

Last lecture: logistic regression
touched overfitting

Today: overfitting & regularization

simplest regression

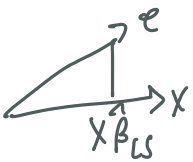
1		
⋮		
i	x_i	y_i
⋮		
n		

model: $y_i = x_i \beta + e_i$ $e_i \sim N(0, \sigma^2)$ iid

hypothesis testing: $H_0: \beta = 0$ simpler model
 $H_1: \beta \neq 0$ more complex

$$\hat{\beta}_{LS} = \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n x_i^2} \quad \stackrel{H_0}{=} \quad \frac{\sum_{i=1}^n x_i e_i}{\sum_{i=1}^n x_i^2} \neq 0$$

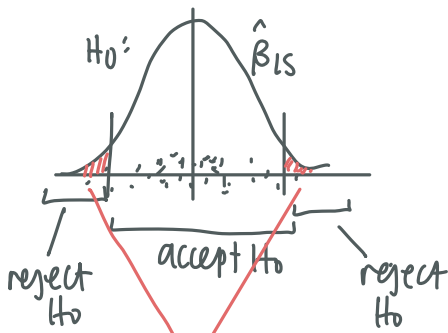
there is always a certain amount of overfitting in $\hat{\beta}$ even if true $\beta = 0$



overfitting: model absorbs noise

ways to avoid overfitting

① hypothesis testing



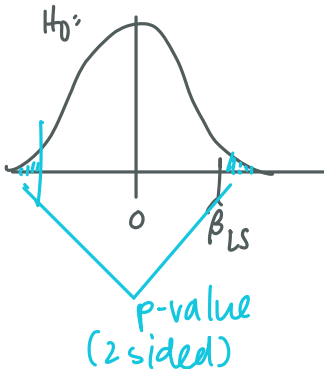
let $t = \text{threshold}$

$$\hat{\beta}_{\text{hypothesis testing}} = \begin{cases} 0 & |\hat{\beta}_{LS}| < t \\ \hat{\beta}_{LS} & |\hat{\beta}_{LS}| > t \end{cases}$$

hard thresholding

error: α , usually taken as 5%. (2.5% in each tail)

p-value: how extreme $\hat{\beta}_{LS}$ is relative to H_0



Logistic regression

$$H_0: \beta = 0 \Rightarrow s_i = x_i \beta = 0$$

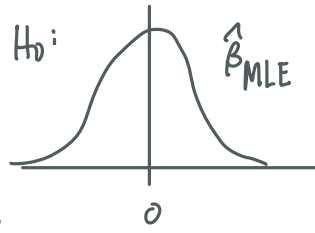
$$p_i = \frac{e^{s_i}}{1 + e^{s_i}} = \frac{1}{2}$$

$$y_i \sim \text{Bernoulli}(\frac{1}{2})$$

nothing more than a coin flip

$$\hat{\beta}_{MLE} \neq 0$$

max likelihood estimate



classification

noise in the logistic regression scenario is overinterpretation of coin flipping, as opposed to gaussian noise observed in linear regression

2) Regularization, ML treatment

1) L2 -> ridge regression

keeps all covariates

simplest regression

$$L(\beta) = \sum_{i=1}^n (y_i - x_i \beta)^2 + \lambda \beta^2$$

ls term penalty term — drags estimate of β towards 0 to help avoiding overfitting

$$L'(\beta) = -2 \sum_{i=1}^n (y_i - x_i \beta) x_i + 2\lambda \beta$$

univariate scenario

$$= \sum_{i=1}^n y_i x_i + \sum_{i=1}^n x_i^2 \beta + \lambda \beta = 0$$

$$\hat{\beta}_{\text{ridge}} = \frac{\sum x_i y_i}{\sum x_i^2 + \lambda}$$

shrinkage estimator:

under $H_0: y_i = e_i$

adding λ will help w/ not absorbing as much noise (however, tradeoff is introduction of some bias)

