Paul Johnson April 25, 2013

# Data Management

```
library(foreign)
library(rockchalk)
i <- 18
dat <- read.dta(paste("../student-test2/student-",i,".dta", sep = ""))
```

The variables pprof and pnet are scored as numeric, but really they are factors. So convert them to prevent future mis-understandings.

```
dat$pprof <- factor(dat$pprof, labels = c("NO", "YES"))
dat$pnet <- factor(dat$pnet, labels = c("NO","YES"))
```

```
datsum <- summarize(dat)
```

Table would need some hand customization

```
library(xtable)
print(xtable(datsum$numeric, caption = "Best Automatic Summary Table for Numerics", label =
    "table1"), "latex")
```

|  | act | harv | ibs | sal1 | sal2 | sal3 | sat |
|---|---|---|---|---|---|---|---|
| 0% | 7.88 | 1095.00 | 66.74 | 3933.00 | 5132.00 | 148700.00 | 1080.00 |
| 25% | 18.69 | 1509.00 | 92.55 | 16630.00 | 19260.00 | 161300.00 | 1487.00 |
| 50% | 22.02 | 1606.00 | 98.96 | 20050.00 | 23130.00 | 165400.00 | 1588.00 |
| 75% | 25.55 | 1719.00 | 105.80 | 24380.00 | 27150.00 | 169000.00 | 1694.00 |
| 100% | 38.91 | 2377.00 | 130.50 | 38350.00 | 45250.00 | 187900.00 | 2341.00 |
| mean | 22.10 | 1612.00 | 99.22 | 20270.00 | 23010.00 | 165200.00 | 1589.00 |
| sd | 5.13 | 162.30 | 9.77 | 5647.00 | 5989.00 | 5751.00 | 162.80 |
| var | 26.28 | 26340.00 | 95.40 | 31890000.00 | 35870000.00 | 33080000.00 | 26510.00 |
| NA's | 19.00 | 56.00 | 0.00 | 7.00 | 0.00 | 0.00 | 22.00 |
| N | 541.00 | 541.00 | 541.00 | 541.00 | 541.00 | 541.00 | 541.00 |

Table 1: Best Automatic Summary Table for Numerics

Let students figure way to beautify this:

```
print(datsum$factors)
```

```
          gender                    major                      pnet
F             :300.0000   S             :190.0000   NO            :393.0000
M             :241.0000   H             :189.0000   YES           :148.0000
NA's          :  0.0000   N             :162.0000   NA's          :  0.0000
entropy       :  0.9914   NA's          :  0.0000   entropy       :  0.8465
normedEntropy :  0.9914   entropy       :  1.5812   normedEntropy :  0.8465
N             :541.0000   normedEntropy :  0.9976   N             :541.0000
                          N             :541.0000
          pprof
NO            :378.0000
YES           :163.0000
NA's          :  0.0000
entropy       :  0.8829
normedEntropy :  0.8829
N             :541.0000
```

# Aptitude Test Variables

There's severe multicollinearity between the variables harv, sat, and act. It seems clear we can't estimate both sat and harv, and several students noticed that since harv is a summary of the other tests, then there's some reason to suppose sat is a better variable. (I know for a fact that harv = sat + act).

   Please find Table 2. I left the Iowa Basic Skills variable in my best model, mainly because I wanted to estimate that coefficient, even though the F test below indicates one can exclude harv and ibs from the "full" model without losing any sleep.

```
m1s <- lm(sal1 ~ sat, data = dat)
m1a <- lm(sal1 ~ act, data = dat)
m1i <- lm(sal1 ~ ibs, data = dat)
m1h <- lm(sal1 ~ harv, data = dat)
m1all <- lm(sal1 ~ sat + act + ibs + harv, data = dat)
m1best <- lm(sal1 ~ sat + act + ibs, data = dat)
```

```
mcDiagnose(m1all)
```

```
The following auxiliary models are being estimated and returned in a list:
sat ~ act + ibs + harv
<environment: 0x2566770>
act ~ sat + ibs + harv
<environment: 0x2566770>
ibs ~ sat + act + harv
<environment: 0x2566770>
harv ~ sat + act + ibs
<environment: 0x2566770>
Drum roll please!

And your R_j Squareds are (auxiliary Rsq)
      sat       act       ibs      harv
0.9998491 0.8806357 0.2582530 0.9998531
The Corresponding VIF, 1/(1-R_j^2)
      sat       act       ibs      harv
6627.063268    8.377717    1.348169 6807.321651
Bivariate Correlations for design matrix
      sat  act  ibs harv
sat   1.00 0.41 0.44 1.00
act   0.41 1.00 0.41 0.44
ibs   0.44 0.41 1.00 0.45
harv  1.00 0.44 0.45 1.00
```

```
niceLabels <- c(act = "ACT", sat = "SAT", harv = "Harvard SS", ibs = "Iowa BS",  majorS = "
    Major: Soc.", majorN = "Major: Nat.", majorH = "Major: Hum.", pnetYES = "Parent Network
    : Yes",  pprofYES="Prof. Parents: Yes", genderM = "Gender: Male", "log(harv)"= "ln(
    Harvard SS)",
    "I(harv * harv)"= "Harvard SS$^2$", major2H = "Major 2: Hum.", major2N = "Major 2: Nat.
    ")
outreg(list(m1s, m1a, m1i, m1h, m1all, m1best), tight = TRUE, modelLabels = c("SAT","ACT","
    IBS","Harvard SS", "All", "Best"), varLabels = niceLabels, title = paste("Regression
    with sal1: Student-",i, sep=""), label = "tab:tab2")
```

