

Statistics 191:  
Introduction  
to Applied  
Statistics

Jonathan  
Taylor  
Department of  
Statistics  
Stanford  
University

# Statistics 191: Introduction to Applied Statistics

Bias-Variance tradeoff: penalized techniques

Jonathan Taylor  
Department of Statistics  
Stanford University

March 10, 2010

# Penalized regression

Statistics 191:  
Introduction  
to Applied  
Statistics

Jonathan  
Taylor  
Department of  
Statistics  
Stanford  
University

## Topics

- Bias-Variance tradeoff.
- Penalized regression.
- Cross-validation.

# Penalized regression

Statistics 191:  
Introduction  
to Applied  
Statistics

Jonathan  
Taylor  
Department of  
Statistics  
Stanford  
University

## Bias-variance tradeoff

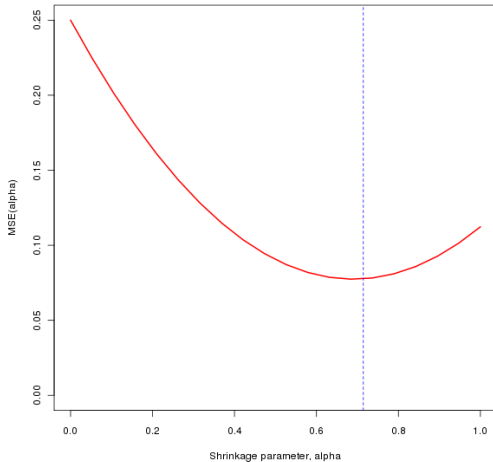
- Arguably, the goal of a regression analysis is to “build” a model that predicts well.
- What does “predict well” mean?

$$\begin{aligned}MSE_{pop}(\mathcal{M}) &= \mathbb{E} \left( (Y_{new} - \hat{Y}_{new, \mathcal{M}})^2 \right) \\&= \text{Var}(Y) + \text{Var}(\hat{Y}_{new, \mathcal{M}}) + \\&\quad \text{Bias}(\hat{Y}_{new, \mathcal{M}})^2.\end{aligned}$$

# Bias-Variance Tradeoff

Statistics 191:  
Introduction  
to Applied  
Statistics

Jonathan  
Taylor  
Department of  
Statistics  
Stanford  
University



# Penalized regression

Statistics 191:  
Introduction  
to Applied  
Statistics

Jonathan  
Taylor  
Department of  
Statistics  
Stanford  
University

## Shrinkage & Penalties

- Shrinkage can be thought of as “constrained” minimization.
- Minimize

$$\sum_{i=1}^n (Y_i - \mu)^2 \quad \text{subject to } \mu^2 \leq C$$

- Lagrange: equivalent to minimizing

$$\sum_{i=1}^n (Y_i - \mu)^2 + \lambda \mu^2$$

for some  $\lambda = \lambda_C$

# Penalized regression

Statistics 191:  
Introduction  
to Applied  
Statistics

Jonathan  
Taylor  
Department of  
Statistics  
Stanford  
University

## Shrinkage & Penalties

- Differentiating:

$$-2 \sum_{i=1}^n (Y_i - \hat{\lambda}_C) + 2\lambda \hat{\mu}_\lambda = 0$$

- Solving

$$\hat{\mu}_\lambda = \frac{\sum_{i=1}^n Y_i}{n + \lambda} = \frac{n}{n + \lambda} \bar{Y}.$$

- As  $\lambda \rightarrow 0$ ,

$$\hat{\mu}_\lambda \rightarrow \bar{Y}.$$

- As  $\lambda \rightarrow \infty$

$$\hat{\mu}_\lambda \rightarrow 0.$$

# Penalized regression

Statistics 191:  
Introduction  
to Applied  
Statistics

Jonathan  
Taylor  
Department of  
Statistics  
Stanford  
University

## Penalties & Priors

- Minimizing

$$\sum_{i=1}^n (Y_i - \mu)^2 + \lambda \mu^2$$

is similar to computing “MLE” of  $\mu$  if the likelihood was proportional to

$$\exp \left( -\frac{1}{2\sigma^2} \left( \sum_{i=1}^n (Y_i - \mu)^2 + \lambda \mu^2 \right) \right).$$

- If  $\lambda = m$ , an integer, then  $\hat{\mu}_\lambda$  is the sample mean of  $(Y_1, \dots, Y_n, 0, \dots, 0) \in \mathbb{R}^{n+m}$ .
- This is equivalent to adding some data with  $Y = 0 \dots$

# Penalized regression

Statistics 191:  
Introduction  
to Applied  
Statistics

Jonathan  
Taylor  
Department of  
Statistics  
Stanford  
University

## Biased regression: penalties

- Not all biased models are better – we need a way to find “good” biased models.
- Inference ( $F$ ,  $\chi^2$  tests, etc) is not quite exact for biased models.
- Generalized one sample problem: penalize large values of  $\beta$ . This should lead to “multivariate” shrinkage of the vector  $\beta$ .
- Heuristically, “large  $\beta$ ” is interpreted as “complex model”. Goal is really to penalize “complex” models, i.e. Occam’s razor.
- If truth really is complex, this may not work! (But, it will then be hard to build a good model anyways ... (statistical lore))



