

# Topic 1: Linear regression and regularization

Eric B. Laber

Department of Statistical Science, Duke University

Statistics 561



*The world needs another penalized regression method.*  
—Nobody (circa 2010).<sup>1</sup>

*Penalization is like the word 'myself,' it's used too much  
(often incorrectly) by people wishing to seem more intelli-  
gent than they are.*  
—Joel Vaughan

---

<sup>1</sup>Relax, I'm 98.5 percent joking.



## Warm-up (5 minutes)

- ▶ Explain to your group
  - ▶ Why is penalization/regularization used with predictive models?
  - ▶ How does penalization fit into the bias-variance trade-off?
  - ▶ What's the connection between ridge-regression and a Bayesian linear model?
- ▶ True or false
  - ▶ Masking occurs when there is a single large outlier
  - ▶ The Gauss-Markov theorem says the least squares estimator minimizes MSE
  - ▶ The owner of "the regulator bookshop" on 9th street thinks it's hilarious when you open the door and holler "REGULATORS!!" like in the Warren G song

# Roadmap

- ▶ Review and reminders
- ▶ All subsets (you don't want the truth)
- ▶ Ridge regression
- ▶ Lasso



# Roadmap

- ▶ **Review and reminders**
- ▶ All subsets (you don't want the truth)
- ▶ Ridge regression
- ▶ Lasso



## Review: fitting a linear model

- ▶ Ordinary least squares estimator as

$$\hat{\beta}_n = \arg \min_{\beta \in \mathbb{R}^p} \mathbb{P}_n(Y - \mathbf{X}^\top \beta)^2 = \arg \min_{\beta \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n (Y_i - \mathbf{X}_i^\top \beta)^2$$

- ▶ Alternatively, view  $\hat{\beta}_n$  as solution to

$$\mathbb{P}_n \mathbf{X} (Y - \mathbf{X}^\top \beta) = 0$$

- ▶ If  $\mathbb{P}_n \mathbf{X} \mathbf{X}^\top$  is invertible  $\hat{\beta}_n = (\mathbb{P}_n \mathbf{X} \mathbf{X}^\top)^{-1} \mathbb{P}_n \mathbf{X} Y$

## Review: fitting a linear model cont'd

- ▶ What to do if  $\mathbb{P}_n \mathbf{X} \mathbf{X}^\top$  is (nearly) singular?



# Estimation and approximation error

- ▶ Suppose  $Y = f(\mathbf{X}) + \epsilon$ 
  - ▶ Let  $\beta^* = \arg \min_{\beta \in \mathbb{R}^p} P(Y - \mathbf{X}^\top \beta)^2$
- ▶ Decompose mean squared error as

$$\begin{aligned}\mathbb{E}(Y - \mathbf{X}^\top \hat{\beta}_n)^2 &= \text{Var}(\epsilon) + \mathbb{E} \{f(\mathbf{X}) - \mathbf{X}^\top \beta^*\}^2 + \mathbb{E} \left\{ \mathbf{X}^\top (\hat{\beta}_n - \beta^*) \right\}^2 \\ &= \text{Noise} + \text{Approx Err} + \text{Est. Err}\end{aligned}$$

- ▶ A complex model reduces the approximation error but increases the estimation error



# Parsimonious models

- ▶ Predictive models used to inform decision making
  - ▶ Generate interesting hypotheses for further study
  - ▶ Forecasts weighed by stakeholders
  - ▶ Drive automatic action, i.e., decide if and what type of push notification to send patient in mHealth
- ▶ Need to build trust with stakeholders<sup>2</sup>
  - ▶ Models must be validated in domain context
  - ▶ Interpretable models often required

---

<sup>2</sup>Note\* this might be you! If you don't trust your own models it will be hard to make progress.

# Parsimonious models cont'd

- ▶ Justification for parsimonious models
  - ▶ Occam's Razor
  - ▶ (Medicine) true optimal decision rule simple
  - ▶ More nuanced mathematical arguments<sup>3</sup>
- ▶ Linear models aren't really that interpretable
  - ▶ Easy when abstracted away from the problem
  - ▶ Notion of one variable moving while all others held fix can be nonsensical to domain experts in some contexts<sup>4</sup>
  - ▶ But the fewer terms they have, the more interpretable they tend to be (we'll explore lists and trees later in this course)

---

<sup>3</sup>Duke's own Cynthia Rudin has some excellent work in this area. Check out <https://impact.duke.edu/story/whats-in-the-box> and her papers on this topic.

