

# Linear Model Selection Regularization

Ryan Finegan  
11/1/2021

```
library(splot2)
library(caret)

## Loading required package: lattice

library(MASS)
library(glmnet)

## Loading required package: Matrix

## Loaded glmnet 4.1-1

library(class)
set.seed(5)
setwd("~/Users/cryanfinegan/Documents") # my working directory
df<-read.csv("10yearforecasting.csv") # file for 10 year prediction
dates<-as.POSIXct(df$Dates, format = "%m/%d/%Y") # converting to get just the year
df$year<-format(dates, format="%Y") # getting the year in dates
df<-df[1:,] # omitting rid of the first row because zeros
df<-subset(df, select = c(thirtyderivative,thirty,movederivative,
move,dxyderivative,dxy)) # splitting at the 2015 mark
train<-df[df$Dates< 2015, ]
test<-df[df$Dates> 2015, ]
train<-subset(train, select = c(Dates)) # getting rid of Dates
test<-subset(test, select = c(Dates)) # getting rid of Dates
# tenderivativevelag1 + thirtylag5 + dxylag1 + movederivative + movederivativevelag1 + tenlag2
dft<-subset(df, select = c(Dates)) # getting rid of Dates
ten.year.lm<-lm(ten~., data = train) # mlr model on training data
summary(ten.year.lm) # regression statistics
```

```
## Call:
## lm(formula = ten ~ ., data = train)
##
##    Min      1Q  Median      3Q     Max
## -0.56025 -0.08878 -0.00526  0.07813  0.56549
##
## Coefficients: (16 not defined because of singularities)
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   3.691e-03  4.018e-02  0.092  0.9268
