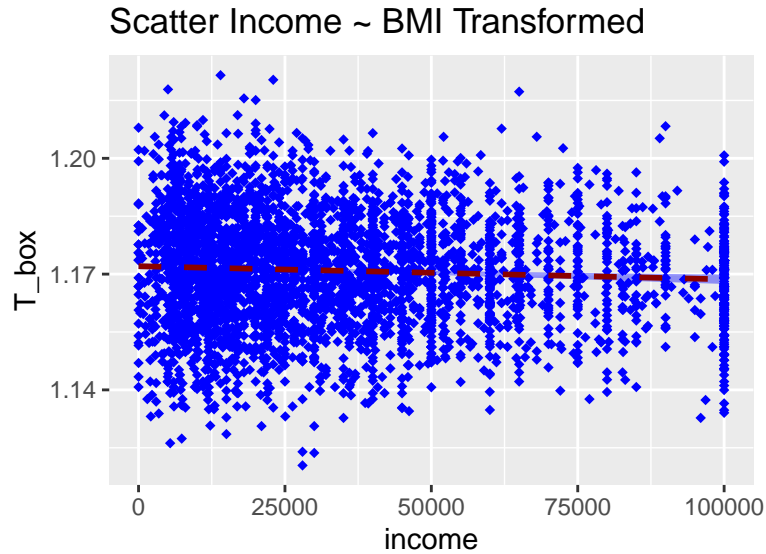
```
## integer(0)
```

### Scatter Plot BMI\_T ~ Age + Age^2

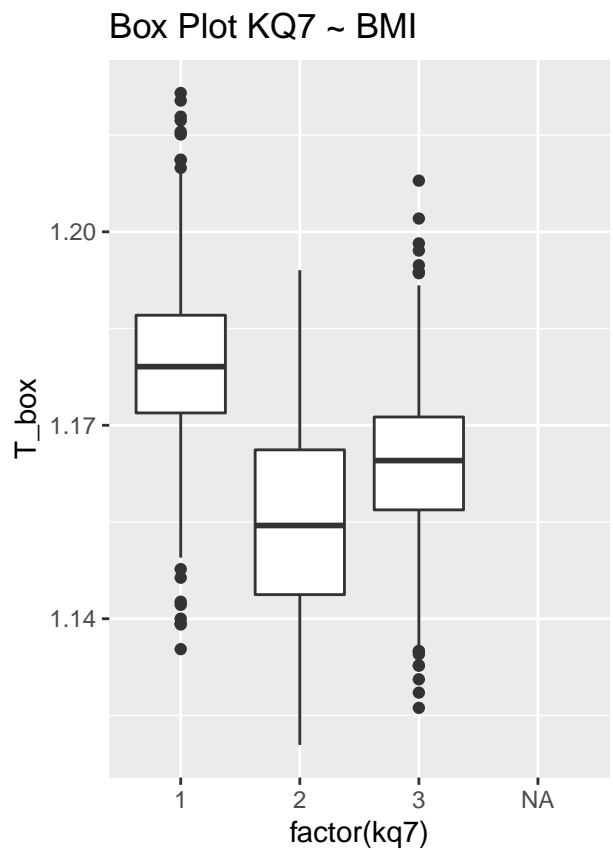
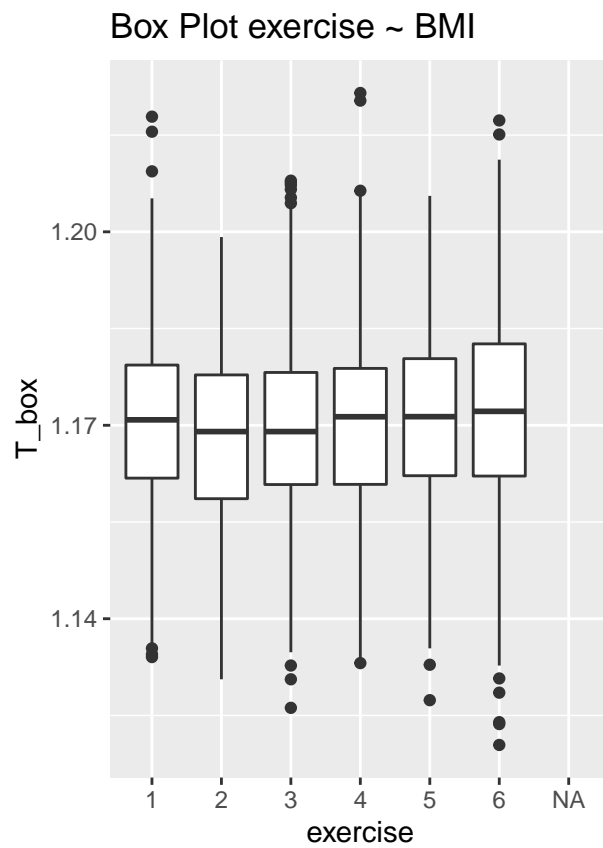


```
## integer(0)
```

Here we can see a scatter plot of the relationship between the income and the income. The trend is pretty linear and subtle. We will suspect that there is not significant relationship between the two variables and will verify this in the following steps.



Here there is a split according to the categorical sets of the Exercise variable. The medians are pretty equal, per all sets. Though by splitting by the kq07, into its 3 categories, we can see a variance within the groups. The right tail observations belong to group 1 which is interesting since it characterizes this group with more erratic answers.



## Part 6 - Univariate Regressions

In this phase we will choose the set of variables ('region','urb','income','age','sex.x','race','grade','exercise','kq7','dt01','dt02','dt03','dt06','dt07','emp\_status') in addition to the scaled variables of the surveys which were too cleaned prior to the analysis. and run univariate regressions of  $BMI \sim X$ , where X will be a univariate variable from the former list. Given the results, we can evaluate which of the following features are not worth consideration in the building of the model for the BMI. I decided to be relaxed with this constraint and set the threshold to 0.05, and let the feature 'urb' enter the following steps. Please note that the results are in the appendix for these tests.