   Could conduct an F test of the hypothesis that $b_{ibs} = b_{harv} = 0$. But which model should I be testing? Test the one with all the variables, to see if *harv* and *ibs* should both be set to 0. To do that, I need to take the data frame used to fit m1all and use it to fit the restricted model. Otherwise, the F test fails.

```
m1alldf <- model.frame(m1all)
m1restricted <- lm(sal1 ~ sat + act, data = m1alldf)
anova(m1restricted, m1all)
```

```
Analysis of Variance Table

Model 1: sal1 ~ sat + act
Model 2: sal1 ~ sat + act + ibs + harv
```

Table 2: Regression with sal1: Student-18

| | SAT Estimate (S.E.) | ACT Estimate (S.E.) | IBS Estimate (S.E.) | Harvard SS Estimate (S.E.) | All Estimate (S.E.) | Best Estimate (S.E.) |
|---|---|---|---|---|---|---|
| (Intercept) | -2892.637 (2236.236) | 13393.584* (1056.889) | 5245.229* (2409.908) | -1508.676 (2324.386) | -4643.679 (2856.416) | -5388.4* (2716.79) |
| SAT | 14.573* (1.399) | . | . | . | 74.371 (120.372) | 11.945* (1.624) |
| ACT | . | 309.022* (46.482) | . | . | 167.947 (133.263) | 120.918* (51.053) |
| Iowa BS | . | . | 151.344* (24.153) | . | 39.399 (28.534) | 40.213 (27.138) |
| Harvard SS | . | . | . | 13.45* (1.434) | -62.647 (120.335) | . |
| N | 513 | 515 | 534 | 478 | 444 | 496 |
| RMSE | 5167.077 | 5402.498 | 5454.873 | 5101.35 | 5038.285 | 5089.915 |
| $R^2$ | 0.175 | 0.079 | 0.069 | 0.156 | 0.185 | 0.196 |
| adj $R^2$ | 0.174 | 0.078 | 0.067 | 0.154 | 0.178 | 0.191 |

$*p \leq 0.05$

```
  Res.Df         RSS Df Sum of Sq       F Pr(>F)
1    441 1.1201e+10
2    439 1.1144e+10  2   56865011 1.1201  0.3272
```

Noticing this sample size problem, I wondered if I should re-do Table 2 so that all are fitted on the exact same data. Since I exclude harv, should those cases that are missing on harv "come back to life" when I exclude harv from the model? I think so. Still, there is something unappetizing about this. Fitting harv causes a loss of cases, no matter how we look at it. So for the best model and the ones for sat and ibs, I use the sample from the best model, but when harv enters the picture, we lose some cases.

```
m1best <- lm(sal1 ~ sat + act + ibs, data = dat)
dat2 <- model.frame(m1best)
m1s <- lm(sal1 ~ sat, data = dat2)
m1a <- lm(sal1 ~ act, data = dat2)
m1i <- lm(sal1 ~ ibs, data = dat2)
m1h <- lm(sal1 ~ harv, data = dat[row.names(dat2), ])
m1all <- lm(sal1 ~ sat + act + ibs + harv, data = dat[row.names(dat2), ])


outreg(list(m1s, m1a, m1i, m1h, m1all, m1best), tight = TRUE, modelLabels = c("SAT","ACT","
    IBS","Harvard SS", "All", "Best"), varLabels = niceLabels)
```

