

# Modern Statistics and Machine Learning: Regularized Regression

37105 Data Science for Marketing Decision Making  
Günter J. Hitsch  
Chicago Booth

Winter 2018

1 / 38

## Overview

1. The bias-variance tradeoff
2. Ridge regression — regularization and shrinkage
3. LASSO — variable selection
4. Cross-validation
5. Elastic nets
6. Classification

2 / 38

## The bias-variance trade-off

The data set:  $\mathcal{D} = ((y_1, \mathbf{x}_1), \dots, (y_n, \mathbf{x}_n))$

To simplify the notation,  $\mathbf{x}_i = (x_{i1}, \dots, x_{ip})$  contains the values for all inputs

$r(\mathbf{x})$  is the regression function,

$$r(\mathbf{x}) = \mathbb{E}[Y|\mathbf{x}]$$

We often propose some model  $f(\mathbf{x}, \boldsymbol{\beta})$  for the regression function and use the data  $\mathcal{D}$  to estimate the model parameters,  $\boldsymbol{\beta}$ .

$\hat{r}(\mathbf{x})$  is the estimated regression function, given the *estimated* model parameters  $\hat{\boldsymbol{\beta}}$ :

$$\hat{r}(\mathbf{x}) = f(\mathbf{x}, \hat{\boldsymbol{\beta}})$$

In the linear regression model:

$$r(\mathbf{x}) = \beta_0 + \sum_{k=1}^p x_k \beta_k$$

3 / 38

The **variance** of the estimated regression function at a specific value of the inputs,  $\mathbf{x}$ , is

$$\begin{aligned} \text{var}(\hat{r}(\mathbf{x})) &= \mathbb{E} \left[ (\hat{r}(\mathbf{x}) - \mathbb{E}[\hat{r}(\mathbf{x})])^2 \right] \\ &= \mathbb{E} \left[ \hat{r}(\mathbf{x})^2 - 2\hat{r}(\mathbf{x})\mathbb{E}[\hat{r}(\mathbf{x})] + \mathbb{E}[\hat{r}(\mathbf{x})]^2 \right] \\ &= \mathbb{E}[\hat{r}(\mathbf{x})^2] - \mathbb{E}[\hat{r}(\mathbf{x})]^2 \end{aligned}$$

The **bias** of the prediction is

$$\begin{aligned} \text{bias}(\hat{r}(\mathbf{x})) &= \mathbb{E}[\hat{r}(\mathbf{x}) - r(\mathbf{x})] \\ &= \mathbb{E}[\hat{r}(\mathbf{x})] - r(\mathbf{x}) \end{aligned}$$

4 / 38

The estimated regression function allows us to predict the output  $\hat{Y} = \hat{r}(\mathbf{x})$  for the input values  $\mathbf{x}$ . The accuracy of this prediction can be measured using the mean-squared error (MSE)

$$\begin{aligned}\mathbb{E}[(r(\mathbf{x}) - \hat{r}(\mathbf{x}))^2] &= \mathbb{E}[r(\mathbf{x})^2 - 2r(\mathbf{x})\hat{r}(\mathbf{x}) + \hat{r}(\mathbf{x})^2] \\ &= \mathbb{E}[r(\mathbf{x})^2 - 2r(\mathbf{x})\hat{r}(\mathbf{x}) + \hat{r}(\mathbf{x})^2 + \mathbb{E}[\hat{r}(\mathbf{x})]^2 - \mathbb{E}[\hat{r}(\mathbf{x})]^2] \\ &= \mathbb{E}[\hat{r}(\mathbf{x})^2] - \mathbb{E}[\hat{r}(\mathbf{x})]^2 + (r(\mathbf{x})^2 - 2r(\mathbf{x})\mathbb{E}[\hat{r}(\mathbf{x})] + \mathbb{E}[\hat{r}(\mathbf{x})]^2) \\ &= (\mathbb{E}[\hat{r}(\mathbf{x})^2] - \mathbb{E}[\hat{r}(\mathbf{x})]^2) + (\mathbb{E}[\hat{r}(\mathbf{x})] - r(\mathbf{x}))^2\end{aligned}$$

Hence, we see that

$$\mathbb{E}[(r(\mathbf{x}) - \hat{r}(\mathbf{x}))^2] = \text{var}(\hat{r}(\mathbf{x})) + \text{bias}(\hat{r}(\mathbf{x}))^2$$

5 / 38

Furthermore, the realized outcome can be written as

$$Y = r(\mathbf{x}) + \epsilon,$$

where  $\mathbb{E}[\epsilon] = 0$  and  $\mathbb{E}[h(\mathbf{x})\epsilon] = 0$  for any function  $h$  of  $\mathbf{x}$ .