# Penalized regression

Statistics 191:  
Introduction  
to Applied  
Statistics

Jonathan  
Taylor  
Department of  
Statistics  
Stanford  
University

## How much to shrink, choosing $\lambda$

- In our one-sample example,

$$\begin{aligned}MSE_{pop}(\lambda) &= \text{Var}(\lambda \bar{Y}) + \text{Bias}(\lambda \bar{Y})^2 \\&= \frac{\lambda^2 \sigma^2}{n} + \mu^2(1 - \lambda)^2\end{aligned}$$

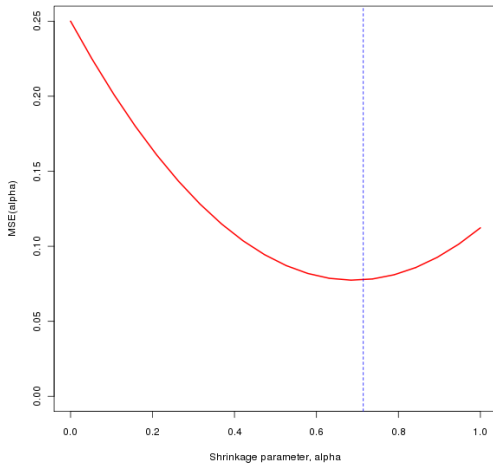
- Differentiating:

$$\begin{aligned}0 &= -2\mu^2(1 - \lambda^*) + 2\frac{\lambda^* \sigma^2}{n} \\ \lambda^* &= \frac{\mu^2}{\mu^2 + \sigma^2/n}\end{aligned}$$

# Bias-Variance Tradeoff

Statistics 191:  
Introduction  
to Applied  
Statistics

Jonathan  
Taylor  
Department of  
Statistics  
Stanford  
University



# Penalized regression

Statistics 191:  
Introduction  
to Applied  
Statistics

Jonathan  
Taylor  
Department of  
Statistics  
Stanford  
University

## Ridge regression

- Assume that columns  $(X_j)_{1 \leq j \leq p-1}$  have zero mean, and length 1 and  $Y$  has zero mean.
- This is called the *standardized model*.
- Minimize

$$SSE_\lambda(\beta) = \sum_{i=1}^n \left( Y_i - \sum_{j=1}^{p-1} X_{ij} \beta_j \right)^2 + \lambda \sum_{j=1}^{p-1} \beta_j^2.$$

- Corresponds (through Lagrange multiplier) to a quadratic constraint on  $\beta$ 's.

# Penalized regression

Statistics 191:  
Introduction  
to Applied  
Statistics

Jonathan  
Taylor  
Department of  
Statistics  
Stanford  
University

## Solving the normal equations

- Normal equations

$$\frac{\partial}{\partial \beta_l} SSE_\lambda(\beta) = -2\langle Y - X\beta, X_l \rangle + 2\lambda\beta_l$$

- 

$$-2\langle Y - X\hat{\beta}_\lambda, X_l \rangle + 2\lambda\hat{\beta}_{l,\lambda} = 0, \quad 1 \leq l \leq p-1$$

- In matrix form

$$-Y^t X + \hat{\beta}_\lambda^t (X^t X + \lambda I) = 0$$

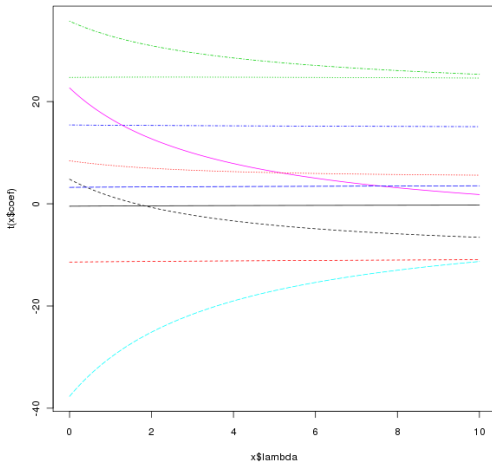
- Or

$$\hat{\beta}_\lambda = (X^t X + \lambda I)^{-1} X^t Y.$$

# Ridge regression

Statistics 191:  
Introduction  
to Applied  
Statistics

Jonathan  
Taylor  
Department of  
Statistics  
Stanford  
University



# Penalized regression

Statistics 191:  
Introduction  
to Applied  
Statistics

Jonathan  
Taylor  
Department of  
Statistics  
Stanford  
University

## LASSO regression

- Another popular penalized regression technique.
- Use the standardized model.
- Minimize

