

Data Management

```
library(foreign)
library(rockchalk)
i <- 26
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%	7.62	1086.00	73.61	2845.00	5831.00	144800.00	1077.00
25%	18.16	1502.00	93.27	17090.00	19790.00	161000.00	1481.00
50%	21.65	1615.00	99.24	20830.00	23700.00	164800.00	1596.00
75%	25.10	1726.00	105.80	24340.00	27450.00	169200.00	1700.00
100%	35.79	2052.00	133.60	35670.00	41580.00	185900.00	2029.00
mean	21.65	1613.00	99.27	20590.00	23470.00	165000.00	1591.00
sd	4.99	164.60	9.93	5492.00	6024.00	6158.00	158.20
var	24.89	27110.00	98.64	30160000.00	36290000.00	37920000.00	25010.00
NA's	12.00	68.00	0.00	10.00	0.00	0.00	32.00
N	558.00	558.00	558.00	558.00	558.00	558.00	558.00

Table 1: Best Automatic Summary Table for Numerics

Let students figure way to beautify this:

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

M	gender	:280	S	major	:190.0000	NO	pnet	:406.0000	NO	pprof	:392
		.0000									
F		:278	H		:187.0000	YES		:152.0000	YES		:166
		.0000									
NA's		: 0	N		:181.0000	NA's		: 0.0000	NA's		: 0
		.0000									
entropy		: 1	NA's		: 0.0000	entropy		: 0.8449	entropy		: 0
		.8782									
normedEntropy		: 1	entropy		: 1.5847	normedEntropy		: 0.8449	normedEntropy		: 0
		.8782									
N		:558	normedEntropy		: 0.9998	N		:558.0000	N		:558
		.0000									
			N		:558.0000						

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: 0x3285118>
act ~ sat + ibs + harv
<environment: 0x3285118>
ibs ~ sat + act + harv
<environment: 0x3285118>
harv ~ sat + act + ibs
<environment: 0x3285118>
Drum roll please!
```

And your R_j Squareds are (auxiliary Rsq)

```
      sat      act      ibs      harv
0.9998508 0.8730606 0.2681290 0.9998552
The Corresponding VIF, 1/(1-Rj2)
      sat      act      ibs      harv
6701.435771  7.877775  1.366361 6904.528134
```

Bivariate Correlations for design matrix

```
      sat  act  ibs harv
sat  1.00 0.46 0.42 1.00
act  0.46 1.00 0.45 0.48
ibs  0.42 0.45 1.00 0.43
harv 1.00 0.48 0.43 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-26

	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)	775.211 (2259.579)	13038.753* (1008.316)	7182.124* (2294.942)	1880.929 (2274.26)	-949.503 (2868.555)	-1290.199 (2764.824)
SAT	12.418* (1.414)	.	.	.	-63.211 (119.342)	9.547* (1.651)
ACT	.	349.214* (45.418)	.	.	83.11 (132.893)	144.633* (53.627)
Iowa BS	.	.	134.913* (22.982)	.	49.749 (28.077)	35.511 (26.566)
Harvard SS	.	.	.	11.652* (1.403)	71.59 (119.374)	.
N	516	536	548	481	446	505
RMSE	5065.087	5236.006	5330.924	5069.266	4977.479	5039.923
R^2	0.131	0.1	0.059	0.126	0.157	0.153
adj R^2	0.129	0.098	0.058	0.124	0.15	0.148

* $p \leq 0.05$

	Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
1	443	1.1018e+10				
2	441	1.0926e+10	2	92003020	1.8567	0.1574

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)	657.689 (2292.366)	13765.195* (1044.387)	7458.973* (2416.284)	1617.497 (2361.866)	-949.503 (2868.555)	-1290.199 (2764.824)
SAT	12.506* (1.434)	.	.	.	-63.211 (119.342)	9.547* (1.651)
ACT	.	314.101* (47.132)	.	.	83.11 (132.893)	144.633* (53.627)
Iowa BS	.	.	131.708* (24.194)	.	49.749 (28.077)	35.511 (26.566)
Harvard SS	.	.	.	11.814* (1.455)	71.59 (119.374)	.
N	505	505	505	446	446	505
RMSE	5092.699	5237.761	5309.925	5042.089	4977.479	5039.923
R^2	0.131	0.081	0.056	0.129	0.157	0.153
adj R^2	0.13	0.079	0.054	0.127	0.15	0.148

* $p \leq 0.05$

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  0.25008684
act  0.11962778
ibs  0.05961271