<sup>4</sup>Imagine increasing the square-footage of a grocery store without changing the products, back-of-house inventory, layout, etc.



*It is our experience and strong belief that better models and a better understanding of ones data result from focused data analysis, guided by substantive theory.*

*–Gwyneth Paltrow<sup>5</sup>*

*Automatic model-building procedures should be avoided at all cost. –D.R. Cox<sup>6</sup>*

---

<sup>5</sup>For the best in psychic vampire repellent see:

<https://goop.com/paper-crane-apothecary-psychic-vampire-repellent/p/>

<sup>6</sup>This is a real one.



# Roadmap

- ▶ Review and reminders
- ▶ **All subsets (you don't want the truth)**
- ▶ Ridge regression
- ▶ Lasso



# Finding the 'true' model

- ▶ Assume that  $\beta_j^* = 0$  for  $j \in \mathcal{J}$ 
  - ▶ Only need to regress  $Y$  on  $X_j$  where  $j \notin \mathcal{J}$
  - ▶ Let  $\mathbf{X}_{\mathcal{J}^c}$  denote relevant predictors
  - ▶ Question: suppose an oracle gave you  $\mathcal{J}$  is  $\hat{\beta}_n^{\mathcal{J}^c} = \arg \min_{\beta} \mathbb{P}_n(Y - \mathbf{X}_{\mathcal{J}^c}^T \beta)^2$  the optimal estimator?
- ▶ R code example: `growingLinearModel.R`<sup>7</sup>

---

<sup>7</sup>We'll mostly use python in this course but sometimes a guy already has some R code written. Geez, get off my back.

# Finding the 'true' model

- ▶ If you want optimal predictions, you may not want to use the 'true' model even if it were available to you
  - ▶ Estimating small effects inflates variance but does little to improve prediction
  - ▶ What constitutes a small effect depends on residual variance and sample size (not an absolute)
- ▶ Exercise: suppose  $X_1, \dots, X_n \sim_{i.i.d.} (\mu, \sigma^2)$ , derive the MSE for the estimators  $\tilde{\mu}_n \equiv 0$  and  $\hat{\mu}_n = \bar{X}_n$ . When is  $\tilde{\mu}_n$  preferable to  $\hat{\mu}_n$ ?

**Blank page for notes**



# Hard- and soft-thresholding

- ▶ Consider our toy example  $X_1, \dots, X_n \sim_{i.i.d.} (\mu, \sigma^2)$
- ▶ Idea: reduce MSE by adaptively shrinking our estimator  $\bar{X}_n$ 
  - ▶ Hard-thresholding  $\hat{\mu}_n^H = \bar{X}_n 1_{f(|\bar{X}_n|) \geq \tau}$  for some function  $f : \mathbb{R}_+ \rightarrow \mathbb{R}_+$
  - ▶ Soft-thresholding  $\hat{\mu}_n^H = \bar{X}_n g(|\bar{X}_n|)$  for some function  $g : \mathbb{R}_+ \rightarrow \mathbb{R}_+$



## Soft-thresholding: try it at home!

- ▶ Consider an estimator of the form  $\tilde{\mu}_n = \alpha \overline{X}_n$ , derive the value of  $\alpha^{\text{opt}}$  that minimizes MSE over  $\alpha$

# Hooray! A new estimator of the mean?

- ▶  $\alpha^{\text{opt}} \bar{X}_n$  resembles an empirical Bayes estimator of the mean <sup>8</sup>
  - ▶ Optimal shrinkage depends on signal-to-noise ratio
  - ▶ Common theme for adaptive shrinkage methods, i.e., adaptive ridge and adaptive lasso
- ▶ Not very exciting for a 1D example, but can be extremely useful/effective in higher dimensions
- ▶ Warning: hard-thresholding estimators are often non-regular which means we cannot estimate their sampling distributions uniformly  $\Rightarrow$  standard inference procedures, i.e., bootstrap or series approximations, perform poorly

---

<sup>8</sup>you will derive the analog for hard-thresholding estimator  $\bar{X}_n 1_{|\bar{X}_n| \geq \alpha}$  in HW2.

## Quick aside: all subsets

- ▶ Before moving on to penalized regression, we should mention the obvious approach of looking at all models and choosing the ‘best one’ based on some criterion
- ▶ An intuitive approach is to examine all possible models
  - ▶ Choose model with lowest BIC/AIC, etc.
  - ▶  $p$  variables  $\Rightarrow 2^p$  possible models
  - ▶ Branch-and-bound algorithms can reduce the search space making all subsets search feasible for  $p \leq 50^9$
- ▶ R code example: `allSubsets.R`

---

<sup>9</sup>For a fun paper on this topic see Furnival, George M., and Robert W. Wilson. "Regressions by leaps and bounds." *Technometrics* 42.1 (2000): 69-79.



