

10. Lagged Correlation

Lagged relationships are characteristic of many natural physical systems. Lagged correlation refers to the correlation between two time series shifted in time relative to one another. Lagged correlation is important in studying the relationship between time series for two reasons. First, one series may have a delayed response to the other series, or perhaps a delayed response to a common stimulus that affects both series. Second, the response of one series to the other series or an outside stimulus may be “smeared” in time, such that a stimulus restricted to one observation elicits a response at multiple observations. For example, because of storage in reservoirs, glaciers, etc., the volume discharge of a river in one year may depend on precipitation in the several preceding years. Or because of changes in crown density and photosynthate storage, the width of a tree-ring in one year may depend on climate of several preceding years. The simple correlation coefficient between the two series properly aligned in time is inadequate to characterize the relationship in such situations. Useful functions we will examine as alternative to the simple correlation coefficient are the cross-correlation function and the impulse response function. The cross-correlation function is the correlation between the series shifted against one another as a function of number of observations of the offset. If the individual series are autocorrelated, the estimated cross-correlation function may be distorted and misleading as a measure of the lagged relationship. We will look at two approaches to clarifying the pattern of cross-correlations. One is to individually remove the persistence from, or prewhiten, the series before cross-correlation estimation. In this approach, the two series are essentially regarded on an “equal footing”. An alternative is the “systems” approach: view the series as a dynamic linear system – one series the input and the other the output – and estimate the impulse response function. The impulse response function is the response of the output at current and future times to a hypothetical “pulse” of input restricted to the current time.

10.1 Cross-Correlation Function

The cross-correlation function (ccf) of two time series is the product-moment correlation as a function of lag, or time-offset, between the series. It is helpful to begin defining the ccf with a definition of the cross-covariance function (ccvf). Consider N pairs of observations on two time series, u_t and y_t . Following Chatfield (2004, p. 158), the sample ccvf is given by

$$c_{uy}(k) = \frac{1}{N} \sum_{t=1}^{N-k} (u_t - \bar{u})(y_{t+k} - \bar{y}) \quad [k = 0, 1, \dots, (N-1)]$$

$$c_{uy}(k) = \frac{1}{N} \sum_{t=1-k}^N (u_t - \bar{u})(y_{t+k} - \bar{y}) \quad [k = -1, -2, \dots, -(N-1)]$$
(1)

where N is the series length, \bar{u} and \bar{y} are the sample means, and k is the lag. The sample cross-correlation function (ccf) is the ccvf scaled by the variances of the two series:

$$r_{uy}(k) = \frac{c_{uy}(k)}{\sqrt{c_{uu}(0)c_{yy}(0)}}$$
(2)

Where $c_{uu}(0)$ and $c_{yy}(0)$ are the sample variances of u_t and y_t .

In a previous chapter we studied the autocovariance function (acvf) and autocorrelation function (acf), and learned that these are *symmetrical* functions (value at lag k equals value at lag $-k$). In contrast, the ccvf and ccf are *asymmetrical* functions. The asymmetry brings about the need for

the two parts of equation (1). The cross-correlation function as described by equation (1) can be described in terms of “lead” and “lag” relationships. The first part of the equation applies to y_t shifted forward relative to u_t . With this direction of shift, u_t is said to “lead” y_t . This is the same saying that y_t “lags” u_t . The second part of equation (1) describes the reverse situation, and summarizes lagged correlations when y_t “leads” u_t .

Consider the two time series of annual tree-ring index plotted in Figure 1. The plots suggest the series are positively correlated, but the detailed year-to-year variations are too noisy to directly judge whether lagged relationships exist or are important. It is obvious, however, that both series are positively autocorrelated, as positive departures from the mean tend to follow positive departures, and negative departures tend to follow negative departures.

