# Regularization

ST552 Lecture 24

Charlotte Wickham

2019-03-06

- Study guide posted along with ~~last year~~'s final. ← *heavy on conceptual*

  a

Three questions worth roughly equal amounts

1. **Calculations, inference and interpretation.** Like Q2 and Q3 on midterm: t-tests, F-tests, confidence intervals, prediction intervals etc. Know your formulas, how to use them and how to interpret the results.

2. **Assumptions and diagnostics.** Examine plots, identify problems, discuss the consequences of the problem, suggest remedies, suggest ways to verify if the suggestions worked. Or suggest ways to diagnose certain problems.

3. Everything else, at a conceptual level.

Some options:

- Revisit a topic (or more depth on a topic). Which topic?
- Exam review session $\longrightarrow$ Concrete task
- No lab

"think about what to do/cover"

#1
#2
#3

"Previous final question that isn't on the practice"

## Today

- The Bias-Variance tradeoff
- Regularized regression: lasso and ridge

For estimates, the mean squared error of an estimate can be broken down into bias and variance terms:

$$MSE(\hat{\theta}) = \mathsf{E}\left((\hat{\theta} - \theta)^2\right) = \mathsf{E}\left(\left(\hat{\theta} - \mathsf{E}\left(\hat{\theta}\right)\right)^2\right) + \left(\mathsf{E}\left(\hat{\theta}\right) - \theta\right)^2$$

$$= \mathsf{Var}\left(\hat{\theta}\right) + \mathsf{Bias}\left(\hat{\theta}\right)^2$$

Often in statistics, we focus on estimates that are unbiased (so the second term is zero), and focus on minimising the variance term.

You might argue that you are willing to introduce a little bias, if it reduces the variance enough to reduce the overall mean squared error in the estimate.

There is a similar breakdown for the mean square error in prediction. Let $\hat{f}(X)$ indicate the regression model for predicting an observation with explanantory values $X$.

*[handwritten: $\hat{\Theta}$]*  *[handwritten: $\Theta$]*

If the true data is generated according to $Y = f(X) + \epsilon$, where $E(\epsilon) = 0$ and $Var(\epsilon) = \sigma^2$, then the MSE for a point $x_0$:

$$MSE(\hat{f}(X)|X = x_0) =$$
$$E\left((Y - \hat{f}(X))^2 \Big| X = x_0\right)$$
$$= E\left(\left(\hat{f}(x_0) - E\left(\hat{f}(x_0)\right)\right)^2\right) + E\left(\left(Y - E\left(\hat{f}(x_0)\right)\right)\right)^2$$
$$= Var\left(\hat{f}(x_0)\right) + Bias\left(\hat{f}(x_0)\right)^2 + \sigma^2$$

*[handwritten: how close is our estimated function to true function on average]*

Bias captures how far our predictions are from the true mean on average (over repeated samples).
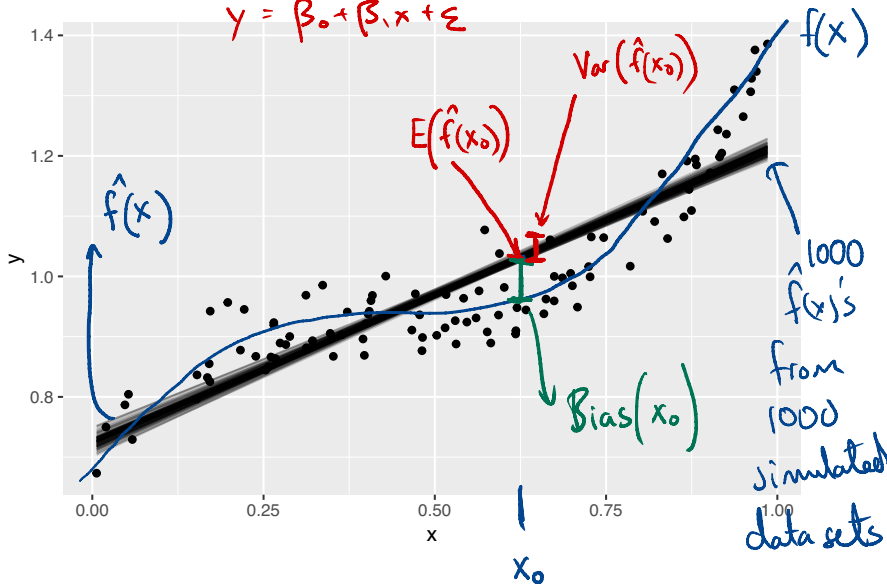
Variance captures how much our predictions vary (over repeated samples).