# Roadmap

- ▶ Review and reminders
- ▶ All subsets (you don't want the truth)
- ▶ **Ridge regression**
- ▶ Lasso



*People say sometimes that Beauty is superficial. That may be so. But at least it is not so superficial as Thought is. To me, Beauty is the wonder of wonders. It is only shallow people who do not judge by appearances. The true mystery of the world is the visible, not the invisible.*  
—Excerpt "What to Expect When You're Expecting."



# Ridge regression: a superficial first look

- ▶ Least squares estimator is unbiased and has minimum variance among all unbiased estimators
  - ▶ As we've seen, this does not mean it minimizes MSE
  - ▶ Recall  $\text{MSE} = \text{bias}^2 + \text{variance} \Rightarrow \text{small increase in bias} + \text{big reduction in variance} = \text{smaller MSE}$
- ▶ Shrink regression coefficients toward zero by solving

$$\hat{\beta}_n^\lambda = \arg \min_{\beta} \mathbb{P}_n(Y - \mathbf{X}^\top)^2 + \lambda \|\beta\|^2,$$

where  $\lambda \geq 0$  is a tuning parameter

## Historical side-note

- ▶ Proposed by Hoerl and Kennard (1970)
  - ▶ Goal was to stabilize OLS estimator
  - ▶ Proved there always exists  $\lambda^*$  s.t.

$$\mathbb{E}||\hat{\beta}_n^{\lambda^*} - \beta^*||^2 < \mathbb{E}||\hat{\beta}_n - \beta^*||^2,$$

but  $\lambda^*$  depends on  $\beta^*$ , suggested looking at ‘ridge trace’ and picking values at each coefficients appear to ‘stabilize’<sup>10</sup> we’ll look at other data-driven tuning methods

- ▶ Fact: Hoerl is a fun name to say, try it out

---

<sup>10</sup>the ridge trace plot is what we might call a solution path today (nothing is new)



## Ridge regression: orthogonal case

- ▶ Ridge regression estimator can give dramatic reduction in MSE, especially when the dimension  $p$  is large
- ▶ Consider first the orthogonal case  $\mathbb{P}_n \mathbf{X} \mathbf{X}^\top = I_p$

$$\hat{\beta}_{n,j}^\lambda = \frac{\hat{\beta}_{n,j}}{1 + \lambda}$$

thus  $\mathbb{E}(\hat{\beta}_n^\lambda) = \beta_j^*/(1 + \lambda)$  and variance is  $\text{Var}(\hat{\beta}_{n,j})/(1 + \lambda)^2$

- ▶ Exercise: suppose  $p = 1$ , find the value  $\lambda$  that minimizes  $\text{MSE}(\hat{\beta}_{n,1}^\lambda)$ . What is the optimal for general  $p$ ?



**Blank page for notes**



## Ridge regression cont'd

- ▶ In the non-orthogonal case

$$\begin{aligned}\hat{\beta}_n^\lambda &= \{\mathbb{P}_n \mathbf{X} \mathbf{X}^\top + n^{-1} \lambda I_p\}^{-1} \mathbb{P}_n \mathbf{X} Y \\&= \left\{ I + n^{-1} \lambda (\mathbb{P}_n \mathbf{X} \mathbf{X}^\top)^{-1} \right\}^{-1} (\mathbb{P}_n \mathbf{X} \mathbf{X}^\top)^{-1} \mathbb{P}_n \mathbf{X} Y \\&= \left\{ I + n^{-1} \lambda (\mathbb{P}_n \mathbf{X} \mathbf{X}^\top)^{-1} \right\}^{-1} \hat{\beta}_n\end{aligned}$$

- ▶ How to choose  $\lambda \geq 0$  in this case?
  - ▶ Information criteria AIC, BIC, etc.
  - ▶ Cross-validation (generalized cross-validation)
  - ▶ You'll explore these criteria in lab

## Tuning $\lambda_n$ (one more criterion)

- ▶ Can use empirical Bayes (EB) to select  $\lambda_n$
- ▶ EB in a nutshell
  - ▶ Posit Bayesian model
  - ▶ Marginal distribution to obtain frequentist estimators of hyper-parameters<sup>11</sup>
- ▶ Pro-tip: if you want to derive an estimator that performs well in practice, posit a Bayesian model, derive posterior mean, call this a frequentist estimator and hide all evidence you ever considered Bayes approach<sup>12</sup>

---

<sup>11</sup>This is a rich area but we don't have the bandwidth to cover it in depth in this class.

