

# Exploratory Data Mining via Search Strategies Lab #2

Ross Jacobucci & Kevin J. Grimm

# Outline

The second lab will go over some more recent techniques in regression –

1. Multivariate Adaptive Regression Splines
2. Regularized Regression

## Ridge and Lasso Regression Including a penalty on the  $\beta$  parameters, and by varying the penalty we can shrink some of the  $\beta$ 's to zero, doing a form of “automatic” subset selection. \ Although there a number of packages to do this, maybe the best is *glmnet*

Note, for glmnet, your data has to be set up in two separate matrices. Doing this can be accomplished by:

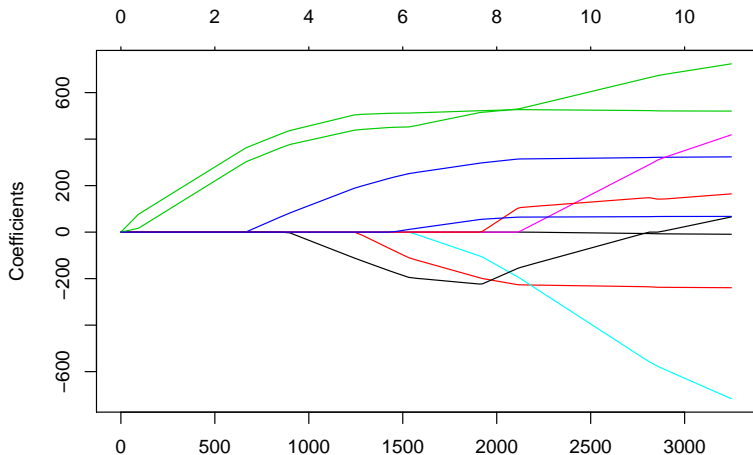
```
library(elasticnet);data(diabetes);  
X <- diabetes$x;Y <- diabetes$y  
diabetes2 <- data.frame(cbind(Y,X))  
YY <- as.matrix(diabetes2$Y)  
XX <- as.matrix(diabetes2[,2:11])
```

Two things to note:

1. Because we are doing regression with a continuous outcome, we specify the family(distribution) as “gaussian”
2. Shrinkage in lasso and ridge is sensitive to the scale of the variables, therefore, it is best to standardize the predictors before entering. glmnet does this by default (look at ?glmnet).

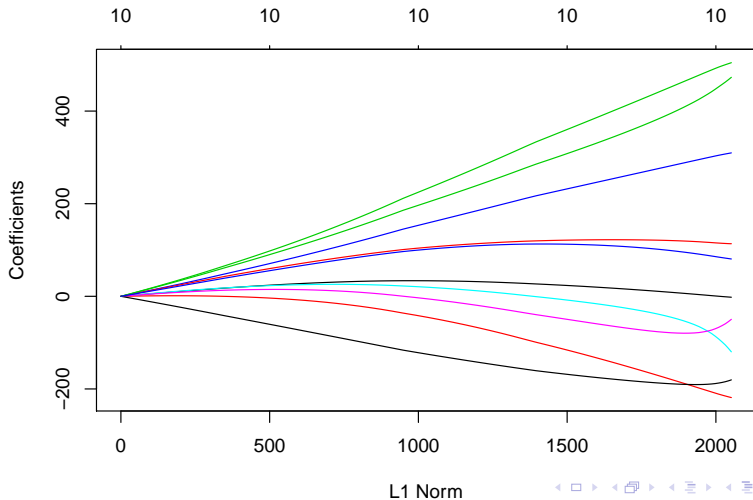
# Lasso

```
library(glmnet)
# ?glmnet
lasso.out <- glmnet(XX,YY,family="gaussian",alpha=1)
plot(lasso.out)
```



# Ridge

```
ridge.out <- glmnet(XX,YY,family="gaussian",alpha=0)  
#plot(ridge.out,type.coef="2norm")  
plot(ridge.out)
```



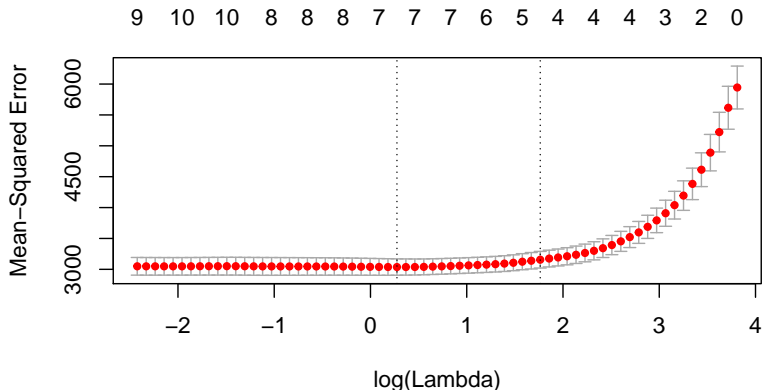
## Regularized Regression Continued

Since ridge regression does not shrink the  $\beta$  coefficients to 0 with increase penalization, it does not do an “automatic” form of subset selection

The problem now becomes, which value of  $\lambda$  (amount of shrinkage) do we choose?

