

Data Management

```
library(foreign)
library(rockchalk)
i <- 8
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%	9.09	1069.00	74.34	5029.00	7886.00	151000.00	1057.00
25%	18.84	1518.00	93.78	16720.00	19300.00	161600.00	1497.00
50%	22.28	1625.00	100.60	20590.00	23270.00	165200.00	1600.00
75%	25.32	1736.00	107.00	24300.00	27400.00	169100.00	1706.00
100%	36.25	2090.00	131.10	39480.00	42930.00	182300.00	2075.00
mean	22.20	1625.00	100.30	20530.00	23430.00	165400.00	1601.00
sd	4.85	159.10	9.89	5616.00	6182.00	5350.00	159.60
var	23.52	25310.00	97.91	31540000.00	38210000.00	28620000.00	25480.00
NA's	15.00	56.00	0.00	10.00	0.00	0.00	16.00
N	555.00	555.00	555.00	555.00	555.00	555.00	555.00

Table 1: Best Automatic Summary Table for Numerics

Let students figure way to beautify this:

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

	gender		major		pnet
M	:284.0000	S	:188.0000	NO	:405.0000
F	:271.0000	H	:185.0000	YES	:150.0000
NA's	: 0.0000	N	:182.0000	NA's	: 0.0000
entropy	: 0.9996	NA's	: 0.0000	entropy	: 0.8419
normedEntropy	: 0.9996	entropy	: 1.5848	normedEntropy	: 0.8419
N	:555.0000	normedEntropy	: 0.9999	N	:555.0000
		N	:555.0000		
	pprof				
NO	:379.0000				
YES	:176.0000				
NA's	: 0.0000				
entropy	: 0.9012				
normedEntropy	: 0.9012				
N	:555.0000				

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: 0x14fec28>
act ~ sat + ibs + harv
<environment: 0x14fec28>
ibs ~ sat + act + harv
<environment: 0x14fec28>
harv ~ sat + act + ibs
<environment: 0x14fec28>
Drum roll please!
```

And your R_j Squareds are (auxiliary Rsq)

```
      sat      act      ibs      harv
0.9998391 0.8583910 0.2319107 0.9998436
The Corresponding VIF, 1/(1-Rj2)
      sat      act      ibs      harv
6215.086662   7.061699   1.301932 6394.060127
```

Bivariate Correlations for design matrix

```
      sat  act  ibs  harv
sat  1.00 0.45 0.39 1.00
act  0.45 1.00 0.43 0.48
ibs  0.39 0.43 1.00 0.40
harv 1.00 0.48 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 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 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-8

	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)	-4187.133 (2201.132)	12235.622* (1088.074)	11025.407* (2414.351)	-3855.203 (2369.098)	-2509.684 (2892.831)	-2619.501 (2706.268)
SAT	15.476* (1.368)	.	.	.	-21.123 (117.231)	13.798* (1.581)
ACT	.	372.9* (47.822)	.	.	182.242 (127.873)	198.474* (52.668)
Iowa BS	.	.	94.843* (23.968)	.	-34.138 (26.739)	-32.669 (25.019)
Harvard SS	.	.	.	15.096* (1.45)	34.763 (117.277)	.
N	529	531	545	490	464	515
RMSE	5024.552	5362.257	5541.948	5102.832	5026.617	5002.909
R^2	0.195	0.103	0.028	0.182	0.218	0.217
adj R^2	0.194	0.101	0.026	0.18	0.211	0.213

* $p \leq 0.05$

	Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
1	461	1.1640e+10				
2	459	1.1597e+10	2	42682529	0.8446	0.4304

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)	-4378.138 (2248.372)	12433.338* (1110.246)	11543.804* (2495.942)	-4824.038* (2419.17)	-2509.684 (2892.831)	-2619.501 (2706.268)
SAT	15.598* (1.398)	.	.	.	-21.123 (117.231)	13.798* (1.581)
ACT	.	367.418* (48.912)	.	.	182.242 (127.873)	198.474* (52.668)
Iowa BS	.	.	90.03* (24.74)	.	-34.138 (26.739)	-32.669 (25.019)
Harvard SS	.	.	.	15.747* (1.482)	34.763 (117.277)	.
N	515	515	515	464	464	515
RMSE	5062.384	5356.326	5571.768	5078.571	5026.617	5002.909
R^2	0.195	0.099	0.025	0.196	0.218	0.217
adj R^2	0.194	0.097	0.023	0.195	0.211	0.213

* $p \leq 0.05$

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)
```

```

      s all
s all -1.00000000
s at  0.36014490
a ct  0.16443389
i bs  -0.05766811