<sup>12</sup>Credit to Derek Bingham for this nugget of wisdom.



## Tuning $\lambda_n$ with EB

- ▶ Assume linear model correct  $Y = \mathbf{X}^\top \boldsymbol{\beta} + \epsilon$  where  $\epsilon \sim \text{Normal}(0, \sigma^2)$  and  $\boldsymbol{\beta} \sim \text{Normal}(\mu, \tau^2)$  and  $\boldsymbol{\beta} \perp \epsilon$
- ▶ Posterior mean is

$$\hat{\boldsymbol{\beta}}_n = (\mathbb{P}_n \mathbf{X} \mathbf{X}^\top + \lambda_n I)^{-1} \mathbb{P}_n (\mathbf{X} Y + \lambda_n \mu),$$

where  $\lambda_n = \sigma^2 / (n \tau^2)$

- ▶ Setting  $\mu = 0$  yields ridge estimator, given estimators  $\hat{\sigma}_n^2$  and  $\hat{\tau}_n^2$  yield a plug-in estimator of  $\lambda_n$ , i.e.,  $\hat{\lambda}_n = \hat{\sigma}_n^2 / (n \hat{\tau}_n^2)$

## Tuning $\lambda_n$ with EB cont'd

- ▶ Marginal distribution<sup>13</sup>  $\mathbb{E}Y^2 = \tau^2 \text{trace}(\mathbf{X}\mathbf{X}^\top) + \sigma^2$
- ▶ MOM estimator matches

$$\mathbb{P}_n Y^2 \approx \tau^2 \text{trace}(\mathbb{P}_n \mathbf{X}\mathbf{X}^\top) + \sigma^2$$

plugging in  $\hat{\sigma}_n^2 = \mathbb{P}_n(Y - \mathbf{X}^\top \hat{\boldsymbol{\beta}}_n)^2$  yields

$$\hat{\tau}_n^2 = \frac{\mathbb{P}_n Y^2 - \hat{\sigma}_n^2}{\text{trace}(\mathbb{P}_n \mathbf{X}\mathbf{X}^\top)}$$

and thus  $\hat{\lambda}_n = \hat{\sigma}_n^2 / (\hat{\tau}_n^2 n)$

---

<sup>13</sup>To align with classic derivations, we're conditioning on  $\mathbf{X}$  and treating it as fixed here. Which makes the expression  $\mathbb{P}_n \mathbf{X}\mathbf{X}^\top$  a bit of an abuse of notation.

## Tuning $\lambda_n$ with EB more notes

- ▶ Note that  $\hat{\lambda}_n$  is a well-defined statistic without reference to the Bayes model
  - ▶ We can use this estimator and analyze it from a frequentist point-of-view
  - ▶ Bayesian connection is just icing
- ▶ Related idea: when we talk about RL, we might use a Bayesian framework for exploration-exploitation but not require these models to be correct when analyzing alg performance

**Blank page for notes**



## Ridge regression: BIC

- ▶ Let  $\mathbb{X}$  denote the design matrix, effective degrees of freedom

$$df(\lambda) = \text{trace} \left\{ \mathbb{X}^T (\mathbb{X}^T \mathbb{X} + \lambda I)^{-1} \mathbb{X} \right\}$$

to compute this efficiently we use the SVD of  $\mathbb{X}$

$$\mathbb{X} = UDV^T \Rightarrow df(\lambda) = \sum_{j=1}^p \frac{d_j^2}{d_j^2 + \lambda}$$

- ▶ BIC for each value of  $\lambda$  is

$$\text{BIC}(\lambda) = \log \{ \text{RSS}(\lambda) \} + df(\lambda) \log(n)/n$$

choose  $\hat{\lambda}^{\text{BIC}} = \arg \min_{\lambda} \text{BIC}(\lambda)$



# Asymptotic behavior

- ▶ Under what conditions on  $\{\lambda_n\}_{n \geq 1}$  will:
  - ▶  $\hat{\beta}_n^{\lambda_n} \rightarrow \beta^*$
  - ▶  $\sqrt{n} \left( \hat{\beta}_n^{\lambda_n} - \beta^* \right)$  converge to a Gaussian limit
- ▶ Before we mathematize anything, what does your intuition say? Should  $\lambda_n$  shrink? Grow? Converge to a constant?