Hence, the mean-squared error between the output  $Y$  and the prediction  $\hat{Y} = \hat{r}(\mathbf{x})$  is

$$\begin{aligned}\text{MSE} &= \mathbb{E}[(Y - \hat{r}(\mathbf{x}))^2] \\ &= \mathbb{E}[(r(\mathbf{x}) + \epsilon - \hat{r}(\mathbf{x}))^2] \\ &= \mathbb{E}[(r(\mathbf{x}) - \hat{r}(\mathbf{x}))^2] + \mathbb{E}[\epsilon^2]\end{aligned}$$

Note that  $\mathbb{E}[\epsilon^2] = \text{var}(\epsilon^2)$  because  $\epsilon$  has mean 0.

Hence, we see that

$$\text{MSE} = \text{var}(\epsilon^2) + \text{var}(\hat{r}(\mathbf{x})) + \text{bias}(\hat{r}(\mathbf{x}))^2$$

6 / 38

$$\text{MSE} = \text{var}(\epsilon^2) + \text{var}(\hat{r}(\mathbf{x})) + \text{bias}(\hat{r}(\mathbf{x}))^2$$

In words, the MSE is given by the variance of the error term plus the variance of the prediction plus the squared bias of the prediction.

The accuracy of the prediction is worse if

- (i) The inherent variance in  $Y$  through  $\epsilon$  (the irreducible error) increases
- (ii) The sampling variability of the estimated regression function  $\hat{r}(\mathbf{x})$  increases
- (iii) The systematic bias in the prediction of  $\hat{Y} = \hat{r}(\mathbf{x})$  increases

Note: We can do nothing about the variance of the error term  $\epsilon$ . Even if we knew the regression function  $r(\mathbf{x})$ , we would still have a prediction error through the variance in  $Y$  through  $\epsilon$ .

7 / 38

To be perfectly clear what the expectation used to characterize the mean-squared error means, we should write

$$\text{MSE} = \mathbb{E}_{\mathcal{D}, Y} [(Y - \hat{r}(\mathbf{x}))^2]$$

The expectation is taken over realized training data sets,  $\mathcal{D}$ , and realized prediction cases,  $Y$ .

Hence, the MSE is an average (expected) measure of the prediction accuracy over many data sets and prediction problems = applications of the estimation method.

8 / 38

## Special case: Linear regression model

Suppose the key condition, mean-independence of the error term is satisfied:

$$\mathbb{E}[\epsilon|\mathbf{x}] = 0$$

We already know that in this case the OLS estimator is unbiased:

$$\begin{aligned}\mathbb{E}[\hat{\beta}] &= \beta \\ \Rightarrow \text{bias}(\hat{r}(\mathbf{x})) &= \mathbb{E}[\hat{r}(\mathbf{x})] - r(\mathbf{x}) = 0\end{aligned}$$

Also, according to the *Gauss-Markov Theorem*, the OLS estimator is BLUE = the **best linear unbiased estimator**, i.e. the estimator with the smallest variance among all linear unbiased estimators.

Implication: If we require our estimator to be unbiased, then the *OLS estimator is the best linear estimator*

9 / 38

The “old school” thinking was that unbiasedness was the most important quality measure of an estimator. Non-linear estimation methods were hard to work with, and hence the OLS estimator was it.

But modern statistics recognizes that the price for an unbiased estimator can be high. If we tolerate some (hopefully small) amount of bias, we can find alternative estimators with lower variance, and the overall quality of the prediction will increase.