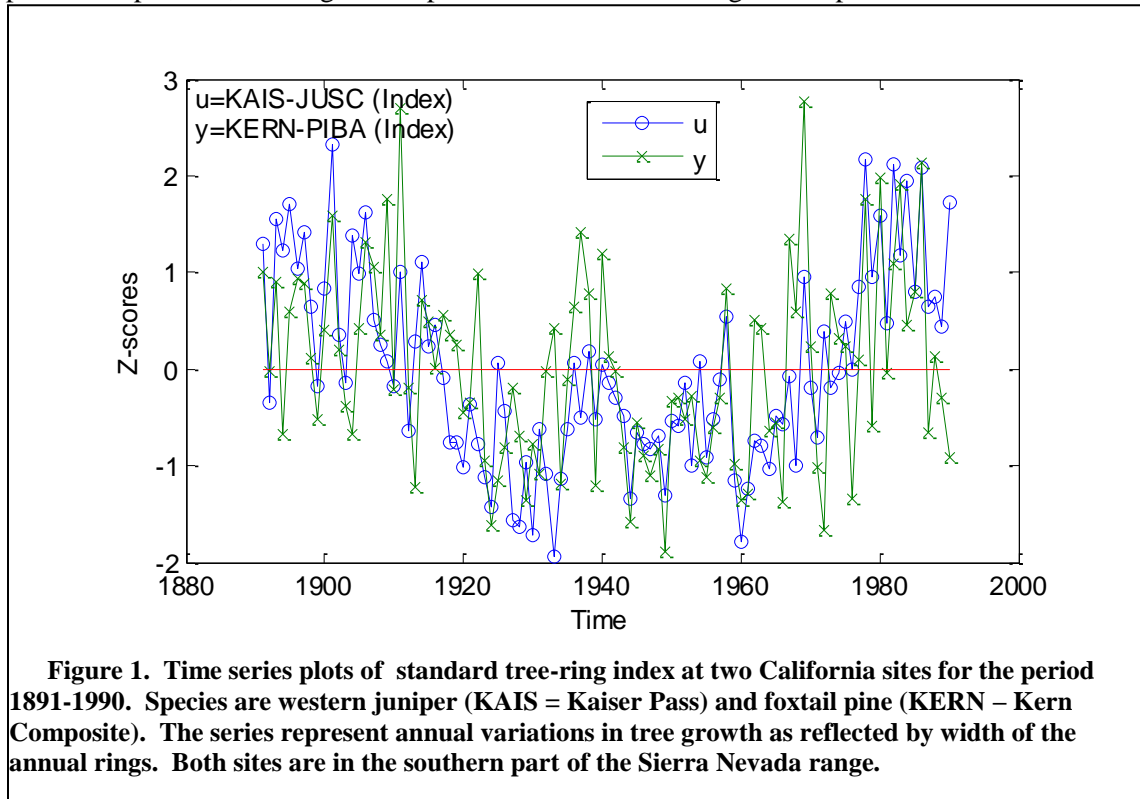
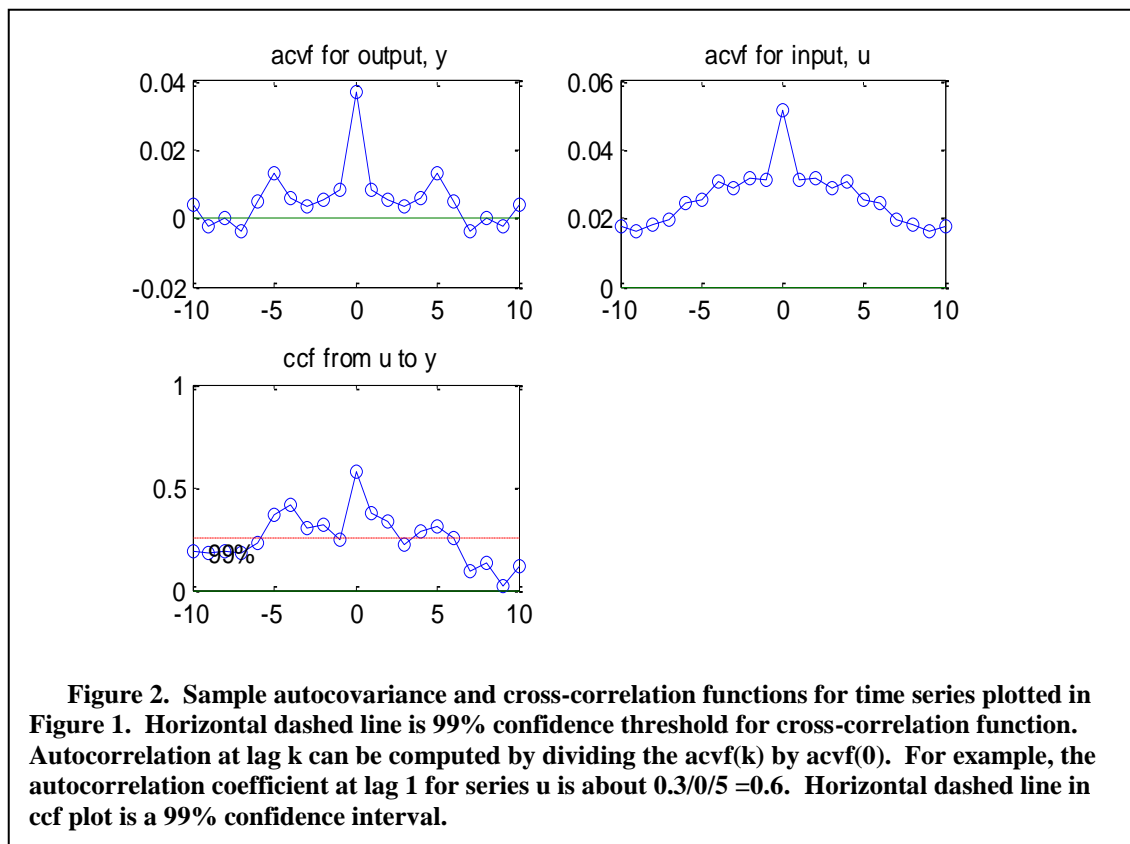


Figure 1. Time series plots of standard tree-ring index at two California sites for the period 1891-1990. Species are western juniper (KAIS = Kaiser Pass) and foxtail pine (KERN – Kern Composite). The series represent annual variations in tree growth as reflected by width of the annual rings. Both sites are in the southern part of the Sierra Nevada range.

The autocovariance functions (acvf's) of the two series and the cross-correlation function (ccf) between them are plotted in Figure 2. The acvf is the autocorrelation function (acf) scaled by multiplying by the variance, and so the acvf and acf have exactly the same shape but a different scaling. The acvf is symmetric, as evidenced by $acvf(k)$ equaling $acvf(-k)$. The plotted acvf's in Figure 2 are consistent with positive autocorrelation in both series, as positive acvf at some lag k would re-scale to positive autocorrelation, but the statistical significance cannot be judged from these plots. The autocovariance at lag $k = 0$ is identically the variance, such that we can see from the top plots in Figure 2 alone that the variance of u is slightly greater than the variance of y .