# Asymptotic behavior notes



# Ridge regression: prostate data example

- ▶ Example: prostate cancer data set
  - ▶  $n = 97$  patients, response prostate-specific antigen (lspa)
  - ▶ Predictors:
    - ▶ lccavol: log (cancer volume)
    - ▶ lweight: log(prostate weight)
    - ▶ age: age
    - ▶ lbph: log (benign prostatic hyperplasia)
    - ▶ svi: seminal vesicle invasion
    - ▶ lcp: log (capsular penetration)
    - ▶ gleason: Gleason score
    - ▶ pgg45: percent of Gleason scores 4 or 5
- ▶ R code example: ridgeExample.R

# Ridge regression and dropout

- ▶ Dropout is popular heuristic to reduce overfitting in machine learning models (esp. nnets)
  - ▶ Randomly 'zero-out' some components in the training input
  - ▶ Lots of heuristic explanations for why this works<sup>14</sup>, with a linear model we can obtain a more rigorous answer
- ▶ Let  $\mathbf{Z} \in \{0, 1\}^p$  be a vector of independent Bernoulli random variables so that  $P(Z_j = 0) = \phi$  (and  $P(Z_j = 1) = 1 - \phi$ )
  - ▶ Suppose  $\mathbf{Z}$  is ind of  $(\mathbf{X}, Y)$
  - ▶ Random dropout  $\mathbf{X} \odot \mathbf{Z} = (X_1 Z_1, \dots, X_p Z_p)^\top$

---

<sup>14</sup>Search reddit or quora for 'why dropout works'.



## Ridge regression and dropout cont'd

- ▶ Imagine for each  $i = 1, \dots, n$  we generate a bajillion<sup>15</sup> values of  $\mathbf{Z}$ , say  $\mathbf{Z}_i^1, \dots, \mathbf{Z}_i^B$ , and subsequently generate new dataset  $\left\{ \left[ Y_i, \mathbf{X}_i \odot \mathbf{Z}_i^k / (1 - \phi) \right] \mid k = 1, \dots, B, i = 1, \dots, n \right\}$
- ▶ Now suppose we fit a linear model using this data

$$\begin{aligned}\tilde{\beta}_n &= \arg \min_{\beta} \sum_{k=1}^B \sum_{i=1}^n \left\{ Y_i - (\mathbf{X}_i \odot \mathbf{Z}_i^k)^{\top} \beta / (1 - \phi) \right\}^2 \\ &\approx \arg \min_{\beta} \mathbb{E}_{\mathbf{Z}} \mathbb{P}_n \left\{ Y - (\mathbf{X} \odot \mathbf{Z})^{\top} \beta / (1 - \phi) \right\}^2,\end{aligned}$$

where we've approximated the sample average with expectation (recall we generate as many  $\mathbf{Z}$ 's as we want)

---

<sup>15</sup>By a bajillion, I mean a lot. Also, this lecture is sponsored by "Bajillion Dollar Properties," watch now on Amazon Prime.



## Ridge regression and dropout cont'd

- ▶ Suppose data has been centered and scaled derive a closed form expression for  $\tilde{\beta}_n$

# Ridge regression and dropout discussion

- ▶ Reducing info available to the model to make a prediction at any given point ( $\mathbf{X} \mapsto \mathbf{X} \odot \mathbf{Z}$ ) has a regularizing effect
  - ▶ General strategy can be used with models that are harder to penalize explicitly, e.g., this is used when building trees in random forests (we'll talk more about this later)
  - ▶ Smoother version of dropout can be obtained by replacing Bernoulli's with any unit-mean r.v.'s

## Ridge regression and noise addition

- ▶ Another way to reduce info available to the model is to replace  $\mathbf{X}$  with  $\mathbf{X} + \mathbf{Z}$ , w/  $\mathbf{Z}$  vector of independent  $(0, \lambda)$  r.v.'s
- ▶ As with dropout, note that this is a general strategy that can be applied with any ML algorithm (for supervised learning)
- ▶ In-class exercise: compute

$$\tilde{\beta} = \arg \min_{\beta} \mathbb{E}_{\mathbf{Z}} \mathbb{P}_n \{ Y - (\mathbf{X} + \mathbf{Z})^\top \beta \}^2$$



**Blank page for notes**



## Ridge regression: discussion

- ▶ Reduces MSE by shrinking regression coefficients, very effective in a wide range of settings (battle-tested)
- ▶ Does not perform variable selection, i.e., all variables are kept in the model, this can be a problem if parsimony is critical but we can always threshold small values to zero
- ▶ Choose amount of penalization using information criteria (BIC, AIC, etc.) or cross-validation

