

Data Analysis Project 2: Multiple Regression Computing Project

Introduction

In this project, we are first given a dataset containing one dependent variable and twenty four independent variables. We regarded variables $E1$ to $E4$ as environmental (or control) variables. We also regarded the variables $G1$ to $G20$ as genetic risk variables. The G variables are indicator variables and the data that we are analyzing is synthetic. Our task is to find the model that the TA used for my data. The background for this project is that there was a paper written by Caspi et al. that reports a finding of a gene-environment interaction. The paper used multiple regression techniques as the methodology for its findings. The author was determined to figure out why stressful experiences lead to depression in some people but not in others? It was hypothesized that the (5-HTT) gene is found to moderate the influence of stressful life events on depression. There were Three lines of experimental research that suggest this hypothesis of a gene-by-environment ($G \times E$) interaction. It was tested in mice, in rhesus macaques and in human neuroim- aging research that suggests that the stress response is mediated by variations in the 5-HTTLPR (Caspi et al). Also, after the hypothesis was made, it was tested “among members of the Dunedin Multidisciplinary Health and Development Study (16). This representative birth cohort of 1037 children (52% male) has been assessed at ages 3, 5, 7, 9, 11, 13, 15, 18, and 21 and was virtually intact (96%) at the age of 26 years” (Caspi et al). The two big questions were to figure out whether there are any associations of the synthetic Y variable with any of the G variables and whether there are any interactions of G and E variables or G and G variables.

Method

The first step that was done was to import my dataset into R Studios and realized that there was no missing data. My dataset is complete. I had 1,253 observations and 25 variables. There was 1 dependent variable which is Y and 24 Independent variables ($E1$ to $E4$ and $G1$ to $G20$). Next what I did was model the Environmental variables and got the Summary portion. Then, it was time to work with the Genetic Variables. I used the lm function to model it and I assumed this model includes all interaction terms up to the 2nd order. A Residual plot was created and saw that it was adequate because all I saw was a flat ellipse. The summary was then created. I felt like the residual plot could be a bit better. Therefore, a transformation was made. I did the boxcox and saw that it's not going to be a log transformation. I used this line of code `best.lam = bc$x[which(bc$y==max(bc$y))]` and got the lambda value to be 1.81 which I rounded it to 2. Therefore, I used the transformation y^2 . The new residual plot was created and the summary. Then, we performed Stepwise Regression and the R Studios produced a proposed model. The model summary table is shown below with adjusted R^2 and the Bayesian Information Criterion (BIC). The code for that is written below. I chose variables selected in the 3rd model as candidates; namely $E3$, $E1$, $E3$, $G9$, and $G14$. That is, I split the $E1E3$ and $G9G14$ interactions. Besides this, I want to also make sure main effects that are significant are in the

model. The code for that is also shown below in the appendix. The code allowed me to get sig coefficients and the variables with all of the p-values. The E1:E3 interaction was eliminated.

Results:

After modeling the Environmental variables, the adjusted R squared value was 0.5899. As for the additional contribution of the Genetic Variables, the residual plot was created and the adjusted R squared value after assuming that this model includes all interaction terms up to the 2nd order in the model is 0.6002644. This was already a good model, however the transformation was found and it was a little bit better. I used the y^2 transformation using boxcox. The lambda value that I got was 1.81 which I rounded it to 2. I got the new adjusted R squared value to be 0.6006907. Then, after doing the stepwise regression, according to the model summary table below we see that, there is an obvious increase in the adjR2 from the 2nd model to the 3rd model. There is only a very small increase from the 3rd to the 4th model, which may not be significant. Another measure that researchers use to assess models is the Bayesian Information Criterion (BIC). The BIC values consistently decrease. The decrease from the 3rd to the 4th and the 5th model is much smaller than the other decreases. I ended up choosing the 3rd model as candidates namely E3, E1, E3, G9, and G14. That is, I split the E1E3 and G9G14 interactions. However, the table titled Sig Coefficient shows variables E1, E2 and G14 that are significant in the main effect model. I then changed my model to include all interaction terms up to the 1st order. The sig coefficient table showed E1, E3, and G14 to be significant in both the 2nd order and the first order interaction. Based on the last table where all the p-values are shown, the E1:E3 interaction was eliminated. We only had G9:G14 interaction. In the end, adding the G variables only made a little improvement since the adjusted R squared value went from 0.5899 to 0.6002644.

