

Stat 804

Lecture 13 Notes

Forecasting: an introduction

Given data X_0, \dots, X_{T-1} our goal will be to guess, or forecast, X_T or more generally X_{T+r} . There are a variety of *ad hoc* methods as well as a variety of statistically derived methods. I illustrate the *ad hoc* methods with the exponentially weighted moving average (EWMA). In this case we simply take

$$\hat{X}_T = (X_{T-1} + aX_{T-2} + a^2X_{T-3} + \dots + a^{T-1}X_0)/c(a, T)$$

where $c(a, T)$ makes it a weighted average: $c(a, T) = (1 - a^T)/(1 - a)$. If we take a near to 1 we are almost using the sample mean while if we take a near 0 we are virtually using X_{T-1} . You are supposed to choose a to trade off the desire to use lots of data against the possibility that the structure of the series has changed over time.

Statistically based methods concentrate on some measure of the size of $X_T - \hat{X}_T$; the **mean squared prediction error** $E[(X_T - \hat{X}_T)^2]$ is the most common.

In general \hat{X}_T must be some function $f(X_0, \dots, X_{T-1})$. The mean squared prediction error can be seen by conditioning on the data to be minimized by

$$\hat{X}_T = E(X_T | X_0, \dots, X_{T-1})$$

For most distributions of the X 's this would be hard to compute but for Gaussian processes the solution is the usual linear regression of X_T on the data, namely

$$\hat{X}_T = \mu_T + a_1(X_{T-1} - \mu_{T-1}) + \dots + a_T(X_0 - \mu_0)$$

where the coefficient vector a is given by

$$a = \text{Cov}(X_T, (X_{T-1}, \dots, X_0)^T) \text{Var}(X_{T-1}, \dots, X_0)^{-1}$$

When T is large the computation of these forecasts is difficult in general. There are some shortcuts, however.

Forecasting AR(p) processes

When the process is an AR the computation of the conditional expectation is easier:

$$\begin{aligned}\hat{X}_T &= E(X_T|X_0, \dots, X_{T-1}) \\ &= E(\epsilon_T + \sum_{i=1}^p a_i X_{T-i} | X_0, \dots, X_{T-1}) \\ &= \sum_{i=1}^p a_i X_{T-i}\end{aligned}$$

For $r > 0$ we have the recursion

$$\begin{aligned}E(X_{T+r}|X_0, \dots, X_{T-1}) &= E(\epsilon_{T+r} + \sum_{i=1}^p a_i X_{T+r-i} | X_0, \dots, X_{T-1}) \\ &= \sum_{i=1}^p a_i \hat{X}_{T+r-i}\end{aligned}$$

Notice the the forecast into the future uses current values where these are available and forecasts already calculated for the other X 's.

Forecasting ARMA(p, q) processes

An ARMA(p, q) can be inverted to be an infinite order AR process. We could then use the method just given for the AR except that now the formula actually mentions values of X_t for $t < 0$. In practice we simply truncate the series and ignore the missing terms in the forecast, assuming that the coefficients of these omitted terms are very small. Remember each term is built up out of a geometric series for $(I - \alpha B)^{-1}$ with $|\alpha| < 1$.

A more direct method goes like this:

$$\begin{aligned}\hat{X}_{T+r} &= E(\epsilon_{T+r}|X) + \sum_{i=1}^p a_i \hat{X}_{T+r-i} \\ &\quad + \sum_{i=1}^q b_i E(\epsilon_{T+r-i}|X)\end{aligned}$$

where now the conditioning “ $|X$ ” means given the observed data.

Whenever the time index on an epsilon is T or more the conditional expectations are 0. For $T+r-i < T$ we need to guess the value of ϵ_{T+r-i} . The

same recursion can be re-arranged to help compute $E(\epsilon_t|X)$ for $0 \leq t \leq T-1$, at least approximately:

$$E(\epsilon_t|X) = X_t - \sum a_i X_{t-i} + \sum b_i E(\epsilon_{t-i}|X)$$

This recursion works you backward but you have to get it started. Generally we start the recursion by putting

$$\hat{\epsilon}_t = 0$$

for negative t and then using the recursion. The coefficients b are such that the effect of getting these values of ϵ wrong is damped out at a geometric rate as we increase t so if we have enough data and the smallest root of the characteristic polynomial for the MA part is not too close to 1 then we will have accurate values for $\hat{\epsilon}_t$ for t near T .