```

```
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
s at  0.116676032
a ct  0.021756093
i bs  0.002612234

```

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%         36.40    34.82   43.21
50%         49.44    46.32   53.36
75%         60.46    57.80   64.01
100%        100.00   100.00  100.00

```

```

mean  48.92  45.94  53.37
sd    18.05  17.51  15.68
var   325.80 306.70 246.00
NA's   0.00   0.00   0.00
N     515.00 515.00 515.00

```

```

$ 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
-11642.9  -3432.4  -136.2   3167.3  16104.6

```

```

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) 11336.49    874.10  12.969 < 2e-16 ***
satpoms      140.51     16.10   8.727 < 2e-16 ***
actpoms       53.11     14.09   3.768 0.000183 ***
ibspoms     -18.53     14.19  -1.306 0.192215

```

```

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

```

```

Residual standard error: 5003 on 511 degrees of freedom
Multiple R2: 0.2171, Adjusted R2: 0.2125
F-statistic: 47.24 on 3 and 511 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 -1.000000000
sat  -0.008409815
act   0.066375064
ibs  -0.059486629
harv  0.013834131

```

```

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  0.00005531884
act  0.00346096248
ibs  0.00277745909
harv 0.00014971189

```

```

options(scipen = 5)

```

Additional Variables

Please see Table 3 for the regressions.

Table 3: Regression with sal2: Student-8

	Test Scores Only	All Predictors
	Estimate	Estimate
	(S.E.)	(S.E.)
(Intercept)	222.867 (3062.626)	-3163.449 (2706.909)
SAT	12.656* (1.785)	13.546* (1.565)
ACT	194.81* (59.645)	226.006* (52.471)
Iowa BS	-13.294 (28.345)	-36.328 (24.912)
Major: Soc.	.	2534.789* (535.921)
Major: Nat.	.	6697.052* (535.529)
Prof. Parents: Yes	.	704.716 (478.052)
Parent Network: Yes	.	788.328 (493.875)
Gender: Male	.	236.346 (438.054)
N	524	524
RMSE	5703.265	4990.373
R^2	0.163	0.365
adj R^2	0.158	0.355

* $p \leq 0.05$

```
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 = -83.6118952038042 Denominator = 648.228315936162"
```

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

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

```
tval
```

```
difference
-0.1289853
```

```
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
0.8974197
```

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
-83.6118952	648.2283159	-0.1289853	515.0000000	0.8974197

```
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     520 16914162390
2     515 12825469846  5 4088692544 32.836 < 2.2e-16 ***
---
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