The ccf as plotted indicates that u_t and y_t are significantly positively correlated at lag $k = 0$. The horizontal dashed line in Figure 2 marks the upper 99% confidence level for significance of the ccf. This confidence interval relies on several simplifying assumptions and can be computed from the sample size alone. For a two-tailed test, the approximate 99% confidence interval ($\alpha = 0.01$) is $0 \pm 2.58/\sqrt{N}$, where N is the sample size. The value 2.58 is the 0.995 probability-point of the cdf of the normal distribution. This confidence interval relies on



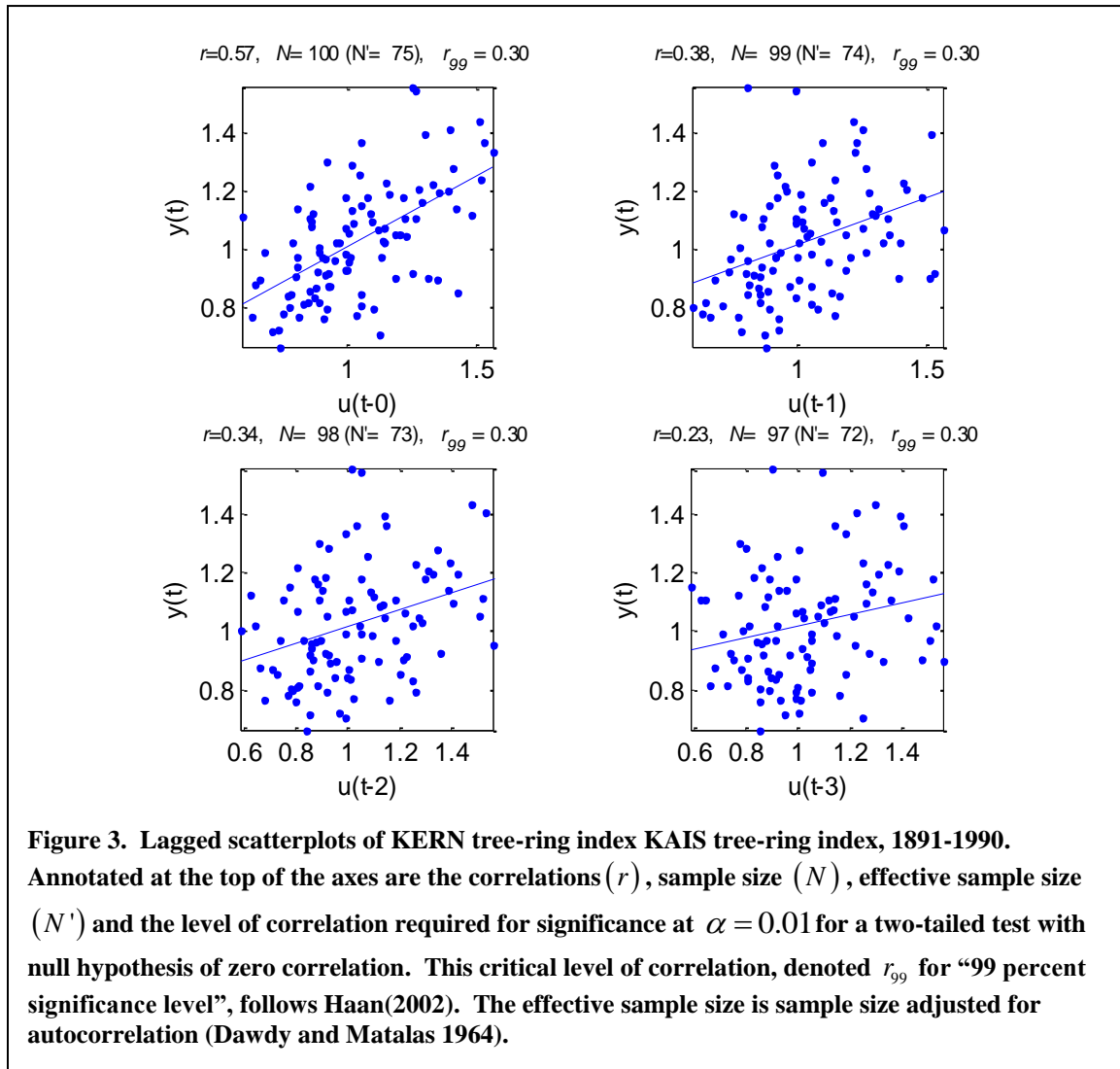
assumptions that the processes generating u_t and v_t are uncorrelated, are not autocorrelated, that the populations are normally distributed, and that the sample size is large. Under those assumptions, the sample cross-correlations are approximately $N(0, 1/N)$, or normally distributed with mean zero and variance $1/N$. A ccf estimate outside the confidence interval is “significant” in that the null hypothesis that the true ccf at that lag is zero must be rejected against the alternative hypothesis that the true ccf is non-zero.

The 99% confidence level for the ccf plotted in Figure 2 is a horizontal line, but that simple interval relies on the stated assumptions. More generally, the confidence interval would vary from lag to lag, and among other factors would depend on the autocorrelation in the individual series. That is why the ccf of autocorrelated time series is of limited use for directly inferring lagged relationships. One must question whether the “significance” of the ccf at some lags as portrayed by the simplified confidence interval is artificial and the result of ignoring the effects of autocorrelation. Modifications aimed at dealing with autocorrelation are discussed later.

