

# Cross-validation and prediction error

Patrick Breheny

February 21

# Introduction

- Today we will discuss the selection of  $\lambda$  as well as the related but somewhat different task of estimating (and obtaining confidence intervals for) the prediction error of a model
- For the lasso, both of these involve tend to revolve around cross-validation, although we will discuss a few different approaches

# Degrees of freedom

- In our discussion of ridge regression, we used information criteria to select  $\lambda$
- All of the criteria we discussed required an estimate of the degrees of freedom of the model
- For linear fitting methods, we saw that  $df = \text{tr}(\mathbf{S})$
- The lasso, however, is not a linear fitting method; there is no exact, closed form solution to  $\text{Cov}(\mathbf{y}, \hat{\mathbf{y}})$

# Degrees of freedom for the lasso

- A natural proposal would be to use  $\text{df}(\lambda) = \|\hat{\beta}(\lambda)\|_0$ , the number of nonzero coefficients
- From one perspective, this might seem to underestimate the true degrees of freedom, as the variables were not prespecified
- For example, in our forward selection example from our first class, we selected 5 features but the true df was  $\approx 19$
- On the other hand, shrinkage reduces the degrees of freedom in an estimator, as we have seen in ridge regression; from this perspective,  $\|\hat{\beta}(\lambda)\|_0$  might seem to overestimate the true degrees of freedom

## Degrees of freedom for the lasso (cont'd)

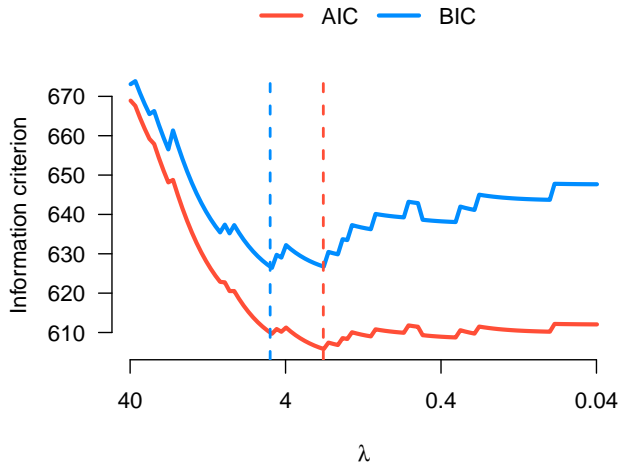
- Surprisingly, it turns out that these two factors exactly cancel and  $\text{df}(\lambda) = \|\widehat{\beta}(\lambda)\|_0$  can be shown to be an unbiased estimate of the lasso degrees of freedom
- Given this estimate, we can then use information criteria such as BIC for the purposes of selecting  $\lambda$

## ncvreg

- To illustrate, we will use the `ncvreg` package to fit the lasso path
- The primary purpose of `ncvreg` is to provide penalties other than the lasso, which we will discuss in our next topic
- However, it provides a `logLik` method, unlike `glmnet`, so it can be used with R's AIC and BIC functions:

```
fit <- ncvreg(X, y, penalty="lasso")  
AIC(fit)  
BIC(fit)
```

# AIC, BIC for pollution data



# Remarks

- As we would expect, BIC applies a stronger penalty for overfitting and chooses a smaller, more parsimonious model than does AIC
- The main advantage of AIC and BIC is that they are computationally convenient: they can be calculated using the fit of the lasso model at very little computational cost
- The primary disadvantage is that both AIC and BIC rely on a number of asymptotic approximations that can be quite inaccurate for high-dimensional data



# Cross-validation: Introduction

- As we have discussed, a reasonable approach to selecting  $\lambda$  in an objective manner is to choose the value of  $\lambda$  that yields the greatest predictive power
- An alternative to the approximations of AIC and BIC is to assess predictive power more directly and empirically through a technique called cross-validation
- Cross-validation is much more reliable, although it comes at an added computation cost

# Sample splitting

- Using the observed agreement between fitted values and the data is too optimistic; we require independent data to test predictive accuracy
- One solution, known as *sample splitting*, is to split the data set into two fractions, a training set and test set, using one portion to estimate  $\hat{\beta}$  (i.e., “train” the model) and the other to evaluate how well  $\mathbf{X}\hat{\beta}$  predicts the observations in the second portion (i.e., “test” the model)
- The problem with this solution is that we rarely have so much data that we can freely part with half of it solely for the purpose of choosing  $\lambda$

# Cross-validation

To finesse this problem, *cross-validation* splits the data into  $K$  folds, fits the data on  $K - 1$  of the folds, and evaluates prediction error on the fold that was left out



