

Predictive Analytics: practical 2 solutions

```
library("caret")
data(FuelEconomy, package = "AppliedPredictiveModeling")
set.seed(25)
```

Cross validation and the bootstrap

- Fit a linear regression model to the cars2010 data set with FE as the response, using EngDispl, NumCyl and NumGears as predictors.

The data set can be loaded
`data("FuelEconomy", package = "AppliedPredictiveModeling")`.

```
mLM = train(FE ~ EngDispl + NumCyl + NumGears, method = "lm", data = cars2010)
```

- What is the training error rate (RMSE) for this model?

Hint: The training error can be found by taking the square root of the average square residuals. The `sqrt` and `resid` functions may be useful.

```
res = resid(mLM)
(trainRMSE = sqrt(mean(res * res)))

## [1] 4.59
```

- Re-train your model using the validation set approach to estimate a test RMSE, make your validation set equivalent to a half of the full available data.

```
## pick an index for samples floor just rounds down so we only try to sample
## a whole number
index = sample(nrow(cars2010), floor(nrow(cars2010)/2))
## set up a train control object
tcVS = trainControl(method = "cv", index = list(Fold1 = (1:nrow(cars2010))[-index]))
## train the model
mLMVS = train(FE ~ EngDispl + NumCyl + NumGears, method = "lm", data = cars2010,
              trControl = tcVS)
```

- How does this compare to the training error that we estimated above?

```
# it's larger, often training error under estimates test error
getTrainPerf(mLMVS)

##   TrainRMSE TrainRsquared method
## 1      4.847         0.616     lm

trainRMSE

## [1] 4.59
```

- Go through the same process using the different methods for estimating test error. That is leave one out and k -fold crossvalidation as well as bootstrapping.

10-fold cross validation can be shown to be a good choice for almost any situation.

```
# set up train control objects
tcL00CV = trainControl(method = "L00CV")
tcKFOLD = trainControl(method = "cv", number = 10)
tcBOOT = trainControl(method = "boot")
```

```
# train the model
mLML00CV = train(FE ~ EngDispl + NumCyl + NumGears, method = "lm", data = cars2010,
  trControl = tcL00CV)
mLMKFOLD = train(FE ~ EngDispl + NumCyl + NumGears, method = "lm", data = cars2010,
  trControl = tcKFOLD)
mLMB00T = train(FE ~ EngDispl + NumCyl + NumGears, method = "lm", data = cars2010,
  trControl = tcB00T)
```

- How do these estimates compare with the validation set approach?

```
getTrainPerf(mLMVS)

##      TrainRMSE TrainRsquared method
## 1      4.847      0.616      lm

getTrainPerf(mLML00CV)

##      TrainRMSE TrainRsquared method
## 1      4.612      0.6214     lm

getTrainPerf(mLMKFOLD)

##      TrainRMSE TrainRsquared method
## 1      4.585      0.6315     lm

getTrainPerf(mLMB00T)

##      TrainRMSE TrainRsquared method
## 1      4.567      0.6326     lm

# all lower than validation set, we mentioned it tended to over estimate
# test error
```

- The object returned by train also contains timing information that can be accessed via the times component of the list. Which of the methods is fastest?

The \$ notation can be used to pick a single list component.

```
mLMVS$times$everything

##      user  system elapsed
## 0.244    0.000    0.247

mLML00CV$times$everything

##      user  system elapsed
## 4.796    0.008    4.823

mLMKFOLD$times$everything

##      user  system elapsed
## 0.280    0.004    0.283

mLMB00T$times$everything

##      user  system elapsed
## 0.344    0.004    0.350
```

