

Data Management

```
library(foreign)
library(rockchalk)
i <- 31
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%	7.28	1182.00	67.01	5383.00	6654.00	147600.00	1163.00
25%	18.19	1512.00	93.58	16750.00	19610.00	161400.00	1483.00
50%	21.94	1622.00	99.92	20290.00	23260.00	165500.00	1594.00
75%	25.60	1717.00	106.80	24130.00	27370.00	169500.00	1696.00
100%	38.24	2088.00	130.50	36370.00	40010.00	182700.00	2052.00
mean	21.90	1620.00	100.10	20360.00	23480.00	165500.00	1596.00
sd	5.12	160.30	10.23	5551.00	5880.00	6024.00	159.00
var	26.21	25680.00	104.60	30820000.00	34570000.00	36280000.00	25270.00
NA's	11.00	52.00	0.00	13.00	0.00	0.00	32.00
N	518.00	518.00	518.00	518.00	518.00	518.00	518.00

Table 1: Best Automatic Summary Table for Numerics

Let students figure way to beautify this:

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

	gender		major		pnet
M	:271.0000	N	:190.0000	NO	:358.0000
F	:247.0000	S	:173.0000	YES	:160.0000
NA's	: 0.0000	H	:155.0000	NA's	: 0.0000
entropy	: 0.9985	NA's	: 0.0000	entropy	: 0.8919
normedEntropy	: 0.9985	entropy	: 1.5800	normedEntropy	: 0.8919
N	:518.0000	normedEntropy	: 0.9969	N	:518.0000
		N	:518.0000		
	pprof				
NO	:362.0000				
YES	:156.0000				
NA's	: 0.0000				
entropy	: 0.8827				
normedEntropy	: 0.8827				
N	:518.0000				

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 $\text{harv} = \text{sat} + \text{act}$).

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

```
mcDiagnose(m1all)
```

The following auxiliary models are being estimated and returned in a list:

```
sat ~ act + ibs + harv
<environment: 0x2d334e8>
act ~ sat + ibs + harv
<environment: 0x2d334e8>
ibs ~ sat + act + harv
<environment: 0x2d334e8>
harv ~ sat + act + ibs
<environment: 0x2d334e8>
Drum roll please!
```

And your R_j Squareds are (auxiliary Rsq)

```
      sat      act      ibs      harv
0.9998438 0.8741502 0.2167947 0.9998478
The Corresponding VIF, 1/(1-Rj2)
      sat      act      ibs      harv
6402.349159  7.945979  1.276804 6571.372296
```

Bivariate Correlations for design matrix

```
      sat  act  ibs  harv
sat  1.00 0.40 0.39 1.00
act  0.40 1.00 0.39 0.42
ibs  0.39 0.39 1.00 0.40
harv 1.00 0.42 0.40 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 SS2", 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 sall: 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(sall ~ sat + act, data = m1alldf)
anova(m1restricted, m1all)
```

Analysis of Variance Table

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

Table 2: Regression with sall: Student-31

	SAT	ACT	IBS	Harvard SS	All	Best
	Estimate	Estimate	Estimate	Estimate	Estimate	Estimate
	(S.E.)	(S.E.)	(S.E.)	(S.E.)	(S.E.)	(S.E.)
(Intercept)	-66.568 (2419.766)	12284.853* (1028.083)	11966.483* (2421.874)	-676.268 (2446.461)	2421.766 (3001.622)	2388.155 (2849.656)
SAT	12.811* (1.51)	.	.	.	-51.307 (125.405)	10.421* (1.701)
ACT	.	369.729* (45.785)	.	.	253.526 (133.19)	280.541* (51.828)
Iowa BS	.	.	83.836* (24.063)	.	-57.283* (27.415)	-47.212 (26.34)
Harvard SS	.	.	.	12.956* (1.503)	61.819 (125.346)	.
N	473	494	505	455	415	463
RMSE	5206.071	5220.714	5490.829	5158.82	5034.047	5076.177
R^2	0.133	0.117	0.024	0.141	0.205	0.185
adj R^2	0.131	0.115	0.022	0.139	0.198	0.18

* $p \leq 0.05$

	Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
1	412	1.0508e+10				
2	410	1.0390e+10	2	117746616	2.3232	0.09925

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

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(sall ~ sat + act + ibs, data = dat)
dat2 <- model.frame(m1best)
m1s <- lm(sall ~ sat, data = dat2)
m1a <- lm(sall ~ act, data = dat2)
m1i <- lm(sall ~ ibs, data = dat2)
m1h <- lm(sall ~ harv, data = dat[row.names(dat2), ])
m1all <- lm(sall ~ 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	ACT	IBS	Harvard SS	All	Best
	Estimate	Estimate	Estimate	Estimate	Estimate	Estimate
	(S.E.)	(S.E.)	(S.E.)	(S.E.)	(S.E.)	(S.E.)
(Intercept)	-186.701 (2458.675)	12243.981* (1063.199)	13111.993* (2567.197)	-1156.941 (2601.581)	2421.766 (3001.622)	2388.155 (2849.656)
SAT	12.91* (1.534)	.	.	.	-51.307 (125.405)	10.421* (1.701)
ACT	.	373.746* (47.416)	.	.	253.526 (133.19)	280.541* (51.828)
Iowa BS	.	.	72.771* (25.506)	.	-57.283* (27.415)	-47.212 (26.34)
Harvard SS	.	.	.	13.291* (1.6)	61.819 (125.346)	.
N	463	463	463	415	415	463
RMSE	5225.001	5268.084	5562.971	5207.867	5034.047	5076.177
R^2	0.133	0.119	0.017	0.143	0.205	0.185
adj R^2	0.131	0.117	0.015	0.141	0.198	0.18

* $p \leq 0.05$

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)
```

```

      sall
sall -1.00000000
sat  0.27487447
act  0.24495782
ibs  -0.08336934