In the following table we can see the results of the univariate analysis. I decided to leave the following features due to their significant p value which indicates their potential contribution to the analysis: 'region','urb','income','age','sex.x','race','grade','exercise','kq7','dt01','dt02','dt03','dt06','dt07','emp\_status'



**Characteristic**	**N**	**Beta**	**95% CI**	**p-value**
region	3,685			0.008
1				
2		0.00	0.00, 0.00	
3		0.00	0.00, 0.00	
4		0.00	0.00, 0.00	
urb	3,685			0.14
1				
2		0.00	0.00, 0.00	
3		0.00	0.00, 0.00	
income	3,685	0.00	0.00, 0.00	<0.001
age	3,685	0.00	0.00, 0.00	0.003
sex.x	3,685			<0.001
1				
2		0.00	0.00, 0.00	
race	3,685			<0.001
1				
2		0.01	0.00, 0.01	
3		-0.01	-0.01, -0.01	
4		0.01	0.00, 0.01	
5		0.00	0.00, 0.00	
grade	3,685	0.00	0.00, 0.00	<0.001
exercise	3,685			<0.001
1				
2		0.00	0.00, 0.00	
3		0.00	0.00, 0.00	
4		0.00	0.00, 0.00	
5		0.00	0.00, 0.00	
6		0.00	0.00, 0.00	
kq7	3,685			<0.001
1				
2		-0.02	-0.02, -0.02	
3		-0.02	-0.02, -0.01	
dt01	3,685			<0.001
1				
2		-0.01	-0.01, -0.01	
dt02	3,685			0.001
1				
2		0.00	0.00, 0.00	
dt03	3,685			<0.001
1				
2		0.00	-0.01, 0.00	
dt06	3,685			0.35
1				
2		0.00	-0.01, 0.00	
dt07	3,685			<0.001
1				
2		-0.01	-0.01, 0.00	
emp_status	3,685			0.039
1				
2		0.00	0.00, 0.00	
3		0.00	0.00, 0.00	
4		0.00	0.00, 0.00	

## Part 7 - Multiple Regression

### Variable Screening

We screened the datasets features using the former runs results , and the pvalue outcomes.

### STEP AIC

Once we have the fileterd set , we can run the Step AIC model , (using both sides) to filter out more features using a more sophisticated method. AIC is found to be pretty useful in this situation, and since the target variable is converted to a normal distribution , we are able to apply a log likelihood model.

We can see that the final run gives us an adj. R squared of 0.07 , F statistic of a significance. Although there are some predictors here which come out to be not as significant in their contribution , I decided to leave them in .

Another thought to be considered is the fact that there are some features that one sub group within the categorical feature has a significant P val , and other members of the domain dont. This could be due to a redundant split of the category, which could result in us merging some of the groups back together.

An example for this could be seen under the category Race , where some sub groups have a quite high p value and low |t-stat| , and could possibly be merged into another one. (Possibly expand Other to be all but White)

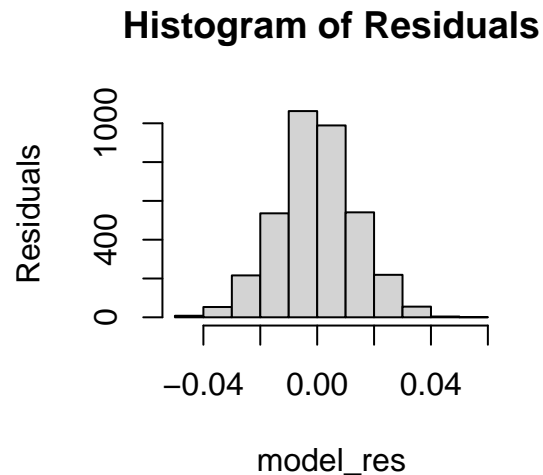
Some additional insights that are worth noting here are that females have a significant negative relation with BMI , derived from the T stat. Grade does as well , giving an indication that potentially , the lower the grade the higher the BMI , and vice verse.

Generally speaking , this model presents several significant negative correlations with some of the predictors. ( where it it noticable that the intercept is a positive number , so there is a compensation represented there against all the negative coeff predictors.)

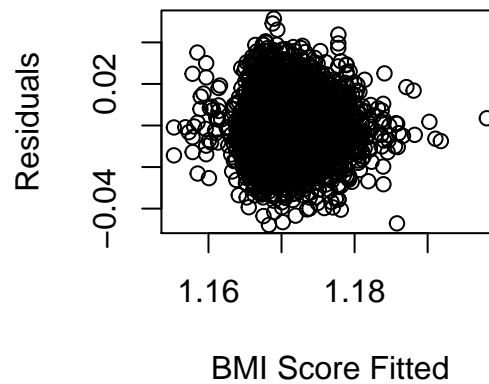
```
##
## Call:
## lm(formula = uvariate_filtered_df$T_box ~ region + sex.x + race +
##      grade + exercise + dt01 + dt02 + dt07, data = uvariate_filtered_df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.047848 -0.008727 -0.000208  0.008744  0.051399
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  1.190e+00  1.990e-03  597.969  < 2e-16 ***
## region2      1.087e-03  6.886e-04   1.578  0.114605
## region3     -3.971e-05  6.504e-04  -0.061  0.951319
## region4     -7.005e-04  7.416e-04  -0.945  0.344946
## sex.x2      -3.552e-03  4.643e-04  -7.650  2.55e-14 ***
## race2        5.630e-03  7.351e-04   7.659  2.39e-14 ***
## race3       -7.840e-03  1.945e-03  -4.030  5.69e-05 ***
## race4        6.983e-03  2.847e-03   2.453  0.014232 *
## race5        6.760e-04  1.206e-03   0.561  0.575080
## grade       -3.725e-04  7.899e-05  -4.716  2.50e-06 ***
## exercise2   -8.998e-04  1.041e-03  -0.864  0.387554
## exercise3    4.710e-04  7.101e-04   0.663  0.507231
```

```
## exercise4      1.287e-03  1.002e-03   1.285 0.198796
## exercise5      1.647e-03  1.169e-03   1.409 0.159006
## exercise6      2.400e-03  6.297e-04   3.811 0.000141 ***
## dt012          -8.838e-03  9.925e-04  -8.904 < 2e-16 ***
## dt022          -1.594e-03  8.138e-04  -1.959 0.050223 .
## dt072          -5.081e-03  1.235e-03  -4.116 3.95e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.01378 on 3667 degrees of freedom
## Multiple R-squared:  0.07671,    Adjusted R-squared:  0.07243
## F-statistic: 17.92 on 17 and 3667 DF,  p-value: < 2.2e-16
```