As we discussed in the section on estimation these computed estimates of the epsilon's can be improved by backcasting the values of ϵ_t for negative t and then forecasting and backcasting, etc.

Forecasting ARIMA(p, d, q) series

If $Z = (I - B)^d X$ and X is ARIMA(p, d, q) then we: compute Z , forecast Z and reconstruct X by undoing the differencing. For $d = 1$ for example we just have

$$\hat{X}_t = \hat{Z}_t + \hat{X}_{t-1}.$$

Forecast standard errors

You should remind yourself that the computations of conditional expectations we have just made used the fact that the a 's and b 's are constants – the true parameter values. In fact we then replace the parameter values with estimates. The quality of our forecasts will be summarized by the forecast standard error:

$$\sqrt{E[(X_t - \hat{X}_t)^2]}.$$

We will compute this ignoring the estimation of the parameters and then discuss how much that might have cost us.

If $\hat{X}_t = E(X_t|X)$ then $E(\hat{X}_t) = E(X_t)$ so that our forecast standard error is just the variance of $X_t - \hat{X}_t$.

Consider first the case of an AR(1) and one step ahead forecasting:

$$X_T - \hat{X}_T = \epsilon_T.$$

The variance of this forecast is σ_ϵ^2 so that the forecast standard error is just σ_ϵ .

For forecasts further ahead in time we have

$$\hat{X}_{T+r} = a\hat{X}_{T+r-1}$$

and

$$X_{T+r} = aX_{T+r-1} + \epsilon_{T+r}$$

Subtracting we see that

$$\text{Var}(X_{T+r} - \hat{X}_{T+r}) = \sigma_\epsilon^2 + \text{Var}(X_{T+r-1} - \hat{X}_{T+r-1})$$

so that we may calculate forecast standard errors recursively. As $r \rightarrow \infty$ we can check that the forecast variance converges to

$$\sigma_\epsilon^2/(1 - a^2)$$

which is simply the variance of individual X s. When you forecast a *stationary* series far into the future the forecast error is just the standard deviation of the series.

Turn now to a general ARMA(p, q). Rewrite the process as the infinite order AR

$$X_t = \sum_{s>0} c_s X_{t-s} + \epsilon_t$$

to see that again, ignoring the truncation of the infinite sum in the forecast we have

$$X_T - \hat{X}_T = \epsilon_T$$

so that the one step ahead forecast standard error is again σ_ϵ .

Parallel to the AR(1) argument we see that

$$X_{T+r} - \hat{X}_{T+r} = \sum_{j=0}^{r-1} a_j (X_{T+j} - \hat{X}_{T+j}) + \epsilon_{T+r}.$$

The errors on the right hand side are not independent of one another so that computation of the variance requires either computation of the covariances or recognition of the fact that the right hand side is a linear combination of $\epsilon_T, \dots, \epsilon_{T+r}$.

A simpler approach is to write the process as an infinite order MA:

$$X_t = \epsilon_t + \sum_{s>0} d_s \epsilon_{t-s}$$

for suitable coefficients d_s . Now if we treat conditioning on the data as being effectively equivalent to conditioning on all X_t for $t < T$ we are effectively conditioning on ϵ_t for all $t < T$. This means that

$$\begin{aligned} E(X_{T+r} | X_{T-1}, X_{T-2}, \dots) &= E(X_{T+r} | \epsilon_{T-1}, \epsilon_{T-2}, \dots) \\ &= \sum_{s>r} d_s \epsilon_{T+r-s} \end{aligned}$$

and the forecast error is just

$$X_{T+r} - \hat{X}_{T+r} = \epsilon_t + \sum_{s=1}^r d_s \epsilon_{T+r-s}$$

so that the forecast standard error is

$$\sigma_\epsilon \sqrt{1 + \sum_{s=1}^r d_s^2}.$$

Again as $r \rightarrow \infty$ this converges to σ_X .