```

```
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  0.066582198
act  0.052002794
ibs  0.005701797

```

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%         35.09    42.16   35.98
50%         46.93    51.94   48.30
75%         59.06    62.66   59.89
100%        100.00   100.00  100.00

```

```

mean    46.97    52.15    48.54
sd      16.70    15.97    17.83
var     278.80   255.10   317.80
NA's    0.00     0.00     0.00
N       463.00   463.00   463.00

```

```

$ factors
NULL

```

```

mlpoms <- lm(sall ~ satpoms + actpoms + ibspoms, data = dat2)
summary(mlpoms)

```

```

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

```

```

Residuals:
    Min       1Q   Median       3Q      Max
-16823.6  -3392.4    94.3   3636.1  14643.3

```

```

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) 13388.49    935.98  14.304 < 2e-16 ***
satpoms       92.62     15.12   6.125 1.95e-09 ***
actpoms       86.86     16.05   5.413 1.00e-07 ***
ibspoms      -29.99     16.73  -1.792  0.0737 .

```

```

Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

```

Residual standard error: 5076 on 459 degrees of freedom
Multiple R2: 0.1854, Adjusted R2: 0.18
F-statistic: 34.81 on 3 and 459 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(mlall)

```

```

           sall
sall -1.00000000
sat  -0.02020139
act   0.09359368
ibs  -0.10264769
harv  0.02434952

```

```

getDeltaRsquare(mlall)

```

```

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  0.0003244579
act  0.0070231607
ibs  0.0084628578
harv 0.0004714731

```

```

options(scipen = 5)

```

Additional Variables

Please see Table 3 for the regressions.

Table 3: Regression with sal2: Student-31

	Test Scores Only	All Predictors
	Estimate	Estimate
	(S.E.)	(S.E.)
(Intercept)	3170.82 (2988.427)	1152.089 (2850.448)
SAT	10.82* (1.81)	10.625* (1.708)
ACT	272.674* (55.148)	265.213* (52.002)
Iowa BS	-28.451 (27.796)	-33.645 (26.214)
Major: Soc.	.	1180.716* (596.195)
Major: Nat.	.	4277.163* (580.528)
Prof. Parents: Yes	.	1163.723* (512.782)
Parent Network: Yes	.	966.662 (510.758)
Gender: Male	.	730.027 (470.959)
N	476	476
RMSE	5442.348	5106.048
R^2	0.176	0.282
adj R^2	0.17	0.27

* $p \leq 0.05$

```
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 = 197.061419359701 Denominator = 731.016848312413"
```

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

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

```
tval
```

```
difference
0.2695717
```

```
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
0.7876089
```