- Using k-fold cross validation to estimate test error investigate how the number of folds effects the resultant estimates and computation time.

```
# a number of trainControl objects
tc2 = trainControl(method = "cv", number = 2)
tc5 = trainControl(method = "cv", number = 5)
tc10 = trainControl(method = "cv", number = 10)
tc15 = trainControl(method = "cv", number = 15)
tc20 = trainControl(method = "cv", number = 20)
# train the model using each
mLM2 = train(FE ~ EngDispl + NumCyl + NumGears, method = "lm", data = cars2010,
             trControl = tc2)
mLM5 = train(FE ~ EngDispl + NumCyl + NumGears, method = "lm", data = cars2010,
             trControl = tc5)
mLM10 = train(FE ~ EngDispl + NumCyl + NumGears, method = "lm", data = cars2010,
              trControl = tc10)
mLM15 = train(FE ~ EngDispl + NumCyl + NumGears, method = "lm", data = cars2010,
              trControl = tc15)
mLM20 = train(FE ~ EngDispl + NumCyl + NumGears, method = "lm", data = cars2010,
              trControl = tc20)
# use a data frame to store all of the relevant information
(info = data.frame(Folds = c(2, 5, 10, 15, 20), Time = c(mLM2$times$everything[1],
                                                         mLM5$times$everything[1], mLM10$times$everything[1], mLM15$times$everything[1],
                                                         mLM20$times$everything[1]), Estimate = c(mLM2$results$RMSE, mLM5$results$RMSE,
                                                         mLM10$results$RMSE, mLM15$results$RMSE, mLM20$results$RMSE)))

##   Folds  Time Estimate
## 1     2 0.264   4.595
## 2     5 0.260   4.597
## 3    10 0.280   4.604
## 4    15 0.300   4.583
## 5    20 0.316   4.551

# as there are more folds it takes longer to compute, not an issue with such
# a small model but something to consider on more complicated models.
# Estimates are going down as the number of folds increases. This is
# because for each held out fold we are using a greater proportion of the
# data in training so expect to get a better model.
```

- Experiment with adding terms to the model, transformations of the predictors and interactions say and use cross validation to estimate test error for each. What is the best model you can find? You can still use the validate and mark functions to look at how your models fair on the unseen data.

Penalised regression

The diabetes data set in the lars package contains measurements of a number of predictors to model a response y , a measure of disease progression. There are other columns in the data set which contain interactions so we will extract just the predictors and the response. The data has already been normalized.

```
## load the data in
data(diabetes, package = "lars")
diabetesdata = cbind(diabetes$x, y = diabetes$y)
```

- Try fitting a lasso, ridge and elastic net model using all of the main effects, pairwise interactions and square terms from each of the predictors.¹

```
modelformula = as.formula(paste("y~(.)^2 + ", paste0("I(", colnames(diabetesdata),
  "^2)", collapse = "+")))
mLASSO = train(modelformula, data = diabetesdata, method = "lasso")
mRIDGE = train(modelformula, data = diabetesdata, method = "ridge")
mENET = train(modelformula, data = diabetesdata, method = "enet")
```

¹ Hint: see notes for shortcut on creating model formula. Also be aware that if the predictor is a factor a polynomial term doesn't make sense

- Try to narrow in on the region of lowest RMSE for each model, don't forget about the tuneGrid argument to the train function.

```
# examine previous output then train over a finer grid near the better results
mLASSOfine = train(modelformula, data = diabetesdata, method = "lasso",
  tuneGrid = data.frame(fraction = seq(0.1, 0.5, by = 0.05)))
mLASSOfine$results
```

##	fraction	RMSE	Rsquared	RMSESD	RsquaredSD
## 1	0.10	16.88	0.9531	1.198	0.005822
## 2	0.15	17.06	0.9513	1.371	0.006878
## 3	0.20	17.36	0.9496	1.456	0.007611
## 4	0.25	17.47	0.9490	1.422	0.007500
## 5	0.30	17.56	0.9485	1.415	0.007553
## 6	0.35	17.63	0.9481	1.423	0.007681
## 7	0.40	17.70	0.9477	1.418	0.007710
## 8	0.45	17.75	0.9474	1.413	0.007728
## 9	0.50	17.80	0.9472	1.413	0.007802

best still right down at the 0.1 end

```
mLASSOfiner = train(modelformula, data = diabetesdata, method = "lasso", tuneGrid = data.frame(fraction = seq(0.1, 0.15, by = 0.01)))
mLASSOfiner$results
```