## tenlag1       9.159e-01  7.889e-02  11.609 <2e-16 ***
## tenlag2       1.390e-01  1.052e-01  1.322  0.1864
## tenlag3      -6.029e-02  1.057e-01 -0.571  0.5684
## tenlag4      -1.211e-01  1.059e-01 -1.146  0.2522
## tenlag5      1.181e-01  7.827e-02  1.490  0.1365
## movederivelag1 7.747e-04  4.424e-04  1.751  0.0802 .
## movederivelag2 -1.235e-03  5.801e-04 -2.128  0.0335 *
## movederivelag3 3.912e-04  5.839e-04  0.670  0.5030
## movederivelag4 -1.197e-04  5.839e-04 -0.205  0.8376
## movederivelag5 8.511e-05  4.456e-04  0.191  0.8486
## thirtyderivelag1 5.457e-03  8.708e-02  0.063  0.9500
## thirtyderivelag2 1.932e-02  1.166e-01  0.166  0.8684
## thirtyderivelag3 2.982e-02  1.171e-01  0.255  0.7990
## thirtyderivelag4 1.588e-01  1.171e-01  1.356  0.1752
## thirtyderivelag5 -2.072e-01  8.750e-02 -2.368  0.0180 *
## dxylag1       6.342e-03  3.596e-03  1.764  0.0780 .
## dxylag2       -8.511e-03  5.048e-03 -1.686  0.0920 .
## dxylag3       9.712e-04  5.049e-03  0.192  0.8475
## dxylag4       3.933e-03  5.038e-03  0.781  0.4351
## dxyderivative1 -2.627e-03  3.592e-03 -0.721  0.4646
## dxyderivative2 NA NA NA NA NA
## dxyderivative3 NA NA NA NA NA
## dxyderivative4 NA NA NA NA NA
## dxyderivative5 NA NA NA NA NA
## movederivative1 NA NA NA NA NA
## movederivative2 NA NA NA NA NA
## movederivative3 NA NA NA NA NA
## movederivative4 NA NA NA NA NA
## movederivative5 -3.497e-04  4.364e-04 -0.801  0.4231
## thirtyderivative1 NA NA NA NA NA
## thirtyderivative2 NA NA NA NA NA
## thirtyderivative3 NA NA NA NA NA
## thirtyderivative4 NA NA NA NA NA
## thirtyderivative5 NA NA NA NA NA
## dxyderivative1 NA NA NA NA NA
## dxyderivative2 NA NA NA NA NA
## dxyderivative3 NA NA NA NA NA
## dxyderivative4 NA NA NA NA NA
## dxyderivative5 7.223e-03  3.599e-03  2.007  0.0450 *
```

```
## ---
## 3 signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## Residual standard error: 0.128 on 1273 degrees of freedom
## Multiple R-squared: 0.9949, Adjusted R-squared: 0.9949
## F-statistic 1.045e+04 on 24 and 1273 DF, p-value: < 2.2e-16
```

```
ols.prediction<-predict(ten.year.lm, test) # getting the MSE for the OLS method
```

```
## Warning in predict.lm(ten.year.lm, test): prediction from a rank-deficient fit
## may be misleading
```

```
(ols.mse<-mean(ols.prediction - test$ten)^2) # the predictions on the test data split vs actual
```

```
## [1] 0.008799786
```

```
library(dplyr)
```

```
##
## Attaching package: 'dplyr'
```

```
## The following object is masked from 'package:MASS':
##
## select
```

```
## The following objects are masked from 'package:stats':
##
## filter, lag
```

```
## The following objects are masked from 'package:base':
##
## intersect, setdiff, setequal, union
```

```
# tenderivativevelag1 + thirtylag5 + dxylag1 + dxylag2 + movederivative + movederivativevelag1 + tenlag2
train.matrix<-dummyVars(ten~., data = train, fullRank = F) %>%
predict(newdata = train) %>%
as.matrix()
# tenderivativevelag1 + thirtylag5 + dxylag1 + dxylag2 + movederivative + movederivativevelag1 + tenlag2
test.matrix<-dummyVars(ten~., data = test, fullRank = F) %>%
predict(newdata = test) %>%
as.matrix()
# ridge has an alpha of zero and the coefficients never shrink to zero
mod.ridge<-cv.glmnet(y = train$ten,
X = train.matrix,
alpha = 0,
lambda = 10^seq(2,-2, length = 100),
standardize = TRUE,
nfolds = 5)
```

```
data.frame(lambda = mod.ridge$lambda,
cv_mse = mod.ridge$cvm) %>%
ggplot(aes(x = lambda, y = cv_mse)) +
geom_point() +
geom_line() +
geom_vline(xintercept = mod.ridge$lambda.min, col = "deepskyblue3") +