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
197.0614194	731.0168483	0.2695717	467.0000000	0.7876089

```
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     472 13980237277
2     467 12175494529  5 1804742748 13.844 1.261e-12 ***
```

```
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

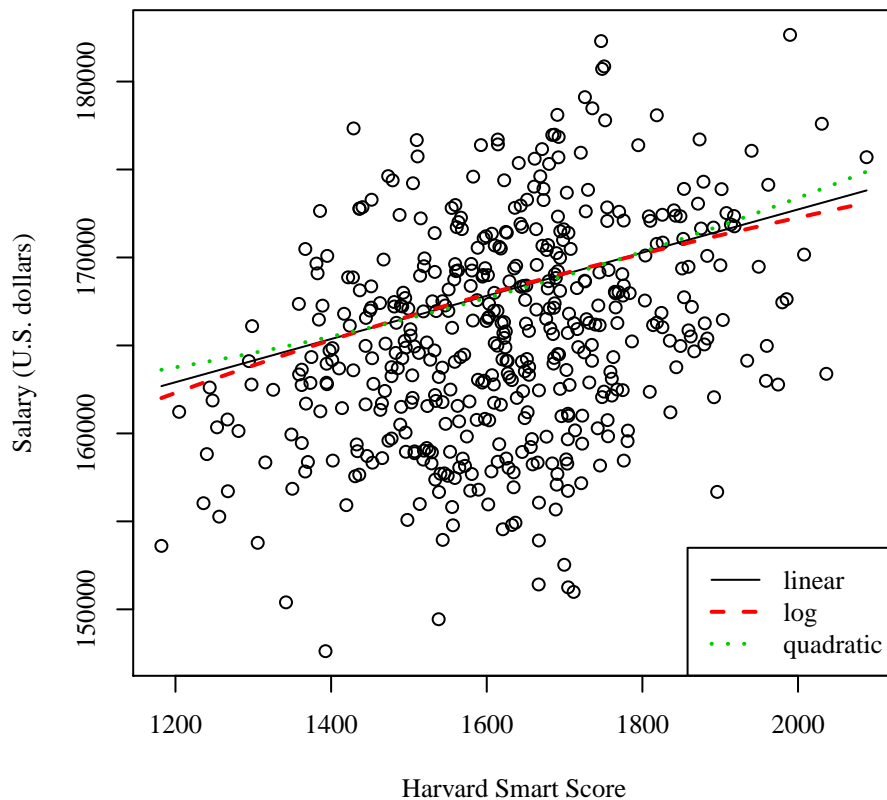
Table 4: Regression with sal3: Student-31

	Linear Estimate (S.E.)	Log Estimate (S.E.)	Quadratic Estimate (S.E.)
(Intercept)	143181.145* (2478.796)	19073.436 (17697.31)	157462.165* (17678.683)
Harvard SS	12.274* (1.495)	.	-5.552 (21.901)
Gender: Male	-762.35 (479.809)	-751.092 (480.342)	-775.583 (480.258)
Major: Soc.	1677.134* (607.353)	1642.847* (608.126)	1724.661* (610.361)
Major: Nat.	5762.184* (589.757)	5727.007* (590.491)	5811.393* (593.047)
Prof. Parents: Yes	149.417 (522.712)	156.441 (523.295)	143.01 (522.961)
Parent Network: Yes	519.752 (518.715)	512.123 (519.285)	526.755 (518.975)
ln(Harvard SS)	.	19499.884* (2396.536)	.
Harvard SS ²	.	.	0.005 (0.007)
N	466	466	466
RMSE	5144.553	5150.622	5146.428
R^2	0.277	0.275	0.278
adj R^2	0.268	0.266	0.267

* $p \leq 0.05$

```
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
N (40%) 25714.72  N
S (30%) 22665.99  S
H (30%) 21647.45  H

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

```
predictOMatic(cm2)
```

```
$major2
      fit major2
N (40%) 25714.72  N
S (30%) 22665.99  S
H (30%) 21647.45  H

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

Table 5: Categorical Regressions: Student-31

	major	major2	major full	major2 full
	Estimate	Estimate	Estimate	Estimate
	(S.E.)	(S.E.)	(S.E.)	(S.E.)
(Intercept)	21647.451*	22665.994*	1152.089	2332.805
	(451.753)	(427.606)	(2850.448)	(2882.48)
Major: Soc.	1018.542	.	1180.716*	.
	(622.035)		(596.195)	
Major: Nat.	4067.271*	.	4277.163*	.
	(608.742)		(580.528)	
Major 2: Hum.	.	-1018.542	.	-1180.716*
		(622.035)		(596.195)
Major 2: Nat.	.	3048.729*	.	3096.447*
		(591.044)		(560.52)
SAT	.	.	10.625*	10.625*
			(1.708)	(1.708)
ACT	.	.	265.213*	265.213*
			(52.002)	(52.002)
Iowa BS	.	.	-33.645	-33.645
			(26.214)	(26.214)
Prof. Parents: Yes	.	.	1163.723*	1163.723*
			(512.782)	(512.782)
Parent Network: Yes	.	.	966.662	966.662
			(510.758)	(510.758)
Gender: Male	.	.	730.027	730.027
			(470.959)	(470.959)
N	518	518	476	476
RMSE	5624.276	5624.276	5106.048	5106.048
R^2	0.089	0.089	0.282	0.282
adj R^2	0.085	0.085	0.27	0.27

* $p \leq 0.05$