# 5.5 Selecting predictors

When there are many possible predictors, we need some strategy for selecting the best predictors to use in a regression model.

A common approach that is *not recommended* is to plot the forecast variable against a particular predictor and if there is no noticeable relationship, drop that predictor from the model. This is invalid because it is not always possible to see the relationship from a scatterplot, especially when the effects of other predictors have not been accounted for.

Another common approach which is also invalid is to do a multiple linear regression on all the predictors and disregard all variables whose $p$-values are greater than 0.05. To start with, statistical significance does not always indicate predictive value. Even if forecasting is not the goal, this is not a good strategy because the $p$-values can be misleading when two or more predictors are correlated with each other (see Section 5.9).

Instead, we will use a measure of predictive accuracy. Five such measures are introduced in this section. They can be calculated using the `CV()` function, here applied to the model for US consumption:

```
CV(fit.consMR)
#>         CV       AIC      AICc       BIC      AdjR2
#>     0.1163 -409.2980 -408.8314 -389.9114    0.7486
```

We compare these values against the corresponding values from other models. For the CV, AIC, AICc and BIC measures, we want to find the model with the lowest value; for Adjusted $R^2$, we seek the model with the highest value.

## Adjusted $R^2$

Computer output for a regression will always give the $R^2$ value, discussed in Section 5.2. However, it is not a good measure of the predictive ability of a model. It measures how well the model fits the historical data, but not how well the model will forecast future data.

In addition, $R^2$ does not allow for "degrees of freedom." Adding *any* variable tends to increase the value of $R^2$, even if that variable is irrelevant. For these reasons, forecasters should not use $R^2$ to determine whether a model will give good predictions, as it will lead to overfitting.

An equivalent idea is to select the model which gives the minimum sum of squared errors (SSE), given by

$$\text{SSE} = \sum_{t=1}^{T} e_t^2.$$

Minimising the SSE is equivalent to maximising $R^2$ and will always choose the model with the most variables, and so is not a valid way of selecting predictors.

An alternative which is designed to overcome these problems is the adjusted $R^2$ (also called "R-bar-squared"):

$$\bar{R}^2 = 1 - (1 - R^2)\frac{T-1}{T-k-1},$$

where $T$ is the number of observations and $k$ is the number of predictors. This is an improvement on $R^2$, as it will no longer increase with each added predictor. Using this measure, the best model will be the one with the largest value of $\bar{R}^2$. Maximising $\bar{R}^2$ is equivalent to minimising the standard error $\hat{\sigma}_e$ given in Equation (5.3).

Maximising $\bar{R}^2$ works quite well as a method of selecting predictors, although it does tend to err on the side of selecting too many predictors.

## Cross-validation

Time series cross-validation was introduced in Section 3.4 as a general tool for determining the predictive ability of a model. For regression models, it is also possible to use classical leave-one-out cross-validation to selection predictors (Bergmeir, Hyndman, & Koo, 2018). This is faster and makes more efficient use of the data. The procedure uses the following steps:

1. Remove observation $t$ from the data set, and fit the model using the remaining data. Then compute the error ($e_t^* = y_t - \hat{y}_t$) for the omitted observation. (This is not the

same as the residual because the $t$th observation was not used in estimating the value of $\hat{y}_t$.)

2. Repeat step 1 for $t = 1, \ldots, T$.
3. Compute the MSE from $e_1^*, \ldots, e_T^*$. We shall call this the **CV**.

Although this looks like a time-consuming procedure, there are fast methods of calculating CV, so that it takes no longer than fitting one model to the full data set. The equation for computing CV efficiently is given in Section 5.7. Under this criterion, the best model is the one with the smallest value of CV.

## Akaike's Information Criterion

A closely-related method is Akaike's Information Criterion, which we define as

$$\text{AIC} = T \log\left(\frac{\text{SSE}}{T}\right) + 2(k + 2),$$

where $T$ is the number of observations used for estimation and $k$ is the number of predictors in the model. Different computer packages use slightly different definitions for the AIC, although they should all lead to the same model being selected. The $k + 2$ part of the equation occurs because there are $k + 2$ parameters in the model: the $k$ coefficients for the predictors, the intercept and the variance of the residuals. The idea here is to penalise the fit of the model (SSE) with the number of parameters that need to be estimated.

The model with the minimum value of the AIC is often the best model for forecasting. For large values of $T$, minimising the AIC is equivalent to minimising the CV value.

## Corrected Akaike's Information Criterion

For small values of $T$, the AIC tends to select too many predictors, and so a bias-corrected version of the AIC has been developed,

$$\text{AIC}_{\text{c}} = \text{AIC} + \frac{2(k + 2)(k + 3)}{T - k - 3}.$$

As with the AIC, the AICc should be minimised.

## Schwarz's Bayesian Information Criterion

A related measure is Schwarz's Bayesian Information Criterion (usually abbreviated to BIC, SBIC or SC):

$$\text{BIC} = T \log\left(\frac{\text{SSE}}{T}\right) + (k+2)\log(T).$$

As with the AIC, minimising the BIC is intended to give the best model. The model chosen by the BIC is either the same as that chosen by the AIC, or one with fewer terms. This is because the BIC penalises the number of parameters more heavily than the AIC. For large values of $T$, minimising BIC is similar to leave-$v$-out cross-validation when $v = T[1 - 1/(\log(T) - 1)]$.

## Which measure should we use?

While $\bar{R}^2$ is widely used, and has been around longer than the other measures, its tendency to select too many predictor variables makes it less suitable for forecasting.