Finally consider forecasting the ARIMA(p, d, q) process $(I - B)^d X = W$ where W is ARMA(p, q). The forecast errors in X can clearly be written as a linear combination of forecast errors for W permitting the forecast error in X to be written as a linear combination of the underlying errors ϵ_t . As an example consider first the ARIMA(0,1,0) process $X_t = \epsilon_t + X_{t-1}$. The forecast of ϵ_{T+r} is just 0 and so the forecast of X_{T+r} is just

$$\hat{X}_{T+r} = \hat{X}_{T+r-1} = \dots = X_{T-1}.$$

The forecast error is

$$\epsilon_{T+r} + \dots + \epsilon_T$$

whose standard deviation is $\sigma\sqrt{r+1}$. Notice that the forecast standard error grows to infinity as $r \rightarrow \infty$. For a general ARIMA($p, 1, q$) we have

$$\hat{X}_{T+r} = \hat{X}_{T+r-1} + \hat{W}_{T+r}$$

and

$$X_{T+r} - \hat{X}_{T+r} = (W_{T+r} - \hat{W}_{T+r}) + \cdots + (W_T - \hat{W}_T)$$

which can be combined with the expression above for the forecast error for an ARMA(p, q) to compute standard errors.

Software

The S-Plus function *arima.forecast* can do the forecasting.

Comments

I have ignored the effects of parameter estimation throughout. In ordinary least squares when we predict the Y corresponding to a new x we get a forecast standard error of

$$\sqrt{\text{Var}(Y - x\hat{\beta})} = \sqrt{\text{Var}(\epsilon + x(\beta - \hat{\beta}))}$$

which is

$$\sigma \sqrt{1 + x(X^T X)^{-1}x^T}.$$

The procedure used here corresponds to ignoring the term $x(X^T X)^{-1}x^T$ which is the variance of the fitted value. Typically this value is rather smaller than the 1 to which it is added. In a 1 sample problem for instance it is simply $1/n$. Generally the major component of forecast error is the standard error of the noise and the effect of parameter estimation is unimportant.

Forecast standard errors

You should remind yourself that the computations of conditional expectations we have made used the fact that the a 's and b 's are constants – the true parameter values. In fact we then replace the parameter values with estimates. The quality of our forecasts will be summarized by the forecast standard error:

$$\sqrt{\text{E}[(X_t - \hat{X}_t)^2]}.$$

We will compute this ignoring the estimation of the parameters and then discuss how much that might have cost us.

If $\hat{X}_t = \text{E}(X_t|X)$ then $\text{E}(\hat{X}_t) = \text{E}(X_t)$ so that our forecast standard error is just the variance of $X_t - \hat{X}_t$.

Consider first the case of an AR(1) and one step ahead forecasting:

$$X_T - \hat{X}_T = \epsilon_T.$$

The variance of this forecast is σ_ϵ^2 so that the forecast standard error is just σ_ϵ .

For forecasts further ahead in time we have

$$\hat{X}_{T+r} = a\hat{X}_{T+r-1}$$

and

$$X_{T+r} = aX_{T+r-1} + \epsilon_{T+r}$$

Subtracting we see that

$$\text{Var}(X_{T+r} - \hat{X}_{T+r}) = \sigma_\epsilon^2 + \text{Var}(X_{T+r-1} - \hat{X}_{T+r-1})$$

so that we may calculate forecast standard errors recursively. As $r \rightarrow \infty$ we can check that the forecast variance converges to

$$\sigma_\epsilon^2/(1 - a^2)$$

which is simply the variance of individual X s. When you forecast a *stationary* series far into the future the forecast error is just the standard deviation of the series.

Turn now to a general ARMA(p, q). Rewrite the process as the infinite order AR

$$X_t = \sum_{s>0} c_s X_{t-s} + \epsilon_t$$

to see that again, ignoring the truncation of the infinite sum in the forecast we have

$$X_T - \hat{X}_T = \epsilon_T$$

so that the one step ahead forecast standard error is again σ_ϵ .

