

Model Comparison + Information Criteria

Mallow's C_p
bias/variance tradeoff
prediction error

Prediction Problem: Given data \underline{x} , $X = (x_1 \cdots x_p)$

- 1) Obtain OLS estimate $\hat{\beta}$
- 2) Construct fitted values $\hat{m} = X\hat{\beta} \in \mathbb{R}^{n \times 1}$
- 3) Use \hat{m} to predict x^* , a new data vector generated under the same conditions as \underline{x} , the original data vector.

this means $E[\underline{x}] = E[x^*] = \underline{m}$, (might not be $X\beta$!)

Q: How well does \hat{m} do at predicting y^* ?

A: Evaluate with expected prediction sum of squares

$$PSS = \|y^* - \hat{m}\|^2, \text{ evaluate } E[PSS]$$

Setup: $\underline{x} = \underline{m} + \underline{\epsilon}$ $V(\underline{\epsilon}) = V(\underline{\epsilon}^*) = \sigma^2 I$, ϵ, ϵ^* uncorrelated.
 $\underline{x}^* = \underline{m} + \underline{\epsilon}^*$

Given X , $\hat{\beta} = (X^T X)^{-1} X^T y$
 $\hat{m} = X\hat{\beta} = X(X^T X)^{-1} X^T y \in Py$ (P is a projection matrix)

$$\text{IDENTITY \#1 : } E[\|y - \hat{m}\|^2] = n\sigma^2 + p\sigma^2 + \|(\mathbf{I} - \mathbf{P})m\|^2$$

↑ ↑ ↑
 sampling variability estimation var bias²

$$\begin{aligned} \text{Proof: } \|y - \hat{m}\|^2 &= \|m + \epsilon^* - \hat{m}\|^2 \\ &= \|\epsilon^* + (m - \hat{m})\|^2 \\ &= \|\epsilon^*\|^2 + \|m - \hat{m}\|^2 + \epsilon^{*T}(m - \hat{m}) \end{aligned}$$

$$\begin{aligned} E["] &= E[\|\epsilon^*\|^2] + E[\|m - \hat{m}\|^2] + E[\epsilon^{*T}(m - \hat{m})] \\ &= n\sigma^2 + E[\|m - \hat{m}\|^2] + 0 \\ &= n\sigma^2 + \underline{E[\|m - \hat{m}\|^2]} \end{aligned}$$

$$\begin{aligned} \|m - \hat{m}\|^2 &= \|m - \mathbf{P}_y\|^2 = \|m - \mathbf{P}m - \mathbf{P}\epsilon\|^2 \\ &= \|(\mathbf{I} - \mathbf{P})m - \mathbf{P}\epsilon\|^2 \\ &= \|(\mathbf{I} - \mathbf{P})m\|^2 + \|\mathbf{P}\epsilon\|^2 - 2\epsilon^{*T}\mathbf{P}(\mathbf{I} - \mathbf{P})m \\ &= \|(\mathbf{I} - \mathbf{P})m\|^2 + \epsilon^{*T}\mathbf{P}\epsilon - 0 \end{aligned}$$

$$E(") = \|(\mathbf{I} - \mathbf{P})m\|^2 + p\sigma^2$$

$$\text{So } E[\|y - \hat{m}\|^2] = n\sigma^2 + p\sigma^2 + \underline{\|(\mathbf{I} - \mathbf{P})m\|^2}$$

① All else being equal, prediction error is necessary in \mathbf{P}

② Suppose $m = X\beta$. Then $(\mathbf{I} - \mathbf{P}_x)m = X\beta - X(X^T X)^{-1}X^T\beta$
 $= \underline{0}$ Bias is zero.

③ Suppose $X = [X_1, X_2]$ but $\hat{m} = X_1\beta$ (too many predictors)

$$X_1 \in \mathbb{R}^{n \times p_1}, X_2 \in \mathbb{R}^{n \times p_2}$$

if we use $\hat{u} = x_1(x_1^T x_1)^{-1} x_1^T y$, then
 $= P_1 y$

$$\begin{aligned} E(\|y^* - \hat{u}\|^2) &= n\sigma^2 + p_1\sigma^2 + \|(\mathbf{I} - P_1) u\|^2 \\ &= n\sigma^2 + p_1\sigma^2 \end{aligned}$$

if we use $\hat{u} = x(x^T x)^{-1} x^T y$ ($x = (x_1, x_2)$, too many predictors)
 $= P_x y$

$$\begin{aligned} E(\|y^* - \hat{u}\|^2) &= n\sigma^2 + (p_1 + p_2)\sigma^2 + \|(\mathbf{I} - P_x) u\|^2 \\ &= n\sigma^2 + (p_1 + p_2)\sigma^2 \end{aligned}$$

So adding unnecessary predictors increases the prediction error.

④ If $u = [x \ x^*] \beta = \tilde{x} \beta$ (not enough predictors)

$$\text{then } E(\|y^* - \hat{u}\|^2) = n\sigma^2 + np^2 + \|(\mathbf{I} - P_x) u\|^2$$

\uparrow
not zero.

$$\text{In general, } E(\|y^* - \hat{u}\|^2) = n\sigma^2 + p\sigma^2 + \|(\mathbf{I} - P) u\|^2$$

|
goes up when
adding predictors

|
goes down when
adding predictors.

Goal: Add predictors that reduce bias
Don't add unnecessary predictors

How to do this? Find a way to estimate the prediction error

$$\text{IDENTITY \#2: } E[\|y^* - \hat{y}\|^2] = E[\text{RSS}] + 2p\sigma^2$$

This means $\text{RSS} + 2p\sigma^2$ is an unbiased est. of $E[\text{PSS}]$

- In particular, this shows that
- 1) RSS is a bad est. of PSS
 - 2) RSS gets worse at est PSS as p gets bigger.

In general, don't know σ^2 , but suppose $\hat{\sigma}^2$ is an unbiased est.

Then $\underline{\text{RSS} + 2p\hat{\sigma}^2}$ is an unbiased est. of $E[\text{PSS}]$.

Model Selection with C_p

Consider two models, for y , $y \sim X_1$ $\text{ncol}(X_1) = p_1$ Model 1
 $y \sim X_2$ $\text{ncol}(X_2) = p_2$ Model 2

- 1) fit M1, get RSS₁
- 2) fit M2, get RSS₂
- 3) obtain estimate $\hat{\sigma}^2$ of σ^2 by fitting a model with all available regressors.

Then, prefer M1 to M2 if $\text{RSS}_1 + 2p_1\hat{\sigma}^2 < \text{RSS}_2 + 2p_2\hat{\sigma}^2$

$\text{RSS}_j + 2p_j\hat{\sigma}^2$ is the "C_p statistic" for model j.

$$C_p(M_j) = \text{RSS}_j + 2p_j\hat{\sigma}^2$$

Sometimes, C_p is defined as $C_p(M_j) = \frac{\text{RSS}_j}{\hat{\sigma}^2} + 2p_j - n$

but this doesn't matter for model comparison.