

maybe?

Lecture 14: Ridge Regression

Ideally: If I have P potential covariates then I could calc. a penalized metric on all potential models and choose model w/ best metric.

Problem: generally there are 2^P possible models

Use a greedy approach

(1) Forward Stepwise Selection

(i) start w/ model just containing intercept

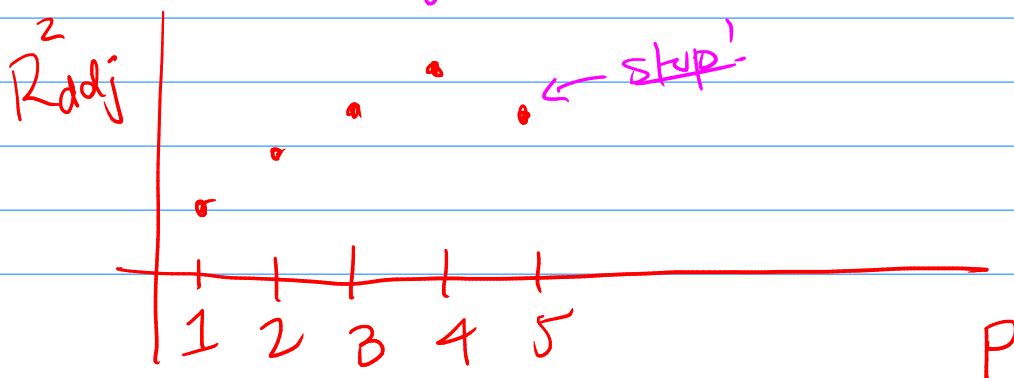
(ii) add best variable to model

(iii) Calc. my penalized metric

choose var. that improves metric most

(iv) Loop to step (ii)

Stop when my penalized metric gets worse



② Backwards Stepwise Selection

- starts w/ a model w/ all vars
 - removes one at a time
-

Trying all possible subsets is a very flexible model building process

Might be better off w/ a more constrained stepwise search.

(think flexibility and Bias-Var tradeoff)

Ridge Regression

Q: Can we deal with ill-conditioning in a continuous way?

Recall for OLS (ordinary least squares)

$$L(\beta) = \text{RSS}(\beta) = \|y - X\beta\|^2$$

and $\hat{\beta}^{(\text{OLS})} = \underset{\beta}{\operatorname{argmin}} L(\beta)$

problem is that if some variables are highly co-linear (highly correlated) the assoc. $\hat{\beta}$ s tend to blow up ($\pm \infty$)

Ridge regression solves this by penalizing large β s

$$\hat{\beta}^{(\text{ridge})} = \arg \min_{\beta} \left[L(\beta) + \lambda \|\beta\|^2 \right]$$

penalty

$\lambda > 0$ = regularization strength

by adding $\lambda \|\beta\|^2$ we penalize/avoid choosing $\hat{\beta}$ w/ really large components

$\lambda = 0 \Rightarrow$ OLS estimates

$\lambda \rightarrow \infty \Rightarrow \hat{\beta} \rightarrow 0$

① Typically: ignore β_0
 $\|\beta\|^2 = \sum_{j=1}^p \beta_j^2$

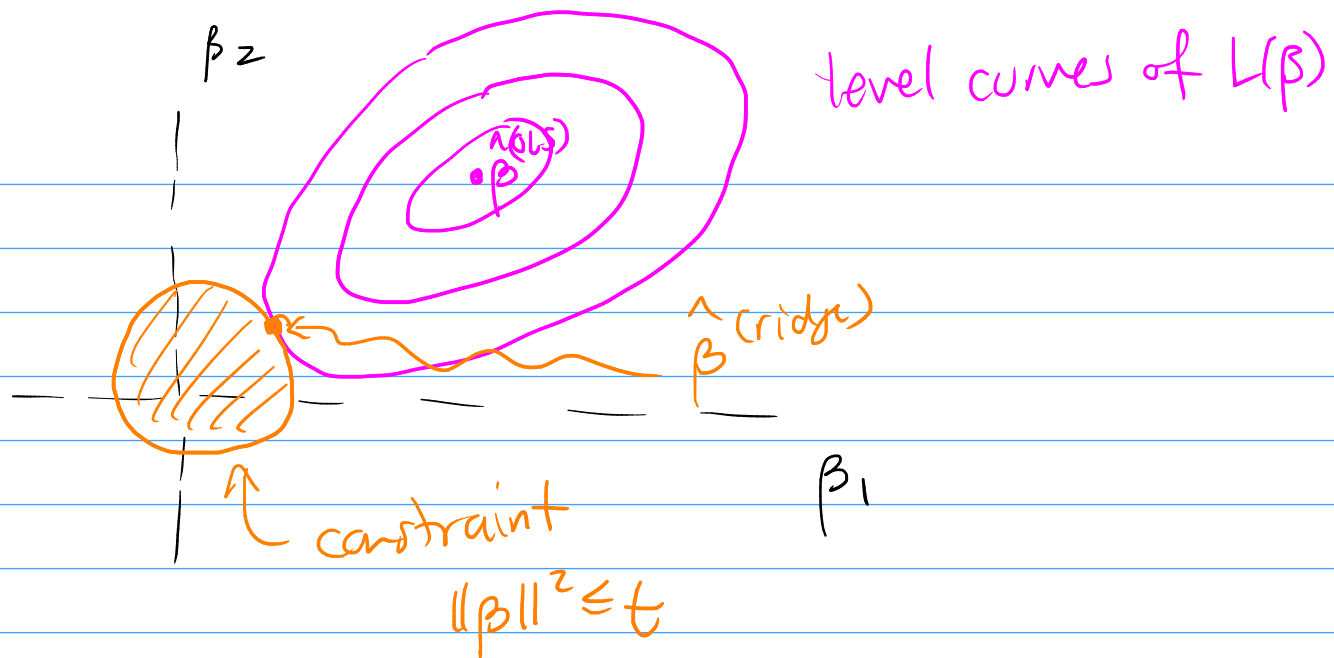
② Typically: standardize X s so that β s are comparable

③ Typically: choose λ using x -validation

Second interpretation: Ridge is equivalent to

$$\hat{\beta}^{(\text{ridge})} = \arg \min_{\beta} L(\beta) \quad \text{s.t.} \quad \|\beta\|^2 \leq t$$

1-1
corresp.
w/ λ



How do I get $\hat{\beta}^{(ridge)}$. Because penalty $\lambda \|\beta\|^2$ is quadratic, there is a closed form soln for $\hat{\beta}^{(ridge)}$.

OLS: had to solve $(X^T X) \beta = X^T y$ may not be invertible
 to get $\hat{\beta}^{(OLS)} = (X^T X)^{-1} X^T y$

Ridge: if we take deriv. and set to zero we get $(X^T X + \lambda I) \beta = X^T y$ $\lambda > 0$ this is always invertible

so $\hat{\beta}^{(ridge)} = (X^T X + \lambda I)^{-1} X^T y$

So I've replaced $X^T X$ w/ $X^T X + \lambda I$

For OLS conditiony depended on $K(X^T X)$
 Ridge " " $K(X^T X + \lambda I)$