The significant cross-correlations in Figure 2 extend over several lags, such that u_t is significantly correlated with y_t shifted several lags forward and backwards. If the ccf in Figure 2 were significant at only zero and positive lags, it could be said that y_t tended to “lag” u_t . Such a pattern with input u_t and output y_t would be consistent with a “causal” system. The plotted ccf in Figure 2 is significant different from zero at negative as well as non-negative lags, and consequently is not consistent with a causal system. This is not surprising for these data, which are two tree-ring series with no physically causal relationship to imply that variations in one series should in any sense “drive” variations in the other.

The ccf at a particular lag can be regarded as a correlation coefficient between two time series, one of which just happens to be the other shifted relative to itself some number of time units. Just as the bivariate relationship for any two variables can be examined with a scatterplot, the bivariate relationship posed by the ccf can be examined with lagged scatterplots. Such plots

might indicate oddities in the data – for example, that the ccf at a particular lag is driven by one or two outlier values of the time series. Lagged scatterplots for the two time series used as an example above are plotted in Figure 3. The plots show (y_t) in year t plotted against (u_t) in years $t, t-1, t-2$, and $t-3$.



Like the ccf, the scatterplots show the strongest relationship at lag 0. Correlation at lag 0 ($r=0.57$) is by definition the value of the cross-correlation function at lag 0 (compare correlation coefficient for lag 0 in Figure 3 with lag-0 cross-correlation in Figure 2). The computed correlation coefficients (r) in Figure 3 can be compared with the threshold correlation required for significance at the 99% confidence level, r_{99} . The scatterplots indicate significant correlations at lags $k = 0, -1, -2$, which is consistent with the ccf plot in Figure 2. It should be noted, that an exact match of lags with significant correlation will not generally be found in comparing the lagged scatterplots with the ccf plot because the significance thresholds in the scatterplots have been adjusted for autocorrelation in the individual series (see previous lecture). For example, the 99% confidence level is at $r_{u,y}(k) = 0.26$ in the ccf (Figure 2), and at $r_{99} = 0.30$ in the scatterplots of Figure 3

10.2 Prewhitening to clarify lagged relationships

As already mentioned, lagged correlations between time series can be misleading as evidence of lagged relationships and dependence, especially if the individual time series are autocorrelated. Jenkins and Watts (1968) suggest that to clarify lagged relationships the individual series first be prewhitened before estimating the ccf. In contrast to the systems approach described later, this approach to lagged relationships treats time series on an “equal footing” (Chatfield 2004, p. 155). The systems approach, on the other hand, treats one series as an input to a hypothetical system and the other series as an output.

Prewhitening in the context of estimation of the ccf is intended to deal with the complicating effects of autocorrelation on the estimated ccf and its standard deviations. Referring to the estimators of the ccvf and ccf ($c_{xy}(k)$ and $r_{xy}(k)$), Chatfield (2004, p. 158) says:

It can be shown that these estimators are asymptotically unbiased and consistent. However it can also be shown that estimators at neighbouring lags are themselves autocorrelated. Furthermore, it can be shown that the variances of sample cross-correlations depend on the autocorrelation functions of the two components. In general, the variances will be inflated. Thus even for moderately large values of N up to about 200, or even higher, it is possible for two series, which are actually unrelated, to give rise to apparently ‘large’ cross-correlation coefficients which are spurious, in that they arise solely from autocorrelations within the two series. Thus, of a test is required for non-zero correlation between two time series, then (at least) one of the series should first be filtered to convert it to (approximate) white noise.

Chatfield goes on to suggest the systems approach, in which one time series, regarded as “input”, is prewhitened by a time series model, the other series is filtered by the same model, and correlations between the prewhitened input and filtered output are plotted and examined. He also mentions the alternative approach of prewhitening *both* time series by fitting them individually to time series models, and then examining the cross-correlations between the prewhitened series. This latter approach is suggested by Jenkins and Watts (1968) and Brockwell and Davis (2002), and consists of a non-systems, or “equal-footing” approach.

10.3 “Equal-Footing” Approach

As in section 10.1, we consider two time series, u_t and y_t . Each series is individually fit to a high-order autoregressive model (say, AR(10)), and each is then prewhitened by its respective model to generate the two prewhitened series α_t and β_t . The ccf is computed for these prewhitened series. Because neither is autocorrelated, we can apply the conventional non-adjusted confidence interval estimation for a correlation coefficient in evaluating the significance of the ccf at various lags. The approximate variance of the cross-correlations is $1/N$, where N is the number of observations. An approximate 95% confidence interval for the ccf is therefore $0 \pm 1.96/N$, and an approximate 99% interval is $0 \pm 2.58/N$.