**Regularization** is a key technique to construct such estimators.

10 / 38

## Ridge regression

In a ridge regression the estimated coefficients,  $\beta_{\text{ridge}}$ , are found by minimizing the objective:

$$\min_{\beta} \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{k=1}^p x_{ik} \beta_k \right)^2 + \lambda \sum_{k=1}^p \beta_k^2$$

The first term is the sum of squared residuals (RSS) — well known from the linear regression model.

The second term is a **penalty term**, which depends on the magnitude of the coefficients ( $\beta_k^2$ ) and the tuning parameter  $\lambda \geq 0$ . The penalty allows for **regularization**.

Note: There is no penalty for the intercept.

11 / 38

If  $\lambda = 0$ , the penalty is zero, and minimization of the RSS yields the OLS (ordinary least-squares) estimator. To minimize the objective, the regression coefficients,  $\beta_{\text{OLS}}$ , are chosen to make the squared residuals as small as possible.

If  $\lambda > 0$ , then the penalty term increases. Therefore, increasing the magnitude of a coefficient (measured by the square  $\beta_k^2$ ) away from 0 has a cost. Minimizing the squared residuals needs to be balanced with minimizing the cost of the coefficient values.

Note: A ridge regression should be performed based on *standardized (scaled) inputs*:

$$x'_{ik} = \frac{x_{ik}}{\text{sd}(x_k)}$$

This makes the magnitudes and hence the cost of all inputs comparable.

12 / 38

## Regularization and shrinkage methods

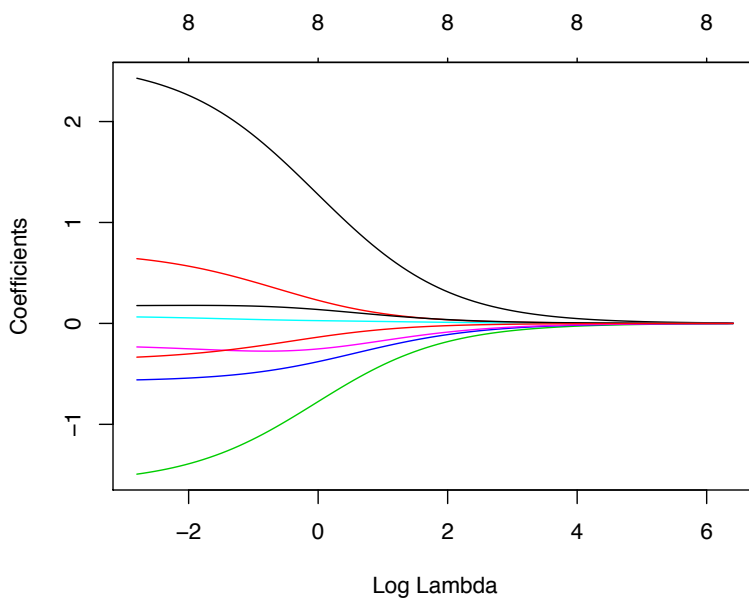
Ridge estimators are **shrinkage estimators**. If  $\lambda$  increases, the regression coefficients  $\beta_k$  are shrunk towards zero.

We thus achieve **regularization**, in the sense that we constrain or stabilize the  $\beta_k$  values.

At the limit,  $\lambda = \infty$ , the price to be paid for the coefficients is so high that they will all be set to zero.

13 / 38

Example: Regression with 8 inputs. As  $(\log) \lambda$  increases, the coefficients are shrunk towards zero.



14 / 38

## Model exploration

Explore how the model works using the examples in the `Regularization-Variable-Selection.Rmd` R Markdown file

Compare the OLS estimates with the ridge regression estimates: Even parameters that do not influence the output,  $\beta_k = 0$ , will sometimes have large magnitudes and be statistically different from zero. This phenomenon is called **overfitting**. Overfitting is more likely to occur if the model is very complex, i.e. includes a large number of inputs that have no (or very little) impact on the output. Regularization reduces the model complexity and ameliorates overfitting.

Compare the predicted MSE (mean-squares error) of OLS with the predict MSE of the ridge regression.

Compare the results under different degrees of model complexity (change number of inputs).