Using cross-validation is one of the better ways, and is implemented the glmnet package

```
cv.lasso <- cv.glmnet(XX,YY,family="gaussian",alpha=1)
plot(cv.lasso)
```



# Choosing the Optimal Lambda (Penalty)

Two-strategies for selecting  $\lambda$ : either pick the lowest CV error, or the best solution within 1 standard error.

I don't think that there is a clear best choice. The one advantage of using the 1SE rule is that you need fewer predictors. In our example 4 instead of 7.

```
#str(cv.lasso)
(lmin <- cv.lasso$lambda.min)
```

```
## [1] 1.316439
```

```
(lminSE <- cv.lasso$lambda.1se)
```

```
## [1] 5.832642
```

```
lasso.out2 = glmnet(XX,YY,family="gaussian",alpha=1,lambda=lminSE)
lasso.out2
```

```
##
```

```
## Call: glmnet(x = XX, y = YY, family = "gaussian", alpha = 1, lambda = lminSE)
```

```
##
```

```
##          Df      %Dev Lambda
```

```
## [1,]    5 0.4823  5.833
```

# Highly correlated predictors

```
library(lavaan)
sim.mod <- '
y ~ 1*x1 + 1*x2
x1~~0.999*x2
'

set.seed(3)
dat <- simulateData(sim.mod, model.type="sem", sample.nobs=100)

out <- lm(y ~ ., data=dat)
summary(out)
```

```
##
## Call:
## lm(formula = y ~ ., data = dat)
##
## Residuals:
```

	Min	1Q	Median	3Q	Max
	-1.9320	-0.7527	-0.1309	0.7976	2.5258

```
##
## Coefficients:
```

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	0.02609	0.10487	0.249	0.804
x1	1.58969	2.30343	0.690	0.492
x2	0.16792	2.30542	0.073	0.942

```
##
```

## Residual standard error: 1.049 on 97 degrees of freedom

# Ridge Regression

```
library(glmnet)
X <- matrix(cbind(dat$x1,dat$x2),100,2)
Y <- data.matrix(dat$y)
ridge <- glmnet(X,Y,family="gaussian",alpha=0)
coef(ridge,s=0.2)
```

```
## 3 x 1 sparse Matrix of class "dgCMatrix"
##              1
## (Intercept) 0.02293853
## V1          0.84523667
## V2          0.82491602
```



# Lasso for correlated variables

```
#library(glmnet)  
lasso <- glmnet(X,Y,family="gaussian",alpha=1) # change alpha  
coef(lasso,s=0.01)
```

```
## 3 x 1 sparse Matrix of class "dgCMatrix"  
##              1  
## (Intercept) 0.02680360  
## V1          1.71635450  
## V2          0.02973456
```

Lasso doesn't have the same properties as ridge for collinear predictors. This is the rationale for the elastic net

# Elastic Net

```
enet1 <- glmnet(X,Y,family="gaussian",alpha=0.5) # mixture  
coef(enet1,s=0.01)
```

```
## 3 x 1 sparse Matrix of class "dgCMatrix"  
##              1  
## (Intercept) 0.0232833  
## V1          1.0197302  
## V2          0.7301989
```

```
enet2 <- glmnet(X,Y,family="gaussian",alpha=0.5) # mixture  
coef(enet2,s=0.2)
```

```
## 3 x 1 sparse Matrix of class "dgCMatrix"  
##              1  
## (Intercept) 0.02326065  
## V1          0.82358307  
## V2          0.78180790
```

```
enet3 <- glmnet(X,Y,family="gaussian",alpha=0.5) # mixture  
coef(enet3,s=2)
```

```
## 3 x 1 sparse Matrix of class "dgCMatrix"  
##              1  
## (Intercept) 0.02767924  
## V1          0.25954355  
## V2          0.25609747
```