|  | SAT Estimate (S.E.) | ACT Estimate (S.E.) | IBS Estimate (S.E.) | Harvard SS Estimate (S.E.) | All Estimate (S.E.) | Best Estimate (S.E.) |
|---|---|---|---|---|---|---|
| (Intercept) | -2948.921 (2253.726) | 13338.625* (1074.709) | 4941.978* (2493.27) | -2399.291 (2372.501) | -4643.679 (2856.416) | -5388.4* (2716.79) |
| SAT | 14.608* (1.411) | . | . | . | 74.371 (120.372) | 11.945* (1.624) |
| ACT | . | 312.386* (47.21) | . | . | 167.947 (133.263) | 120.918* (51.053) |
| Iowa BS | . | . | 154.421* (25.006) | . | 39.399 (28.534) | 40.213 (27.138) |
| Harvard SS | . | . | . | 14.033* (1.465) | -62.647 (120.335) | . |
| N | 496 | 496 | 496 | 444 | 444 | 496 |
| RMSE | 5135.115 | 5429.508 | 5458.241 | 5062.237 | 5038.285 | 5089.915 |
| $R^2$ | 0.178 | 0.081 | 0.072 | 0.172 | 0.185 | 0.196 |
| adj $R^2$ | 0.177 | 0.08 | 0.07 | 0.17 | 0.178 | 0.191 |

$*p \leq 0.05$

Deciding what's "important"? We have lots of ways. If I've settled on a "best" model, it seems like I should be confined to the variables in that model. And the diagnostics should not depend on harv. Here are the partial and semi-partial correlations.

```
getPartialCor(m1best)
```

```
            sal1
sal1  −1.00000000
sat    0.31468948
act    0.10617604
ibs    0.06665687
```

```
getDeltaRsquare(m1best)
```

```
The deltaR−square values: the change in the R−square
      observed when a single term is removed.
Same as the square of the 'semi−partial correlation coefficient'
    deltaRsquare
sat   0.088371431
act   0.009167159
ibs   0.003588242
```

I admit, it is tough to conceptualize the scales of the different variables. I suppose I could scale the sat, act, and ibs scores so that they are all on the same 0-100 scale. Then I'll re-run the model. (This is called "percent of maximum" scoring (POMS)). Since we KNOW from previous work that re-scaling a variable has absolutely no substantive impact on the fit, and it is just for convenience of interpretation, this is an innocuous change.

```
dat2$satpoms <− 100*(dat2$sat − min(dat2$sat))/(max(dat2$sat) − min(dat2$sat))
dat2$actpoms <− 100*(dat2$act − min(dat2$act))/(max(dat2$act) − min(dat2$act))
dat2$ibspoms <− 100*(dat2$ibs − min(dat2$ibs))/(max(dat2$ibs) − min(dat2$ibs))
summarize(dat2[, c("satpoms","actpoms","ibspoms")])
```

```
$numerics
     actpoms  ibspoms  satpoms
0%      0.00     0.00     0.00
25%    34.97    40.71    32.20
50%    45.84    50.30    40.28
75%    57.12    61.39    48.63
100%  100.00   100.00   100.00
```

```
mean     46.05    50.93    40.38
sd       16.66    15.38    12.97
var     277.50   236.60   168.20
NA's      0.00     0.00     0.00
N       496.00   496.00   496.00

$factors
NULL
```

```
m1poms <- lm(sal1 ~ satpoms + actpoms + ibspoms, data = dat2)
summary(m1poms)
```

```
Call:
lm(formula = sal1 ~ satpoms + actpoms + ibspoms, data = dat2)

Residuals:
     Min       1Q    Median       3Q       Max
 -14035.4  -3702.2    -98.4   3448.3   14262.9

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  11145.74     926.69  12.027  < 2e-16 ***
satpoms        150.66      20.49   7.354 8.11e-13 ***
actpoms         37.52      15.84   2.368   0.0182 *
ibspoms         25.65      17.31   1.482   0.1390
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 5090 on 492 degrees of freedom
Multiple R^2: 0.196 , Adjusted R^2: 0.1911
F-statistic: 39.98 on 3 and 492 DF,  p-value: < 2.2e-16
```