Many statisticians like to use the BIC because it has the feature that if there is a true underlying model, the BIC will select that model given enough data. However, in reality, there is rarely, if ever, a true underlying model, and even if there was a true underlying model, selecting that model will not necessarily give the best forecasts (because the parameter estimates may not be accurate).

Consequently, we recommend that one of the AICc, AIC, or CV statistics be used, each of which has forecasting as their objective. If the value of $T$ is large enough, they will all lead to the same model. In most of the examples in this book, we use the AICc value to select the forecasting model.

### Example: US consumption

In the multiple regression example for forecasting US consumption we considered four predictors. With four predictors, there are $2^4 = 16$ possible models. Now we can check if all four predictors are actually useful, or whether we can drop one or more of them. All 16 models were fitted and the results are summarised in Table 5.1. A "1" indicates

that the predictor was included in the model, and a "0" means that the predictor was not included in the model. Hence the first row shows the measures of predictive accuracy for a model including all four predictors.

The results have been sorted according to the AICc. Therefore the best models are given at the top of the table, and the worst at the bottom of the table.

Table 5.1: All 16 possible models for forecasting US consumption with 4 predictors.

| Income | Production | Savings | Unemployment | CV | AIC | AICc | BIC | AdjR2 |
|--------|-----------|---------|--------------|-----|-----|------|-----|-------|
| 1 | 1 | 1 | 1 | 0.116 | −409.3 | −408.8 | −389.9 | 0.749 |
| 1 | 0 | 1 | 1 | 0.116 | −408.1 | −407.8 | −391.9 | 0.746 |
| 1 | 1 | 1 | 0 | 0.118 | −407.5 | −407.1 | −391.3 | 0.745 |
| 1 | 0 | 1 | 0 | 0.129 | −388.7 | −388.5 | −375.8 | 0.716 |
| 1 | 1 | 0 | 1 | 0.278 | −243.2 | −242.8 | −227.0 | 0.386 |
| 1 | 0 | 0 | 1 | 0.283 | −237.9 | −237.7 | −225.0 | 0.365 |
| 1 | 1 | 0 | 0 | 0.289 | −236.1 | −235.9 | −223.2 | 0.359 |
| 0 | 1 | 1 | 1 | 0.293 | −234.4 | −234.0 | −218.2 | 0.356 |
| 0 | 1 | 1 | 0 | 0.300 | −228.9 | −228.7 | −216.0 | 0.334 |
| 0 | 1 | 0 | 1 | 0.303 | −226.3 | −226.1 | −213.4 | 0.324 |
| 0 | 0 | 1 | 1 | 0.306 | −224.6 | −224.4 | −211.7 | 0.318 |
| 0 | 1 | 0 | 0 | 0.314 | −219.6 | −219.5 | −209.9 | 0.296 |
| 0 | 0 | 0 | 1 | 0.314 | −217.7 | −217.5 | −208.0 | 0.288 |
| 1 | 0 | 0 | 0 | 0.372 | −185.4 | −185.3 | −175.7 | 0.154 |
| 0 | 0 | 1 | 0 | 0.414 | −164.1 | −164.0 | −154.4 | 0.052 |
| 0 | 0 | 0 | 0 | 0.432 | −155.1 | −155.0 | −148.6 | 0.000 |

The best model contains all four predictors. However, a closer look at the results reveals some interesting features. There is clear separation between the models in the first four rows and the ones below. This indicates that Income and Savings are both more important variables than Production and Unemployment. Also, the first two rows have almost identical values of CV, AIC and AICc. So we could possibly drop the Production variable and get similar forecasts. Note that Production and Unemployment are highly (negatively) correlated, as shown in Figure 5.5, so most of the predictive information in Production is also contained in the Unemployment variable.

# Best subset regression

Where possible, all potential regression models should be fitted (as was done in the example above) and the best model should be selected based on one of the measures discussed. This is known as "best subsets" regression or "all possible subsets" regression.

## Stepwise regression

If there are a large number of predictors, it is not possible to fit all possible models. For example, 40 predictors leads to $2^{40} > 1$ trillion possible models! Consequently, a strategy is required to limit the number of models to be explored.

An approach that works quite well is *backwards stepwise regression*:

- Start with the model containing all potential predictors.
- Remove one predictor at a time. Keep the model if it improves the measure of predictive accuracy.
- Iterate until no further improvement.

If the number of potential predictors is too large, then the backwards stepwise regression will not work and *forward stepwise regression* can be used instead. This procedure starts with a model that includes only the intercept. Predictors are added one at a time, and the one that most improves the measure of predictive accuracy is retained in the model. The procedure is repeated until no further improvement can be achieved.

Alternatively for either the backward or forward direction, a starting model can be one that includes a subset of potential predictors. In this case, an extra step needs to be included. For the backwards procedure we should also consider adding a predictor with each step, and for the forward procedure we should also consider dropping a predictor with each step. These are referred to as *hybrid* procedures.

It is important to realise that any stepwise approach is not guaranteed to lead to the best possible model, but it almost always leads to a good model. For further details see James, Witten, Hastie, & Tibshirani (2014).

## Beware of inference after selecting predictors

We do not discuss statistical inference of the predictors in this book (e.g., looking at $p$-values associated with each predictor). If you do wish to look at the statistical significance of the predictors, beware that *any* procedure involving selecting predictors first will invalidate the assumptions behind the $p$-values. The procedures we recommend for selecting predictors are helpful when the model is used for forecasting; they are not helpful if you wish to study the effect of any predictor on the forecast variable.