# More P's than People

```
set.seed(1)
N <- 30; P <- 100
X <- matrix(rnorm(N*P),N,P)
Y <- rnorm(N)

out <- lm(Y ~ X)
head(summary(out)$coefficients)
```

##		Estimate	Std. Error	t value	Pr(> t )
##	(Intercept)	-1.37504720	NaN	NaN	NaN
##	X1	-0.08138609	NaN	NaN	NaN
##	X2	4.79719809	NaN	NaN	NaN
##	X3	-6.98146901	NaN	NaN	NaN
##	X4	-4.93789272	NaN	NaN	NaN
##	X5	1.50237557	NaN	NaN	NaN

Other parts of the summary list the errors and non-singularity of the information matrix.  
Can't invert a matrix that is wider than long.

# Regularized Regression for $P > N$

Both Ridge and Lasso (& Enet) can handle this case

```
# just use lasso  
lasso2 <- glmnet(X,Y,alpha=1)  
head(coef(lasso2,0.0001))
```

```
## 6 x 1 sparse Matrix of class "dgCMatrix"  
##              1  
## (Intercept) -0.16476453  
## V1          -0.04102075  
## V2          .  
## V3          .  
## V4          .  
## V5          .
```

The addition of penalties effectively reduces the dimensionality of the parameter space. In this case, don't need much for penalty because a lot of coefficients are set immediately to zero.

# P-values for Lasso

The traditional lasso does not output p-values. Only really assessing “importance” in the sense of what relationships we think with generalize through the use of cross-validation.

**Need two new packages**

```
library(lars)  
library(covTest)
```

## Lasso P-Values Continued

Exploratory, but only current method that accounts for adaptive nature without having to split the sample.

```
X <- diabetes$x; Y <- diabetes$y
lars.out <- lars(X,Y)
cov.out <- covTest(lars.out,X,Y)
cov.out
```

```
## $results
## Predictor_Number Drop_in_covariance P-value
##           3           19.5084  0.0000
##           9           50.8005  0.0000
##           4            5.5744  0.0041
##           7            6.0161  0.0026
##           2            4.8769  0.0080
##          10            0.1565  0.8552
##           5            3.2893  0.0382
##           8            0.6105  0.5436
##           6            0.1393  0.8700
##           1            0.0161  0.9841
##
## $sigma
## [1] 54.0915
##
## $null.dist
## [1] "F(2,432)"
```

## Lasso P-Values Continued

```
sig <- cov.out$results[, "P-value"] < 0.05
vars <- cov.out$results[sig, "Predictor_Number"]
colnames(X[, vars])
```

```
## [1] "bmi" "ltg" "map" "hdl" "sex" "tc"
```

Compare to normal lasso procedure

```
glmnet.out <- cv.glmnet(X, Y)
coef(glmnet.out, s = glmnet.out$lambda.1se)
```

```
## 11 x 1 sparse Matrix of class "dgCMatrix"
```

```
##              1
## (Intercept) 152.13348
## age          .
## sex          -13.59236
## bmi          506.66309
## map          199.02672
## tc           .
## ldl           .
## hdl          -124.16053
## tch          .
## ltg          441.69986
## glu          .
```

# Relaxed Lasso

So the lasso has found to be biased in that it shrinks non-zero coefficients too much. To compensate for this, the relaxed lasso is a two step procedure in that the steps include:

1. Fit lasso, select non-zero coefficients
2. Re-fit linear regression with only non-zero coefficients included

Step 1

```
y <- data.matrix(mtcars$mpg)
x <- as.matrix(mtcars[,2:11],nrow(mtcars),10)
lasso.out1 <- cv.glmnet(x,y)
coef(lasso.out1,lasso.out1$lambda.1se)
```

```
## 11 x 1 sparse Matrix of class "dgCMatrix"
##              1
## (Intercept) 33.940487812
## cyl        -0.843038418
## disp         .
## hp          -0.006965929
## drat         .
## wt          -2.365917425
## qsec         .
## vs           .
## am           .
## gear         .
## carb         .
```



## Relaxed Lasso Step 2

```
lm.out <- lm(mpg ~ cyl + hp + wt,mtcars)
coef(lm.out)
```

```
## (Intercept)          cyl          hp          wt
##  38.7517874  -0.9416168  -0.0180381  -3.1669731
```

The coefficients are larger in step 2.