##	fraction	RMSE	Rsquared	RMSESD	RsquaredSD
## 1	0.01	53.80	0.9541	14.736	0.003887
## 2	0.02	38.43	0.9546	18.834	0.004405
## 3	0.03	30.59	0.9545	17.109	0.005125
## 4	0.04	25.97	0.9540	14.843	0.005017
## 5	0.05	23.24	0.9540	12.779	0.004980
## 6	0.06	21.58	0.9532	10.932	0.005182
## 7	0.07	20.60	0.9527	9.168	0.005150
## 8	0.08	19.90	0.9522	7.476	0.005188
## 9	0.09	19.28	0.9517	5.914	0.005419
## 10	0.10	18.72	0.9513	4.526	0.005481
## 11	0.11	18.25	0.9509	3.331	0.005565
## 12	0.12	17.88	0.9505	2.334	0.005592
## 13	0.13	17.62	0.9502	1.615	0.005573
## 14	0.14	17.45	0.9501	1.218	0.005595
## 15	0.15	17.37	0.9499	1.123	0.005772

best is

```
mLASSOfiner$bestTune
```

fraction = 0 is the same as the null model.

$y \sim (.)^2$ is short hand for a model that includes pairwise interactions for each predictor, so if we use this we should only need to add the square terms

```
## fraction
## 15 0.15

mRIDGEfine = train(modelformula, data = diabetesdata, method = "ridge", tuneGrid = data.frame(lambda = seq(0,
  0.1, by = 0.01)))
mRIDGEfine$results

##      lambda  RMSE Rsquared RMSESD RsquaredSD
## 1 0.00 18.54 0.9436 1.5934 0.010896
## 2 0.01 17.16 0.9514 0.7828 0.005701
## 3 0.02 17.07 0.9519 0.7510 0.005412
## 4 0.03 17.09 0.9518 0.7575 0.005467
## 5 0.04 17.18 0.9512 0.7813 0.005677
## 6 0.05 17.31 0.9505 0.8134 0.005966
## 7 0.06 17.49 0.9494 0.8494 0.006297
## 8 0.07 17.69 0.9483 0.8867 0.006651
## 9 0.08 17.91 0.9470 0.9239 0.007016
## 10 0.09 18.16 0.9456 0.9602 0.007385
## 11 0.10 18.41 0.9441 0.9953 0.007756

mRIDGEfiner = train(modelformula, data = diabetesdata, method = "ridge", tuneGrid = data.frame(lambda = seq(0.00
  0.03, by = 0.001)))
mRIDGEfiner$results

##      lambda  RMSE Rsquared RMSESD RsquaredSD
## 1 0.005 16.98 0.9524 1.055 0.005344
## 2 0.006 16.95 0.9526 1.061 0.005314
## 3 0.007 16.92 0.9527 1.067 0.005295
## 4 0.008 16.90 0.9528 1.073 0.005285
## 5 0.009 16.89 0.9529 1.079 0.005282
## 6 0.010 16.88 0.9529 1.084 0.005282
## 7 0.011 16.87 0.9530 1.090 0.005287
## 8 0.012 16.86 0.9530 1.095 0.005294
## 9 0.013 16.86 0.9530 1.101 0.005304
## 10 0.014 16.86 0.9530 1.106 0.005316
## 11 0.015 16.86 0.9530 1.111 0.005330
## 12 0.016 16.86 0.9530 1.116 0.005345
## 13 0.017 16.86 0.9530 1.121 0.005362
## 14 0.018 16.87 0.9529 1.126 0.005380
## 15 0.019 16.87 0.9529 1.131 0.005399
## 16 0.020 16.88 0.9529 1.136 0.005420
## 17 0.021 16.89 0.9528 1.140 0.005441
## 18 0.022 16.90 0.9528 1.145 0.005464
## 19 0.023 16.90 0.9527 1.150 0.005487
## 20 0.024 16.91 0.9527 1.154 0.005511
## 21 0.025 16.92 0.9526 1.159 0.005536
## 22 0.026 16.94 0.9525 1.163 0.005561
## 23 0.027 16.95 0.9525 1.168 0.005587
## 24 0.028 16.96 0.9524 1.172 0.005614
## 25 0.029 16.97 0.9523 1.176 0.005641
## 26 0.030 16.99 0.9522 1.181 0.005669

# the best one
mRIDGEfiner$bestTune

##      lambda
## 10 0.014
```

```

mENETfine = train(modelformula, data = diabetesdata, method = "enet", tuneGrid = expand.grid(lambda = c(0.001,
  0.01, 0.1), fraction = c(0.4, 0.5, 0.6)))
mENETfine$results

##   lambda fraction  RMSE Rsquared RMSESD RsquaredSD
## 1  0.001      0.4 16.26   0.9565 0.8848   0.003724
## 4  0.010      0.4 16.18   0.9581 1.0260   0.003675
## 7  0.100      0.4 22.94   0.9553 2.1321   0.004127
## 2  0.001      0.5 16.64   0.9546 0.8795   0.003647
## 5  0.010      0.5 15.89   0.9584 0.8946   0.003591
## 8  0.100      0.5 17.05   0.9559 0.9928   0.003703
## 3  0.001      0.6 16.84   0.9535 0.9578   0.004200
## 6  0.010      0.6 16.23   0.9566 0.8379   0.003280
## 9  0.100      0.6 16.74   0.9536 0.8435   0.004589

mENETfiner = train(modelformula, data = diabetesdata, method = "enet", tuneGrid = expand.grid(lambda = seq(0.001,
  0.1, length.out = 10), fraction = 0.5))
mENETfiner$results

##   lambda fraction  RMSE Rsquared RMSESD RsquaredSD
## 1  0.001      0.5 16.85   0.9535 0.8808   0.005661
## 2  0.012      0.5 15.98   0.9578 0.7835   0.003957
## 3  0.023      0.5 16.03   0.9576 0.8601   0.003964
## 4  0.034      0.5 16.18   0.9570 0.9269   0.004037
## 5  0.045      0.5 16.36   0.9566 0.9742   0.004042
## 6  0.056      0.5 16.53   0.9562 1.0125   0.003977
## 7  0.067      0.5 16.70   0.9558 1.0340   0.003958
## 8  0.078      0.5 16.84   0.9555 1.0587   0.004075
## 9  0.089      0.5 16.96   0.9552 1.0920   0.004301
## 10 0.100      0.5 17.06   0.9549 1.1178   0.004488

mENETfiner$bestTune

##   fraction lambda
## 2      0.5  0.012