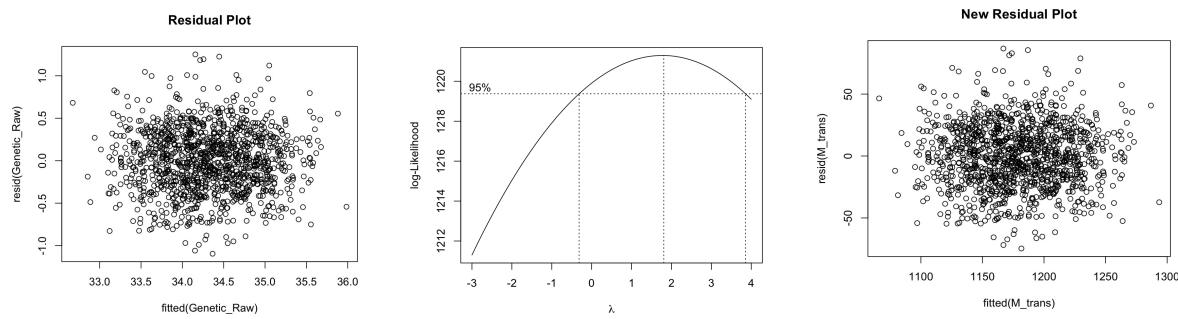


Table: Model Summary

| model | adjR2 | IC |
|---|-------------------|-------------------|
| (Intercept)+E1:E3 | 0.577745343507559 | -1067.00441622411 |
| (Intercept)+E1:E3+G9:G14 | 0.590059481964278 | -1097.95754934016 |
| (Intercept)+E3+E1:E3+G9:G14 | 0.601626601448881 | -1127.69081955919 |
| (Intercept)+E3+E1:E3+G2:G8+G9:G14 | 0.603590146881052 | -1127.75231891966 |
| (Intercept)+E3+E1:E3+G2:G8+G9:G14+G17:G18 | 0.605381471055271 | -1127.29840653493 |

Table: Sig Coefficients

| | Estimate | Std. Error | t value | Pr(> t) |
|-------------|-------------|-------------|------------|-----------|
| (Intercept) | 410.0854061 | 25.94876561 | 15.8036581 | 0.0000001 |
| E1 | 14.5466411 | 0.56201621 | 25.8829561 | 0.0000001 |
| E3 | 19.4175781 | 0.56870371 | 34.1435741 | 0.0000001 |
| G14 | 6.7229981 | 1.81804011 | 3.6979371 | 0.0002271 |

Conclusion

Our analysis found associations with genetic variables after the environmental variables had been controlled. Specifically the G9-G14 interaction with t-value 4.4014975 significantly associated with the χ^2 . Including these variables significantly increased the r-squared value from 0.5912 to 0.6044. The model that the TA used for my data is

$y^2 = \beta_0 + \beta_1 E_3 + \beta_2 E_1 + \beta_3 G_9 G_{14} + \epsilon$. The final model with estimated parameters is:

$$y^2 = 513.311353 + 15.319715E + 9.908984E + 17.808464G G \quad . \text{ This is shown in}$$

the table below.

| | Estimate | Std. Error | t value | Pr(> t) |
|-------------|------------|------------|------------|--------------|
| (Intercept) | 513.311353 | 200.292591 | 2.5628075 | 1.049997e-02 |
| E3 | 15.319715 | 8.869358 | 1.7272630 | 8.436905e-02 |
| E1 | 9.908984 | 8.883250 | 1.1154683 | 2.648656e-01 |
| G9 | 12.824424 | 40.920885 | 0.3133956 | 7.540328e-01 |
| G14 | -26.413922 | 39.205027 | -0.6737381 | 5.006032e-01 |
| G9:G14 | 17.808464 | 4.046001 | 4.4014975 | 1.167364e-05 |

Citation:

Caspi, Avshalom, et al. "Influence of Life Stress on Depression: Moderation by a Polymorphism in the 5-HTT Gene." *Science*, American Association for the Advancement of Science, 18 July 2003, science.sciencemag.org/content/301/5631/386/tabs-figures-data.