# Multivariate Adaptive Splines

Get ECLS dataset

```
ecls <- read.table("C:/Users/RJacobucci/Documents/GitHub/SearchWkshp_labs16/ec1
```

```
names(ecls) = c('gender','kage',  
  'k_read_irt','k_read1','k_read2','k_read3','k_read4',  
  'k_print','k_read_tht',  
    'k_math_irt','k_math1','k_math2','k_math3','k_math4',  
    'k_math_tht',  
    'k_gk_irt','k_gk_tht',  
    'f_mtr','g_mtr',  
    'P1LEARN','P1CONTRO','P1SOCIAL','P1SADLON','P1IMPULS',  
    'ars_lit','ars_mth','ars_gk',  
    'T1LEARN','T1CONTRO','T1INTERP','T1EXTERN','T1INTERN',  
    'height','weight','bmi',  
    'hisp','na_amer','asian','black','pac_isl','white','m_race',  
    'ses_c','ses_cat','poor','income',  
    'g8_read','g8_read_tht','g8_math','g8_math_tht',  
    'g8_sci','g8_sci_tht')
```

```
myvars = c('gender','kage',  
  'k_read_irt','k_print',  
    'k_math_irt',  
    'k_gk_irt',  
    'f_mtr','g_mtr',  
    'P1LEARN','P1CONTRO','P1SOCIAL','P1SADLON','P1IMPULS',
```

# Use the earth package

Predicting science scores at grade 8

```
library(earth)
```

```
## Loading required package: plotmo
```

```
## Loading required package: plotrix
```

```
## Loading required package: TeachingDemos
```

```
mars.1 = earth(g8_sci ~ ., degree=1, data=ecls.1)
```

```
mars.1
```

```
## Selected 17 of 22 terms, and 11 of 28 predictors
```

```
## Termination condition: RSq changed by less than 0.001 at 22 terms
```

```
## Importance: k_math_irt, k_gk_irt, ses_c, kage, f_mtr, gender, T1LEARN, ...
```

```
## Number of terms at each degree of interaction: 1 16 (additive model)
```

```
## GCV 116.3739    RSS 522546.8    GRSq 0.4635932    RSq 0.4711035
```

```
summary(mars.1)
```

```
## Call: earth(formula=g8_sci~., data=ecls.1, degree=1)
```

```
##
```

```
##               coefficients
```

```
## (Intercept)      92.271391
```

```
## gender           -2.966051
```

```
## h(78.53-kage)      0.465144
```

```
## h(kage-78.53)     -18.785860
```

## Second Mars Model

```
mars.2 = earth(g8_sci ~ ., data = ecl.1,  
              degree = 1, nfold = 10, pmethod = 'cv')
```

```
summary(mars.2)
```

```
## Call: earth(formula=g8_sci~., data=ecl.1, pmethod="cv", degree=1,  
##           nfold=10)
```

```
##               coefficients
```

```
## (Intercept)      91.368480
```

```
## gender          -2.953725
```

```
## h(78.53-kage)    0.458340
```

```
## h(kage-78.53)   -18.858686
```

```
## h(29.88-k_math_irt) -0.699111
```

```
## h(k_math_irt-29.88) 0.152866
```

```
## h(17.05-k_gk_irt) -1.367930
```

```
## h(k_gk_irt-17.05)  0.518428
```

```
## h(6-f_mtr)        -1.084918
```

```
## h(f_mtr-6)         0.629474
```

```
## h(g_mtr-1)         -0.367090
```

```
## h(2.33-P1LEARN)   -13.862618
```

```
## h(P1LEARN-2.33)    1.055330
```

```
## h(2.33-P1SOCIAL)   4.815832
```

```
## h(P1SOCIAL-2.33)   -1.182026
```

```
## h(3.5-P1IMPULS)    0.553900
```

```
## h(P1IMPULS-3.5)    7.566074
```

## Second Model Continued

```
# Examine plot of one predictor
plot(ecls.1$k_math_irt, ecl.1$g8_sci, ylab = '8th Grade Science', xlab = 'Math'
      xlim=c(min(ecls.1$k_math_irt), max(ecls.1$k_math_irt)), ylim=c(min(ecls.1$g8_sci), max(ecls.1$g8_sci))
#mars.2$coefficients
b0=mars.2$coefficients[1]; b1=mars.2$coefficients[5]; b2=mars.2$coefficients[4]

curve(b0 + b1*pmax(29.88-x,0) + b2*pmax(x-29.88,0),
      from = min(ecls.1$k_math_irt), to = max(ecls.1$k_math_irt),
      n = 100, col = "red", lwd = 2, ann = F, add = T)
```