15 / 38

## Selection of tuning parameter $\lambda$

Can we use the training data to directly estimate both the regression coefficients  $\beta$  and the tuning parameter  $\lambda \geq 0$ ?

$$\min_{\beta, \lambda \geq 0} \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{k=1}^p x_{ik} \beta_k \right)^2 + \lambda \sum_{k=1}^p \beta_k^2$$

What estimate of  $\lambda$  do you expect if you try this approach?

16 / 38



## Cross-validation

Take the training data and *randomly* split them into  $K$  **folds**.



Most commonly used:  $K = 10$



Choose some value for the tuning parameter  $\lambda$ .

Pick one of the folds,  $k$ , and set it aside as a validation data set.

Estimate the model only using the data in the other folds,

$1, \dots, k-1, k+1, \dots, 10$

Then predict the output  $y_i$  in fold  $k$  based on the model estimates, and record the MSE (mean-squared error).

Repeat for all folds  $k$ , and compute the average MSE over all folds. We thus obtain an **out-of-sample prediction error** for the chosen tuning parameter  $\lambda$ .

17 / 38

## Selection of $\lambda$

Select a range of tuning parameter  $\lambda$  values,  $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_L)$ , and split the training data into  $K$  folds.

Iterate over all folds  $k = 1, \dots, K$

► For each value of  $\lambda_l$  in  $\boldsymbol{\lambda}$ :

1. Train (estimate) the model using all folds but fold  $k$
2. Predict the output,  $y_{\lambda_l, -k}(\mathbf{x}_i)$ , for all observations  $\mathbf{x}_i$  in fold  $k$ .  
Then calculate the total prediction error in fold  $k$ ,

$$e_k(\lambda_l) = \sum_{i \in \text{fold}(k)} (y_i - y_{\lambda_l, -k}(\mathbf{x}_i))^2$$

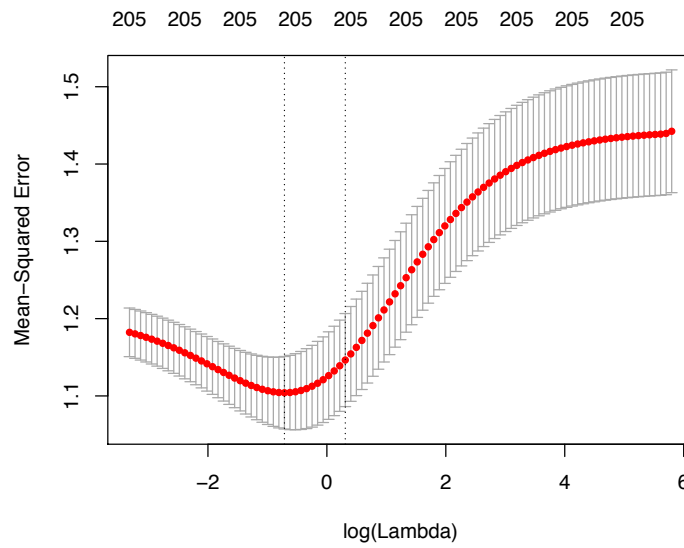
Calculate the cross-validation error for each  $\lambda_l$ ,

$$CV(\lambda_l) = \frac{1}{n} \sum_{k=1}^K e_k(\lambda_l)$$

$CV(\lambda_l)$  is the out-of-sample MSE given the tuning parameter  $\lambda_l$

18 / 38

Select tuning parameter  $\lambda$  based on  $CV(\lambda)$ , the **cross-validation error curve**



Plot indicates  $\lambda_{\min}$  that minimizes  $CV(\lambda)$  and  $\lambda_{1se}$  for most regularized model such that  $CV(\lambda_{1se})$  is within one standard error of the minimum.

19 / 38

The ridge regression is able to achieve a smaller out-of-sample prediction error by imposing a penalty on the regression coefficients through the tuning parameter  $\lambda$  — regularization.