Common choices for  $K$  are 5, 10, or  $n$  (also known as leave-one-out cross-validation)

# Cross-validation: Details

- (1) Specify a grid of regularization parameter values  
 $\Lambda = \{\lambda_1, \dots, \lambda_K\}$
- (2) Divide the data into  $V$  roughly equal parts  $D_1, \dots, D_V$
- (3) For each  $v = 1, \dots, V$ , compute the lasso solution path using the observations in  $\{D_u, u \neq v\}$
- (4) For each  $\lambda \in \Lambda$ , compute the mean squared prediction error

$$\text{MSPE}_v(\lambda) = \frac{1}{n_v} \sum_{i \in D_v} \{y_i - \mathbf{x}_i^\top \hat{\boldsymbol{\beta}}_{-v}(\lambda)\}^2,$$

where  $n_v$  is the number of observations in  $D_v$ , as well as

$$\text{CV}(\lambda) = \frac{1}{V} \sum_{v=1}^V \text{MSPE}_v(\lambda).$$

## Cross-validation: Details (cont'd)

- Then  $\hat{\lambda}$  is taken to be the value that minimizes  $CV(\lambda)$  and  $\hat{\beta} \equiv \hat{\beta}(\hat{\lambda})$  the estimator of the regression coefficients
- Note that
  - $MSPE_v(\lambda)$  is the mean squared prediction error for the model based on the training data  $\{D_u, u \neq v\}$  in predicting the response variables in  $D_v$
  - $CV(\lambda)$  is an estimate of ... actually, that turns out to be a good question; what *exactly* is CV estimating?

# What is CV estimating?

- An attractive target would be

$$\mathbb{E}_{\tilde{y}} \left\{ \frac{1}{n} \sum [\tilde{y}_i - \hat{f}(\mathbf{x}_i)]^2 \right\},$$

where  $\tilde{y}$  is an independent copy of  $\mathbf{y}$  based on the same  $\mathbf{X}$  (drawn from the same conditional distribution  $y_i | \mathbf{x}_i$ )

- This is, in words, the expected prediction error over new random outcomes but conditional on the specific data set  $\mathbf{X}, \mathbf{y}$  that we collected
- CV does not estimate this quantity well (as it turns out, this quantity is rather challenging to estimate)

## What is CV estimating? (cont'd)

- Another possibility is the “Same-X” prediction error:

$$\mathbb{E}_{\mathbf{X}, \mathbf{y}, \tilde{\mathbf{y}}} \left\{ \frac{1}{n} \sum [\tilde{y}_i - \hat{f}(\mathbf{x}_i)]^2 \right\};$$

this is what Cp and AIC estimate

- Several papers have shown, however, that CV doesn't estimate this well either
- CV lies much closer to the “Random-X” prediction error:

$$\mathbb{E}_{\mathbf{X}, \mathbf{y}, \mathbf{x}^*, \mathbf{y}^*} \left\{ [\mathbf{y}^* - \hat{f}(\mathbf{x}^*)]^2 \right\},$$

where  $(\mathbf{x}^*, \mathbf{y}^*)$  are drawn from the same distribution that  $\mathbf{X}$  and  $\mathbf{y}$  came from

# Variability of CV estimates

- With this in mind as the estimand (estimation target) of CV, we can think about constructing confidence intervals
- To begin, note that regardless of the number of cross-validation folds, each observation in the data appears exactly once in a test set
- Letting  $\hat{\mu}_i(\lambda) = \mathbf{x}_i^\top \hat{\boldsymbol{\beta}}_{u(i)}(\lambda)$ , the mean of  $\{y_i - \hat{\mu}_i(\lambda)\}_{i=1}^n$  is equal to  $\text{CV}(\lambda)$
- Its variability, however, is useful for estimating the accuracy with which  $\mathbb{E}(\text{MSPE}(\lambda))$  is estimated



## CV standard errors

- Letting  $\text{SD}_{\text{CV}}(\lambda)$  denote the sample standard deviation of the  $\{y_i - \hat{\mu}_i(\lambda)\}^2$  values, the standard error of  $\text{CV}(\lambda)$  is

$$\text{SE}_{\text{CV}}(\lambda) = \frac{\text{SD}_{\text{CV}}(\lambda)}{\sqrt{n}},$$