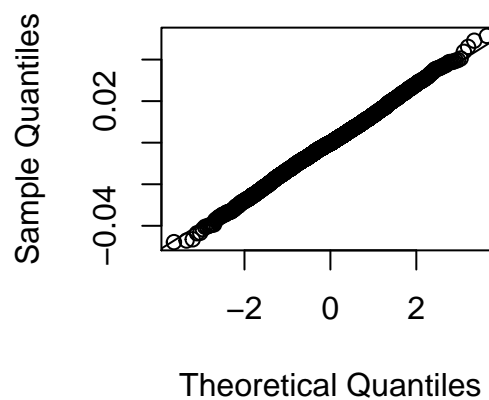
Lets review the residuals and break them down to understand their behavior. In the following plots we can see the histogram of the residuals and the scatter plot of the fitted vs the residuals. We can see that the scatter surrounds 0 which is encouraging. The QQ plot confirms the normality of the residuals. We can also see that looking at the Durbin Watson metric, we can see that the score is ~2.05 which given that it is close to 2 , indicates that there is close to no auto correlation within the residuals - meaning that the residuals are independant.



## BMI Prediction (residuals)



## Normal Q-Q Plot

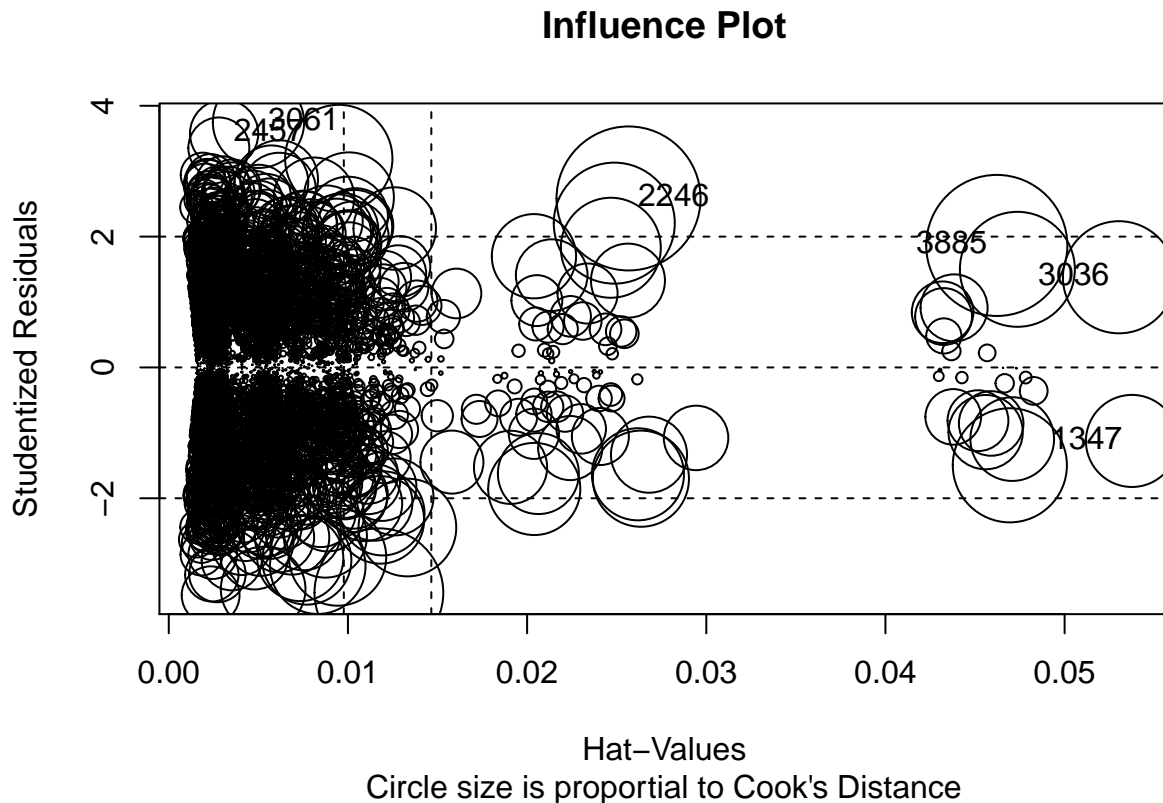


## Influence Plot

Here we will test out the influence of the different predictors , extracting the Hat values , and the Cooks distance. In the following plot we can see that there are two clusters which are differed by the Hat values (above and below 0.04).

There are also two observations with very high hat values , but significantly small cooks distance sizes. (far right of the graph).since the Cooks distance is small , this is not too alarming to our analysis.

Value 3885 has a pretty high Cooks distance and Hat value. This could be considered the most influential observant.



## Interactions Effects

Lets consider using interactions between variables. For this , we can leverage the previos architechture of the AIC model , just run it including the interactions of the different predictors. We will test this using the forwards direction. The results were quite encouraging. While adding the additional interactions we manage to filter out most of them using the AIC , and remain with multiple interactions (and single variates): please see appendix for full report.

To summarize this test , we will address some of the features that were included in the model:

Age is considered a significant predictor , having a strong negative T stat , indicating that the younger the age , the higher potentially the BMI.

Age along with Low calorie diets are negatively correlated with BMI, with a T stat  $\sim -3$ .