6

$y = \beta_0 + \beta_1 x + \varepsilon$

$f(x)$

$\text{Var}(\hat{f}(x_0))$

$E(\hat{f}(x_0))$

$\hat{f}(x)$

$\text{Bias}(x_0)$

1000 $\hat{f}(x)'s$ from 1000 simulated data sets

$x_0$

8

$$Y = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \varepsilon$$

$f(x)$

variance similar to linear fits

$x_0$

bigger variance

In general, more complex models decrease bias and increase variance. We hunt for the sweet spot where MSE is minimized.

There aren't nice partitions into bias and variance for other metrics, but the pattern is usually the same. Increasing complexity only improves performance to a point, then it decreases performance.

11

## Regularized regression

One approach that introduces bias into the coefficient estimates, is regularized (a.k.a. penalized) regression. Instead of minimising

$$\sum_{i=1}^{n} (y_i - \hat{y}_i)^2$$

minimise

$$\sum_{i=1}^{n} (y_i - \hat{y}_i)^2 + \lambda \sum_{j=1}^{p} f(\beta_j)$$

*→ "tuning" parameter*

*penalty*
*big coefficients =>*
*big penalty*

When $f(\beta_j) = \beta_j^2$, the method is called **ridge regression**, and when $f(\beta_j) = |\beta_j|$ the method is called **lasso**.

**The general idea:** the first term rewards good fit to the data, the second penalizses for large values on the coefficients.

The result: estimates shrink toward zero (introducing bias) and have smaller variance.

You can also view ridge and lasso as constrained minimisation, where we minimise
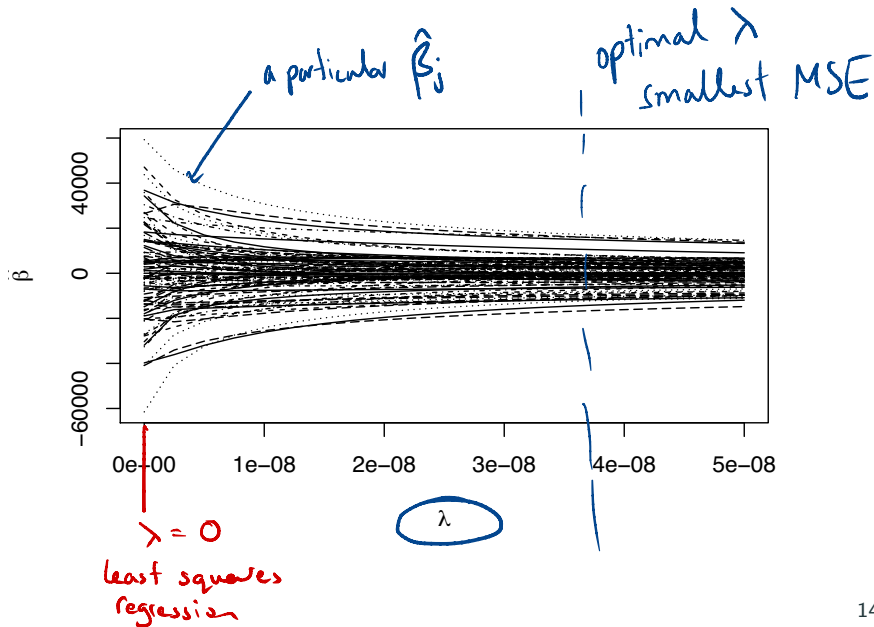
$$\sum_{i=1}^{n} (y_i - \hat{y}_i)^2$$