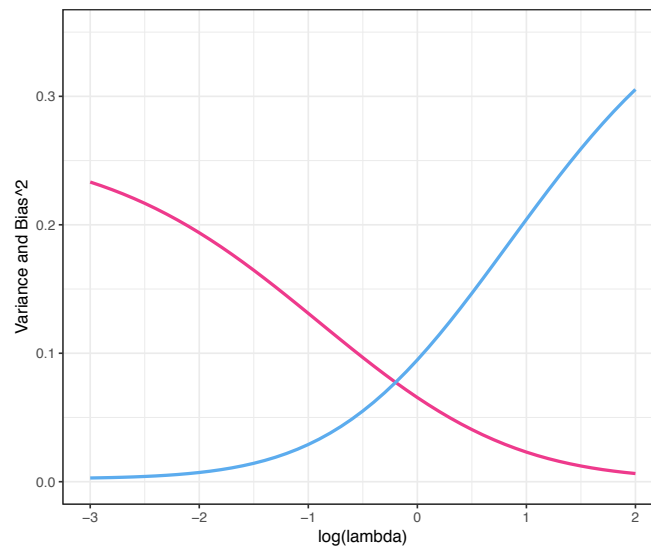
Hence, a ridge regression is able to predict better than the OLS estimator, which is unbiased under sufficient conditions.

Why can a biased estimator predict better than a biased estimator?

20 / 38

## Bias-variance tradeoff

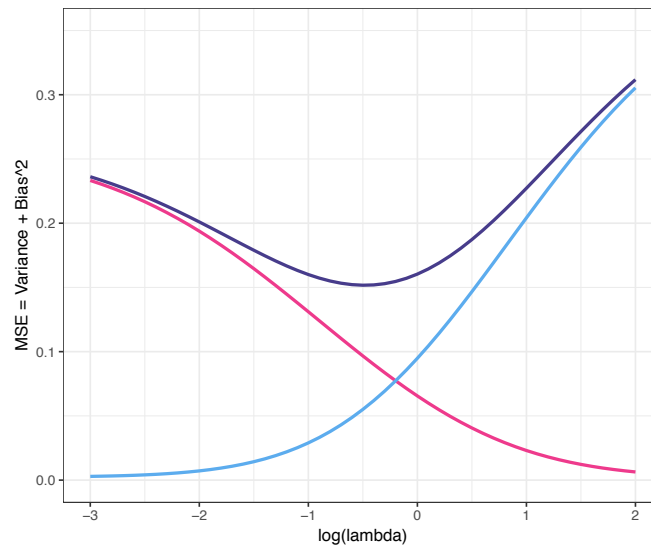
Plot shows the squared bias,  $\text{bias}(\hat{r}(\mathbf{x}))^2$ , and variance,  $\text{var}(\hat{r}(\mathbf{x}))$ , conditional on  $(\log) \lambda$  (averaged over many data sets)



Pink:  $\text{var}(\hat{r}(\mathbf{x}))$ , blue:  $\text{bias}(\hat{r}(\mathbf{x}))^2$

21 / 38

Ignoring the irreducible error variance  $\text{var}(\epsilon)$ , the mean-squared error is  $\text{MSE} = \text{var}(\hat{r}(\mathbf{x})) + \text{bias}(\hat{r}(\mathbf{x}))^2$ . By trading off an increase in the bias in exchange for a reduction in the variance of the prediction we can lower the prediction error.



Dark blue: MSE

22 / 38

## Variable selection

Especially when  $p$  is large compared to  $n$ , or if there are good reasons to believe that many of the coefficients are zero,  $\beta_k = 0$ , we are interested in variable selection methods.

Variable selection: Estimator sets some of the estimated coefficients to zero.

Can a ridge regression perform variable selection? — Examine the marginal cost of increasing  $\beta_k$  in the ridge regression penalty factor:

$$\frac{d}{d\beta_k} \lambda \sum_{k=1}^p \beta_k^2 = 2\lambda\beta_k$$

Hence, the marginal cost of increasing  $\beta_k$  at  $\beta_k = 0$  is  $2\lambda \cdot 0 = 0$ . Therefore, even though a ridge regression shrinks the regression coefficients, it will never set all coefficients to exactly zero.