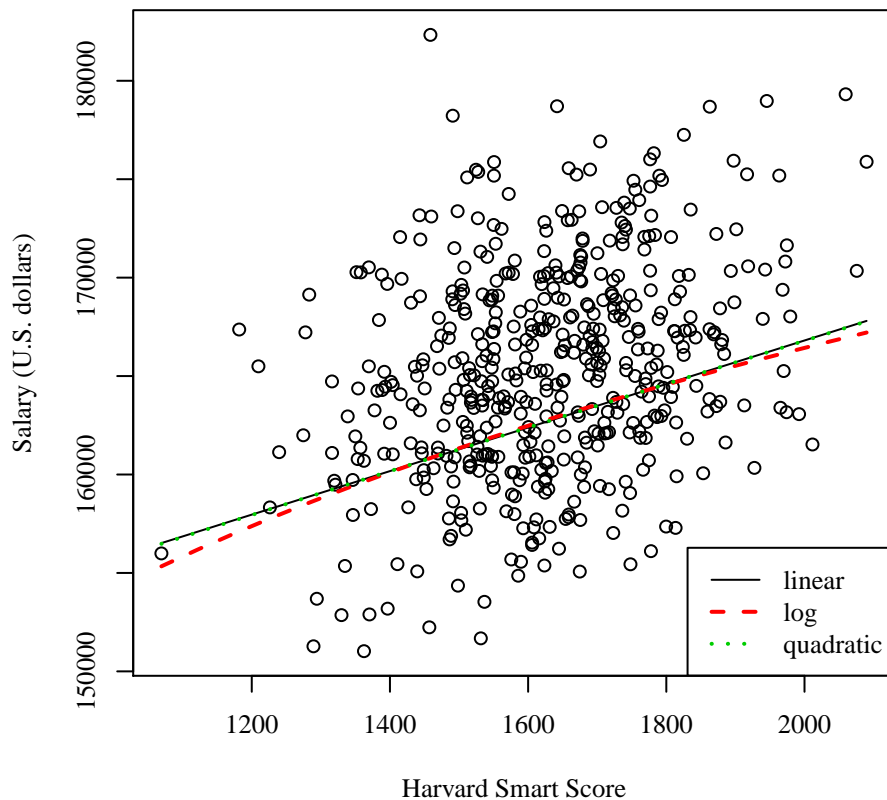
Table 4: Regression with sal3: Student-8

	Linear Estimate (S.E.)	Log Estimate (S.E.)	Quadratic Estimate (S.E.)
(Intercept)	144653.849* (2182.525)	31709.121* (15650.601)	144318.252* (15148.691)
Harvard SS	11.058* (1.32)	.	11.475 (18.67)
Gender: Male	37.762 (417.893)	42.563 (417.989)	37.874 (418.348)
Major: Soc.	2171.902* (509.696)	2177.717* (509.798)	2172.076* (510.274)
Major: Nat.	5221.248* (512.219)	5226.975* (512.382)	5221.498* (512.861)
Prof. Parents: Yes	1139.465* (451.731)	1129.346* (451.902)	1139.103* (452.479)
Parent Network: Yes	-311.979 (474.562)	-328.813 (474.819)	-312.636 (475.95)
ln(Harvard SS)	.	17719.017* (2118.732)	.
Harvard SS ²	.	.	0 (0.006)
N	499	499	499
RMSE	4652.611	4653.686	4657.344
R^2	0.26	0.259	0.26
adj R^2	0.251	0.25	0.249

* $p \leq 0.05$

```
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 (30%) 22953.78   S
H (30%) 20619.26   H
N (30%) 26764.37   N

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

```
predictOMatic(cm2)
```

```
$major2
      fit major2
S (30%) 22953.78   S
H (30%) 20619.26   H
N (30%) 26764.37   N

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

Table 5: Categorical Regressions: Student-8

	major	major2	major full	major2 full
	Estimate	Estimate	Estimate	Estimate
	(S.E.)	(S.E.)	(S.E.)	(S.E.)
(Intercept)	20619.262*	22953.782*	-3163.449	-628.66
	(415.638)	(412.308)	(2706.909)	(2727.109)
Major: Soc.	2334.52*	.	2534.789*	.
	(585.451)		(535.921)	
Major: Nat.	6145.111*	.	6697.052*	.
	(590.218)		(535.529)	
Major 2: Hum.	.	-2334.52*	.	-2534.789*
		(585.451)		(535.921)
Major 2: Nat.	.	3810.591*	.	4162.263*
		(587.878)		(539.376)
SAT	.	.	13.546*	13.546*
			(1.565)	(1.565)
ACT	.	.	226.006*	226.006*
			(52.471)	(52.471)
Iowa BS	.	.	-36.328	-36.328
			(24.912)	(24.912)
Prof. Parents: Yes	.	.	704.716	704.716
			(478.052)	(478.052)
Parent Network: Yes	.	.	788.328	788.328
			(493.875)	(493.875)
Gender: Male	.	.	236.346	236.346
			(438.054)	(438.054)
N	555	555	524	524
RMSE	5653.284	5653.284	4990.373	4990.373
R^2	0.167	0.167	0.365	0.365
adj R^2	0.164	0.164	0.355	0.355

* $p \leq 0.05$