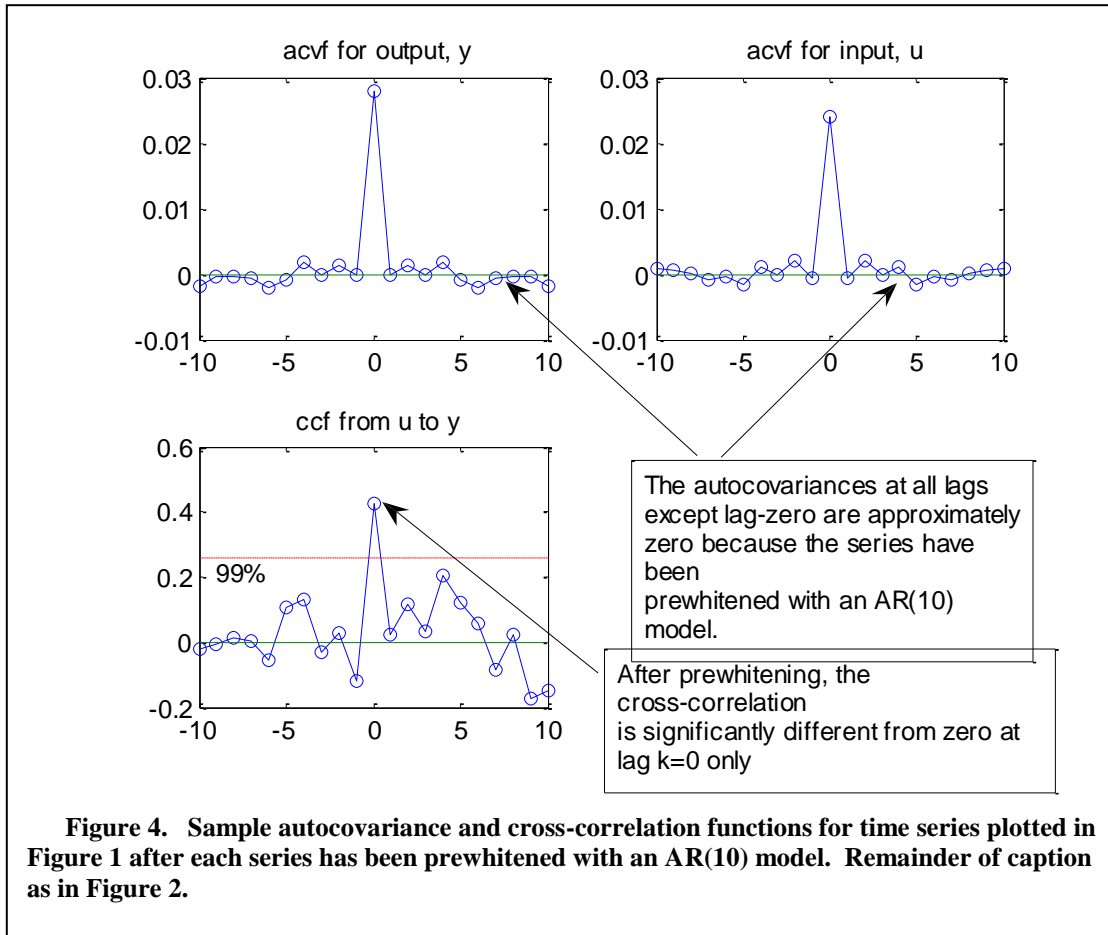
For the tree-ring example discussed previously, the ccf of the prewhitened series indicates that the series are significantly correlated at lag zero only (Figure 4). This result is quite different from the ccf on the original series (Figure 2), and suggests no lag in the relationship between the two tree-ring series. The significant ccf values at non-zero lags in Figure 2 might therefore be considered artifacts resulting from autocorrelation in the individual series

10.4 Systems Approach -- Impulse Response Function (IRF)

As discussed above, autocorrelation in the individual time series makes direct use of the ccf to study lagged relationships problematic. Essentially, the questions asked of the ccf are 1) how strongly is one series related to another, 2) is the relationship simultaneous or distributed over

several time steps, and 3) if distributed, how many lags are involved and what is their relative importance?

The prewhitening approach described in the preceding section treated each series on an equal footing in dealing with autocorrelation. These questions can also be addressed by a systems



approach in which the series are regarded as input to and output from a *linear dynamic system* (Figure 5). *Dynamic* refers to the possible dependence of the output at time t on the input signal at many previous times. For such a system, the hypothetical response to a unit pulse of input at time $t = 0$ is given by the *impulse response function* (IRF).

For the tree-ring example, the IRF gives the hypothetical response of series KERN at times $t, t+1, t+2, \dots$ to a pulse of anomalous (above or below mean) “input” from series KAIS at time t . (Of course in a tree-ring context this is an artificial conceptualized system, and is used here only to illustrate methods. A more natural input-output system might have precipitation as input and tree-growth as output). The system ideally has the following properties:

- *time invariant*: response to an input signal does not depend on absolute time
- *linear*: output response to a linear combination of inputs is the same as the linear combination of the output responses to the individual inputs
- *causal*: output at a certain time depends on the input up to that time only

The system sketched in Figure 5 is described by the equation

$$y(t) = \sum_{k=1}^{\infty} g(k)u(t-k) + v(t)$$

where

$$y(t) = \text{output in year } t \quad (3)$$

$$u(t-k) = \text{input in year } t-k$$

$$g(k) = \text{impulse response function at lag } k$$

$$v(t) = \text{noise term in year } t$$

The summation as written, following the notation in Ljung (1987), does not allow a contemporaneous response, or a response at time t to an input at time t . This can be remedied with no loss of generality by shifting the input one time step relative to the output (re-aligning the series) so that the summation effectively starts with $k=0$. Realignment can also insure that for a system with nominally assigned input and output, the output response does not precede the input stimulus.