Code for Modeling the Environmental Variables

```
> library(readr)
> SP21_P2_949881 <- read_csv("~/Downloads/Project 2 Data Set SP2021/SP21_P2_949881.csv")

--- Column specification ---
cols(
  .default = col_double()
)
i Use `spec()` for the full column specifications.

> View(SP21_P2_949881)
> Environmental_fit <- lm(Y ~ E1+E2+E3+E4+E5, data=SP21_P2_949881)
Error in eval(predvars, data, env) : object 'E5' not found
> Environmental_fit <- lm(Y ~ E1+E2+E3+E4, data=SP21_P2_949881)
> summary(Environmental_fit)

Call:
lm(formula = Y ~ E1 + E2 + E3 + E4, data = SP21_P2_949881)

Residuals:
    Min      1Q  Median      3Q     Max 
-1.36359 -0.30496 -0.00475  0.29390  1.51875 

Coefficients:
            Estimate Std. Error t value Pr(>|t|)    
(Intercept) 23.144364   0.375954  61.562   <2e-16 ***
E1          0.212730   0.008210  25.911   <2e-16 ***
E2          0.002782   0.008264   0.337    0.736    
E3          0.281853   0.008274  34.066   <2e-16 ***
E4         -0.001355   0.008354  -0.162    0.871    
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.425 on 1248 degrees of freedom
Multiple R-squared:  0.5912,    Adjusted R-squared:  0.5899 
F-statistic: 451.3 on 4 and 1248 DF,  p-value: < 2.2e-16
```

Code for Modeling the Additional Contribution of the Genetic Variables, Getting the Residual Plot and Finding the Transformation

```
> Genetic_Raw <- lm( Y ~ (E1+E2+E3+E4+G1+G2+G3+G4+G5+G6+G7+G8+G9+G10+G11+G12+G13+G14+G15+G16+G17+G18+G19+G20)^2, data=SP21_P2_949881 )
> plot(resid(Genetic_Raw) ~ fitted(Genetic_Raw), main='Residual Plot')
> boxcox(Genetic_Raw)
> Genetic_trans <- lm( I(log(Y)) ~ .)^2, data=SP21_P2_949881 )
> summary(Genetic_Raw)$adj.r.square
[1] 0.6002644
> summary(Genetic_trans)$adj.r.square
[1] 0.5996779
> Genetic_trans <- lm( I(Y)^(1/2) ~ .)^2, data=SP21_P2_949881 )
> summary(Genetic_trans)$adj.r.square
[1] 0.5999911
> Genetic_trans <- lm( I(Y)^(1/2) ~ .)^2, data=SP21_P2_949881 )
> Genetic_trans <- lm( I(Y)^(2) ~ .)^2, data=SP21_P2_949881 )
> summary(Genetic_trans)$adj.r.square
[1] 0.6006907
```

Code for Getting the Summary of the Additional Contribution of the Genetic Variables Without Transformation

MD Mahmudur Rahman

```

> summary(Genetic_Raw)

Call:
lm(formula = Y ~ (E1 + E2 + E3 + E4 + G1 + G2 + G3 + G4 + G5 +
G6 + G7 + G8 + G9 + G10 + G11 + G12 + G13 + G14 + G15 + G16 +
G17 + G18 + G19 + G20)^2, data = SP21_P2_949881)

Residuals:
    Min      1Q  Median      3Q     Max 
-1.09452 -0.25063  0.00228  0.25494  1.25136 

Coefficients:
            Estimate Std. Error t value Pr(>|t|)    
(Intercept) 24.2032327  8.1042194  2.986  0.00289 **  
E1          0.1632109  0.2492668  0.655  0.51278    
E2         -0.0893265  0.2613321 -0.342  0.73257    
E3          0.2630954  0.2563671  1.026  0.38504    
E4          0.1035604  0.2583491  0.401  0.68862    
G1          0.1899178  0.9572110  0.198  0.84277    
G2          0.0654399  0.9429396  0.069  0.94469    
G3          -0.2766230  0.9837082 -0.281  0.77861    
G4         -1.4893894  0.9591759 -1.553  0.12081    
G5          -0.2022686  0.9660415 -0.209  0.83420    
G6          -0.7066605  0.9995759 -0.707  0.47976    
G7         -2.4758661  0.9933794 -2.492  0.01286 *  
G8          -0.9661723  0.9898797 -0.976  0.32929    
G9          1.4265073  1.0083603  1.415  0.15749    
G10         0.1276429  0.9379250  0.136  0.89178    
G11         1.7731926  0.9442782  1.878  0.06071 .  
G12         0.8261814  1.0050599  0.822  0.41127    
G13         -0.2717844  0.9791398 -0.278  0.78140    
G14         -0.0492792  0.9608174 -0.051  0.95911    
G15         0.0109222  0.9534526  0.011  0.99086    
G16         0.7283826  1.0085192  0.722  0.47033    
G17        -0.3275600  0.9915678 -0.330  0.74121    
G18         -0.5959492  0.9535185 -0.625  0.53212    
G19        -1.7462114  0.9731426 -1.794  0.07307 .  
G20         0.0095157  0.9178548  0.010  0.99173    
E1:E2        0.0036022  0.0064525  0.558  0.57680    
E1:E3        0.0053444  0.0066567  0.803  0.42226    
E1:E4       -0.0071957  0.0065432 -1.100  0.27173    
E1:G1       -0.0210337  0.0207295 -1.015  0.31052    
E1:G2        0.0338809  0.0200545  1.689  0.09146 .  
E1:G3        0.0111566  0.0210796  0.529  0.59675    
E1:G4        0.0433236  0.0211074  2.053  0.04039 *  