```

```
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 0.056533933
act 0.012302767
ibs 0.003022056

```

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%     37.34    32.90   42.51
50%     49.73    42.98   54.47
75%     61.95    53.54   65.28
100%    100.00   100.00  100.00

```

```

mean  49.63  42.95  53.97
sd    17.57  16.29  16.61
var   308.80 265.20 275.80
NA's   0.00   0.00   0.00
N     505.00 505.00 505.00

```

```

$ 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
-17412.6  -3348.4   156.9   3477.4  14198.3

```

```

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) 12705.13    864.95  14.689 < 2e-16 ***
satpoms       90.93     15.73   5.781  1.3e-08 ***
actpoms       40.74     15.11   2.697  0.00723 **
ibspoms       21.32     15.95   1.337  0.18193

```

```

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

```

```

Residual standard error: 5040 on 501 degrees of freedom
Multiple R2: 0.1526, Adjusted R2: 0.1475
F-statistic: 30.08 on 3 and 501 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  -0.02521380
act   0.02976725
ibs   0.08407633
harv  0.02854615

```

```

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  0.0005361715
act  0.0007475036
ibs  0.0060003876
harv 0.0006873848

```

```

options(scipen = 5)

```

Additional Variables

Please see Table 3 for the regressions.

Table 3: Regression with sal2: Student-26

	Test Scores Only	All Predictors
	Estimate	Estimate
	(S.E.)	(S.E.)
(Intercept)	1110.763 (3024.582)	-2189.352 (2758.461)
SAT	10.1* (1.81)	9.571* (1.642)
ACT	129.454* (58.829)	156.072* (53.906)
Iowa BS	34.75 (29.049)	40.884 (26.465)
Major: Soc.	.	1510.475* (540.643)
Major: Nat.	.	5242.979* (547.116)
Prof. Parents: Yes	.	2107.245* (487.799)
Parent Network: Yes	.	751.917 (496.534)
Gender: Male	.	-174.662 (449.719)
N	515	515
RMSE	5576.203	5030.607
R^2	0.131	0.3
adj R^2	0.126	0.289

* $p \leq 0.05$

```
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 = 1355.32839671917 Denominator = 673.815215517464"
```

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

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

```
tval
```

```
difference
2.011424
```

```
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
0.04481037
```

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
1355.32839672	673.81521552	2.01142445	506.00000000	0.04481037

```
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	511	15889055693				
2	506	12805346521	5	3083709172	24.37	< 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

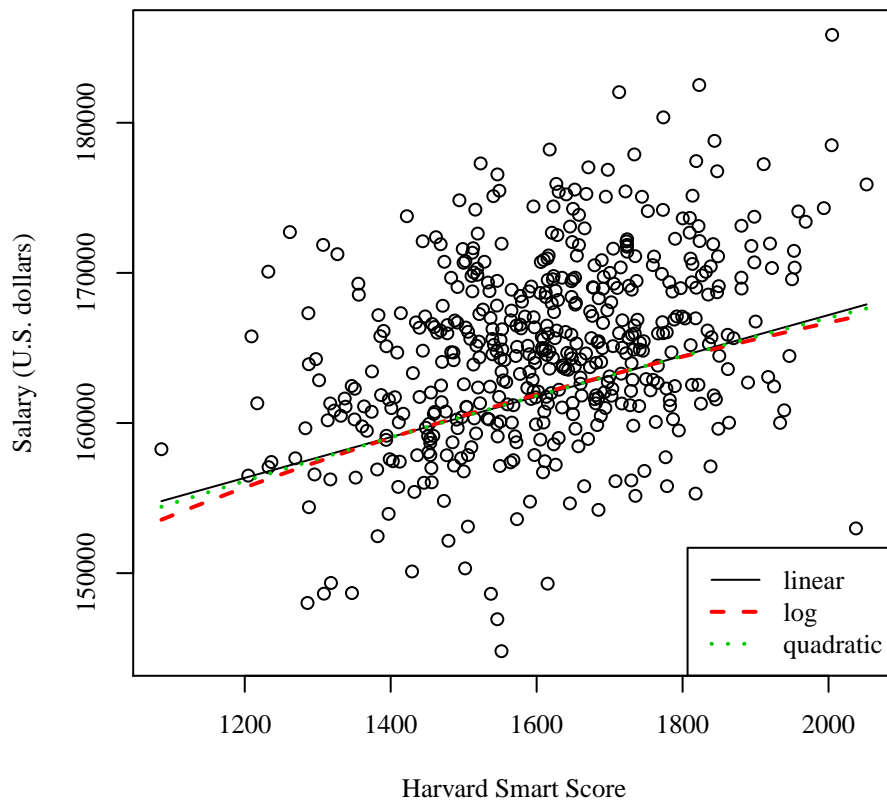
Table 4: Regression with sal3: Student-26