The model represented by (3) gives the output as a linear combination of past (and possibly current) input. The numbers $\{g(k)\}$ are called the impulse response function (IRF). Generally, the IRF is unknown, and must be estimated from the data -- the input signal

$$u(t), t = 1, \dots, N \quad (4)$$

and the output signal

$$y(t), t = 1, \dots, N. \quad (5)$$

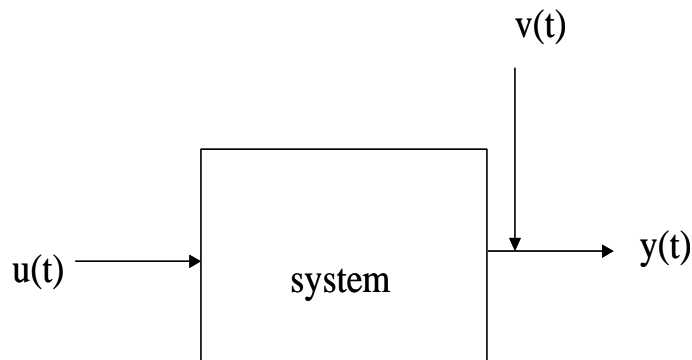


Figure 5. Input-output system with superposed noise. For example, input $u(t)$ might be rainfall, output $y(t)$ a tree-ring index, and $v(t)$ a disturbance variable incorporating all other influences on the tree-ring index. In the example used in class, one tree-ring chronology (KAIS) is used as the nominal input, and another chronology (KERN) as the nominal output (see Figure 1).

10.5 Estimation of the IRF by correlation analysis

The first “systems” method for estimating the IRF to be described is based on reducing one of the series to white noise before computing its correlation with the other series. Unlike in the “equal footing” approach described in section 10.3, only the input series is explicitly prewhitened.

The systems method to clarifying the lagged relationships in the presence of autocorrelation amounts to passing the input and output series through a filter before computing the cross-correlation function. The filter is chosen such that it reduces the input series to white noise (removes the autocorrelation). The filtered input series is therefore called the “prewhitened” input. The output is passed through the same filter, but the filtered output will generally not be white noise because the filter has been designed specifically to prewhiten the input, not the output. The ccf between the prewhitened input and filtered output is an estimate of the IRF of the system.

The filter for this operation is a high-order (e.g., order 10) autoregressive (AR) model fit to the input series. Figure 6 illustrates estimation of the IRF by this method for the tree-ring example. Before computing the ccf, we have prewhitened the input series with a 10th order autoregressive model and filtered the output series with that same model. The 10th order AR