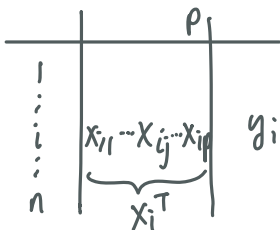
$\lambda \rightarrow 0$

$$\hat{\beta}_{\text{ls}} = \frac{\sum x_i y_i}{\sum x_i^2}$$

interpolation stwn

H_0 & H_1 , as opposed to hard thresholding

multivariate scenario



$$L(\beta) = \sum_{i=1}^n (y_i - x_i^T \beta)^2 + \lambda \|\beta\|_{l2}^2 \rightarrow \sum_{j=1}^p \beta_j^2$$

$$L'(\beta) = -2 \sum_{i=1}^n x_i (y_i - x_i^T \beta) + 2\lambda \beta = 0$$

$$= -\sum_{i=1}^n x_i y_i + \sum_{i=1}^n x_i x_i^T \beta + \lambda \beta = 0$$

$$\left(\sum_{i=1}^n x_i x_i^T + \lambda I \right) \beta = \sum_{i=1}^n x_i y_i$$

$$\hat{\beta}_{\text{ridge}} = \left(\sum x_i x_i^T + \lambda I \right)^{-1} \left(\sum x_i y_i \right) = \left(X^T X + \lambda I \right)^{-1} X^T Y$$

raw input	features p	
x_i	$\sin x_i, \cos x_i, \dots$	y_i
n	$p \gg n$	

$$y_i = \beta_{1,1} \sin x_i + \beta_{1,2} \cos x_i + \beta_{2,1} \sin 2x_i + \beta_{2,2} \cos 2x_i + \dots + e_i$$

$|\beta|_{\ell_2}^2$ small \rightarrow smaller curve

otherwise, training error = 0 \rightarrow model absorbs too much noise

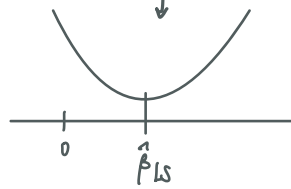
② L1 regularization \rightarrow LASSO
(least absolute shrinkage selection operator)

simplest regression

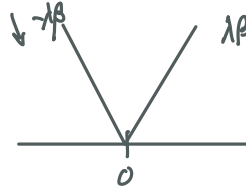
i	x_i	y_i
-----	-------	-------

$$L(\beta) = \frac{1}{2} \sum_{i=1}^n (y_i - x_i \beta)^2 + \lambda |\beta|$$

shrinkage + select

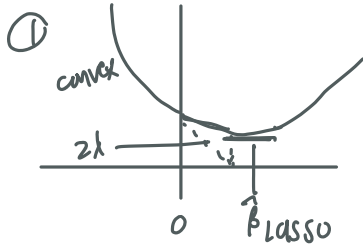


smooth, differentiable @ 0

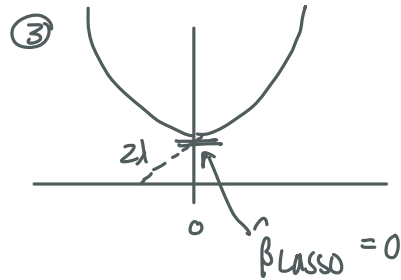


sharp, not differentiable @ 0

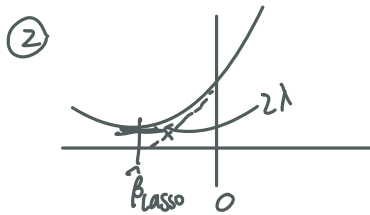
Three scenarios:



change in slope 2λ
 $\beta > 0$



accept H_0 in this scenario



change in slope 2λ

scenario ① $\hat{\beta}_{\text{lasso}} > 0 \longrightarrow \hat{\beta}_{\text{LS}} > \frac{\lambda}{\sum x_i^2}$ for $\hat{\beta}_{\text{lasso}} > 0$