subject to the constraints
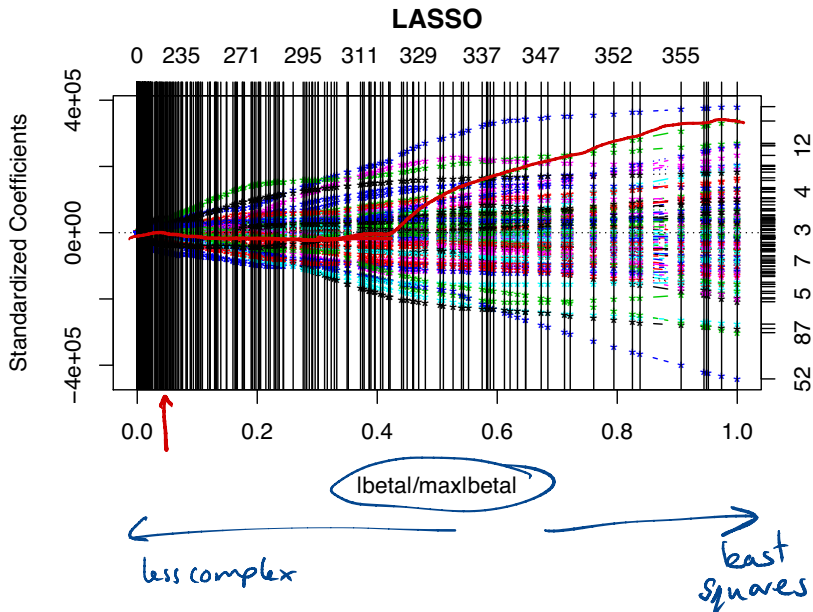
$$\sum_{j=1}^{p} \beta_j^2 \leq t \quad \text{for ridge}$$

$$\sum_{j=1}^{p} |\beta_j| \leq s \quad \text{for lasso}$$
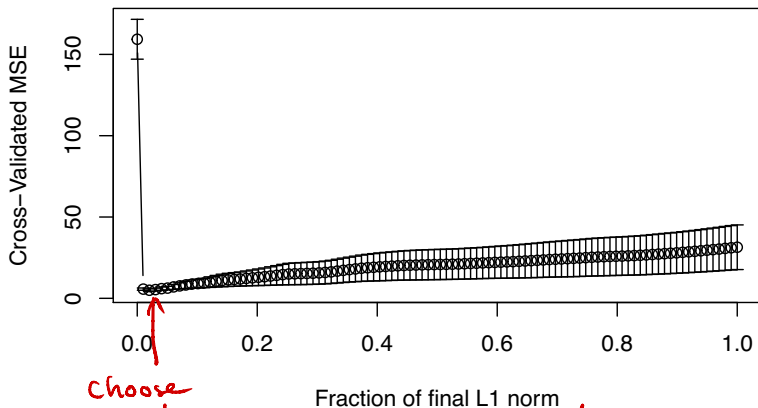
*'tuning parameters*

⇒ *instead of* $\lambda$

a particular $\hat{\beta}_j$

optimal $\lambda$
smallest MSE

$\lambda = 0$
least squares
regression

$\lambda$

14

LASSO

15

```
cvout <- cv.lars(as.matrix(trainmeat[ , -101]), trainmeat$fat)
```



Lasso

Choose
this value
of tuning parameter · tuning parameter

## Tuning parameter

```
(best_s <- cvout$index[which.min(cvout$cv)])
```
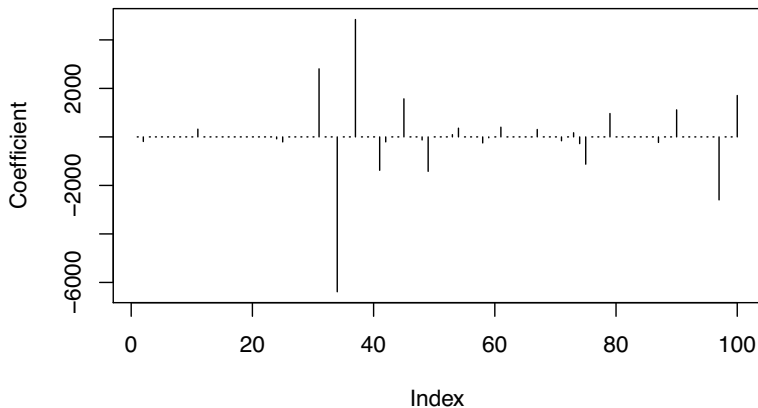
```
## [1] 0.02020202
```

RMSE for this value:

```
testx <- as.matrix(testmeat[,-101])

predlars <- predict(fit_lasso, testx, s=best_s,
  mode="fraction")
sqrt(mean((testmeat$fat - predlars$fit)^2))
```

```
## [1] 2.062124
```

```
## [1] 27
```
# $\hat{\beta}_j$'s are not zero

## Key points

- Regularized/Penalized regression models have a tuning parameter that controls the degree of penalization/shrinkage
- The tuning parameter may be chosen to optimize some kind of criterion
- Lasso estimates can be exactly zero (so it performs model selection as well)