# 1 Model Selection

## 1.1 Introduction

- It won't always be obvious **which predictors** should be **included** in a multiple regression model and **which** should be **omitted**.

  When the predictors are *uncorrelated*, it's safe to simply drop from the model those whose coefficients aren't statistically significant according to the **$t$ tests**.

  More often, though, there will be some degree of multicollinearity among the predictors, and in this case special ***model selection*** procedures should be used.

  The goal is to find a model that accomplishes *both* of two *competing* objectives:

  1. The model should **fit** the data **well**.
  2. The model should be **parsimonious** (contain only a small number of predictors).

  The challenge is that there's a **tradeoff** – the more parsimonious the model, the less well it fits the data.

- We'll look at several **model selection criteria** for comparing models:

  1. $R_p^2$ and $\mathrm{SSE}_p$
  2. $R_{a,p}^2$ and $\mathrm{MSE}_p$
  3. $\mathrm{AIC}_p$ and $\mathrm{BIC}_p$ (or $\mathrm{SBC}_p$)
  4. $\mathrm{PRESS}_p$

- We'll use the following notation:

$$P - 1 \;=\; \text{The total number of predictors } \textbf{available} \text{ for inclusion in a model.}$$
$$p - 1 \;=\; \text{The number of predictors } \textbf{in} \text{ the model } \textbf{currently being considered} \text{ (so } p - 1 \;\leq\; P - 1).$$

## 1.2 $R_p^2$ and $\mathrm{SSE}_p$

- $R_p^2$ and $\mathrm{SSE}_p$ are just the usual **coefficient of multiple determination** and **error sum of squares** (Class Notes 11), i.e.

$$\mathrm{SSE}_p \;=\; \sum_{i=1}^{n}(Y_i - \hat{Y}_i)^2$$

and

$$R_p^2 = 1 - \frac{\text{SSE}_p}{\text{SSTO}},$$

except now, because they'll be used to compare models with **different** numbers of parameters, the number of parameters in the model ($p$) is explicitly represented in the notation.

- Recall that a **small** $\text{SSE}_p$ and **large** $R_p^2$ indicate that the model **fits** the data **well**.

- But recall also that $\text{SSE}_p$ *always increases* and $R_p^2$ *decreases* when a predictor is dropped from a model. So **neither** of these is useful for comparing two models that have **different** numbers of predictors.

## 1.3 $R_{a,p}^2$ and $\text{MSE}_p$

- $R_{a,p}^2$ and $\text{MSE}_p$ are just the **adjusted coefficient of multiple determination** and usual **mean squared error** (Class Notes 11), i.e.

$$\text{MSE}_p = \frac{\text{SSE}_p}{n - p}$$

and

$$R_{a,p}^2 = 1 - \frac{\text{SSE}_p/(n - p)}{\text{SSTO}/(n - 1)},$$

except now the number of parameters in the model ($p$) is explicitly represented in the notation.

- Recall that a **small** $\text{MSE}_p$ and **large** $R_{a,p}^2$ indicate that the model **fits** the data **well**.

- These criteria take into account the number of predictors in the model, so that they're **useful** for comparing two models that have **different** numbers of predictors.

  Using these criteria, the model that has **larger** $R_{a,p}^2$ or, equivalently, **smaller** $\text{MSE}_p$ is **preferred**.

## 1.4 $\text{AIC}_p$ and $\text{BIC}_p$

- *Akaike's Information Criterion*, denoted **AIC$_p$**, is defined as

> **Akaike's Information Criterion:**
>
> $$\text{AIC}_p = n \log \text{SSE}_p - n \log n + 2p$$

- The *__Bayesian Information Criterion__*, denoted **BIC$_p$** is defined as

> **Bayesian Information Criterion**:
>
> $$\text{BIC}_p \;=\; n \log \text{SSE}_p - n \log n + (\log n)\, p$$

- For both $\text{AIC}_p$ and $\text{BIC}_p$:

  1. The first term $\boldsymbol{n \log \textbf{SSE}_p}$ will be **small** if the model **fits** the data **well**.
  2. The second term $\boldsymbol{n \log n}$ is **constant** for fixed $n$ (i.e. it doesn't depend on the how many predictors are in the model or on how well the model fits the data).
  3. The last term $\boldsymbol{2p}$ or $\boldsymbol{(\log n)\, p}$ will be **small** if the model is **parsimonious** (i.e. if the number of predictors in the model, $\boldsymbol{p - 1}$, is small).

  Using these criteria, the model that has **smaller** $\text{AIC}_p$ (or $\text{BIC}_p$) is **preferred**, and accomplishes better the *two* competing objectives of Subsection 1.1.

  Note that the term $\boldsymbol{2p}$ in $\text{AIC}_p$ (and $\boldsymbol{(\log n)\, p}$ in $\text{BIC}_p$) acts as a *__penalty__* for including too many predictors in the model.

## 1.5 PRESS$_p$

- The idea behind **PRESS$_p$** is to successively **delete one observation** (row) at a time from the data set, **fit** a given **model** to the **remaining $\boldsymbol{n - 1}$ observations**, and for each fitted model calculate the *__delete-one prediction error__*

$$\text{Delete-one Prediction Error} \;=\; Y_i - \hat{Y}_{i(i)}\,,$$

  where $\hat{Y}_{i(i)}$ is the predicted value for the deleted $Y_i$ based on the model fitted to the other $n - 1$ observations.

  **PRESS$_p$** is the **sum** of **squared delete-one prediction errors**:

> **PRESS$_p$**:
>
> $$\text{PRESS}_p \;=\; \sum_{i=1}^{n} (Y_i - \hat{Y}_{i(i)})^2$$

  Using this criteria, the model that has **smaller PRESS$_p$** is **preferred**.