

Lec 7

Thursday, September 19, 2019 11:16

Recap:

- Bias-var tradeoff
- Subset selection for OLS (feature) gives a range of complexities
- Cross-val as a way to
  - a) Estimate the risk of any algo
  - b) Choose a complexity param (we often use the 1-std-err rule)

Shrinkage

Subset selection is very discrete

- good for interpretation
- potentially bad for prediction (not necessarily)

Shrinkage is a more cts way to trade off bias & var

$$\hat{\beta}^{ridge} = \underset{\beta \in \mathbb{R}^{p+1}}{\text{argmin}} \left\| Y - X\beta \right\|_2^2 + \lambda \sum_{j=1}^p \beta_j^2$$

excluding the intercept  
 after standard the data so all  $\beta_j$  on same scale

Shrink(s)  $\hat{\beta}^{OLS}$  toward simplest model  
 most complex model

At one extreme  $\lambda=0$  :  $\hat{\beta}^{ridge} = \hat{\beta}^{OLS}$

— || —  $\lambda=\infty$  :  $\hat{\beta}^{ridge} = \text{intercept model at sample mean of } \bar{Y} = \frac{1}{n}$

In (w): shrinking toward 0 prevent  $\beta$

from trying to reach far off observations with extreme slopes

$$\hat{\beta}^{\text{ridge}} = (\mathbf{X}^T \mathbf{X} + \Lambda)^{-1} \mathbf{X}^T \mathbf{Y}$$

$$\Lambda = \begin{pmatrix} 0 & & & 0 \\ & \lambda & & \\ & & \ddots & \\ 0 & & & \lambda \end{pmatrix}$$

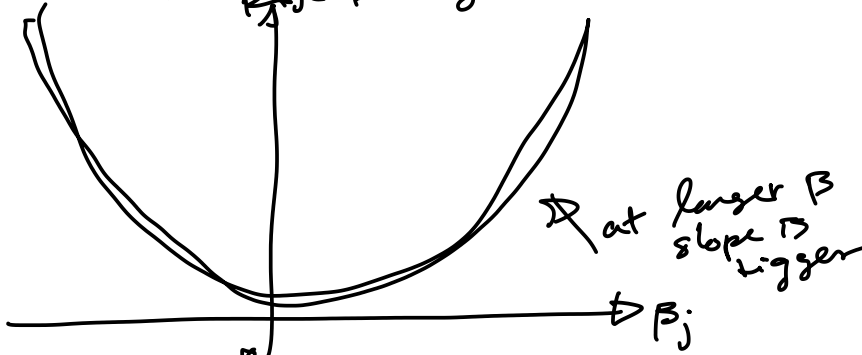
### Lasso

Combine shrinkage & subset selection

$$\hat{\beta}^{\text{lasso}} = \underset{\beta \in \mathbb{R}^{p+1}}{\text{argmin}} \left( \|\mathbf{Y} - \mathbf{X}\beta\|_2^2 + \lambda \sum_{j=1}^p |\beta_j| \right)$$

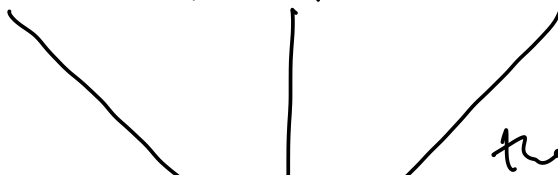
excluding intercept

Idea:  
ridge penalty

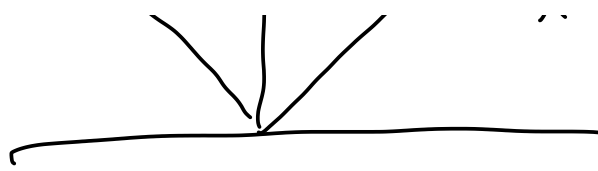


at 0, have slope 0  
 $\Rightarrow$  pay essentially nothing to go from  $\beta_j = 0$  to  $\beta_j = \epsilon$

lasso penalty



absolute val of the slope is +1 everywhere


 $\Rightarrow$  We start paying immediately for making  $\beta_j$  nonzero  
 This encourages 0 coeffs

No closed form for lasso

But can use grad descent methods (after fall break)

In Sciklearn, this is implemented in `sklearn.linear_model.lasso` - path

compute lasso for all  $\lambda$ 's simultaneously

## Naive Bayes

Another ex of classification algo

$$\begin{aligned}
 P(Y=j|X=x) &= \frac{P(X=x|Y=j) P(Y=j)}{\sum_{k} P(X=x|Y=k) P(Y=k)} && \text{(Bayes' rule)} \\
 &= \frac{f_j(x) \pi_j}{\sum_{k} f_k(x) \pi_k}
 \end{aligned}$$

$\pi_j$  — fraction of label- $j$  examples

$f_j(x)$  — prob/density of  $X=x$  in the population of label- $j$  example

$f_j(x) = f_j(x_1, \dots, x_p)$  a fu from  $\mathbb{R}^p$  to  $\mathbb{R}$

hard to estimate for large (even moderate)  $p$

## Naive assumption

$$f_j(x) = f_{j1}(x_1) f_{j2}(x_2) \cdots f_{jp}(x_p)$$

$$= \prod_{k=1}^p f_{jk}(x_k)$$

meaning: Assuming that  $X_1, \dots, X_p$  are statistically independent given  $Y$

$$\Rightarrow P(Y=j | X=x) = \frac{f_j(x) \pi_j}{\sum_k f_k(x) \pi_k}$$

$$= \frac{\pi_j \prod_{l=1}^p f_{jl}(x_l)}{\sum_k \pi_k \prod_{l=1}^p f_{kl}(x_l)}$$

## Naive Bayes classifier:

estimate each of  $\pi_j, f_{jl} \forall j, l$   
 plug in above to get  $\hat{P}(Y=j | X=x)$   
 & mimic the Bayes classifier  
 by maximizing  $\hat{P}(Y=j | X=x)$

For class frequencies, use empirical frequencies

$$\hat{\pi}_j = \frac{1}{n} \sum_{i=1}^n \mathbb{I}[Y_i=j]$$

What about  $f_{jl}$ ?

Discrete features:

just use empirical frequencies

$$\hat{f}_{j\ell}(x_\ell) = \frac{\sum_{i=1}^n \mathbb{I}[y_i = j, x_{i\ell} = x_\ell]}{\sum_{i=1}^n \mathbb{I}[y_i = j]}$$

What if  $\hat{f}_{j\ell}(x_\ell) = 0$  while  $\hat{f}_{i\ell}(x_\ell) \neq 0$ ?

Get zero prob no matter

What happens w/ other  $p_i$  features

One soln: Laplace smoothing

just pad the data

$$\hat{f}_{j\ell}(x_\ell) = \frac{\left( \sum_{i=1}^n \mathbb{I}[y_i = j, x_{i\ell} = x_\ell] \right) + \alpha}{\left( \sum_{i=1}^n \mathbb{I}[y_i = j] \right) + \alpha}$$

$$\alpha \geq 0$$

## Continuous features

Several options:

- kernel density estimation

next time

- parametric density estimate

Eg. Assume  $f_{j\ell}(x_\ell) = \mathcal{Q}\left(\frac{x_\ell - \mu}{\sigma}\right)$   
 $\mathcal{Q}$ : pdf of std normal

Set  $\hat{\mu}, \hat{\sigma}$  to the  
sample mean & std dev

$$\text{of } \left\{ \begin{array}{l} X_{ik} : i=1, \dots, h \\ Y_i = j \end{array} \right\}$$