This can be a disadvantage if many coefficient are in fact zero. Also, the regression results can be harder to **interpret** if no variable selection is performed.

23 / 38

## The LASSO

The LASSO estimator,  $\beta_{\text{LASSO}}$ , minimizes the objective:

$$\min_{\beta} \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{k=1}^p x_{ik} \beta_k \right)^2 + \lambda \sum_{k=1}^p |\beta_k|$$

This appears very similar to a ridge regression. Ridge regression uses the penalty  $\sum \beta_k^2$ , called an  $\ell_2$  norm, while the LASSO uses the  $\ell_1$  norm  $\sum_k |\beta_k|$ .

LASSO stands for *Least Absolute Selection and Shrinkage Operator*.

Notes:

- ▶ The intercept is not penalized
- ▶ If the tuning parameter  $\lambda = 0$  we obtain the OLS estimator
- ▶ To make the inputs comparable they need to be *standardized (scaled)*:

$$x'_{ik} = \frac{x_{ik}}{\text{sd}(x_k)}$$

24 / 38

## Regularization and shrinkage methods

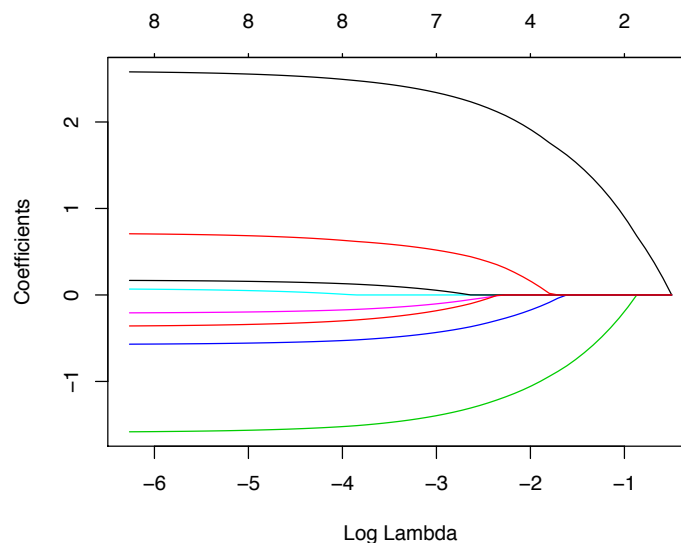
The LASSO is qualitatively different from a ridge regression. In particular, the marginal cost of increasing  $\beta_k$  (in absolute value) is  $\lambda$  and does not vary with the level of  $\beta_k$ . Hence, even at  $\beta_k = 0$  the marginal cost is  $\lambda$ .

Like the ridge regression estimator, the LASSO is a shrinkage estimator: As  $\lambda$  increases, the regression coefficients  $\beta_k$  are shrunk towards zero (**regularization**).

However, because the cost of the coefficient values is always constant at  $\lambda$  (even for tiny values of  $\beta_k$ ) the LASSO will set the coefficients exactly to zero if  $\lambda$  is large enough.

25 / 38

Example: LASSO regression with 8 inputs. Once  $(\log) \lambda$  is large enough, all coefficients are set to exactly zero.



The variable-selection property of the LASSO can be useful to **interpret** the regression results.

26 / 38

## Model exploration

Explore how the model works using the examples in `Regularization-Variable-Selection.Rmd`.

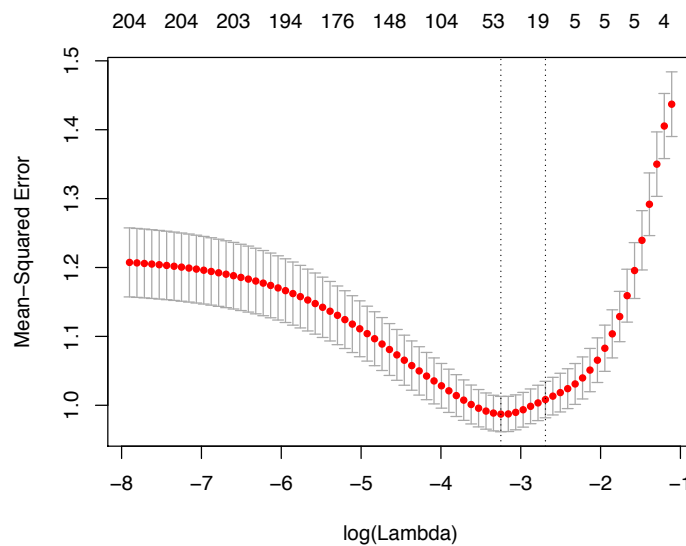
Compare the LASSO with the ridge regression estimates.