geom_hline(yintercept = min(mod.ridge$cvm), col = "deepskyblue3") +
scale_x_continuous(trans = "log10", breaks = c(0.01, 0.1, 1, 10, 100), labels = c(0.01, 0.1, 1, 10, 100)) +
scale_y_continuous(labels = scales::comma_format()) +
theme(legend.position = "bottom") +
labs(x = "Lambda",
y = "CV MSE",
col = "Coefficients", # Coefficients - Can't be zero because Ridge Regression
title = "Ridge Regression") # Lambda Selection with 5 CV
```



```
## Using lambda selection model above
mod.ridge.best<-glmnet(y = train$ten,
X = train.matrix,
alpha = 0,
lambda = 10^seq(2,-2, length = 100))
```

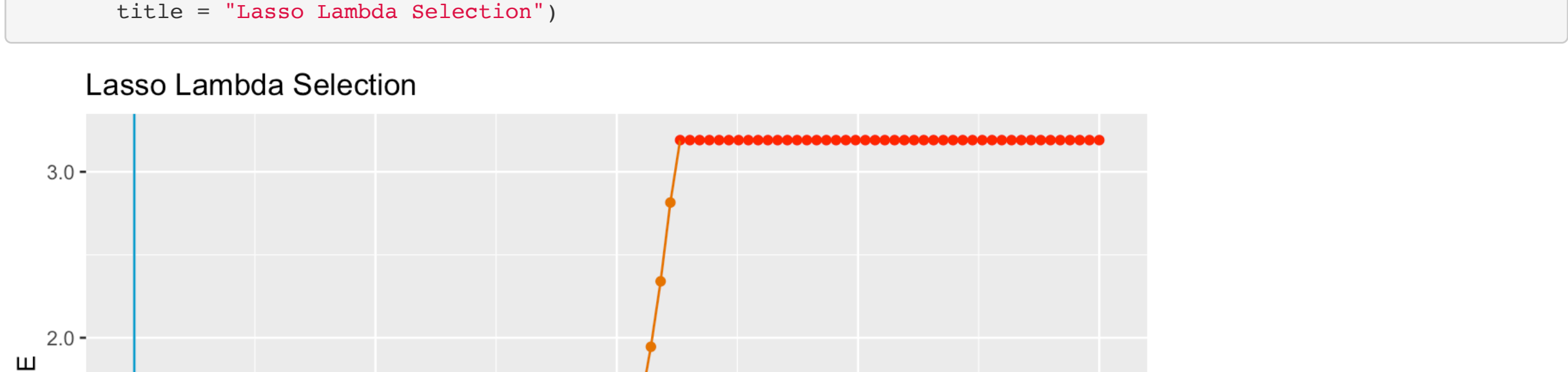
```
ridge.prediction<-predict(mod.ridge.best, s = mod.ridge$lambda.min, newx = test.matrix)
(ridge.coef<-predict(mod.ridge.best, type = "coefficients", s = mod.ridge$lambda.min))
```

```
## 41 x 1 sparse Matrix of class "dgCMatrix"
##
## (Intercept) -9.644693e-02
## tenlag1 2.305408e-01
## tenlag2 1.998801e-01
## tenlag3 1.807223e-01
## tenlag4 1.678480e-01
## tenlag5 1.592517e-01
## movederivelag1 4.055365e-05
## movederivelag2 -6.999793e-05
## movederivelag3 -4.785335e-05
## movederivelag4 -6.147246e-05
## movederivelag5 -5.507020e-05
## thirtyderivelag1 -2.713635e-02
## thirtyderivelag2 2.561110e-03
## thirtyderivelag3 3.980840e-02
## thirtyderivelag4 3.155623e-02
## thirtyderivelag5 3.974564e-02
## dxylag1 -1.627403e-04
## dxylag2 1.251479e-05
## dxylag3 1.625469e-04
## dxylag4 2.601627e-04
## dxylag5 3.001955e-04
## tenderivativevelag1 6.287978e-01
## tenderivativevelag2 5.667343e-01
## tenderivativevelag3 3.408576e-01
## tenderivativevelag4 7.322606e-02
## tenderivativevelag5 -8.111355e-02
## movederivativevelag1 7.338204e-04
## movederivativevelag2 -4.280920e-04
## movederivativevelag3 1.105378e-06
## movederivativevelag4 -7.012943e-05
## movederivativevelag5 -2.986203e-04
## thirtyderivativevelag1 9.171952e-02
## thirtyderivativevelag2 1.083925e-01
## thirtyderivativevelag3 1.032688e-01
## thirtyderivativevelag4 2.092852e-01
## thirtyderivativevelag5 2.543348e-02
## dxyderivativevelag1 7.489715e-03
## dxyderivativevelag2 -9.587432e-04
## dxyderivativevelag3 -3.140408e-04
## dxyderivativevelag4 3.187494e-03
## dxyderivativevelag5 7.406089e-03
```

```
(ridge.mse<-mean((ridge.prediction - test$ten)^2))
```

```
## [1] 0.008894181
```

```
model.lasso<-cv.glmnet(y = train$ten, x = train.matrix, alpha = 1, lambda = 10^seq(2,-2,length = 100), standar
dize = TRUE, nfolds = 5, thresh = 1e-12)
data.frame(lambda = model.lasso$lambda,
cv_mse = model.lasso$cvm,
nonzero_coef = model.lasso$nonzero) %>%
ggplot(aes(x = lambda, y = cv_mse, col = nonzero_coef)) +
geom_point() +
geom_line() +
geom_vline(xintercept = model.lasso$lambda.min, col = "deepskyblue3") +
geom_hline(yintercept = min(model.lasso$cvm), col = "deepskyblue3") +
scale_x_continuous(trans = "log10", breaks = c(0.01, 0.1, 1, 10, 100), labels = c(0.01, 0.1, 1, 10, 100)) +
scale_y_continuous(labels = scales::comma_format()) +
theme(legend.position = "bottom") +
scale_color_gradient(low = "red", high = "green") +
labs(x = "Lambda",
y = "CV MSE",
col = "Coefficients",
title = "Lasso Lambda Selection")
```



```
## Lambda Selection for Lasso Regression
# al
pha at one is a lasso regression
lasso.prediction<-predict(model.lasso.best, s = model.lasso$lambda.min, newx = test.matrix)
(lasso.mse<-mean((lasso.prediction - test$ten)^2)) # MSE for Lasso Regression
```

```
## [1] 0.008903584
```

```
# getting the coefficients below
lasso.weights<-predict(model.lasso.best, type = "coefficients", s = model.lasso$lambda.min)
lasso.weights
```