```

We can view the coefficients via

```

coef = predict(mLASSO$finalModel,
  mode = "fraction",
  s = mLASSO$bestTune$fraction, # which ever fraction was chosen as best
  type = "coefficients"
)

```

- How many features have been chosen by the lasso and enet models?

```

# use predict to find the coefficients
coefLASSO = predict(mLASSO$finalModel, mode = "fraction", type = "coefficient",
  s = mLASSO$bestTune$fraction, )
sum(coefLASSO$coefficients != 0)

## [1] 57

coefENET = predict(mENETfiner$finalModel, mode = "fraction", type = "coefficient",
  s = mENET$bestTune$fraction)
sum(coefENET$coefficients != 0)

## [1] 24

```

- How do these models compare to principal components and partial least squares regression?

```

mPCR = train(modelformula, data = diabetesdata, method = "pcr", tuneGrid = data.frame(ncomp = 1:7))
mPLS = train(modelformula, data = diabetesdata, method = "pls", tuneGrid = data.frame(ncomp = 1:7))
mPLS2 = train(modelformula, data = diabetesdata, method = "pls", tuneGrid = data.frame(ncomp = 5:15))
getTrainPerf(mLASSOfiner)

##   TrainRMSE TrainRsquared method
## 1      17.37         0.9499  lasso

getTrainPerf(mRIDGEfiner)

##   TrainRMSE TrainRsquared method
## 1      16.86         0.953   ridge

getTrainPerf(mENETfiner)

##   TrainRMSE TrainRsquared method
## 1      15.98         0.9578  enet

getTrainPerf(mPCR)

##   TrainRMSE TrainRsquared method
## 1      16.53         0.9557   pcr

getTrainPerf(mPLS2)

##   TrainRMSE TrainRsquared method
## 1      16.05         0.9587   pls

# The elastic net model has the lowest estimated test error, all are fairly
# similar. The elastic net model suggests only 21 non-zero coefficients out
# of all of those included in the model.

```