# Principal components regression

- ▶ A closely related alternative to ridge regression is principal components regression
- ▶ Review: principal components<sup>16</sup>
  - ▶ Write  $\mathbb{X}^T \mathbb{X} / n = V D^2 V^T$  where  $V = [\mathbf{v}_1, \dots, \mathbf{v}_p]$
  - ▶  $Z_j = \mathbf{X}^T \mathbf{v}_j$  is called the  $j$ th principal component of  $\mathbf{X}$
  - ▶  $\text{Var}(Z_1) \geq \text{Var}(Z_2) \geq \dots \geq \text{Var}(Z_p)$
- ▶ Idea: regression  $Y$  on  $Z_1, \dots, Z_q$  for some  $q \leq p$  instead of regressing on  $\mathbf{X}$ 
  - ▶ Capture important features of  $\mathbf{X}$
  - ▶ Reduce dimension  $\Rightarrow$  bias-variance trade-off

---

<sup>16</sup> Assume predictors have been centered

## Principal components regression cont'd

- ▶ Matrix of principal components is  $\mathbb{Z} = \mathbb{X}V$ 
  - ▶ PR regression design matrix is orthogonal  $\mathbb{Z}^T\mathbb{Z} = nD^2$  (why?)
  - ▶ Compute regression of  $Y$  on  $Z_1, \dots, Z_q$  using series of univariate regressions
  - ▶ Tuning parameter  $q$  can be chosen using BIC/AIC, cross-validation, etc.
- ▶ R code example: `pcr.R`

# Principal components regression cont'd

- ▶ Pros:
  - ▶ Reduce MSE relative to standard least squares
  - ▶ Sometimes interpret principal components (see pcr.R)
  - ▶ Parallel computation possible
- ▶ Cons:
  - ▶ Interpretation of principal components subjective
  - ▶ Doesn't involve  $Y$  in construction of features (pc's)

# Roadmap

- ▶ Review and reminders
- ▶ All subsets (you don't want the truth)
- ▶ Ridge regression
- ▶ **Lasso**



# A nice relaxing quiz

- ▶ Explain to your group
  - ▶ What does the acronym LASSO stand for?
  - ▶ What is the problem of inference after model selection?
  - ▶ What goes wrong in regression when  $p \gg n$ ?
- ▶ True or false
  - ▶ Maximum likelihood is limited to parametric models
  - ▶ Large coefficient std errors can be sign of near collinearity
  - ▶ The world's most expensive Donkey cheese sells for several thousand dollars a pound

# Lasso: superficial overview

- ▶ Simultaneous estimation and model selection via penalization

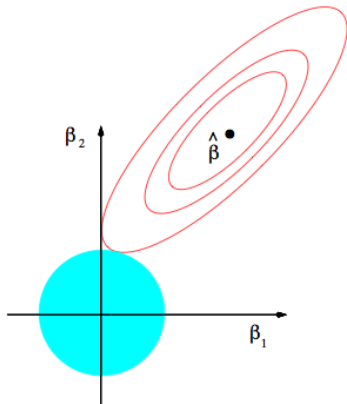
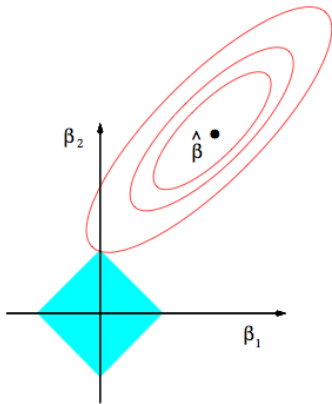
$$\hat{\beta}_n^\tau = \arg \min_{\beta} \mathbb{P}_n(Y - \mathbf{X}^\top \beta)^2 + \tau \|\beta\|_1,$$

where  $\|\beta\|_1 = \sum_{j=1}^p |\beta_j|$  and  $\tau > 0$  is a tuning parameter

