## 6.1.3  Choosing the Optimal Model

Best subset selection, forward selection, and backward selection result in
the creation of a set of models, each of which contains a subset of the $p$
predictors. To apply these methods, we need a way to determine which of
these models is *best*. As we discussed in Section 6.1.1, the model containing
all of the predictors will always have the smallest RSS and the largest $R^2$,
since these quantities are related to the training error. Instead, we wish to
choose a model with a low test error. As is evident here, and as we show
in Chapter 2, the training error can be a poor estimate of the test error.
Therefore, RSS and $R^2$ are not suitable for selecting the best model among
a collection of models with different numbers of predictors.

In order to select the best model with respect to test error, we need to
estimate this test error. There are two common approaches:

1. We can indirectly estimate test error by making an *adjustment* to the
   training error to account for the bias due to overfitting.

2. We can *directly* estimate the test error, using either a validation set
   approach or a cross-validation approach, as discussed in Chapter 5.

We consider both of these approaches below.

### $C_p$, AIC, BIC, and Adjusted $R^2$

We show in Chapter 2 that the training set MSE is generally an under-
estimate of the test MSE. (Recall that MSE = RSS/$n$.) This is because
when we fit a model to the training data using least squares, we specifi-
cally estimate the regression coefficients such that the training RSS (but
not the test RSS) is as small as possible. In particular, the training error
will decrease as more variables are included in the model, but the test error
may not. Therefore, training set RSS and training set $R^2$ cannot be used
to select from among a set of models with different numbers of variables.

However, a number of techniques for *adjusting* the training error for the
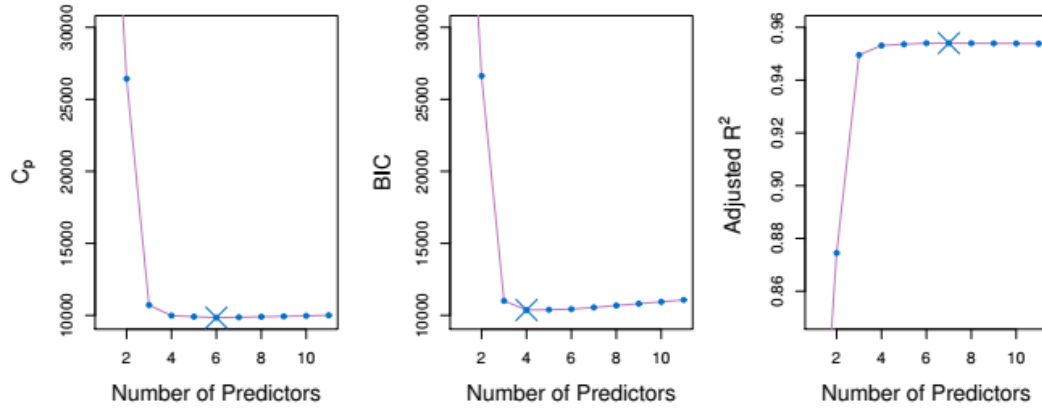model size are available. These approaches can be used to select among a set

**FIGURE 6.2.** $C_p$, BIC, and adjusted $R^2$ are shown for the best models of each size for the Credit data set (the lower frontier in Figure 6.1). $C_p$ and BIC are estimates of test MSE. In the middle plot we see that the BIC estimate of test error shows an increase after four variables are selected. The other two plots are rather flat after four variables are included.

of models with different numbers of variables. We now consider four such approaches: $C_p$, *Akaike information criterion* (AIC), *Bayesian information criterion* (BIC), and *adjusted $R^2$*. Figure 6.2 displays $C_p$, BIC, and adjusted $R^2$ for the best model of each size produced by best subset selection on the Credit data set.

For a fitted least squares model containing $d$ predictors, the $C_p$ estimate of test MSE is computed using the equation

$$C_p = \frac{1}{n}\left(\text{RSS} + 2d\hat{\sigma}^2\right), \tag{6.2}$$

$C_p$
Akaike information criterion
Bayesian information criterion
adjusted $R^2$

2

where $\hat{\sigma}^2$ is an estimate of the variance of the error $\epsilon$ associated with each response measurement in (6.1).[3] Typically $\hat{\sigma}^2$ is estimated using the full model containing all predictors. Essentially, the $C_p$ statistic adds a penalty of $2d\hat{\sigma}^2$ to the training RSS in order to adjust for the fact that the training error tends to underestimate the test error. Clearly, the penalty increases as the number of predictors in the model increases; this is intended to adjust for the corresponding decrease in training RSS. Though it is beyond the scope of this book, one can show that if $\hat{\sigma}^2$ is an unbiased estimate of $\sigma^2$ in (6.2), then $C_p$ is an unbiased estimate of test MSE. As a consequence, the $C_p$ statistic tends to take on a small value for models with a low test error, so when determining which of a set of models is best, we choose the model with the lowest $C_p$ value. In Figure 6.2, $C_p$ selects the six-variable model containing the predictors `income`, `limit`, `rating`, `cards`, `age` and `student`.