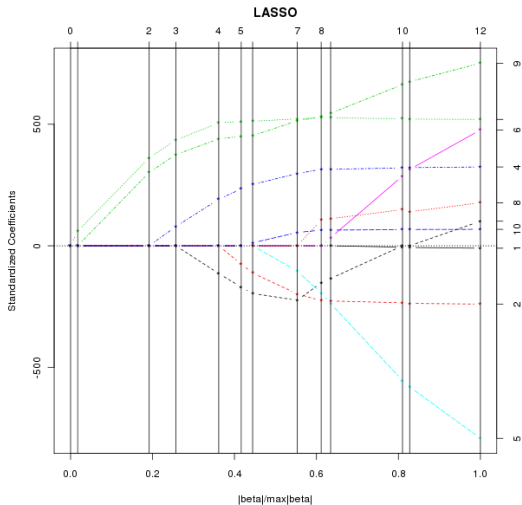
$$SSE_{\lambda}(\beta) = \sum_{i=1}^n \left( Y_i - \sum_{j=1}^{p-1} X_{ij} \beta_j \right)^2 + \lambda \sum_{j=1}^p |\beta_j|.$$

- Corresponds (through Lagrange multiplier) to an  $\ell^1$  constraint on  $\beta$ 's. In theory, it works well when many  $\beta_j$ 's are 0 and gives “sparse” solutions unlike ridge.

# LASSO

Statistics 191:  
Introduction  
to Applied  
Statistics

Jonathan  
Taylor  
Department of  
Statistics  
Stanford  
University



# Penalized regression

Statistics 191:  
Introduction  
to Applied  
Statistics

Jonathan  
Taylor  
Department of  
Statistics  
Stanford  
University

## Choosing $\lambda$ : cross-validation

- If we knew  $MSE$  as a function of  $\lambda$  then we would simply choose the  $\lambda$  that minimizes  $MSE$ .
- To do this, we need to estimate  $MSE$ .
- A popular method is “cross-validation.” Breaks the data up into smaller groups and uses part of the data to predict the rest.
- We saw this in diagnostics: i.e. Cook’s distance measured the fit with and without each point in the data set.



# Penalized regression

Statistics 191:  
Introduction  
to Applied  
Statistics

Jonathan  
Taylor  
Department of  
Statistics  
Stanford  
University

## $K$ -fold cross-validation

- Fix a model (i.e. fix  $\lambda$ ). Break data set into  $K$  approximately equal sized groups  $(G_1, \dots, G_K)$ .
- for (i in 1:K) Use all groups except  $G_i$  to fit model, predict outcome in group  $G_i$  based on this model  $\hat{Y}_{j(i),\lambda}, j \in G_i$ .
- Estimate

$$CV(\lambda) = \frac{1}{n} \sum_{i=1}^K \sum_{j \in G_i} (Y_j - \hat{Y}_{j(i),\lambda})^2.$$

# Penalized regression

Statistics 191:  
Introduction  
to Applied  
Statistics

Jonathan  
Taylor  
Department of  
Statistics  
Stanford  
University

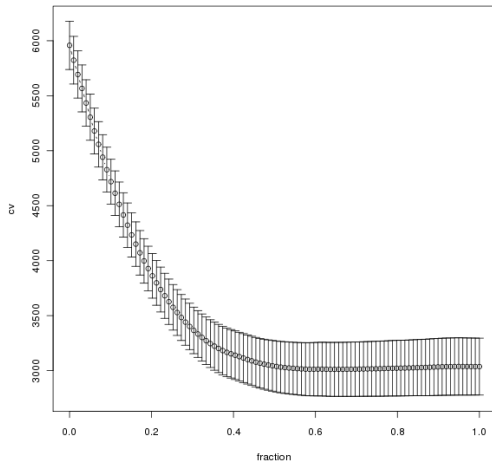
## *K*-fold cross-validation (continued)

- It is a general principle that can be used in other situations, not just for Ridge.
- Pros (partial list): “objective” measure of a model.
- Cons (partial list): inference is, strictly speaking, “out the window” (also true for other model selection procedures in theory).
- If goal is not really inference about certain specific parameters, it is a reasonable way to compare models.

# LASSO: cross-validation

Statistics 191:  
Introduction  
to Applied  
Statistics

Jonathan  
Taylor  
Department of  
Statistics  
Stanford  
University



# Penalized regression

Statistics 191:  
Introduction  
to Applied  
Statistics

Jonathan  
Taylor  
Department of  
Statistics  
Stanford  
University

## Generalized Cross Validation

- A computational shortcut for  $n$ -fold cross-validation (also known as leave-one out cross-validation).

- Let

$$S_\lambda = (X^t X + \lambda I)^{-1} X^t$$

be the matrix in ridge regression.

- Then

$$GCV(\lambda) = \frac{\|Y - S_\lambda Y\|^2}{n - \text{Tr}(S_\lambda)}.$$

- The quantity  $\text{Tr}(S_\lambda)$  is the *effective degrees of freedom*.

# Ridge regression

Statistics 191:  
Introduction  
to Applied  
Statistics

Jonathan  
Taylor  
Department of  
Statistics  
Stanford  
University