```
## 41 x 1 sparse Matrix of class "dgCMatrix"
##
## (Intercept) 0.0357731056
## tenlag1 0.9800292598
## tenlag2 0.0110089311
## tenlag3 .
## tenlag4 .
## tenlag5 .
## movederivelag1 .
## movederivelag2 .
## movederivelag3 .
## movederivelag4 .
## movederivelag5 .
## thirtyderivelag1 0.0006813501
## thirtyderivelag2 .
## thirtyderivelag3 .
## thirtyderivelag4 .
## thirtyderivelag5 .
## dxylag1 .
## dxylag2 .
## dxylag3 .
## dxylag4 .
## dxylag5 .
## tenderivativevelag1 .
## tenderivativevelag2 .
## tenderivativevelag3 .
## tenderivativevelag4 .
## tenderivativevelag5 .
## movederivativevelag1 .
## movederivativevelag2 .
## movederivativevelag3 .
## movederivativevelag4 .
## movederivativevelag5 .
## thirtyderivativevelag1 .
## thirtyderivativevelag2 .
## thirtyderivativevelag3 .
## thirtyderivativevelag4 0.0027976196
## thirtyderivativevelag5 .
## dxyderivativevelag1 .
## dxyderivativevelag2 .
## dxyderivativevelag3 .
## dxyderivativevelag4 .
## dxyderivativevelag5 .
```

```
# principal components regression
library(pls)
```

```
##
## Attaching package: 'pls'
```

```
## The following object is masked from 'package:caret':
##
## R2
```

```
## The following object is masked from 'package:stats':
##
## loadings
```

```
# tenderivativevelag1 + thirtylag5 + dxylag1 + dxylag2 + movederivative + movederivativevelag1 + tenlag2
model.princ<-prc(train~., data = train, scale = T, validation = "CV")
model.princ.mse<-MSEP(model.princ, estimate = "CV")$val %>%
reshape2::melt() %>%
mutate(M = 0:(nrow(.)-1)) %>%
select(M, value) %>%
rename(CV_MSE = value)
model.princ.mse
```

```
## M CV_MSE
## 0 3.18845019
## 1 0.25573058
## 2 0.16155876
## 3 0.05979789
## 4 0.05974783
## 5 0.05805221
## 6 0.05114113
## 7 0.04160151
## 8 0.03263453
## 9 0.03027377
## 10 9.0.03819785
## 11 10.0.03830370
## 12 11.0.03836403
## 13 12.0.03836403
## 14 13.0.03834401
## 15 14.0.03807190
## 16 15.0.03782049
## 17 16.0.03773081
## 18 17.0.03762043
## 19 18.0.03759524
## 20 19.0.03805554
## 21 20.0.03674128
## 22 21.0.02453695
## 23 22.0.02303149
## 24 23.0.02020085
## 25 24.0.01712792
## 26 25.0.01714076
## 27 26.0.01714102
## 28 27.0.01707050
## 29 28.0.01712820
## 30 29.0.01707121
## 31 30.0.01704061
## 32 31.0.01709226
## 33 32.0.01714851
## 34 33.0.01719928
## 35 34.0.01739254
## 36 35.0.01778277
## 37 36.0.01789999
## 38 37.0.01811921
## 39 38.0.01820499
## 40 39.0.01855466
## 41 40.0.01849195
```

```
model.princ.mse %>%
mutate(min_CV_MSE = as.numeric(min(CV_MSE) == CV_MSE)) %>%
ggplot(aes(x = M, y = CV_MSE)) +
geom_line(col = "gray55") +
geom_point(size = 2, aes(col = factor(min_CV_MSE))) +
scale_y_continuous(labels = scales::comma_format()) +
scale_color_manual(values = c("deepskyblue3", "green")) +
theme(legend.position = "none") +
labs(x = "M",
y = "Cross-Validation MSE",
col = "Non-Zero Coefficients",
title = "PCR - M Selection (Using 10-Fold Cross-Validation)")