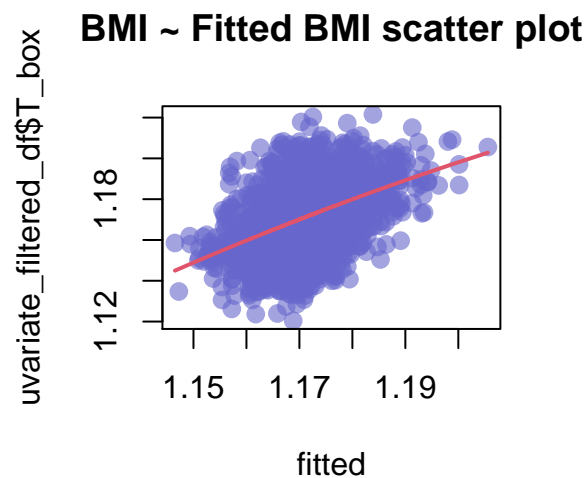
Race interacted with grade have a strong negative correlation , haveing the effect that Black people with a lower age have higher BMIs and vice versa. It is interesting that for some groups of race this is significant but for others it isnt neccesarrily.

There are several more insights that can be derived which can be elaborated in the future.

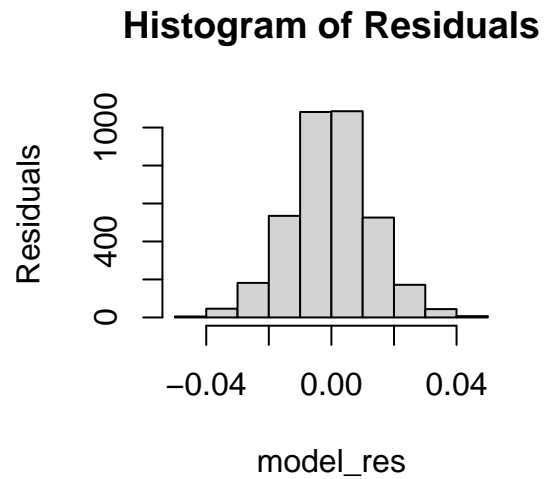
Our R squared increases to 0.11 , which means that there is a significant amount of variance explained from the interactions. Our F statistic is still significant though not as strong as expected , and P value stays significant. Using these results I decided to continue with the added interactions (post feature filtering).

For the sake of the inspection , I plotted out the fitted preditcions vs the actual ones. The intuition here was to validate whether a linear model fits the case we are trying to solve, or whether we might need to address

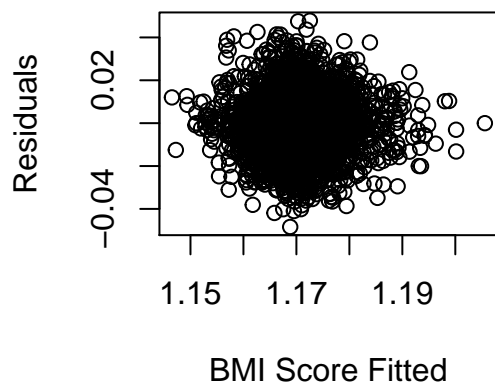
this with a polynomial curve. Given the plot , we can assume that a linear model generalizes the problem well enough.



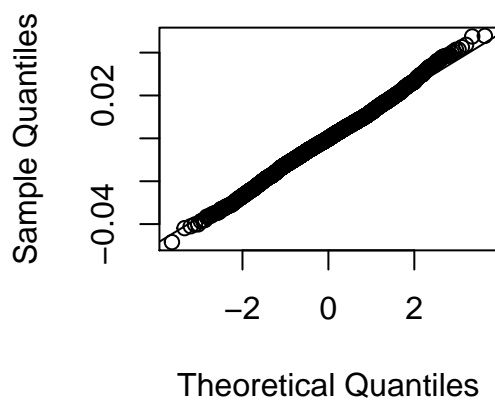
And by re-iterating the residuals evaluation again , we get:



### BMI Prediction (Residuals)



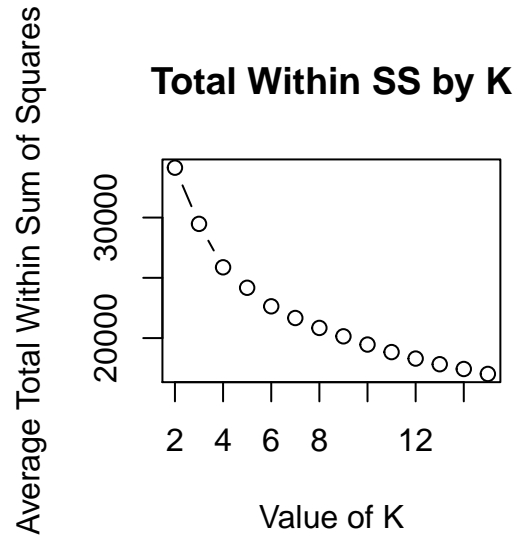
### Normal Q-Q Plot



## Additional Analysis

This part will conduct some additional analysis on the dataset to characterize the data just a bit more and attempt to aid us with fully understanding the behavior and caveats of the dataset.

In the first part I performed a clustering method of K means to try and classify the observants within predefined clusters. We use the elbow method to determine that there should be 4 clusters.



Once we have the optimal number of clusters , we can re-run the kmeans on the set, and assign to each observation its label.

The following plot , describes performance of the Kmeans , reducing the dimension to 2 while still holding the optimal amount of variance using PCA.

We can see that the clusters have an interesting fit , and the number of elements in each cluster is pretty significant with 894, 762 , 958, 1071 obs respectfully.

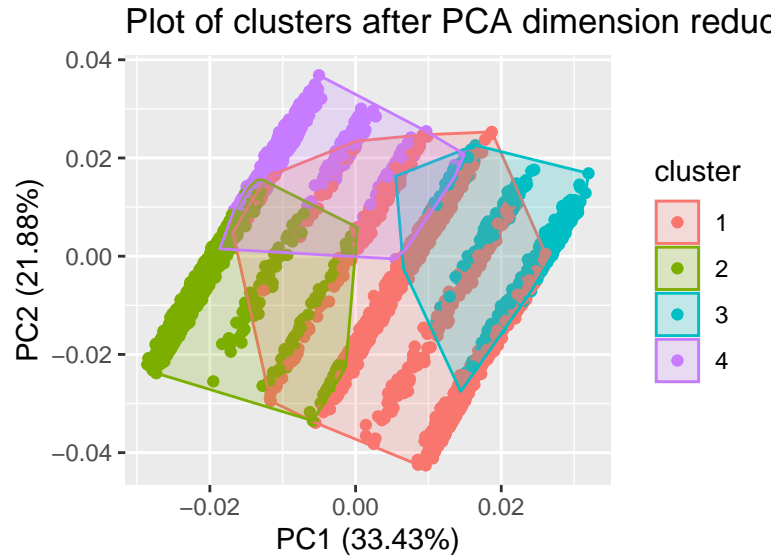
To conclude this analysis, we assign to each observation its label , and group the dataframe by the labels assigned -> and calculate the BMI mean per group.