- ▶ Looks like ridge but use of L1 norm yields sparse solutions



## Lasso vs. ridge (a stolen picture)



## Lasso: orthogonal predictors

- Suppose  $\mathbf{X}^T \mathbf{X} = I$  then

$$\mathbb{P}_n(Y - \mathbf{X}^T \boldsymbol{\beta})^2 + \tau \|\boldsymbol{\beta}\|_1 = \mathbb{P}_n Y^2 - \|\hat{\boldsymbol{\beta}}_n\|^2 + \|\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}\|^2 + \tau \|\boldsymbol{\beta}\|_1,$$

we can re-write the lasso solution as

$$\begin{aligned}\hat{\boldsymbol{\beta}}_n^\tau &= \arg \min_{\boldsymbol{\beta}} \|\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}\|^2 + \tau \|\boldsymbol{\beta}\|_1 \\ &= \arg \min_{\boldsymbol{\beta}} \sum_{j=1}^p \left\{ (\hat{\beta}_{n,j} - \beta_j)^2 + \tau |\beta_j| \right\},\end{aligned}$$

we can look at each component separately

- R code example: `lasso.R`

## Lasso: orthogonal predictors

- In orthogonal case we can compute  $\hat{\beta}_{n,j}^\tau$  explicitly

$$\hat{\beta}_{n,j}^\tau = \text{sgn}(\hat{\beta}_{n,j})(|\hat{\beta}_{n,j}| - \tau/2)_+,$$

where  $(u)_+ = \max(0, u)$

- Can give us further insight into lasso soln (back to lasso.R)

# Computing the lasso solution

- ▶ Can reformulate lasso objective as quadratic program, however, faster iterative algorithms exist
- ▶ As with ridge, select tuning parameter using BIC/AIC, cross-validation, etc.
  - ▶ Problem:  $\widehat{\beta}_n^\tau$  is not linear estimator, how to define degrees of freedom?
  - ▶ Approximate degrees of freedom is the number of nonzero components of  $\widehat{\beta}_n^\tau$

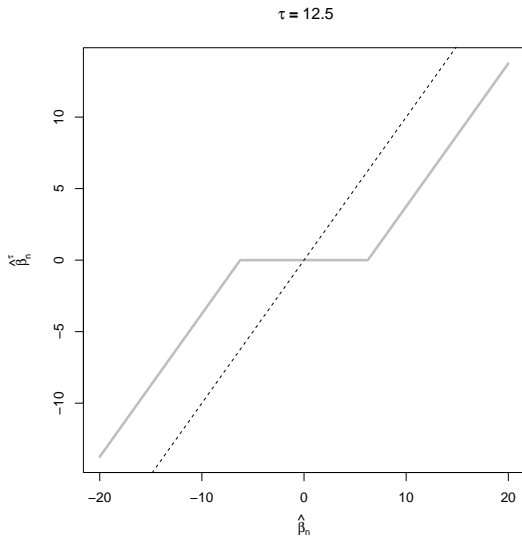
$$df(\tau) = \sum_{j=1}^p 1_{\widehat{\beta}_{n,j}^\tau \neq 0}$$

# Lasso discussion

- ▶ Tool for simultaneous variable selection and estimation
  - ▶ Can be solved very quickly for large problems
  - ▶ Leads to parsimonious models  $\Rightarrow$  interpretable
  - ▶ Variants exist, i.e., for glm, cox-ph, etc. (more on this later)
- ▶ Potential problems
  - ▶ Shrinks all the coefficients, even those that are not close to zero, can lead to excessive bias
  - ▶ Select at most  $\min(n, p)$  variables, problematic in some settings

# Adaptive lasso

- Recall nature of shrinkage for lasso



## Adaptive lasso cont'd

- ▶ Lasso shrinks all coefficients toward zero
  - ▶ Introduces excessive bias when true regression coef is large
  - ▶ Potentially better strategy is to shrink more aggressively when coefficients are small and less aggressively if they are large
  - ▶ How can we do this if true coefficients are unknown?
- ▶ Idea: use ordinary least squares estimates as surrogates

## Adaptive lasso cont'd

- ▶ Adaptive lasso estimator

$$\hat{\beta}_n^\delta = \arg \min_{\beta} \mathbb{P}_n(Y - \mathbf{X}^\top \beta)^2 + \delta \sum_{j=1}^p |\beta_j| / |\hat{\beta}_{n,j}|$$

- ▶ If  $|\hat{\beta}_{n,j}|$  is small  $\Rightarrow$  more shrinkage is applied
- ▶ Note\* if  $p$  is large can use ridge estimator as weights



## Adaptive lasso: orthogonal case

- In the orthogonal case the adaptive lasso estimator is

$$\hat{\beta}_{n,j}^{\delta} = \text{sgn}(\hat{\beta}_{n,j}) \left( |\hat{\beta}_{n,j}| - \frac{\delta}{2|\hat{\beta}_{n,j}|} \right)_{+}$$

for comparison

$$\text{Lasso: } \hat{\beta}_{n,j}^{\tau} = \text{sgn}(\hat{\beta}_{n,j}) \left( |\hat{\beta}_{n,j}| - \frac{\tau}{2} \right)_{+}$$

$$\text{Ridge: } \hat{\beta}_{n,j}^{\lambda} = \frac{\hat{\beta}_{n,j}}{1 + \lambda}$$