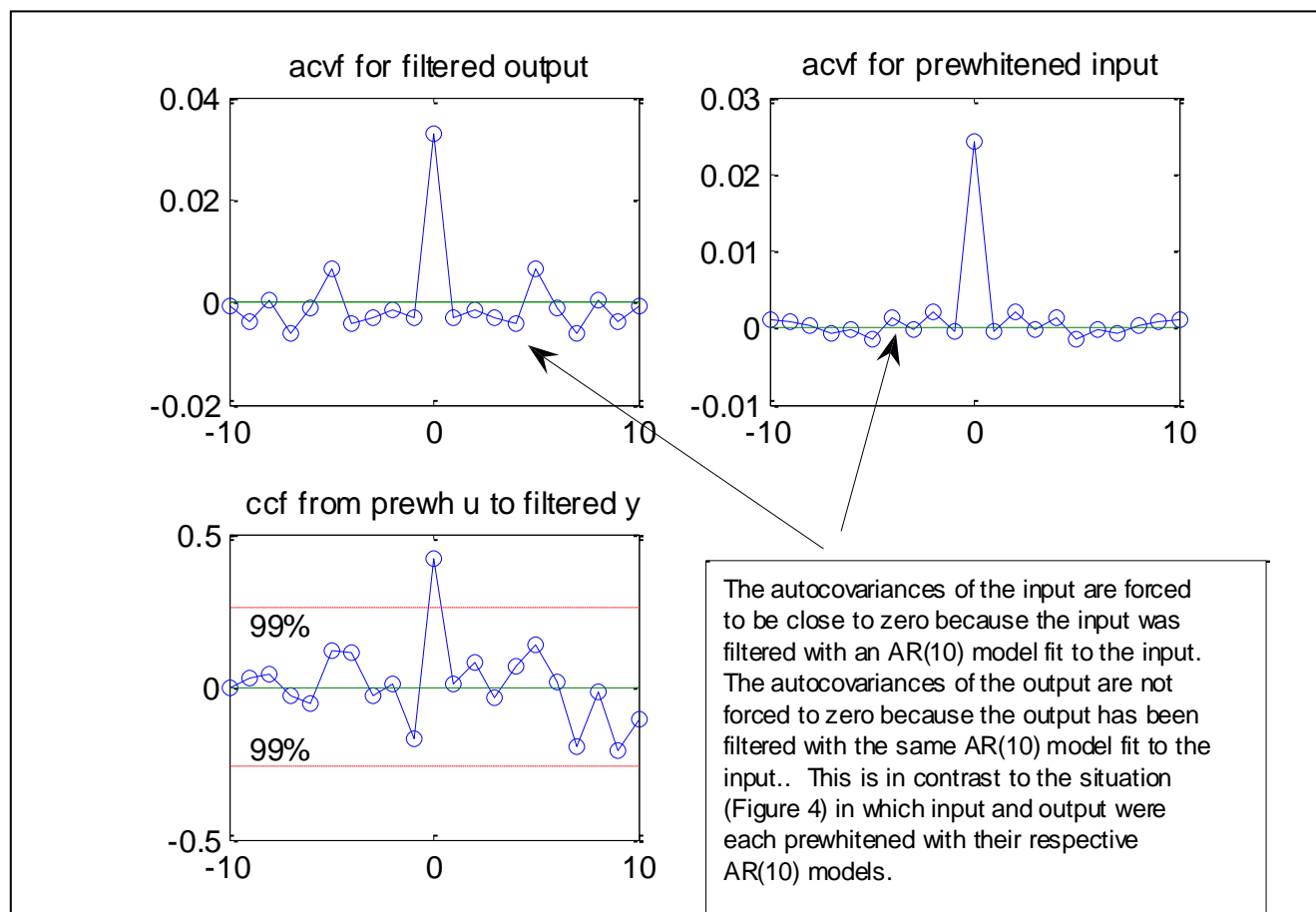


Figure 6. Impulse response function estimated by correlation analysis. Series for analysis same as in Figure 2. Top right is acvf of the the input prewhitened with an AR(10) model. Top left is the acvf of the output (tree-ring series) filtered with that same model. Lower plot is the ccf between prewhitened input and filtered output.

model has of course modified the input series so that its autocovariance is essentially zero through lag 10. The acvf of the output filtered by this same model is clearly affected, but has not been reduced to zero. In other words, the filtered output is not white noise.

The IRF, or the ccf between the prewhitened input and filtered output describes the lagged correlation structure disentangled from the influence of autocorrelation. The IRF in Figure 6 indicates significant response is restricted to lag 0. This result is in agreement with the results of the “equal footing” approach in section 10.3.

The 99% confidence interval for the IRF is once again (as in Figure 2) a horizontal line at $0 \pm 2.58/\sqrt{N}$. A constant CI is deemed applicable this estimated IRF because one of the series is now approximately white (Ljung 1987).

10.6 Estimation of the IRF by linear least-squares

The second “systems” method for estimating the IRF is essentially a minor extension to the method described above, and includes additional steps of 1) specifying maximum positive and negative lags for the possible dependence of output on input, and 2) regressing the filtered output on the prewhitened input for those lags. Recall that the model for the linear dynamic system (equation (3)) is

$$y(t) = \sum_{k=1}^{\infty} g(k)u(t-k) + v(t)$$

Assuming the earliest response occurs in same time step as the input, we can begin the summation at $k=0$ rather than $k=1$ and define a high-order finite impulse response (FIR) model

$$y(t) = g(0)u(t) + g(1)u(t-1) + g(2)u(t-2) + \dots + g(n)u(t-n) \quad (6)$$

where the g 's are the model coefficients and n is some high lag. To check for possible non-causal effects, such as feedback from output to input, negative lags may be included in the model:

$$y(t) = g(-m)u(t+m) + \dots + g(-1)u(t+1) + g(0)u(t) + g(1)u(t-1) + \dots + g(n)u(t-n) \quad (7)$$

The coefficients $g(-m), \dots, g(-1), g(0), g(1), \dots, g(n)$ are estimated by linear least squares, and if $u(t)$ is white will be estimated correctly. Equation (7) describes a multiple linear regression in which the predictand is $y(t)$, the predictors are the values of u at the various lags, and the regression coefficients are the g 's. If we estimate the regression equation by least squares, we obtain confidence intervals for the regression coefficients. The confidence intervals for the regression coefficients are the confidence intervals for the IRF. In summary the steps for the regression-based procedure for the IRF are:

1. Fit input time series u with a high-order AR model (default order 10)
2. Prewhiten u with this AR model
3. Filter the output y with the same AR model
4. Choose a maximum negative lag (m) and positive lag (n) for the regression
5. Compute the least-squares regression weights between the prewhitened input at time $t=0$ and the filtered output at times $-m \leq t \leq n$, where m is the maximum allowed lag of output preceding input, and n is the maximum lag of input preceding output. Reasonable values of m and n can be set depending on the maximum likely lag in effect. For example, n might be set to 10 if the input at time t is believed to have no likely effect on output after time

$t+10$, and m might be set to 3 to check for possible feedback effects. The estimation procedure also gives standard deviations of the estimated regression weights. These are derived by the usual regression approach, assuming a normal distribution, and are applied to compute confidence intervals for the IRF.

A plot of the estimated IRF as a function of lag summarizes the lagged response in the relationship. A single significant weight at lag 0 indicates a no-lag relationship, in which the response of the output to a pulse of input at time t is restricted to time t . A single significant weight at some lag $k > 0$ indicates a simple delay in response. Significant weights at several positive lags indicate a spreading out of the response over several time units. Significant weights at negative lags might indicate feedback in the system (e.g., tree-rings “affecting” precipitation). More often such results are likely to be spurious. The IRF pattern of weights is likely to differ from that of the ccf estimated by cross-correlation only (previous method) because regression optimizes the estimates such that the filtered output is “best” predicted in a least-squares sense from the prewhitened input, while the correlation approach does not.

The IRF computed by linear least squares for the tree-ring example is plotted in Figure 7 along with an approximate 99% confidence interval. The plotted IRF has highly significant weights at lag 0 only, again suggesting that the relationship between the two tree-ring series is contemporaneous only. It should be emphasized that for purposes of illustration, the tree-ring series have been specified as nominal input and output. The intent is not to imply a direct causal relationship between the two variables. The significant lag-0 relationship in this case is probably due to correlation of both tree-ring series with a common climatic variable (e.g., precipitation).

