

## Data Management

```
library(foreign)
library(rockchalk)
i <- 15
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%	8.68	1185.00	73.16	6368.00	8114.00	148200.00	1167.00
25%	18.80	1515.00	92.86	16850.00	19610.00	161500.00	1492.00
50%	21.79	1626.00	99.43	20610.00	23610.00	165500.00	1600.00
75%	24.72	1735.00	106.30	24240.00	27480.00	169000.00	1707.00
100%	35.83	2083.00	137.20	38310.00	41270.00	186500.00	2051.00
mean	21.89	1620.00	99.40	20430.00	23460.00	165500.00	1598.00
sd	4.67	167.00	9.99	5277.00	5663.00	5909.00	165.00
var	21.86	27890.00	99.82	27840000.00	32070000.00	34920000.00	27230.00
NA's	11.00	70.00	0.00	22.00	0.00	0.00	29.00
N	546.00	546.00	546.00	546.00	546.00	546.00	546.00

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>	
F	:277.0000	N	:185.0000	NO	:380.0000
M	:269.0000	S	:181.0000	YES	:166.0000
NA's	: 0.0000	H	:180.0000	NA's	: 0.0000
entropy	: 0.9998	NA's	: 0.0000	entropy	: 0.8862
normedEntropy	: 0.9998	entropy	: 1.5849	normedEntropy	: 0.8862
N	:546.0000	normedEntropy	: 0.9999	N	:546.0000
		N	:546.0000		
<b>pprof</b>					
NO	:369.0000				
YES	:177.0000				
NA's	: 0.0000				
entropy	: 0.9089				
normedEntropy	: 0.9089				
N	:546.0000				

## 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: 0x1f5ffe8>
act ~ sat + ibs + harv
<environment: 0x1f5ffe8>
ibs ~ sat + act + harv
<environment: 0x1f5ffe8>
harv ~ sat + act + ibs
<environment: 0x1f5ffe8>
Drum roll please!
```

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

```
      sat      act      ibs      harv
0.9998592 0.8476091 0.2655526 0.9998628
The Corresponding VIF, 1/(1-Rj2)
      sat      act      ibs      harv
7100.844460  6.562072  1.361568 7286.427066
```

Bivariate Correlations for design matrix

```
      sat  act  ibs  harv
sat  1.00 0.47 0.42 1.00
act  0.47 1.00 0.46 0.49
ibs  0.42 0.46 1.00 0.42
harv 1.00 0.49 0.42 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-15

	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)	4162.534 (2175.457)	11644.22* (1042.554)	9827.385* (2239.564)	3289.836 (2289.192)	2097.156 (2795.046)	3162.859 (2598.904)
SAT	10.212* (1.355)	.	.	.	184.302 (122.286)	6.248* (1.539)
ACT	.	403.829* (46.615)	.	.	447.874* (134.447)	282.357* (56.194)
Iowa BS	.	.	106.787* (22.438)	.	17.535 (26.835)	12.094 (24.854)
Harvard SS	.	.	.	10.65* (1.406)	-177.498 (122.098)	.
N	497	513	524	455	422	488
RMSE	4991.082	4944.46	5170.851	5024.814	4900.25	4865.991
$R^2$	0.103	0.128	0.042	0.112	0.159	0.155
adj $R^2$	0.101	0.126	0.04	0.11	0.151	0.15

\* $p \leq 0.05$ 

	Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
1	419	1.0073e+10				
2	417	1.0013e+10	2	60056911	1.2505	0.2874

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)	4261.814 (2190.037)	11840.576* (1078.225)	9681.537* (2297.984)	3563.657 (2391.198)	2097.156 (2795.046)	3162.859 (2598.904)
SAT	10.185* (1.363)	.	.	.	184.302 (122.286)	6.248* (1.539)
ACT	.	396.726* (48.133)	.	.	447.874* (134.447)	282.357* (56.194)
Iowa BS	.	.	109.192* (23.002)	.	17.535 (26.835)	12.094 (24.854)
Harvard SS	.	.	.	10.545* (1.467)	-177.498 (122.098)	.
N	488	488	488	422	422	488
RMSE	5004.278	4949.166	5165.363	5024.858	4900.25	4865.991
$R^2$	0.103	0.123	0.044	0.11	0.159	0.155
adj $R^2$	0.101	0.121	0.042	0.107	0.151	0.15

\* $p \leq 0.05$

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  0.18143241
act  0.22266100
ibs  0.02211322