- R code example: `adaptiveLasso.R`

## Adaptive lasso: fitting non-orthogonal case

- ▶ Adaptive lasso objective can be recast as quadratic program
  - ▶ We can use existing lasso software by modifying design matrix
  - ▶ Write

$$\mathbb{P}_n(Y - \mathbf{X}^\top \beta)^2 + \delta \sum_{j=1}^p |\beta_j| / |\hat{\beta}_{n,j}| = \mathbb{P}_n(Y - \tilde{\mathbf{X}}^\top \tilde{\beta})^2 + \delta \sum_{j=1}^p |\tilde{\beta}_j|,$$

where  $\tilde{X}_j = \hat{\beta}_{n,j} X_j$  and  $\tilde{\beta}_j = \beta_j / \hat{\beta}_{n,j}$

# Adaptive lasso: fitting non-orthogonal case cont'd

► Pseudo-code (assumes centered, scaled data)

1. Fit ordinary least squares, say  $\hat{\beta}_n$
2. Create scaled design matrix  $\tilde{\mathbf{X}} = \mathbf{X}\text{diag}(\hat{\beta}_n)$
3. For each  $\delta$  under consideration compute lasso estimator

$$\tilde{\beta}_n^\delta = \arg \min_{\beta} \mathbb{P}_n(Y - \tilde{\mathbf{X}}^\top \beta)^2 + \delta \|\beta\|_1,$$

compute adaptive lasso estimator  $\hat{\beta}_{n,j}^\delta = \hat{\beta}_{n,j} \tilde{\beta}_{n,j}^\delta$

4. Select  $\delta$  using BIC/AIC, etc.

# Adaptive lasso discussion

- ▶ Reduce excessive bias due to overshrinkage of large coefficients
  - ▶ An alternative is to refit the model selected by lasso using least squares, this generally works well in practice
- ▶ Can use existing statistical software, generally much more computationally efficient than alternative proposals (SCAD, etc.)

## Robustness and the lasso

- ▶ Let  $c_1, \dots, c_p$  be positive constants and define

$$\mathcal{M} = \left\{ \mathbf{\Delta} \in \mathbb{R}^{n \times p} : \sqrt{\sum_{i=1}^n \Delta_{ij}^2} \leq c_j, j = 1, \dots, p \right\},$$

i.e., the set of  $n \times p$  matrices where norm of  $j$ th column bounded by  $c_j$  for  $j = 1, \dots, p$

- ▶ Robust lin regression estimator under perturbations  $\mathcal{M}$  is

$$\tilde{\beta}_n = \arg \min_{\beta} \left\{ \max_{\mathbf{\Delta} \in \mathcal{M}} \|\mathbb{Y} - (\mathbb{X} + \mathbf{\Delta})\beta\|_2 \right\}$$

## Robustness and lasso equivalence thm

- Claim: lasso and robust regression soln are equivalent in that

$$\begin{aligned}\tilde{\beta}_n &= \arg \min_{\beta} \left\{ \max_{\Delta \in \mathcal{M}} \|\mathbb{Y} - (\mathbb{X} + \Delta)\beta\|_2 \right\} \\ &= \arg \min_{\beta} \left\{ \|\mathbb{Y} - \mathbb{X}\beta\|_2 + \sum_{j=1}^p c_j |\beta_j| \right\}\end{aligned}$$

**Blank page for notes**



# Robustness and lasso discussion

- ▶ Showed equivalence between so-called square-root lasso and robust regression
  - ▶ Set  $c_j \equiv \tau$  yields standard (square-root) lasso
  - ▶ Setting  $c_j \equiv 1/|\hat{\beta}_{n,j}|$  yields adaptive (square-root) lasso
- ▶ The square-root lasso has some desirable properties in terms of tuning (see papers by A. Belloni here at Duke), however, the solution paths are the same



# Penalization and regularization discussion

- ▶ Bias-variance trade-off  $\Rightarrow$  regularization needed to improve predictive performance
- ▶ Ridge regression smoothly (soft) penalizes coefficients but can combine with thresholding to obtain sparse solutions
- ▶ Lasso automatically yields sparse solutions but may perform poorly if true signal is dense and comprised of weak signals
- ▶ Showed that several information deletion/distortion methods are equivalent to penalized regression  $\Rightarrow$  more general strategy for regularizing complex estimators

Thank you.

`eric.laber@duke.edu`

`laber-labs.com`