	Linear Estimate (S.E.)	Log Estimate (S.E.)	Quadratic Estimate (S.E.)
(Intercept)	139939.879* (2364.247)	3179.992 (16887.424)	136009.254* (16870.309)
Harvard SS	13.56* (1.443)	.	18.503 (21.054)
Gender: Male	133.569 (475.518)	143.811 (475.556)	136.284 (476.123)
Major: Soc.	2477.47* (579.364)	2466.349* (579.482)	2473.621* (580.162)
Major: Nat.	5312.386* (582.173)	5302.937* (582.262)	5308.818* (582.94)
Prof. Parents: Yes	714.207 (517.715)	717.015 (517.799)	716.114 (518.285)
Parent Network: Yes	1254.321* (532.951)	1276.141* (532.927)	1261.991* (534.468)
ln(Harvard SS)	.	21492.506* (2289.76)	.
Harvard SS ²	.	.	-0.002 (0.007)
N	490	490	490
RMSE	5240.665	5241.556	5245.797
R^2	0.274	0.274	0.275
adj R^2	0.265	0.265	0.264

* $p \leq 0.05$

```
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%) 23123.31  S
H (30%) 20974.57  H
N (30%) 26419.56  N

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

```
predictOMatic(cm2)
```

```
$major2
      fit major2
S (30%) 23123.31  S
H (30%) 20974.57  H
N (30%) 26419.56  N

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

Table 5: Categorical Regressions: Student-26

	major	major2	major full	major2 full
	Estimate	Estimate	Estimate	Estimate
	(S.E.)	(S.E.)	(S.E.)	(S.E.)
(Intercept)	20974.574*	23123.311*	-2189.352	-678.877
	(410.057)	(406.807)	(2758.461)	(2752.008)
Major: Soc.	2148.737*	.	1510.475*	.
	(577.615)		(540.643)	
Major: Nat.	5444.986*	.	5242.979*	.
	(584.695)		(547.116)	
Major 2: Hum.	.	-2148.737*	.	-1510.475*
		(577.615)		(540.643)
Major 2: Nat.	.	3296.249*	.	3732.503*
		(582.42)		(549.431)
SAT	.	.	9.571*	9.571*
			(1.642)	(1.642)
ACT	.	.	156.072*	156.072*
			(53.906)	(53.906)
Iowa BS	.	.	40.884	40.884
			(26.465)	(26.465)
Prof. Parents: Yes	.	.	2107.245*	2107.245*
			(487.799)	(487.799)
Parent Network: Yes	.	.	751.917	751.917
			(496.534)	(496.534)
Gender: Male	.	.	-174.662	-174.662
			(449.719)	(449.719)
N	558	558	515	515
RMSE	5607.448	5607.448	5030.607	5030.607
R^2	0.137	0.137	0.3	0.3
adj R^2	0.134	0.134	0.289	0.289

* $p \leq 0.05$