

## Data Management

```
library(foreign)
library(rockchalk)
i <- 29
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%	4.65	969.30	72.75	4849.00	5747.00	149800.00	944.10
25%	18.20	1509.00	92.76	16840.00	19440.00	161900.00	1487.00
50%	22.09	1614.00	100.40	20840.00	23390.00	165300.00	1584.00
75%	25.86	1722.00	107.30	23860.00	27250.00	169700.00	1699.00
100%	39.09	2154.00	130.30	36510.00	42330.00	183800.00	2121.00
mean	21.99	1615.00	100.40	20540.00	23340.00	165600.00	1591.00
sd	5.33	164.20	10.11	5228.00	5754.00	5673.00	161.10
var	28.38	26960.00	102.20	27330000.00	33110000.00	32190000.00	25950.00
NA's	15.00	51.00	0.00	15.00	0.00	0.00	22.00
N	542.00	542.00	542.00	542.00	542.00	542.00	542.00

Table 1: Best Automatic Summary Table for Numerics

Let students figure way to beautify this:

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

	<b>gender</b>		<b>major</b>		<b>pnet</b>
M	:276.0000	H	:185.0000	NO	:387.0000
F	:266.0000	S	:182.0000	YES	:155.0000
NA's	: 0.0000	N	:175.0000	NA's	: 0.0000
entropy	: 0.9998	NA's	: 0.0000	entropy	: 0.8635
normedEntropy	: 0.9998	entropy	: 1.5846	normedEntropy	: 0.8635
N	:542.0000	normedEntropy	: 0.9998	N	:542.0000
		N	:542.0000		
	<b>pprof</b>				
NO	:378.0000				
YES	:164.0000				
NA's	: 0.0000				
entropy	: 0.8844				
normedEntropy	: 0.8844				
N	:542.0000				

## 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: 0x2e680a0>
act ~ sat + ibs + harv
<environment: 0x2e680a0>
ibs ~ sat + act + harv
<environment: 0x2e680a0>
harv ~ sat + act + ibs
<environment: 0x2e680a0>
Drum roll please!
```

And your R<sub>j</sub> Squareds are (auxiliary Rsq)

```
      sat      act      ibs      harv
0.9998455 0.8796857 0.2457630 0.9998494
The Corresponding VIF, 1/(1-Rj2)
      sat      act      ibs      harv
6472.534373  8.311563  1.325843 6641.290365
```

Bivariate Correlations for design matrix

```
      sat  act  ibs  harv
sat  1.00 0.36 0.43 1.00
act  0.36 1.00 0.37 0.39
ibs  0.43 0.37 1.00 0.44
harv 1.00 0.39 0.44 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-29

	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)	-947.028 (2154.676)	13727.262* (922.381)	6597.593* (2188.298)	-1044.075 (2173.226)	-3763.583 (2593.184)	-3052.82 (2524.572)
SAT	13.515* (1.348)	.	.	.	252.213* (110.23)	10.704* (1.533)
ACT	.	309.287* (40.734)	.	.	404.707* (118.789)	177.403* (43.75)
Iowa BS	.	.	138.832* (21.681)	.	32.145 (24.764)	26.397 (24.111)
Harvard SS	.	.	.	13.367* (1.339)	-241.255* (110.271)	.
N	505	512	527	478	443	490
RMSE	4822.563	4939.025	5039.816	4748.09	4647.833	4705.028
$R^2$	0.167	0.102	0.072	0.173	0.219	0.203
adj $R^2$	0.165	0.1	0.071	0.171	0.212	0.199

\* $p \leq 0.05$

	Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
1	440	9590139876				
2	438	9461830549	2	128309326	2.9698	0.05235

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)	-1104.155 (2186.628)	13740.387* (948.822)	7040.679* (2290.375)	-1612.101 (2239.934)	-3763.583 (2593.184)	-3052.82 (2524.572)
SAT	13.6* (1.367)	.	.	.	252.213* (110.23)	10.704* (1.533)
ACT	.	308.89* (41.882)	.	.	404.707* (118.789)	177.403* (43.75)
Iowa BS	.	.	134.257* (22.668)	.	32.145 (24.764)	26.397 (24.111)
Harvard SS	.	.	.	13.713* (1.379)	-241.255* (110.271)	.
N	490	490	490	443	443	490
RMSE	4797.091	4990.153	5081.447	4737.473	4647.833	4705.028
$R^2$	0.169	0.1	0.067	0.183	0.219	0.203
adj $R^2$	0.167	0.098	0.065	0.181	0.212	0.199

\* $p \leq 0.05$

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  0.30190956
act  0.18089933
ibs  0.04960104