Compare the prediction MSE (mean square error) of OLS, ridge regression, and the LASSO.

Compare the results under different degrees of model complexity (number of inputs) and different parameter vectors.

27 / 38

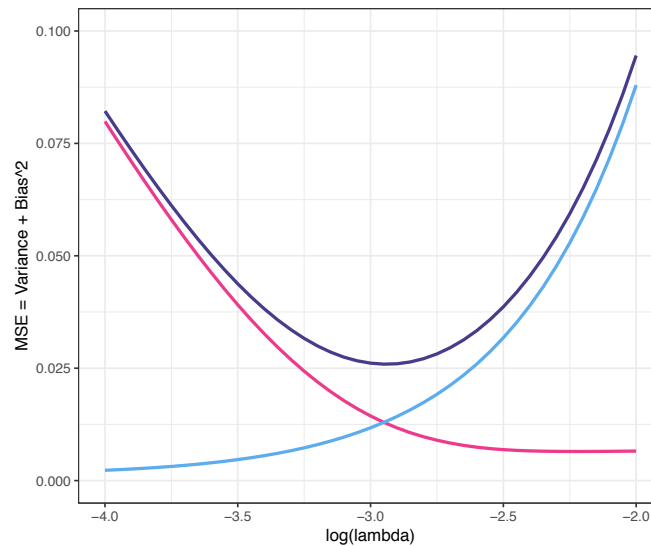
## Tuning $\lambda$ using the cross-validation error curve



Plot indicates  $\lambda_{\min}$  that minimizes  $CV(\lambda)$  and  $\lambda_{1se}$  for most regularized model such that  $CV(\lambda_{1se})$  is within one standard error of the minimum.

28 / 38

## Bias-variance tradeoff



29 / 38

## Other variable selection methods: Best subset selection

Algorithm:

1. Fit the null model  $\mathcal{M}_0$  that contains no inputs—fit based on sample mean of all observations
2. For all values  $k = 1, \dots, p$ :
  - (i) Fit all models that contain exactly  $k$  inputs
  - (ii) Find the model with  $k$  inputs that has the smallest RSS, and call it  $\mathcal{M}_k$
3. Compare all models  $\mathcal{M}_0, \mathcal{M}_1, \dots, \mathcal{M}_p$  and choose the one that is “best” based on one of the following:
  - ▶ Cross-validated error
  - ▶ AIC (Akaike information criterion) or BIC (Bayesian information criterion)

30 / 38

Limitations: Need to evaluate fit for  $2^p$  subsets

With  $p = 25$  inputs,  $2^{25} = 33,554,432$  (OK)

$p = 50$ :  $2^{50} = 1,125,899,906,842,624$  (more than 1,125 trillion) subsets

As  $p$  becomes large, best subset selection becomes computationally infeasible

31 / 38

## Alternatives to best subset selection

Forward stepwise selection

- ▶ Start with the null model  $\mathcal{M}_0$  that contains no inputs
- ▶ Successively add predictors to improve model fit

Backward stepwise selection

- ▶ Start with the full model that contains all inputs,  $\mathcal{M}_p$
- ▶ Successively remove predictors

These approaches are computationally much less intensive, but they are not guaranteed to find the best subset of predictors.

For details, consult *An Introduction to Statistical Learning* (James et al.) or *The Elements of Statistical Learning* (Hastie et al.)

32 / 38



## Variable selection based on $p$ -values

Why not select variables based on  $p$ -values, i.e. if an input is statistically different from 0 at the level  $\alpha$  (e.g.  $\alpha = 0.05$ )?

