

Homework 7

Due Thursday, 10/31/19

1. (Leave one out cross validation) Let $\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix} \in \mathbb{R}^{n \times p}$ be a full rank design matrix, $\mathbf{H} =$

$$\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T, h_{ii} \text{ be the } i\text{th diagonal element of } \mathbf{H} \text{ and } \mathbf{Y} = \begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix} \in \mathbb{R}^n \text{ be the response.}$$

For each $i = 1, \dots, n$, let $\mathbf{X}_{(-i)} \in \mathbb{R}^{(n-1) \times p}$ be the design matrix with the i th row removed and $\mathbf{Y}_{(-i)} \in \mathbb{R}^{n-1}$ be the response with the i th element removed. Define

$$\hat{\beta}_{(-i)} = (\mathbf{X}_{(-i)}^T \mathbf{X}_{(-i)})^{-1} \mathbf{X}_{(-i)}^T \mathbf{Y}_{(-i)}$$

to be the ordinary least squares estimator that ignores the i th sample.

- (a) Show that $(\mathbf{X}_{(-i)}^T \mathbf{X}_{(-i)})^{-1} = (\mathbf{X}^T \mathbf{X})^{-1} + (1 - h_{ii})^{-1} \mathbf{x}_i \mathbf{x}_i^T (\mathbf{X}^T \mathbf{X})^{-1}$.
 (b) Show that the leave one out cross validation error,

$$PRESS = \sum_{i=1}^n (Y_i - \mathbf{x}_i^T \hat{\beta}_{(-i)})^2,$$

can be written as

$$PRESS = \sum_{i=1}^n \left(\frac{Y_i - \hat{Y}_i}{1 - h_{ii}} \right)^2,$$

where \hat{Y}_i is the i th element of $\mathbf{H}\mathbf{Y} \in \mathbb{R}^n$.

2. (AIC and BIC) Assume $\mathbf{X} \in \mathbb{R}^{n \times p}$ is a full rank design matrix and $\mathbf{Y} \sim N(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_n)$. Let $L(\mathbf{b}, v) = (2\pi v)^{-n/2} \exp\left\{-\frac{1}{2v} (\mathbf{Y} - \mathbf{X}\mathbf{b})^T (\mathbf{Y} - \mathbf{X}\mathbf{b})\right\}$ be the likelihood for the normal distribution and $\hat{L} = L(\hat{\boldsymbol{\beta}}_{(MLE)}, \hat{\sigma}_{(MLE)}^2)$, where $\hat{\boldsymbol{\beta}}_{(MLE)}$ and $\hat{\sigma}_{(MLE)}^2$ are the maximum likelihood estimates for $\boldsymbol{\beta}$ and σ^2 . Show that up to constants that do not depend on \mathbf{Y} , \mathbf{X} or p ,

$$AIC = n \log(SSE) + 2p, \quad BIC = n \log(SSE) + \log(n)p,$$

where AIC and BIC are defined on slide 15 of lecture 15. Use this to show that for Gaussian data and n suitably large, the model chosen by AIC is always at least as large as that chosen by BIC.

3. Consider the data Boston.txt that we looked at in lecture on 10/24 and regress mvalue onto all predictors. Using residual \times covariate plots, suggest at least one way to improve the mean model. Use AIC, BIC or adjusted R2 (or all 3) to show that your new mean model is superior to the one we considered in class.

4. Consider the data “Steam.txt” where you are trying to predict steam usage using fat, glycerine, wind, freezday and temp. Perform forward and backward selection with $\alpha_1 = 0.1$ and $\alpha_2 = 0.2$, respectively. Which covariates did you include in the full model? How does this compare when you perform best subset regression using AIC to rank models? Does the model chosen with best subset regression change when you use BIC to rank models?