```

```
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  0.07988763
act  0.02694896
ibs  0.00196458

```

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%     39.31   35.05   46.12
50%     50.74   48.24   54.39
75%     61.45   60.21   64.02
100%    100.00  100.00  100.00

```

```

mean  50.39  48.28  54.99
sd    15.64  17.62  13.48
var   244.80 310.40 181.70
NA's  0.00   0.00   0.00
N     490.00 490.00 490.00

```

```

$ 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
-14179.1  -3220.2   160.8   2931.5  14694.6

```

```

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  9797.89    987.91   9.918 < 2e-16 ***
satpoms       125.98     18.04   6.982 9.63e-12 ***
actpoms        61.10     15.07   4.055 5.84e-05 ***
ibspoms        15.19     13.87   1.095  0.274

```

```

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

```

```

Residual standard error: 4705 on 486 degrees of freedom
Multiple R2: 0.2034, Adjusted R2: 0.1985
F-statistic: 41.37 on 3 and 486 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  0.10868035
act  0.16067492
ibs  0.06190421
harv -0.10397257

```

```

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  0.009334056
act  0.020694909
ibs  0.003004112
harv 0.008534265

```

```

options(scipen = 5)

```

## Additional Variables

Please see Table 3 for the regressions.

Table 3: Regression with sal2: Student-29

	Test Scores Only	All Predictors
	Estimate	Estimate
	(S.E.)	(S.E.)
(Intercept)	-645.459 (2741.472)	-2559.696 (2535.688)
SAT	10.739* (1.653)	10.048* (1.506)
ACT	206.819* (47.915)	180.834* (43.822)
Iowa BS	23.288 (26.369)	30.379 (24.054)
Major: Soc.	.	2422.31* (515.377)
Major: Nat.	.	5443.171* (518.659)
Prof. Parents: Yes	.	904.669* (459.893)
Parent Network: Yes	.	406.115 (469.816)
Gender: Male	.	-160.315 (424.972)
N	505	505
RMSE	5216.452	4725.612
$R^2$	0.185	0.338
adj $R^2$	0.18	0.327

\* $p \leq 0.05$ 

```
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 = 498.553340770017 Denominator = 662.911440166179"
```

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

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

```
tval
```

```
difference
0.7520663
```

```
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
0.4523679
```

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
498.5533408	662.9114402	0.7520663	496.0000000	0.4523679

```
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     501 13632897874
2     496 11076378000   5 2556519874 22.896 < 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-29