```

MD Mahmudur Rahman

| | | | | |
|--------|------------|-----------|--------|-----------|
| G3:G7 | 0.0489207 | 0.0693335 | 0.706 | 0.48062 |
| G3:G8 | -0.0762264 | 0.0703958 | -1.083 | 0.27916 |
| G3:G9 | -0.0820956 | 0.0672419 | -1.221 | 0.22243 |
| G3:G10 | 0.0771642 | 0.0692639 | 1.114 | 0.26553 |
| G3:G11 | 0.0053330 | 0.0675484 | 0.079 | 0.93709 |
| G3:G12 | -0.0795181 | 0.0714112 | -1.114 | 0.26576 |
| G3:G13 | -0.0009792 | 0.0669437 | -0.015 | 0.98833 |
| G3:G14 | 0.0887053 | 0.0689812 | 1.286 | 0.19878 |
| G3:G15 | -0.0893085 | 0.0671225 | -1.331 | 0.18366 |
| G3:G16 | -0.0182270 | 0.0695817 | -0.262 | 0.79342 |
| G3:G17 | 0.0317662 | 0.0705801 | 0.450 | 0.65276 |
| G3:G18 | 0.0124986 | 0.0671750 | 0.186 | 0.85244 |
| G3:G19 | 0.0580081 | 0.0676586 | 0.857 | 0.39146 |
| G3:G20 | 0.0631692 | 0.0703168 | 0.898 | 0.36923 |
| G4:G5 | 0.0019520 | 0.0667221 | 0.029 | 0.97667 |
| G4:G6 | 0.1101162 | 0.0680457 | 1.618 | 0.10594 |
| G4:G7 | 0.0766008 | 0.0686973 | 1.115 | 0.26511 |
| G4:G8 | 0.1250608 | 0.0721622 | 1.733 | 0.08341 |
| G4:G9 | -0.0270300 | 0.0687989 | -0.393 | 0.69449 |
| G4:G10 | 0.0213548 | 0.0663559 | 0.322 | 0.74766 |
| G4:G11 | -0.0639321 | 0.0673599 | -0.949 | 0.34280 |
| G4:G12 | 0.0648314 | 0.0723671 | 0.896 | 0.37055 |
| G4:G13 | 0.0490973 | 0.0677568 | 0.725 | 0.46887 |
| G4:G14 | 0.0857092 | 0.0658415 | 1.302 | 0.19332 |
| G4:G15 | 0.0179232 | 0.0665023 | 0.270 | 0.78759 |
| G4:G16 | 0.0134290 | 0.0697239 | 0.193 | 0.84731 |
| G4:G17 | 0.0254842 | 0.0722091 | 0.353 | 0.72422 |
| G4:G18 | 0.0297874 | 0.0705311 | 0.422 | 0.67288 |
| G4:G19 | -0.0333188 | 0.0684265 | -0.487 | 0.62642 |
| G4:G20 | 0.1358131 | 0.0664229 | 2.045 | 0.04116 * |
| G5:G6 | 0.0642481 | 0.0703443 | 0.913 | 0.36130 |
| G5:G7 | -0.0298923 | 0.0707585 | -0.422 | 0.67279 |
| G5:G8 | 0.0737206 | 0.0734092 | 1.004 | 0.31552 |
| G5:G9 | 0.0448837 | 0.0696234 | 0.645 | 0.51930 |
| G5:G10 | 0.0591099 | 0.0669099 | 0.883 | 0.37723 |
| G5:G11 | -0.0314655 | 0.0697970 | -0.451 | 0.65223 |
| G5:G12 | -0.0198348 | 0.0707792 | -0.280 | 0.77936 |
| G5:G13 | 0.0666785 | 0.0691935 | 0.964 | 0.33546 |
| G5:G14 | -0.0139349 | 0.0679491 | -0.205 | 0.83756 |
| G5:G15 | -0.0380960 | 0.0669689 | -0.569 | 0.56959 |
| G5:G16 | 0.0187042 | 0.0713402 | 0.262 | 0.79324 |
| G5:G17 | 0.0912430 | 0.0704469 | 1.295 | 0.19556 |
| G5:G18 | 0.1115888 | 0.0702722 | 1.588 | 0.11263 |
| G5:G19 | -0.0450347 | 0.0689877 | -0.653 | 0.51405 |
| G5:G20 | -0.0356792 | 0.0669387 | -0.533 | 0.59415 |

```

G6:G7      0.0700004  0.0676090  1.035  0.30076
G6:G8      0.0715260  0.0716130  0.999  0.31815
G6:G9      -0.1101537  0.0702140  -1.569  0.11702
G6:G10     0.1076522  0.0696508  1.546  0.12253
[ reached getOption("max.print") -- omitted 101 rows ]
---
```