---

[3]Mallow's $C_p$ is sometimes defined as $C_p' = \text{RSS}/\hat{\sigma}^2 + 2d - n$. This is equivalent to the definition given above in the sense that $C_p = \frac{1}{n}\hat{\sigma}^2(C_p' + n)$, and so the model with smallest $C_p$ also has smallest $C_p'$.

The AIC criterion is defined for a large class of models fit by maximum likelihood. In the case of the model (6.1) with Gaussian errors, maximum likelihood and least squares are the same thing. In this case AIC is given by

$$\text{AIC} = \frac{1}{n}\left(\text{RSS} + 2d\hat{\sigma}^2\right),$$

where, for simplicity, we have omitted irrelevant constants.[4] Hence for least squares models, $C_p$ and AIC are proportional to each other, and so only $C_p$ is displayed in Figure 6.2.

BIC is derived from a Bayesian point of view, but ends up looking similar to $C_p$ (and AIC) as well. For the least squares model with $d$ predictors, the BIC is, up to irrelevant constants, given by

$$\text{BIC} = \frac{1}{n}\left(\text{RSS} + \log(n)d\hat{\sigma}^2\right). \tag{6.3}$$

Like $C_p$, the BIC will tend to take on a small value for a model with a low test error, and so generally we select the model that has the lowest BIC value. Notice that BIC replaces the $2d\hat{\sigma}^2$ used by $C_p$ with a $\log(n)d\hat{\sigma}^2$ term, where $n$ is the number of observations. Since $\log n > 2$ for any $n > 7$, the BIC statistic generally places a heavier penalty on models with many variables, and hence results in the selection of smaller models than $C_p$. In Figure 6.2, we see that this is indeed the case for the `Credit` data set; BIC chooses a model that contains only the four predictors `income`, `limit`, `cards`, and `student`. In this case the curves are very flat and so there does not appear to be much difference in accuracy between the four-variable and six-variable models.

The adjusted $R^2$ statistic is another popular approach for selecting among a set of models that contain different numbers of variables. Recall from Chapter 3 that the usual $R^2$ is defined as $1 - \text{RSS}/\text{TSS}$, where $\text{TSS} = \sum(y_i - \bar{y})^2$ is the *total sum of squares* for the response. Since RSS always decreases as more variables are added to the model, the $R^2$ always increases as more variables are added. For a least squares model with $d$ variables, the adjusted $R^2$ statistic is calculated as

$$\text{Adjusted } R^2 = 1 - \frac{\text{RSS}/(n-d-1)}{\text{TSS}/(n-1)}. \tag{6.4}$$

Unlike $C_p$, AIC, and BIC, for which a *small* value indicates a model with a *low* test error, a *large* value of adjusted $R^2$ indicates a model with a

---

[4]There are two formulas for AIC for least squares regression. The formula that we provide here requires an expression for $\sigma^2$, which we obtain using the full model containing all predictors. The second formula is appropriate when $\sigma^2$ is unknown and we do not want to explicitly estimate it; that formula has a log(RSS) term instead of an RSS term. Detailed derivations of these two formulas are outside of the scope of this book.

small test error. Maximizing the adjusted $R^2$ is equivalent to minimizing $\frac{\text{RSS}}{n-d-1}$. While RSS always decreases as the number of variables in the model increases, $\frac{\text{RSS}}{n-d-1}$ may increase or decrease, due to the presence of $d$ in the denominator.

The intuition behind the adjusted $R^2$ is that once all of the correct variables have been included in the model, adding additional *noise* variables will lead to only a very small decrease in RSS. Since adding noise variables leads to an increase in $d$, such variables will lead to an increase in $\frac{\text{RSS}}{n-d-1}$, and consequently a decrease in the adjusted $R^2$. Therefore, in theory, the model with the largest adjusted $R^2$ will have only correct variables and no noise variables. Unlike the $R^2$ statistic, the adjusted $R^2$ statistic *pays a price* for the inclusion of unnecessary variables in the model. Figure 6.2 displays the adjusted $R^2$ for the `Credit` data set. Using this statistic results in the selection of a model that contains seven variables, adding `own` to the model selected by $C_p$ and AIC.

$C_p$, AIC, and BIC all have rigorous theoretical justifications that are beyond the scope of this book. These justifications rely on asymptotic arguments (scenarios where the sample size $n$ is very large). Despite its popularity, and even though it is quite intuitive, the adjusted $R^2$ is not as well motivated in statistical theory as AIC, BIC, and $C_p$. All of these measures are simple to use and compute. Here we have presented their formulas in the case of a linear model fit using least squares; however, AIC and BIC can also be defined for more general types of models.