Oh, one more thing. Recall my point that partial and semi-partial correlations are completely worthless when 1) there is multicollinearity and 2) we are uncertain which variables should be in consideration. Notice how crazy your conclusions would be if you based them on the "full" model.

```
options(scipen = 10)
getPartialCor(m1all)
```

```
             sal1
sal1  -1.00000000
sat    0.02947527
act    0.06004081
ibs    0.06575769
harv  -0.02483929
```

```
getDeltaRsquare(m1all)
```

```
The deltaR-square values: the change in the R-square
      observed when a single term is removed.
Same as the square of the 'semi-partial correlation coefficient'
     deltaRsquare
sat   0.0007084979
act   0.0029478606
ibs   0.0035385116
harv  0.0005030277
```

```
options(scipen = 5)
```

# Additional Variables

Please see Table 3 for the regressions.

Table 3: Regression with sal2: Student-18

| | Test Scores Only Estimate (S.E.) | All Predictors Estimate (S.E.) |
|---|---|---|
| (Intercept) | -2127.231 | -4835.422 |
| | (2903.705) | (2753.496) |
| SAT | 11.759* | 11.877* |
| | (1.736) | (1.613) |
| ACT | 146.112* | 128.346* |
| | (54.403) | (50.748) |
| Iowa BS | 32.472 | 39.298 |
| | (29.029) | (27.151) |
| Major: Soc. | . | 1510.397* |
| | | (546.414) |
| Major: Nat. | . | 4943.763* |
| | | (564.78) |
| Prof. Parents: Yes | . | 560.202 |
| | | (492.372) |
| Parent Network: Yes | . | 991.785 |
| | | (516.592) |
| Gender: Male | . | -474.161 |
| | | (455.713) |
| N | 502 | 502 |
| RMSE | 5462.866 | 5071.223 |
| $R^2$ | 0.175 | 0.296 |
| adj $R^2$ | 0.17 | 0.285 |

$*p \leq 0.05$

```
m2small <- lm(sal2 ~ sat + act + ibs, data = dat)
m2all <- lm(sal2 ~ sat + act + ibs + major + pprof + pnet + gender, data = dat)
outreg(list(m2small, m2all), tight = TRUE, title = paste("Regression with sal2: Student-",i
    , sep=""),modelLabels = c("Test Scores Only","All Predictors"), varLabels = niceLabels,
        label = "table3")
```

Fancy T test. Lets use the big model to find out if $b_{pnetYES} = b_{pprofYES}$.

```
m2allc <- coef(m2all)
m2allv <- vcov(m2all)
numer <- m2allc["pprofYES"] - m2allc["pnetYES"]
names(numer) <- "difference"
denom <- sqrt(m2allv["pprofYES","pprofYES"] + m2allv["pnetYES","pnetYES"] - 2 * m2allv["
    pprofYES","pnetYES"])
print(paste("Fancy T: ", "Numerator = ", numer, "Denominator = ", denom))
```

```
[1] "Fancy T:   Numerator =   -431.583194999368 Denominator =   729.19417767533"
```

```
tval <- numer/denom
print("T ratio is")
```

```
[1] "T ratio is"
```

```
tval
```

```
difference
-0.5918632
```

```
print("The two-tailed test would have p value")
```

```
[1] "The two-tailed test would have p value"
```

```
2 * pt(abs(tval), df = m2all$df, lower.tail = FALSE)
```

```
difference
0.5542135
```

Could I make a function that "just" gets that right and would I be damaging students by ruining their educational experience? This would be very easy if the output had the variable names "pprof" and "pnet", but because I've made them factors, they are now pprofYES and pnetYES, and thus either my function has to be clever or the user's have to be clever in naming their request.

```
fancyT <- function(model, parm1, parm2){
    mc <- coef(model)
    mv <- vcov(model)
    numer <- mc[parm1] - mc[parm2]
    denom <- sqrt(mv[parm1, parm1]
        + mv[parm2, parm2] - 2 * mv[parm1, parm2])
    tval <- numer/denom
    tdf <- model$df
    tvalp <- 2 * pt(abs(tval), df = tdf, lower.tail = FALSE)
    res <- c(numer, denom, tval, tdf, tvalp)
    names(res) <- c("parm1 - parm2", "SE(parm1 - parm2)", "T", "df", "p-value")
    res
}
fancyT(m2all, parm1 = "pprofYES", parm2 = "pnetYES")
```

```
  parm1 - parm2 SE(parm1 - parm2)              T             df         p-value
   -431.5831950       729.1941777     -0.5918632    493.0000000       0.5542135
```

```
m2all <- lm(sal2 ~ sat + act + ibs + major + pprof + pnet + gender, data = dat)
m2alldf <- model.frame(m2all)
m2small <- lm(sal2 ~ sat + act + ibs, data = m2alldf)
anova(m2small, m2all)
```

```
Analysis of Variance Table

Model 1: sal2 ~ sat + act + ibs
Model 2: sal2 ~ sat + act + ibs + major + pprof + pnet + gender
  Res.Df        RSS Df  Sum of Sq      F    Pr(>F)
1    498 14861765797
2    493 12678628362  5 2183137434 16.978 1.679e-15 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

# Nonlinear