$$L'(\beta) = -\sum (y_i - x_i \beta) x_i + \lambda = 0$$

$$= -\sum x_i y_i + \sum x_i^2 \beta + \lambda = 0$$

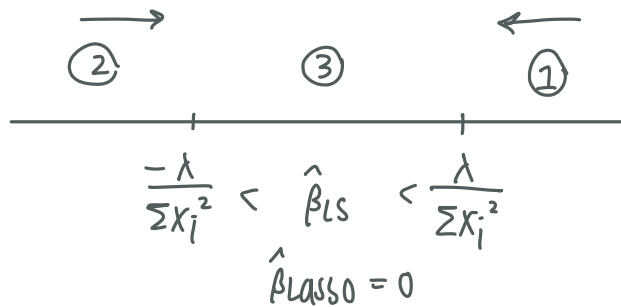
$$\hat{\beta}_{\text{lasso}} = \frac{\sum x_i y_i - \lambda}{\sum x_i^2} = \underbrace{\frac{\sum x_i y_i}{\sum x_i^2}}_{\hat{\beta}_{\text{LS}}} - \frac{\lambda}{\sum x_i^2} \quad \text{shrinkage}$$

scenario ② $\hat{\beta}_{Lasso} < 0 \longrightarrow \hat{\beta}_{LS} < \frac{-\lambda}{\sum x_i^2}$ for $\hat{\beta}_{Lasso} < 0$

$$L(\beta) = -\sum (y_i - x_i \beta) x_i - \lambda = 0$$

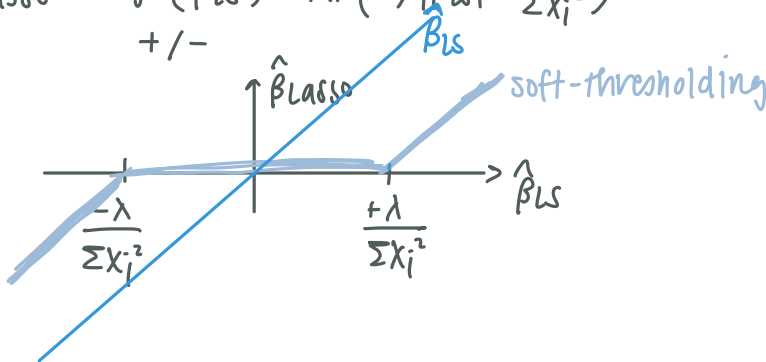
$$= -\sum x_i y_i + \sum x_i^2 \beta - \lambda = 0$$

$$\hat{\beta}_{Lasso} = \frac{\sum x_i y_i + \lambda}{\sum x_i^2} = \underbrace{\frac{\sum x_i y_i}{\sum x_i^2}}_{\hat{\beta}_{LS}} + \frac{\lambda}{\sum x_i^2} = \hat{\beta}_{LS} + \frac{\lambda}{\sum x_i^2}$$



selection in scenario ③

$$\hat{\beta}_{Lasso} = \text{sign}(\hat{\beta}_{LS}) \max\left(0, |\hat{\beta}_{LS}| - \frac{\lambda}{\sum x_i^2}\right)$$



multivariate scenario

$$L(\beta) = \frac{1}{2} \sum_{i=1}^n (y_i - x_i^T \beta)^2 + \lambda \|\beta\|_{L1} \rightarrow \sum_{j=1}^p |\beta_j|$$

conducts selection for covariates in final model as well as minimizing the β s

$$= \frac{1}{2} \left| Y - \sum_{i=1}^p x_i \beta_i \right|_{n \times 1}^2 + \lambda \|\beta\|_{L1}$$

for k in 1 to p
fix $\beta_j, j \neq k$

coordinate descent

solve β_k

$$L(\beta_k) = \frac{1}{2} \left| Y - \sum_{j=1, j \neq k}^p x_j \beta_j - x_k \beta_k \right|^2 + \lambda \sum_{j \neq k} |\beta_j| + \lambda |\beta_k|$$

$$\downarrow$$

$$\left| \hat{Y} - x_k \beta_k \right|^2 + \lambda |\beta_k|$$

start from big λ [β s starts at 0]
gradually reduce λ [β s increase]
solution path.