The motivation here is to see whether the averages are different per group (reminder that BMI was not )

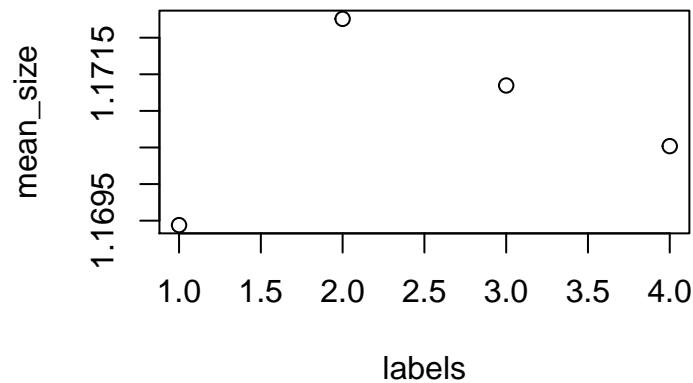
What we can see in the figure below the plot of the 4 averages of BMI , per group. We can see that visually there is a distinctive measurement per group that we can consider while characterizing the BMI set in the future.

Another recommendation here would be to categorize the BMI into 4 groups , and maybe analyze accordingly. We might be able to create stronger relationships while performing the analysis on discrete groups rather than a continuous variable.





#### 4 Cluster BMI Means



Another note here that should be taken is the addition of all the rest of the features. I have added these features in the assessment (the survey question features) and they have been found to be significant, especially using interactions between them and the demographic data. This would conclude a second potential phase to the analysis by adding these features and increasing the R squared to 0.47 (proven in the code).

## Conclusion

To conclude this analysis, we can see that the BMI parameter can be explained and described by the current set of features. There are several more options to enhance this analysis, such as additional features (new features, or transformations of current ones.) We found that there is a clear connection and relationship between the BMI and different elements of the patients such as demographic information, different surveys etc.

Another next step that may be interesting to check is testing additional models (consider decision tree based models) and perform a comparison of predictions between the models . Compare feature importance , and residuals variance.

## Appendix

Per our univariate regression , the full report is the following:

Here we can see the full report of the interaction based AIC.

```
##
## Call:
## lm(formula = uvariate_filtered_df$T_box ~ (region + income +
##       age + sex.x + race + grade + exercise + dt01 + dt02 + dt03 +
##       dt07 + emp_status)^2, data = uvariate_filtered_df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.048403 -0.008326 -0.000013  0.008187  0.047808
##
## Coefficients: (5 not defined because of singularities)
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      1.257e+00  1.685e-02  74.594 < 2e-16 ***
## region2          -1.047e-04  7.243e-03  -0.014  0.988472
## region3          -1.018e-02  6.981e-03  -1.458  0.144801
## region4          -1.293e-02  8.573e-03  -1.508  0.131618
## income           -2.054e-07  1.206e-07  -1.703  0.088590 .
## age              -1.142e-03  1.923e-04  -5.940  3.13e-09 ***
## sex.x2           -2.721e-03  5.302e-03  -0.513  0.607830
## race2            -4.303e-03  7.827e-03  -0.550  0.582506
## race3            -1.603e-02  2.334e-02  -0.687  0.492423
## race4            -8.644e-03  5.032e-02  -0.172  0.863617
## race5            -1.096e-02  1.758e-02  -0.623  0.533160
## grade            -6.773e-04  8.053e-04  -0.841  0.400371
## exercise2        -9.409e-03  9.680e-03  -0.972  0.331087
## exercise3        -7.792e-03  7.796e-03  -0.999  0.317658
## exercise4         2.660e-03  1.298e-02   0.205  0.837621
## exercise5        -8.597e-03  1.617e-02  -0.532  0.594935
## exercise6         8.549e-03  6.895e-03   1.240  0.215116
## dt012            -1.928e-02  1.018e-02  -1.893  0.058441 .
## dt022            -2.177e-02  8.638e-03  -2.520  0.011791 *
## dt032             5.004e-03  1.081e-02   0.463  0.643523
## dt072            -3.499e-02  1.154e-02  -3.031  0.002452 **
## emp_status2       -6.797e-04  9.853e-03  -0.069  0.945003
## emp_status3        6.293e-03  1.501e-02   0.419  0.675139
## emp_status4        1.844e-02  7.099e-03   2.597  0.009439 **
## region2:income     1.535e-08  3.118e-08   0.492  0.622578
## region3:income     2.425e-08  2.995e-08   0.809  0.418343
## region4:income     1.463e-08  3.263e-08   0.449  0.653797
## region2:age        -9.116e-05  5.057e-05  -1.802  0.071567 .
## region3:age        -4.620e-05  4.765e-05  -0.970  0.332317
## region4:age        -3.187e-05  5.427e-05  -0.587  0.557073
## region2:sex.x2     -2.351e-04  1.443e-03  -0.163  0.870559
```