Signif. codes: 0 ‘***’ 0.001 ‘**’ 0.01 ‘*’ 0.05 ‘.’ 0.1 ‘ ’ 1

Residual standard error: 0.4196 on 952 degrees of freedom
 Multiple R-squared: 0.696, Adjusted R-squared: 0.6003
 F-statistic: 7.267 on 300 and 952 DF, p-value: < 2.2e-16

Code for Getting the Summary of the Additional Contribution of the Genetic Variables With Transformation

```
> summary(Genetic_trans)

Call:
lm(formula = I(Y)^4 ~ .)^2, data = SP21_P2_949881)

Residuals:
    Min      1Q  Median      3Q     Max 
-173736 -40998   -285   40959  208627 

Coefficients:
            Estimate Std. Error t value Pr(>|t|)    
(Intercept) 142594.56 1310970.64  0.109  0.91341  
E1          6724.35  40322.38  0.167  0.86759  
E2         -14526.01  42274.12 -0.344  0.73121  
E3          23982.64  41470.95  0.578  0.56320  
E4          20555.12  41791.58  0.492  0.62294  
G1          34096.54  154842.24  0.220  0.82576  
G2          16626.00  152533.54  0.109  0.91323  
G3         -39422.70  159128.54 -0.248  0.80439  
G4         -255969.90  155160.09 -1.650  0.09933 .  
G5          -33651.39  156270.70 -0.215  0.82955  
G6         -122302.16  161695.35 -0.756  0.44961  
G7         -401252.49  160692.99 -2.497  0.01269 *  
G8          -146146.67  160126.86 -0.913  0.36164  
G9          227360.44  163116.36  1.394  0.16369  
G10         10860.97  151722.46  0.072  0.94295  
G11         309905.52  152750.18  2.029  0.04275 *  
G12         134648.20  162582.47  0.828  0.40777  
G13         -38796.09  158389.53 -0.245  0.80655  
G14         -19953.96  155425.64 -0.128  0.89787  
G15          7068.01  154234.17  0.046  0.96346  
G16         130965.19  163142.07  0.803  0.42231  
G17         -55156.46  160399.94 -0.344  0.73102  
G18         -85075.41  154244.93 -0.552  0.58138  
G19        -305062.63  157419.30 -1.938  0.05293 .  
G20          809.35  148475.82  0.005  0.99565  
E1:E2        630.01  1043.78  0.604  0.54627  
E1:E3       1714.15  1076.82  1.592  0.11175  
E1:E4      -1192.83  1058.45 -1.127  0.26004  
E1:G1      -3516.97  3353.28 -1.049  0.29453  
E1:G2        5248.93  3244.09  1.618  0.10599  
E1:G3       1881.15  3409.92  0.552  0.58131  
E1:G4       7265.49  3414.42  2.128  0.03360 *  
E1:G5      -3483.29  3354.95 -1.038  0.29942  
E1:G6       6763.10  3422.17  1.976  0.04841 * 
```