Parallel to the AR(1) argument we see that

$$X_{T+r} - \hat{X}_{T+r} = \sum_{j=0}^{r-1} a_j (X_{T+j} - \hat{X}_{T+j}) + \epsilon_{T+r}.$$

The errors on the right hand side are not independent of one another so that computation of the variance requires either computation of the covariances or recognition of the fact that the right hand side is a linear combination of $\epsilon_T, \dots, \epsilon_{T+r}$.

A simpler approach is to write the process as an infinite order MA:

$$X_t = \epsilon_t + \sum_{s>0} d_s \epsilon_{t-s}$$

for suitable coefficients d_s . Now if we treat conditioning on the data as being effectively equivalent to conditioning on all X_t for $t < T$ we are effectively conditioning on ϵ_t for all $t < T$. This means that

$$\begin{aligned} E(X_{T+r} | X_{T-1}, X_{T-2}, \dots) &= E(X_{T+r} | \epsilon_{T-1}, \epsilon_{T-2}, \dots) \\ &= \sum_{s>r} d_s \epsilon_{T+r-s} \end{aligned}$$

and the forecast error is just

$$X_{T+r} - \hat{X}_{T+r} = \epsilon_t + \sum_{s=1}^r d_s \epsilon_{T+r-s}$$

so that the forecast standard error is

$$\sigma_\epsilon \sqrt{1 + \sum_{s=1}^r d_s^2}.$$

Again as $r \rightarrow \infty$ this converges to σ_X .

Finally consider forecasting the ARIMA(p, d, q) process $(I - B)^d X = W$ where W is ARMA(p, q). The forecast errors in X can clearly be written as a linear combination of forecast errors for W permitting the forecast error in X to be written as a linear combination of the underlying errors ϵ_t . As an example consider first the ARIMA(0,1,0) process $X_t = \epsilon_t + X_{t-1}$. The forecast of ϵ_{T+r} is just 0 and so the forecast of X_{T+r} is just

$$\hat{X}_{T+r} = \hat{X}_{T+r-1} = \dots = X_{T-1}.$$

The forecast error is

$$\epsilon_{T+r} + \dots + \epsilon_T$$

whose standard deviation is $\sigma\sqrt{r+1}$. Notice that the forecast standard error grows to infinity as $r \rightarrow \infty$. For a general ARIMA($p, 1, q$) we have

$$\hat{X}_{T+r} = \hat{X}_{T+r-1} + \hat{W}_{T+r}$$

and

$$X_{T+r} - \hat{X}_{T+r} = (W_{T+r} - \hat{W}_{T+r}) + \cdots + (W_T - \hat{W}_T)$$

which can be combined with the expression above for the forecast error for an ARMA(p, q) to compute standard errors.

Software

The S-Plus function *arima.forecast* can do the forecasting.

Comments

I have ignored the effects of parameter estimation throughout. In ordinary least squares when we predict the Y corresponding to a new x we get a forecast standard error of

$$\sqrt{\text{Var}(Y - x\hat{\beta})} = \sqrt{\text{Var}(\epsilon + x(\beta - \hat{\beta}))}$$

which is

$$\sigma \sqrt{1 + x(X^T X)^{-1}x^T}.$$

The procedure used here corresponds to ignoring the term $x(X^T X)^{-1}x^T$ which is the variance of the fitted value. Typically this value is rather smaller than the 1 to which it is added. In a 1 sample problem for instance it is simply $1/n$. Generally the major component of forecast error is the standard error of the noise and the effect of parameter estimation is unimportant.

In regression we sometimes compute prediction intervals

$$\hat{Y} \pm c\hat{\sigma}_{\hat{Y}}$$

The multiplier c is adjusted to make the coverage probability $P(\frac{|Y - \hat{Y}|}{c} \leq 1)$ close to a desired coverage probability such as 0.95. If the errors are normal then we can get c by taking $t_{0.025, n-p} s \sqrt{1 + x(X^T X)^{-1}x^T}$. When the errors are not normal, however, the error in $Y - \hat{Y}$ is dominated by ϵ which is not normal so that the coverage probability can be radically different from the nominal. Moreover, there is no particular theoretical justification for the use of t critical points. However, even for non-normal errors the prediction standard error is a useful summary of the accuracy of a prediction.