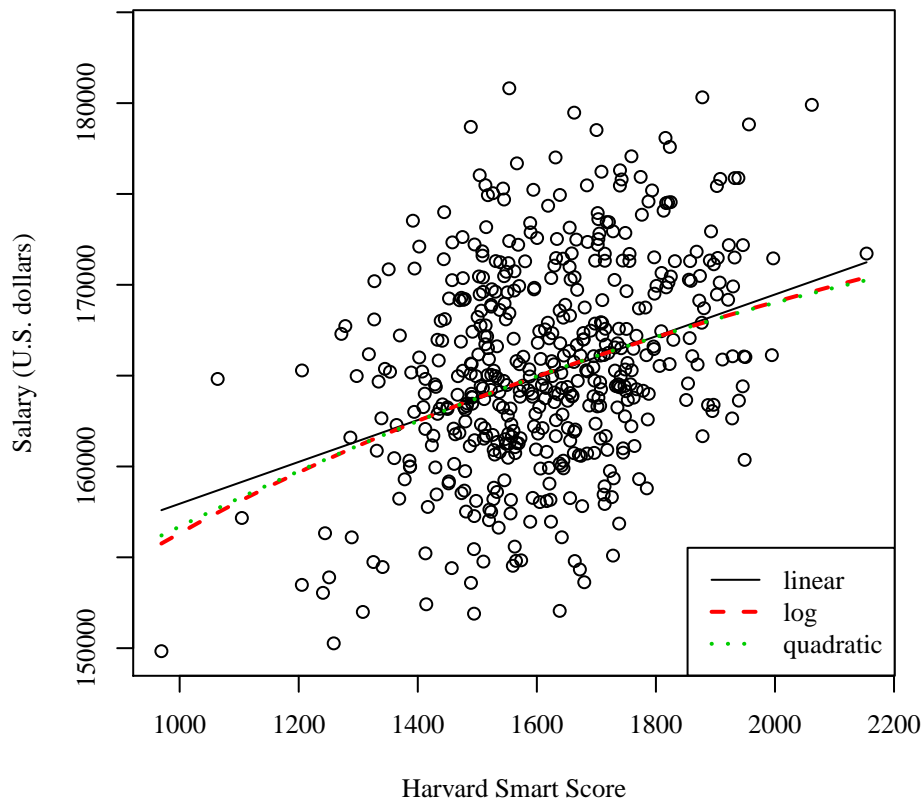
	Linear	Log	Quadratic
	Estimate	Estimate	Estimate
	(S.E.)	(S.E.)	(S.E.)
(Intercept)	144788.25*	27853.595	135432.352*
	(2223.786)	(15688.232)	(13664.778)
Harvard SS	11.518*	.	23.249
	(1.344)		(16.958)
Gender: Male	-155.004	-157.968	-157.445
	(443.628)	(443.129)	(443.88)
Major: Soc.	1796.296*	1818.44*	1815.445*
	(539.726)	(539.134)	(540.72)
Major: Nat.	5214.641*	5237.35*	5232.707*
	(544.638)	(543.898)	(545.552)
Prof. Parents: Yes	177.036	182.093	183.98
	(482.269)	(481.75)	(482.632)
Parent Network: Yes	-183.76	-174.533	-176.067
	(488.939)	(488.42)	(489.327)
ln(Harvard SS)	.	18358.542*	.
		(2123.907)	
Harvard SS <sup>2</sup>	.	.	-0.004
			(0.005)
N	491	491	491
RMSE	4870.982	4865.684	4873.593
$R^2$	0.269	0.271	0.27
adj $R^2$	0.26	0.262	0.26

\* $p \leq 0.05$ 

```
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
H (30%) 20445.60  H
S (30%) 23320.45  S
N (30%) 26435.11  N

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

```
predictOMatic(cm2)
```

```
$major2
      fit major2
H (30%) 20445.60  H
S (30%) 23320.45  S
N (30%) 26435.11  N

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

Table 5: Categorical Regressions: Student-29

	major	major2	major full	major2 full
	Estimate	Estimate	Estimate	Estimate
	(S.E.)	(S.E.)	(S.E.)	(S.E.)
(Intercept)	20445.603*	23320.453*	-2559.696	-137.385
	(383.774)	(386.924)	(2535.688)	(2527.215)
Major: Soc.	2874.85*	.	2422.31*	.
	(544.97)		(515.377)	
Major: Nat.	5989.508*	.	5443.171*	.
	(550.437)		(518.659)	
Major 2: Hum.	.	-2874.85*	.	-2422.31*
		(544.97)		(515.377)
Major 2: Nat.	.	3114.658*	.	3020.861*
		(552.638)		(523.54)
SAT	.	.	10.048*	10.048*
			(1.506)	(1.506)
ACT	.	.	180.834*	180.834*
			(43.822)	(43.822)
Iowa BS	.	.	30.379	30.379
			(24.054)	(24.054)
Prof. Parents: Yes	.	.	904.669*	904.669*
			(459.893)	(459.893)
Parent Network: Yes	.	.	406.115	406.115
			(469.816)	(469.816)
Gender: Male	.	.	-160.315	-160.315
			(424.972)	(424.972)
N	542	542	505	505
RMSE	5219.888	5219.888	4725.612	4725.612
$R^2$	0.18	0.18	0.338	0.338
adj $R^2$	0.177	0.177	0.327	0.327

\* $p \leq 0.05$