MD Mahmudur Rahman

Signif. codes: 0 ‘***’ 0.001 ‘**’ 0.01 ‘*’ 0.05 ‘.’ 0.1 ‘ ’ 1

Code for the Stepwise Regression, Modeling the Summary and the Table for Sig Coefficient

```

> library(leaps)
> M <- regsubsets( model.matrix(M_trans)[,-1], I(SP21_P2_949881$Y^(2))), 
Error: unexpected ',' in "M <- regsubsets( model.matrix(M_trans)[,-1], I(SP21_P2_949881$Y^(2))),"
>           nbest = 1 , nvmax=5,
Error: unexpected ',' in "           nbest = 1 ,"
>           method = 'forward', intercept = TRUE )
Error: unexpected ',' in "           method = 'forward',"
> M <- regsubsets( model.matrix(Genetic_trans)[,-1], I(SP21_P2_949881$Y^2)), nbest = 1 , nvmax=5, method = 'forward',
intercept = TRUE )
Error: unexpected ',' in "M <- regsubsets( model.matrix(Genetic_trans)[,-1], I(SP21_P2_949881$Y^2)),"
> M <- regsubsets( model.matrix(Genetic_trans)[,-1], I(SP21_P2_949881$Y^2),nbest = 1 , nvmax=5,
+           method = 'forward', intercept = TRUE )
> temp <- summary(M)
> Var2 <- colnames(model.matrix(M_trans))
> M_select2 <- apply(temp$which, 1,
+           function(x) paste0(Var2[x], collapse='+'))
> kable(data.frame(cbind( model = M_select2, adjR2 = temp$adjr2, BIC = temp$bic)),
+       caption='Model Summary')

```

Table: Model Summary

| model | adjR2 | BIC |
|--|-------------------|-------------------|
| (Intercept)+E1:E3 | 0.578146731094707 | -1068.19606167482 |
| (Intercept)+E1:E3+G9:G14 | 0.590578080300737 | -1099.54366989543 |
| (Intercept)+E3+E1:E3+G9:G14 | 0.602374779809626 | -1130.04627028489 |
| (Intercept)+E3+E1:E3+G2:G8+G9:G14 | 0.604331544130127 | -1130.09797341833 |
| (Intercept)+E3+E1:E3+G2:G8+G9:G14+G17:G18 | 0.606129517267023 | -1129.67587069904 |
| > M_main <- lm(I(Y^2) ~ ., data=SP21_P2_949881) | | |
| > temp <- summary(M_main) | | |
| > kable(temp\$coefficients[abs(temp\$coefficients[,4]) <= 0.001,], caption='Sig Coefficients') | | |

Table: Sig Coefficients

| | Estimate | Std. Error | t value | Pr(> t) |
|-------------|-------------|-------------|------------|-----------|
| (Intercept) | 410.0854061 | 25.94876561 | 15.8036581 | 0.000000 |
| E1 | 14.5466411 | 0.56201621 | 25.8829561 | 0.000000 |
| E3 | 19.4175781 | 0.56870371 | 34.1435741 | 0.000000 |
| G14 | 6.7229981 | 1.81804011 | 3.6979371 | 0.0002271 |

MD Mahmudur Rahman

```
> M_1st <- lm( I(Y^2) ~ ., data=SP21_P2_949881)
> temp <- summary(M_1st)
> kable(temp$coefficients[ abs(temp$coefficients[,4]) <= 0.001, ], caption='1st Interaction')
```

Table: 1st Interaction

| | Estimate | Std. Error | t value | Pr(> t) |
|-------------|-------------|-------------|------------|-----------|
| (Intercept) | 410.0854061 | 25.94876561 | 15.8036581 | 0.0000001 |
| E1 | 14.5466411 | 0.56201621 | 25.8829561 | 0.0000001 |
| E3 | 19.4175781 | 0.56870371 | 34.1435741 | 0.0000001 |
| G14 | 6.7229981 | 1.81804011 | 3.6979371 | 0.0002271 |

| | Estimate | Std. Error | t value | Pr(> t) |
|-------------|--------------|--------------|--------------|--------------|
| M_2stage | 1.780846e+01 | 4.046001e+00 | 4.401497e+00 | 1.167364e-05 |
| (Intercept) | 513.311353 | 200.292591 | 2.5628075 | 1.049997e-02 |
| E3 | 15.319715 | 8.869358 | 1.7272630 | 8.436905e-02 |
| E1 | 9.908984 | 8.883250 | 1.1154683 | 2.648656e-01 |
| G9 | 12.824424 | 40.920885 | 0.3133956 | 7.540328e-01 |
| G14 | -26.413922 | 39.205027 | -0.6737381 | 5.006032e-01 |
| G9:G14 | 17.808464 | 4.046001 | 4.4014975 | 1.167364e-05 |