```

```
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 0.0287496218
act 0.0440593033
ibs 0.0004132208

```

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%     36.99   30.34   36.71
50%     48.20   41.21   48.95
75%     58.90   51.45   61.02
100%    100.00  100.00  100.00

```

```

mean  48.74  40.94  48.71
sd    17.16  15.89  18.79
var   294.50 252.40 353.20
NA's  0.00   0.00   0.00
N     488.00 488.00 488.00

```

```

$ 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
-14998.0  -3192.9    90.4   3328.5  15570.8

```

```

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) 13787.059    779.976   17.676 < 2e-16 ***
satpoms      55.290     13.622    4.059 5.75e-05 ***
actpoms      76.660     15.257    5.025 7.11e-07 ***
ibspoms      7.746     15.919    0.487  0.627

```

```

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

```

```

Residual standard error: 4866 on 484 degrees of freedom
Multiple R2: 0.1554, Adjusted R2: 0.1501
F-statistic: 29.68 on 3 and 484 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  0.07360495
act  0.16100257
ibs  0.03198274
harv -0.07101014

```

```

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  0.0045798064
act  0.0223741031
ibs  0.0008608941
harv 0.0042609846

```

```

options(scipen = 5)

```

## Additional Variables

Please see Table 3 for the regressions.

Table 3: Regression with sal2: Student-15

	Test Scores Only	All Predictors
	Estimate	Estimate
	(S.E.)	(S.E.)
(Intercept)	5239.951 (2726.229)	3358.382 (2574.746)
SAT	7.357* (1.601)	6.813* (1.483)
ACT	316.07* (58.11)	270.006* (54.153)
Iowa BS	-3.556 (26.165)	8.667 (24.412)
Major: Soc.	.	1047.421* (527.811)
Major: Nat.	.	3902.394* (521.193)
Prof. Parents: Yes	.	625.544 (458.78)
Parent Network: Yes	.	2279.755* (463.205)
Gender: Male	.	-87.155 (427.671)
N	508	508
RMSE	5185.797	4796.496
$R^2$	0.162	0.29
adj $R^2$	0.157	0.279

\* $p \leq 0.05$ 

```
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 = -1654.21027827797 Denominator = 652.131773818115"
```

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

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

```
tval
```

```
difference
-2.53662
```

```
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
0.01149649
```

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
-1654.21027828	652.13177382	-2.53661966	499.00000000	0.01149649

```
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	504	13553816150				
2	499	11480178591	5	2073637559	18.027	< 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-15