## region3:sex.x2	7.607e-04	1.362e-03	0.558	0.576634
## region4:sex.x2	-7.631e-04	1.573e-03	-0.485	0.627555
## region2:race2	1.881e-03	2.381e-03	0.790	0.429653
## region3:race2	5.394e-03	2.094e-03	2.576	0.010050 *
## region4:race2	5.029e-03	3.228e-03	1.558	0.119407
## region2:race3	2.313e-03	8.983e-03	0.258	0.796796
## region3:race3	9.523e-03	8.590e-03	1.109	0.267635
## region4:race3	4.454e-03	7.576e-03	0.588	0.556640
## region2:race4	-2.007e-02	1.766e-02	-1.136	0.255837
## region3:race4	7.079e-04	2.982e-02	0.024	0.981066
## region4:race4	1.199e-02	2.482e-02	0.483	0.629101
## region2:race5	5.617e-03	5.191e-03	1.082	0.279305
## region3:race5	3.622e-03	4.357e-03	0.831	0.405881
## region4:race5	6.743e-03	3.699e-03	1.823	0.068417 .
## region2:grade	-2.916e-04	2.897e-04	-1.007	0.314142
## region3:grade	-2.374e-05	2.556e-04	-0.093	0.926006
## region4:grade	7.493e-05	2.982e-04	0.251	0.801637
## region2:exercise2	-1.048e-03	3.390e-03	-0.309	0.757245
## region3:exercise2	2.124e-03	2.989e-03	0.711	0.477395
## region4:exercise2	-3.095e-03	3.407e-03	-0.909	0.363662
## region2:exercise3	1.894e-03	2.211e-03	0.857	0.391753
## region3:exercise3	1.201e-03	2.110e-03	0.569	0.569327
## region4:exercise3	-1.375e-03	2.394e-03	-0.574	0.565935
## region2:exercise4	1.916e-03	3.036e-03	0.631	0.528015
## region3:exercise4	7.808e-04	2.989e-03	0.261	0.793945
## region4:exercise4	-1.037e-03	3.315e-03	-0.313	0.754401
## region2:exercise5	1.469e-03	4.022e-03	0.365	0.714923
## region3:exercise5	3.018e-03	3.879e-03	0.778	0.436659
## region4:exercise5	-4.434e-03	4.242e-03	-1.045	0.296030
## region2:exercise6	2.327e-03	1.959e-03	1.188	0.235043
## region3:exercise6	2.123e-03	1.840e-03	1.153	0.248848
## region4:exercise6	-6.294e-04	2.176e-03	-0.289	0.772448
## region2:dt012	-1.966e-03	3.293e-03	-0.597	0.550444
## region3:dt012	1.221e-03	3.114e-03	0.392	0.695113
## region4:dt012	-3.597e-03	3.715e-03	-0.968	0.333001
## region2:dt022	2.293e-03	2.585e-03	0.887	0.375067
## region3:dt022	3.159e-04	2.523e-03	0.125	0.900364
## region4:dt022	3.020e-03	2.944e-03	1.026	0.305017
## region2:dt032	-2.147e-03	3.490e-03	-0.615	0.538473
## region3:dt032	-8.443e-04	3.329e-03	-0.254	0.799803
## region4:dt032	5.710e-03	4.059e-03	1.407	0.159529
## region2:dt072	7.276e-03	3.621e-03	2.009	0.044581 *
## region3:dt072	7.952e-03	3.576e-03	2.224	0.026220 *
## region4:dt072	7.117e-03	4.787e-03	1.487	0.137220
## region2:emp_status2	5.415e-03	2.347e-03	2.307	0.021114 *
## region3:emp_status2	3.631e-03	2.284e-03	1.590	0.112001
## region4:emp_status2	2.681e-03	2.518e-03	1.064	0.287188
## region2:emp_status3	-1.621e-03	4.735e-03	-0.342	0.732185
## region3:emp_status3	-2.143e-03	4.373e-03	-0.490	0.624144
## region4:emp_status3	2.428e-03	4.676e-03	0.519	0.603691
## region2:emp_status4	4.161e-03	1.921e-03	2.167	0.030332 *
## region3:emp_status4	1.670e-03	1.780e-03	0.938	0.348298
## region4:emp_status4	1.431e-03	2.015e-03	0.710	0.477694
## income:age	1.044e-09	7.607e-10	1.373	0.169901

## income:sex.x2	-4.633e-08	2.209e-08	-2.098	0.036015	*
## income:race2	7.207e-08	4.306e-08	1.674	0.094251	.
## income:race3	-4.276e-08	8.107e-08	-0.527	0.597892	
## income:race4	5.678e-08	2.575e-07	0.221	0.825484	
## income:race5	3.203e-09	6.558e-08	0.049	0.961042	
## income:grade	-2.146e-09	3.786e-09	-0.567	0.570827	
## income:exercise2	1.267e-08	4.593e-08	0.276	0.782638	
## income:exercise3	-2.437e-08	3.137e-08	-0.777	0.437311	
## income:exercise4	7.458e-08	4.295e-08	1.736	0.082619	.
## income:exercise5	2.614e-08	5.208e-08	0.502	0.615802	
## income:exercise6	3.033e-08	2.975e-08	1.020	0.308002	
## income:dt012	-2.000e-08	4.610e-08	-0.434	0.664464	
## income:dt022	7.119e-08	4.278e-08	1.664	0.096186	.
## income:dt032	-9.139e-09	6.404e-08	-0.143	0.886523	
## income:dt072	1.514e-07	7.136e-08	2.121	0.033981	*
## income:emp_status2	-1.776e-08	3.245e-08	-0.547	0.584156	
## income:emp_status3	-1.265e-07	5.806e-08	-2.179	0.029424	*
## income:emp_status4	-4.215e-08	2.599e-08	-1.622	0.104951	
## age:sex.x2	9.571e-05	3.512e-05	2.725	0.006457	**
## age:race2	5.226e-05	5.162e-05	1.012	0.311436	
## age:race3	1.001e-04	1.700e-04	0.589	0.555998	
## age:race4	-6.322e-05	4.741e-04	-0.133	0.893919	
## age:race5	1.440e-05	9.668e-05	0.149	0.881649	
## age:grade	2.175e-05	6.030e-06	3.607	0.000314	***
## age:exercise2	1.265e-04	7.896e-05	1.602	0.109196	
## age:exercise3	4.657e-05	5.197e-05	0.896	0.370196	
## age:exercise4	7.245e-05	7.535e-05	0.961	0.336376	
## age:exercise5	1.327e-05	8.750e-05	0.152	0.879438	
## age:exercise6	-9.566e-05	4.561e-05	-2.098	0.036016	*
## age:dt012	2.725e-04	7.821e-05	3.484	0.000499	***
## age:dt022	2.556e-04	8.228e-05	3.107	0.001908	**
## age:dt032	7.856e-05	1.045e-04	0.752	0.452211	
## age:dt072	3.429e-04	1.257e-04	2.727	0.006415	**
## age:emp_status2	7.073e-05	5.364e-05	1.319	0.187343	
## age:emp_status3	-3.421e-05	1.207e-04	-0.283	0.776838	
## age:emp_status4	-1.444e-04	4.208e-05	-3.432	0.000605	***
## sex.x2:race2	6.883e-03	1.570e-03	4.384	1.20e-05	***
## sex.x2:race3	2.477e-03	4.835e-03	0.512	0.608492	
## sex.x2:race4	-9.988e-04	1.264e-02	-0.079	0.937014	
## sex.x2:race5	2.336e-03	2.721e-03	0.858	0.390704	
## sex.x2:grade	-1.663e-04	1.876e-04	-0.887	0.375271	
## sex.x2:exercise2	1.573e-03	2.277e-03	0.691	0.489637	
## sex.x2:exercise3	1.437e-03	1.521e-03	0.945	0.344780	
## sex.x2:exercise4	5.623e-03	2.175e-03	2.585	0.009770	**
## sex.x2:exercise5	2.014e-03	2.513e-03	0.801	0.422935	
## sex.x2:exercise6	3.208e-03	1.343e-03	2.389	0.016955	*
## sex.x2:dt012	2.593e-04	2.368e-03	0.110	0.912795	
## sex.x2:dt022	4.520e-03	1.982e-03	2.280	0.022658	*
## sex.x2:dt032	-2.357e-03	2.611e-03	-0.903	0.366852	
## sex.x2:dt072	-7.398e-03	2.784e-03	-2.658	0.007907	**
## sex.x2:emp_status2	4.153e-04	1.605e-03	0.259	0.795862	
## sex.x2:emp_status3	9.111e-04	2.789e-03	0.327	0.743965	
## sex.x2:emp_status4	-6.595e-04	1.297e-03	-0.508	0.611162	
## race2:grade	-1.192e-04	2.927e-04	-0.407	0.683893	