```
nm1 <- lm(sal3 ~ harv + gender + major + pprof + pnet, data = dat)
nm2 <- lm(sal3 ~ log(harv) + gender + major + pprof + pnet, data = dat)
nm3 <- lm(sal3 ~ harv + I(harv*harv) + gender + major + pprof + pnet, data = dat)
library(rockchalk)
nd <- rockchalk::newdata(nm1, predVals = list(harv = plotSeq(dat$harv, 20)))
nd$m1fit <- predict(nm1, newdata = nd)
nd$m2fit <- predict(nm2, newdata = nd)
nd$m3fit <- predict(nm3, newdata = nd)
```

For the regression table, please see Table 4
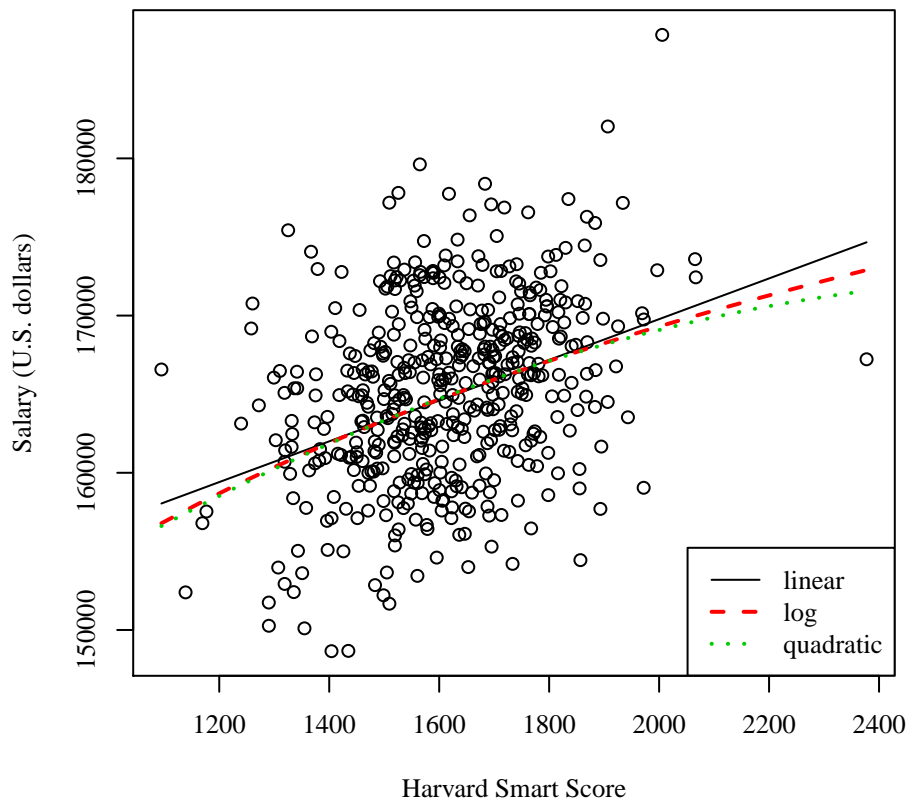
Table 4: Regression with sal3: Student-18

| | Linear Estimate (S.E.) | Log Estimate (S.E.) | Quadratic Estimate (S.E.) |
|---|---|---|---|
| (Intercept) | 141613.567* | 9168.805 | 127002.199* |
| | (2275.276) | (16229.898) | (13534.282) |
| Harvard SS | 12.954* | . | 31.131 |
| | (1.378) | | (16.654) |
| Gender: Male | 186.798 | 202.021 | 204.006 |
| | (447.741) | (447.365) | (447.924) |
| Major: Soc. | 2249.916* | 2243.639* | 2247.918* |
| | (534.879) | (534.4) | (534.771) |
| Major: Nat. | 5529.609* | 5537.845* | 5540.753* |
| | (557.047) | (556.576) | (557.024) |
| Prof. Parents: Yes | 1441.594* | 1471.28* | 1485.426* |
| | (491.466) | (491.081) | (492.991) |
| Parent Network: Yes | -662.997 | -649.854 | -651.272 |
| | (504.01) | (503.465) | (504.019) |
| ln(Harvard SS) | . | 20773.065* | . |
| | | (2197.639) | |
| Harvard SS$^2$ | . | . | -0.006 |
| | | | (0.005) |
| N | 485 | 485 | 485 |
| RMSE | 4906.704 | 4902.299 | 4905.681 |
| $R^2$ | 0.288 | 0.289 | 0.289 |
| adj $R^2$ | 0.279 | 0.28 | 0.279 |

$*p \leq 0.05$