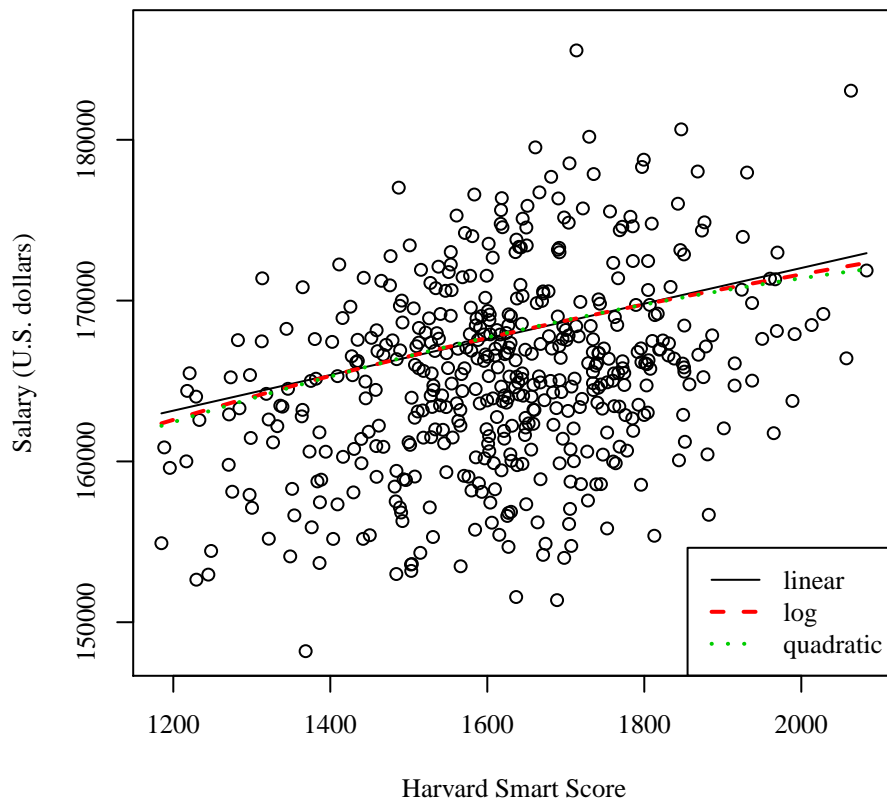
	Linear Estimate (S.E.)	Log Estimate (S.E.)	Quadratic Estimate (S.E.)
(Intercept)	144120.702* (2289.784)	31417.905 (16110.539)	131288.342* (15157.339)
Harvard SS	11.089* (1.374)	.	27.194 (18.855)
Gender: Male	-56.086 (459.458)	-57.224 (459.114)	-57.55 (459.592)
Major: Soc.	2428.282* (563.705)	2442.238* (563.325)	2447.038* (564.291)
Major: Nat.	5783.523* (562.346)	5784.028* (561.936)	5787.65* (562.526)
Prof. Parents: Yes	1222.876* (489.878)	1242.502* (489.58)	1251.818* (491.181)
Parent Network: Yes	416.372 (501.485)	412.126 (501.141)	412.519 (501.647)
ln(Harvard SS)	.	17692.925* (2179.416)	.
Harvard SS <sup>2</sup>	.	.	-0.005 (0.006)
N	476	476	476
RMSE	4984.851	4981.22	4986.268
$R^2$	0.281	0.282	0.282
adj $R^2$	0.272	0.273	0.271

\* $p \leq 0.05$ 

```
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 (30%) 25987.85  N
S (30%) 22760.23  S
H (30%) 21570.18  H

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

```
predictOMatic(cm2)
```

```
$major2
      fit major2
N (30%) 25987.85  N
S (30%) 22760.23  S
H (30%) 21570.18  H

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

Table 5: Categorical Regressions: Student-15

	major	major2	major full	major2 full
	Estimate	Estimate	Estimate	Estimate
	(S.E.)	(S.E.)	(S.E.)	(S.E.)
(Intercept)	21570.178*	22760.226*	3358.382	4405.803
	(399.08)	(397.976)	(2574.746)	(2569.138)
Major: Soc.	1190.047*	.	1047.421*	.
	(563.604)		(527.811)	
Major: Nat.	4417.673*	.	3902.394*	.
	(560.558)		(521.193)	
Major 2: Hum.	.	-1190.047*	.	-1047.421*
		(563.604)		(527.811)
Major 2: Nat.	.	3227.625*	.	2854.973*
		(559.772)		(523.503)
SAT	.	.	6.813*	6.813*
			(1.483)	(1.483)
ACT	.	.	270.006*	270.006*
			(54.153)	(54.153)
Iowa BS	.	.	8.667	8.667
			(24.412)	(24.412)
Prof. Parents: Yes	.	.	625.544	625.544
			(458.78)	(458.78)
Parent Network: Yes	.	.	2279.755*	2279.755*
			(463.205)	(463.205)
Gender: Male	.	.	-87.155	-87.155
			(427.671)	(427.671)
N	546	546	508	508
RMSE	5354.217	5354.217	4796.496	4796.496
$R^2$	0.109	0.109	0.29	0.29
adj $R^2$	0.106	0.106	0.279	0.279

\* $p \leq 0.05$