## race3:grade	-5.468e-04	9.446e-04	-0.579	0.562702
## race4:grade	8.160e-04	2.079e-03	0.393	0.694655
## race5:grade	8.878e-04	3.845e-04	2.309	0.020995 *
## race2:exercise2	-2.670e-03	3.752e-03	-0.712	0.476758
## race3:exercise2	-5.517e-03	1.075e-02	-0.513	0.607720
## race4:exercise2	-2.791e-02	2.787e-02	-1.001	0.316736
## race5:exercise2	2.721e-03	6.171e-03	0.441	0.659243
## race2:exercise3	-7.954e-04	2.466e-03	-0.323	0.747078
## race3:exercise3	-5.151e-03	8.024e-03	-0.642	0.520915
## race4:exercise3	-1.135e-02	2.111e-02	-0.538	0.590715
## race5:exercise3	4.102e-03	4.057e-03	1.011	0.311956
## race2:exercise4	-3.564e-03	4.070e-03	-0.876	0.381335
## race3:exercise4	-1.512e-03	9.252e-03	-0.163	0.870220
## race4:exercise4	-1.661e-02	1.658e-02	-1.001	0.316691
## race5:exercise4	-4.941e-03	5.528e-03	-0.894	0.371525
## race2:exercise5	-2.113e-03	4.051e-03	-0.522	0.602000
## race3:exercise5	4.438e-03	1.046e-02	0.424	0.671283
## race4:exercise5	-1.036e-02	3.390e-02	-0.305	0.760045
## race5:exercise5	-6.566e-03	5.831e-03	-1.126	0.260241
## race2:exercise6	-2.498e-03	2.075e-03	-1.204	0.228656
## race3:exercise6	-6.459e-03	7.338e-03	-0.880	0.378804
## race4:exercise6	-7.455e-03	1.600e-02	-0.466	0.641209
## race5:exercise6	1.155e-03	3.487e-03	0.331	0.740507
## race2:dt012	1.727e-03	3.539e-03	0.488	0.625565
## race3:dt012	1.462e-02	1.093e-02	1.338	0.181073
## race4:dt012	-1.399e-02	1.409e-02	-0.993	0.320789
## race5:dt012	2.664e-06	6.596e-03	0.000	0.999678
## race2:dt022	-2.588e-03	4.145e-03	-0.624	0.532421
## race3:dt022	8.187e-03	1.026e-02	0.798	0.424741
## race4:dt022	-1.234e-02	3.002e-02	-0.411	0.681051
## race5:dt022	3.876e-03	6.525e-03	0.594	0.552549
## race2:dt032	-1.773e-03	4.002e-03	-0.443	0.657854
## race3:dt032	-5.235e-03	1.230e-02	-0.426	0.670298
## race4:dt032	3.326e-02	2.482e-02	1.340	0.180227
## race5:dt032	-1.560e-02	1.363e-02	-1.144	0.252690
## race2:dt072	3.358e-03	3.912e-03	0.858	0.390841
## race3:dt072	NA	NA	NA	NA
## race4:dt072	NA	NA	NA	NA
## race5:dt072	6.161e-03	8.006e-03	0.770	0.441617
## race2:emp_status2	1.412e-03	2.718e-03	0.520	0.603389
## race3:emp_status2	-2.553e-03	6.216e-03	-0.411	0.681286
## race4:emp_status2	-1.741e-02	1.699e-02	-1.024	0.305749
## race5:emp_status2	2.023e-03	4.823e-03	0.419	0.674957
## race2:emp_status3	-2.234e-03	4.224e-03	-0.529	0.596915
## race3:emp_status3	8.746e-03	1.820e-02	0.481	0.630884
## race4:emp_status3	NA	NA	NA	NA
## race5:emp_status3	-1.880e-04	5.393e-03	-0.035	0.972201
## race2:emp_status4	1.210e-04	2.017e-03	0.060	0.952187
## race3:emp_status4	-1.012e-02	5.549e-03	-1.824	0.068206 .
## race4:emp_status4	9.359e-04	1.352e-02	0.069	0.944823
## race5:emp_status4	1.963e-04	3.194e-03	0.061	0.950979
## grade:exercise2	-3.521e-05	4.266e-04	-0.083	0.934221
## grade:exercise3	-1.028e-04	2.896e-04	-0.355	0.722712
## grade:exercise4	-1.026e-03	4.123e-04	-2.487	0.012912 *

```

## grade:exercise5      -3.710e-04  4.825e-04  -0.769  0.442087
## grade:exercise6      -2.676e-04  2.421e-04  -1.105  0.269037
## grade:dt012          -1.218e-05  4.419e-04  -0.028  0.978018
## grade:dt022          -3.952e-05  3.925e-04  -0.101  0.919802
## grade:dt032          -4.217e-04  4.768e-04  -0.884  0.376592
## grade:dt072          9.758e-06  4.827e-04   0.020  0.983872
## grade:emp_status2     3.503e-04  3.315e-04   1.057  0.290655
## grade:emp_status3     6.019e-04  5.318e-04   1.132  0.257743
## grade:emp_status4    -3.089e-04  2.477e-04  -1.247  0.212488
## exercise2:dt012       8.823e-04  4.357e-03   0.203  0.839521
## exercise3:dt012      -1.096e-03  3.196e-03  -0.343  0.731705
## exercise4:dt012       5.153e-03  4.838e-03   1.065  0.286854
## exercise5:dt012       4.671e-03  5.830e-03   0.801  0.423039
## exercise6:dt012       1.400e-04  3.118e-03   0.045  0.964187
## exercise2:dt022       3.932e-03  4.048e-03   0.971  0.331418
## exercise3:dt022      -6.073e-04  2.629e-03  -0.231  0.817337
## exercise4:dt022       8.102e-03  5.729e-03   1.414  0.157378
## exercise5:dt022      -1.102e-02  5.401e-03  -2.040  0.041383 *
## exercise6:dt022      -1.677e-03  2.635e-03  -0.637  0.524398
## exercise2:dt032      -3.256e-04  5.614e-03  -0.058  0.953753
## exercise3:dt032       8.500e-03  3.729e-03   2.279  0.022709 *
## exercise4:dt032      -2.705e-03  6.159e-03  -0.439  0.660548
## exercise5:dt032       9.293e-03  8.372e-03   1.110  0.267099
## exercise6:dt032      -6.096e-05  3.319e-03  -0.018  0.985347
## exercise2:dt072       NA          NA          NA          NA
## exercise3:dt072       8.933e-04  3.960e-03   0.226  0.821538
## exercise4:dt072      -4.046e-03  9.006e-03  -0.449  0.653306
## exercise5:dt072       1.009e-02  1.261e-02   0.800  0.423596
## exercise6:dt072       1.725e-04  3.571e-03   0.048  0.961477
## exercise2:emp_status2 -4.489e-03  3.283e-03  -1.367  0.171571
## exercise3:emp_status2 -4.163e-03  2.372e-03  -1.755  0.079398 .
## exercise4:emp_status2 -5.717e-03  3.376e-03  -1.693  0.090519 .
## exercise5:emp_status2 -2.695e-03  3.570e-03  -0.755  0.450234
## exercise6:emp_status2 -1.449e-03  2.208e-03  -0.656  0.511781
## exercise2:emp_status3 -3.069e-03  6.343e-03  -0.484  0.628568
## exercise3:emp_status3 -7.382e-05  3.993e-03  -0.018  0.985252
## exercise4:emp_status3 -8.741e-03  5.775e-03  -1.514  0.130187
## exercise5:emp_status3 -1.989e-04  9.063e-03  -0.022  0.982494
## exercise6:emp_status3 -7.375e-04  3.843e-03  -0.192  0.847830
## exercise2:emp_status4 -3.137e-03  2.971e-03  -1.056  0.291149
## exercise3:emp_status4  1.105e-03  1.917e-03   0.576  0.564327
## exercise4:emp_status4 -4.667e-03  2.797e-03  -1.669  0.095265 .
## exercise5:emp_status4  9.016e-04  3.361e-03   0.268  0.788521
## exercise6:emp_status4  1.190e-03  1.687e-03   0.706  0.480495
## dt012:dt022           1.506e-04  2.758e-03   0.055  0.956451
## dt012:dt032          -5.199e-05  3.446e-03  -0.015  0.987962
## dt012:dt072           1.388e-04  5.656e-03   0.025  0.980420
## dt012:emp_status2     5.580e-04  3.383e-03   0.165  0.868986
## dt012:emp_status3     -1.586e-03  6.811e-03  -0.233  0.815834
## dt012:emp_status4     -5.411e-03  2.663e-03  -2.032  0.042234 *
## dt022:dt032          -4.345e-03  2.596e-03  -1.673  0.094351 .
## dt022:dt072           6.643e-03  4.532e-03   1.466  0.142785
## dt022:emp_status2     -6.423e-03  3.517e-03  -1.826  0.067895 .
## dt022:emp_status3     6.958e-03  5.530e-03   1.258  0.208408

```

```

## dt022:emp_status4      -3.991e-03  2.560e-03  -1.559  0.119073
## dt032:dt072            -4.403e-03  4.604e-03  -0.956  0.338984
## dt032:emp_status2      4.287e-04  4.693e-03   0.091  0.927223
## dt032:emp_status3     -1.171e-02  9.388e-03  -1.248  0.212231
## dt032:emp_status4      1.308e-03  3.506e-03   0.373  0.709066
## dt072:emp_status2     -3.558e-03  6.244e-03  -0.570  0.568773
## dt072:emp_status3              NA          NA          NA          NA
## dt072:emp_status4     -3.389e-04  4.547e-03  -0.075  0.940593
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.01347 on 3435 degrees of freedom
## Multiple R-squared:  0.1737, Adjusted R-squared:  0.1138
## F-statistic: 2.901 on 249 and 3435 DF,  p-value: < 2.2e-16

```