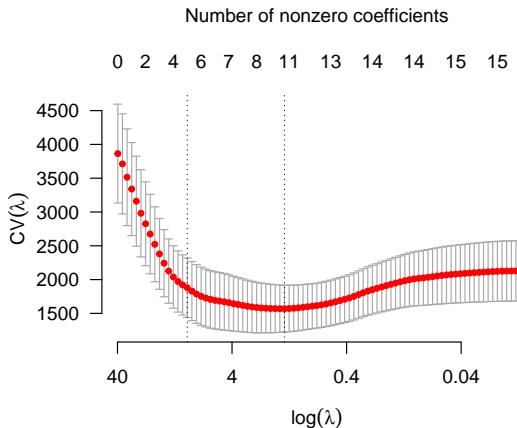
which, in turn, can be used to construct confidence intervals

- The cross-validation procedure described in this section, along with the estimates of  $\text{CV}(\lambda)$  and its standard error, are implemented in `glmnet` and can be carried out using

```
cvfit <- cv.glmnet(X, y)
plot(cvfit)
```

By default, `cv.glmnet` uses  $V = 10$  folds, but this can be changed through the `nfolds` option.

# CV plot for lasso: Pollution data



Intervals are  $\pm 1SE$

## Remarks

- The value  $\lambda = 1.84$  minimizes the cross-validation error, at which point 9 variables are selected
- However, as the confidence intervals show, there is substantial uncertainty about this minimum value
- A fairly wide range of  $\lambda$  values ( $\lambda \in [0.12, 9.83]$ ) yield  $CV(\lambda)$  estimates falling within  $\pm 1SE_{CV}$  of the minimum
- This is almost always the case in model selection: a large number of models could reasonably be considered the “best” model, subject to random variability

# Repeated cross-validation

- Note that  $CV(\lambda)$ , and hence  $\hat{\beta}$ , will change somewhat depending on the random folds
- To avoid this, some people carry out *repeated cross-validation*, and select  $\lambda$  according to the average CV error
- Another option is to carry out  $n$ -fold cross-validation, in which there is only one way to select the fold assignments
- It is important to realize, however, that neither of these approaches does anything to eliminate actual uncertainty with respect to the selection of  $\lambda$

# Nested cross-validation and conformal prediction

- Some recent work has shown that these simple SE calculations have a tendency to be too small, and the confidence intervals for the true prediction error have lower than advertised coverage
- Various solutions to this problem have been proposed, including a nested cross-validation scheme as well as a very different approach altogether called conformal prediction, although both these methods are much more computationally intensive than ordinary CV
- Furthermore, regardless of what exactly CV is estimating, or how accurately it estimates it, picking the  $\lambda$  value that minimizes CV is usually reasonable (which is usually the primary concern in high-dimensional regression)

# Coefficient of determination

- A related goal is estimating the proportion of variance in the outcome that can be explained by the model
- This quantity, familiar from classical regression, is known as the *coefficient of determination* and denoted  $R^2$
- The coefficient of determination is given by

$$R^2 = 1 - \frac{\text{Var}(Y|\mathbf{X})}{\text{Var}(Y)}$$

- Estimation of  $\text{Var}(Y)$  is straightforward
- Estimation of  $\text{Var}(Y|\mathbf{X})$  is (more or less) what CV estimates

## $R^2$ : Calculation in R

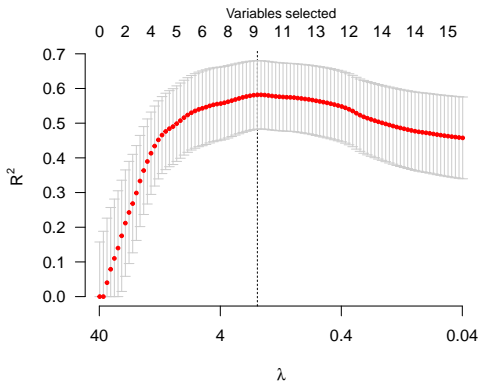
- Once cross-validation is done, calculating  $R^2$  is straightforward
- With `glmnet`:

```
cvfit <- cv.glmnet(X, y)
rsq <- 1-cvfit$cvm/var(y)
```

- Also, the coefficient of determination is available as a plot type in `ncvreg`:

```
cvfit <- cv.ncvreg(X, y, penalty="lasso")
plot(cvfit, type="rsq")
```

# $R^2$ plot: Pollution data



It is worth noting that only a small amount of the explained variability comes from the pollution variables:  $\max R^2 = 0.58$  with the pollution variables;  $\max R^2 = 0.56$  without them



## summary.cv.ncvreg

ncvreg also provides a `summary()` method for its cross-validation object that reports all of this information:

```
> summary(cvfit)
lasso-penalized linear regression with n=60, p=15
At minimum cross-validation error (lambda=1.9762):
-----
Nonzero coefficients: 9
Cross-validation error (deviance): 1591.57
R-squared: 0.58
Signal-to-noise ratio: 1.39
Scale estimate (sigma): 39.895
```