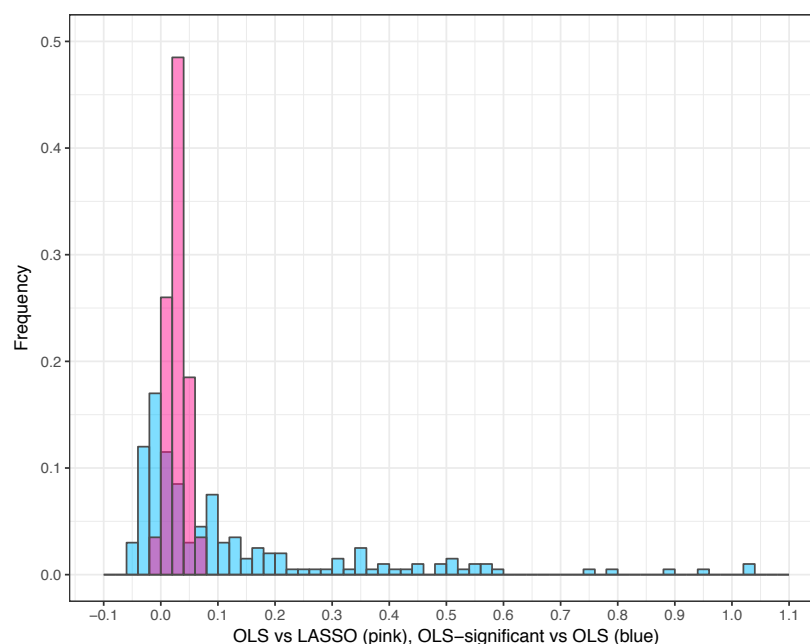
This is a very tempting approach.

However, variable selection based on  $p$ -values is a bad idea:

- ▶ A hypothesis test,  $H_0 : \beta_k = 0$ , is not intended to directly assess how well input  $x_k$  predicts the output.
- ▶ We select variables based on an arbitrary significance level  $\alpha$ , not a tuning parameter  $\lambda$  that is chosen to minimize the out-of-sample prediction error.
- ▶ We include inputs without predictive power if the estimated value happens to be large (i.e. when the inclusion is particularly damaging) and we exclude inputs if the estimated value happens to be small, even though out-of-sample the input explains much of the variation in the data.

33 / 38

Example in `Regularization-Variable-Selection.Rmd`: Compare difference in MSE between (i) OLS and LASSO and (ii) OLS with  $p$ -value based selection and OLS across 200 simulated data sets



34 / 38

The LASSO does best, and often the  $p$ -value selection method does very poorly compared to OLS:

	MSE
OLS	1.107
OLS/ $p$ -value sel.	1.229
LASSO	1.079

35 / 38

## Elastic net

An elastic net can be thought of as a combination of a ridge regression and a LASSO. The penalty factor in an elastic net has the form:

$$\lambda \left( \sum_{k=1}^p [\alpha |\beta_k| + (1 - \alpha) \beta_k^2] \right)$$

for a parameter  $\alpha$  such that  $0 \leq \alpha \leq 1$ .

Special cases:  $\alpha = 1$  (LASSO),  $\alpha = 0$  (ridge regression). The prediction quality of the model may be improved if  $0 < \alpha < 1$  — test this in `Regularization-Variable-Selection.Rmd`.

The  $\alpha$  parameter can be supplied using the `alpha` option in `glmnet` or in `cv.glmnet`.

36 / 38

## Regularized regression for classification

Regularized regression models can easily be applied to classification problems.

To estimate a regularized logistic regression model in `glmnet` (LASSO, ridge regression, or an elastic net) supply the option `family = "binomial"`

37 / 38

## Summary

- ▶ Regularization techniques stabilize or constrain the parameter estimates
  - ▶ Introduces more bias in prediction
- ▶ Bias-variance tradeoff: Allow for some bias if the reduction in the variance of the prediction is sufficiently high — improve overall MSE of prediction
- ▶ Key techniques:
  - ▶ LASSO — variable selection
  - ▶ Ridge regression
  - ▶ Elastic net
- ▶ Selection of tuning parameters using cross-validation

38 / 38