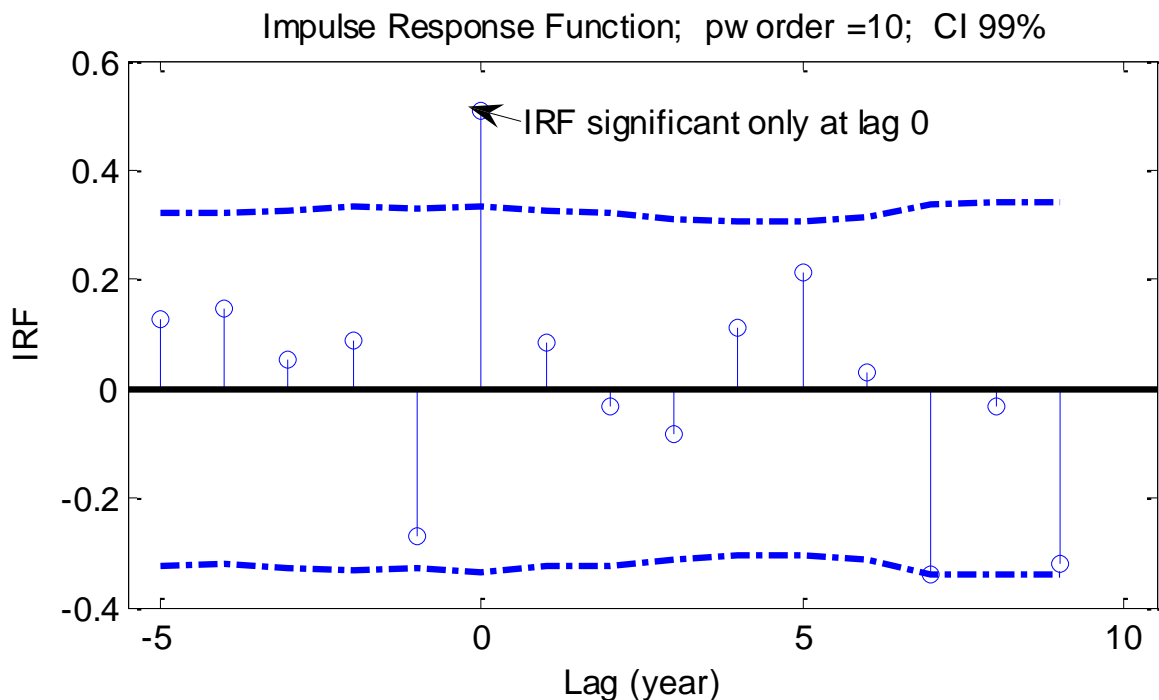


Figure 7. Estimated impulse-response function estimated by regression method for two time series plotted in Figure 1. Dot-dash line is approximate 99% confidence interval.

10.7 Summary Comments

Several alternative methods of examining lagged bivariate relationships have been presented. The simplest is the lagged scatterplot, but this method may become cumbersome when many scatterplots must be examined to cover the possibility of relationship at higher lags. Moreover, for autocorrelated series, the lagged correlations at various lags are interdependent. The cross-correlation function is a the basic tool for studying correlation as a function of lag, but the pattern of weights in the ccf of the original time series is again distorted when the individual series are autocorrelated. The autocorrelation can be dealt with by removing it prior to computation of the ccf. The “equal footing” approach involved separately prewhitening the two series. The “systems” approach involves prewhitening just one of the series and filtering the other by the same model used to prewhiten the first. The systems approach can also be modified such that least-squares regression yields estimates of the lagged response.

An example with two tree-ring chronologies from the same geographical region was used to illustrate the methods. These chronologies are highly autocorrelated, and it is apparent from the results that the autocorrelation distorts the estimated “raw” cross-correlation function sufficiently that a cross-correlation analysis on the original series is not suitable for identifying lagged relationships between the series. The “equal footing” and “systems” approaches to studying the lagged relationship agree for this example in point to a contemporaneous (no lag) relationship.

Prewhitening and filtering are potentially more important in studying relationships between time series the more autocorrelated the series. For example, in the “systems” approach, if the input series is already white noise (non-autocorrelated), an AR model fit to the input will have little effect, and the ccf of the original data will be virtually the same as the ccf and FIR (except for scaling) after filtering the data with the AR model. Likewise, if both series are white noise, the prewhitened series in the “equal footing” method will be approximately the same as the original series, and the ccf of the original data will differ negligibly from the ccf of the prewhitened data.

10.8 References

- Brockwell, P. J., and Davis, P. J., 2002, Introduction to time series and forecasting, 2nd edition. Springer, New York.
- Chatfield, C., 2004, The analysis of time series, an introduction, sixth edition. Chapman & Hall/CRC, New York. 333 pp.
- Dawdy, D.R., and Matalas, N.C., 1964, Statistical and probability analysis of hydrologic data, part III: Analysis of variance, covariance and time series, in Ven Te Chow, ed., Handbook of applied hydrology, a compendium of water-resources technology: New York, McGraw-Hill Book Company, p. 8.68-8.90. (sample size adjustment for first-order autocorrelation)
- Haan, C.T., 2002, Statistical methods in Hydrology, second edition: Ames, Iowa, Iowa State University Press.
- Ljung, L., 1987, System Identification: Theory for the User, Prentice-Hall, Inc., Englewood Cliffs, NJ, 07632, 519 pp. [*Impulse response function*]