```
outreg(list(nm1, nm2, nm3), tight = TRUE, title = paste("Regression with sal3: Student-",i,
    sep=""), modelLabels = c("Linear", "Log", "Quadratic"), varLabels = niceLabels, label
    = "table4")
```

```
plot(sal3 ~ harv, data = dat, xlab = "Harvard Smart Score", ylab = "Salary (U.S. dollars)")
lines(m1fit ~ harv, data = nd, lty = 1, col = 1)
lines(m2fit ~ harv, data = nd, lty = 2, col = 2, lwd = 2)
lines(m3fit ~ harv, data = nd, lty = 3, col = 3, lwd = 2)
legend("bottomright", legend = c("linear", "log", "quadratic"), lty = c(1,2,3), col = c
    (1,2,3), lwd = c(1,2,2))
```

```
cm1 <- lm(sal2 ~ major, data = dat)
dat$major2  <- relevel(dat$major, ref = "S")
cm2 <- lm(sal2 ~ major2, data = dat)
cm3 <- lm(sal2 ~ sat + act + ibs + major + pprof + pnet + gender, data = dat)
cm4 <- lm(sal2 ~ sat + act + ibs + major2 + pprof + pnet + gender, data = dat)
```

```
outreg(list(cm1, cm2, cm3, cm4), tight = TRUE, title = paste("Categorical Regressions:
    Student-", i, sep=""), modelLabels = c("major", "major2", "major full", "major2 full"),
    varLabels = niceLabels)
```

```
predictOMatic(cm1)
```

```
$major
            fit  major
S (40%) 22689.10      S
H (30%) 20951.66      H
N (30%) 25786.66      N

attr(,"flnames")
[1] "major"
```

```
predictOMatic(cm2)
```

```
$major2
            fit  major2
S (40%) 22689.10       S
H (30%) 20951.66       H
N (30%) 25786.66       N

attr(,"flnames")
[1] "major2"
```

Table 5: Categorical Regressions: Student-18

| | major Estimate (S.E.) | major2 Estimate (S.E.) | major full Estimate (S.E.) | major2 full Estimate (S.E.) |
|---|---|---|---|---|
| (Intercept) | 20951.661* (412.493) | 22689.098* (411.406) | -4835.422 (2753.496) | -3325.024 (2747.485) |
| Major: Soc. | 1737.437* (582.585) | . | 1510.397* (546.414) | . |
| Major: Nat. | 4835.003* (607.173) | . | 4943.763* (564.78) | . |
| Major 2: Hum. | . | -1737.437* (582.585) | . | -1510.397* (546.414) |
| Major 2: Nat. | . | 3097.566* (606.435) | . | 3433.366* (565.267) |
| SAT | . | . | 11.877* (1.613) | 11.877* (1.613) |
| ACT | . | . | 128.346* (50.748) | 128.346* (50.748) |
| Iowa BS | . | . | 39.298 (27.151) | 39.298 (27.151) |
| Prof. Parents: Yes | . | . | 560.202 (492.372) | 560.202 (492.372) |
| Parent Network: Yes | . | . | 991.785 (516.592) | 991.785 (516.592) |
| Gender: Male | . | . | -474.161 (455.713) | -474.161 (455.713) |
| N | 541 | 541 | 502 | 502 |
| RMSE | 5670.838 | 5670.838 | 5071.223 | 5071.223 |
| $R^2$ | 0.107 | 0.107 | 0.296 | 0.296 |
| adj $R^2$ | 0.104 | 0.104 | 0.285 | 0.285 |

$*p \leq 0.05$