```

```
PCR - M Selection (Using 10-Fold Cross-Validation)
```



```
## Cross Validation picked M = 7
```

```
princ.pred<-predict(model.princ, test, ncomp = 24)
(princ.mse<-mean((princ.pred - test$ten)^2))
```

```
## [1] 0.008799786
```

```
# tenderivativevelag1 + thirtylag5 + dxylag1 + dxylag2 + movederivative + movederivativevelag1 + tenlag2
mod.partial<-pls(train~., data = train, scale = T, validation = "CV")
mod.partial.mse<-MSEP(mod.partial, estimate = "CV")$val %>%
reshape2::melt() %>%
mutate(M = 0:(nrow(.)-1)) %>%
select(M, value) %>%
rename(CV_MSE = value)
mod.partial.mse
```

```
## M CV_MSE
## 0 3.18845019
## 1 0.25573058
## 2 0.16155876
## 3 0.05979789
## 4 0.05974783
## 5 0.05805221
## 6 0.05114113
## 7 0.04160151
## 8 0.03263453
## 9 0.03027377
## 10 9.0.03819785
## 11 10.0.03830370
## 12 11.0.03836403
## 13 12.0.03836403
## 14 13.0.03834401
## 15 14.0.03807190
## 16 15.0.03782049
## 17 16.0.03773081
## 18 17.0.03762043
## 19 18.0.03759524
## 20 19.0.03805554
## 21 20.0.03674128
## 22 21.0.02453695
## 23 22.0.02303149
## 24 23.0.02020085
## 25 24.0.01712792
## 26 25.0.01714076
## 27 26.0.01714102
## 28 27.0.01707050
## 29 28.0.01712820
## 30 29.0.01707121
## 31 30.0.01704061
## 32 31.0.01709226
## 33 32.0.01714851
## 34 33.0.01719928
## 35 34.0.01739254
## 36 35.0.01778277
## 37 36.0.01789999
## 38 37.0.01811921
## 39 38.0.01820499
## 40 39.0.01855466
## 41 40.0.01849195
```

```
mod.partial.mse %>%
mutate(min_CV_MSE = as.numeric(min(CV_MSE) == CV_MSE)) %>%
ggplot(aes(x = M, y = CV_MSE)) +
geom_line(col = "gray55") +
geom_point(size = 2, aes(col = factor(min_CV_MSE))) +
scale_y_continuous(labels = scales::comma_format()) +
scale_color_manual(values = c("deepskyblue3", "green")) +
theme(legend.position = "none") +
labs(x = "M",
y = "Cross-Validation MSE",
col = "Non-Zero Coefficients",
title = "PLS - M Selection (Using 10-Fold Cross-Validation)")
```

```
PLS - M Selection (Using 10-Fold Cross-Validation)
```



```
partial.pred<-predict(mod.partial, test, ncomp = 14)
(partial.mse<-mean((partial.pred - test$ten)^2))
```

```
## [1] 0.0087995
```

```
## Model Comparison
tss<-sum((test$ten - mean(test$ten))^2) # total sum of squares
## data frame with the five models used before
data.frame(method = c("ols", "ridge", "lasso", "pcr", "pls"),
test.mean.squared.errors = c(ols.mse, ridge.mse, lasso.mse, princ.mse, partial.mse),
test.r2 = c(1 - sum((test$ten - ols.prediction)^2) / tss,
1 - sum((test$ten - ridge.prediction)^2) / tss,
1 - sum((test$ten - lasso.prediction)^2) / tss,
1 - sum((test$ten - princ.pred)^2) / tss,
1 - sum((test$ten - partial.pred)^2) / tss)) %>%
arrange(test.mean.squared.errors)
```

```
## method test.mean.squared.errors test.r2
## 1 PLS 0.008794501 0.9823851
## 2 PCR 0.008799786 0.9823785
## 3 OLS 0.008799786 0.9823785
## 4 ridge 0.008894181 0.9821895
## 5 lasso 0.008893584 0.9821707
```