

Updating Formulas for Correlations

May 16, 2016

Here, we collect updating formulas for correlations between time series:

1 General Time Series

The standard case is to compute the correlation between two time series $x_t(i)$, $t = 1, \dots, T$, $i = 1, \dots, N$, and $y_t(i)$, $t = 1, \dots, T$, $i = 1, \dots, N$. Additionally, it is possible that weights are given for each time step, i.e. there are non-negative number w_t , $t = 1, \dots, T$. Our goal then is to compute the (unnormalized) correlation

$$C(i, j) = \sum_{t=1}^T w_t (x_t(i) - \bar{x}(i)) (y_t(j) - \bar{y}(j)),$$

where $\bar{x}(i)$, $\bar{y}(j)$ denote the weighted mean values of the time series, i.e.

$$\begin{aligned}\bar{x}(i) &= \frac{1}{W_T} \sum_{t=1}^T w_t x_t(i), \\ W_T &= \sum_{t=1}^T w_t.\end{aligned}$$

We are interested in computing the correlation $C(i, j)$ in chunks. That means we split the data into, say, two blocks $x_t(i)$, $t = 1, \dots, T_1$, and $x_t(i)$, $t = T_1 + 1, \dots, T_2 = T$, and the same for y_t . We would then like to compute the correlation of each chunk separately, sum them up and add a correction term. Let us introduce the following notation

$$\bar{x}_{T_1}(i) = \frac{1}{w_{T_1}} \sum_{t=1}^{T_1} w_t x_t, \quad (1.1)$$

$$\bar{x}_{T_2}(i) = \frac{1}{W_{T_2}} \sum_{t=T_1+1}^{T_2} w_t x_t \quad (1.2)$$

$$W_{T_1} = \sum_{t=1}^{T_1} w_t \quad (1.3)$$

$$W_{T_2} = \sum_{t=T_1+1}^{T_2} w_t \quad (1.4)$$

$$S_{T_1}(i, j) = \sum_{t=1}^{T_1} (x_t(i) - \bar{x}_{T_1}(i)) (y_t(j) - \bar{y}_{T_1}(j)) \quad (1.5)$$

$$S_{T_2}(i, j) = \sum_{t=T_1+1}^{T_2} (x_t(i) - \bar{x}_{T_2}(i)) (y_t(j) - \bar{y}_{T_2}(j)). \quad (1.6)$$

Now, the calculations from section 5 show that the full correlation $C(i, j)$ can be computed as

$$C(i, j) = S_{T_1}(i, j) + S_{T_2}(i, j) + \frac{W_{T_1} W_{T_2}}{W_T} (\bar{x}_{T_2}(i) - \bar{x}_{T_1}(i)) (\bar{y}_{T_2}(j) - \bar{y}_{T_1}(j))$$

2 Symmetrization

In some cases, a symmetric correlation matrix is desired, for example if y_t is a time-lagged version of x_t . This can be achieved by redefining the means

$$\bar{x}(i) = \frac{1}{2T} \left[\sum_{t=1}^T x_t(i) + \sum_{t=1}^T y_t(i) \right],$$

and defining the symmetrized correlation by

$$C_s(i, j) = \sum_{t=1}^T (x_t(i) - \bar{x}(i)) (y_t(j) - \bar{x}(j)) + \sum_{t=1}^T (y_t(i) - \bar{x}(i)) (x_t(j) - \bar{x}(j)).$$

Please note that we do not allow time-dependent weights here, as this would lead to problems in the derivations of section 5. Also, it is questionable if a

symmetrization makes sense for weighted correlations. Using the analogues of Eqs. (1.1)-(1.6), where all w_t are set to one and consequently, $W_{T_1} = T_1$, $W_{T_2} = T_2$, we arrive at the updating formula

$$C_s(i, j) = S_{T_1}(i, j) + S_{T_2}(i, j) + \frac{2T_1T_2}{T_1 + T_2} (\overline{x_{T_2}}(i) - \overline{x_{T_1}}(i)) (\overline{x_{T_2}}(j) - \overline{x_{T_1}}(j))$$

see again section 5. Please note that for time-lagged data, T_1 and T_2 must be changed to $T_1 - \tau$ and $T_2 - \tau$, such that the first τ steps of every chunk only appear in x_t , while the last τ steps only appear in y_t .

3 Time-lagged data

If time-lagged data is used without symmetrization, it seems more consistent to define the means only using the $x_t(i)$ time series. Here, weights can again be allowed. Thus, we now assume to be given a time-series $\tilde{x}_t(i)$, $t = 1, \dots, T + \tau$, and define the time-lagged time-series $x_t(i) = \tilde{x}_t(i)$, $t = 1, \dots, T$ and $y_t(i) = \tilde{x}_{t+\tau}$, $t = 1, \dots, T$. The mean is defined by

$$\bar{x}(i) = \frac{1}{W_T} \sum_{t=1}^T w_t x_t(i)$$

and the asymmetric correlation then becomes

$$C_a(i, j) = \sum_{t=1}^T w_t (x_t(i) - \bar{x}(i)) (y_t(j) - \bar{x}(j)).$$

Using the analogues of Eqs. (1.1)-(1.6), we find the updating formula for time-lagged data:

$$\begin{aligned} C_a(i, j) &= S_{T_1}(i, j) + S_{T_2}(i, j) + \frac{W_{T_1}W_{T_2}}{W_T} (\overline{x_{T_2}}(i) - \overline{x_{T_1}}(i)) (\overline{x_{T_2}}(j) - \overline{x_{T_1}}(j)) \\ &\quad + \gamma_1(i) \sum_{t=1}^{T_1} w_t (y_t(j) - \overline{x_{T_1}}(j)) + \gamma_2(i) \sum_{t=T_1+1}^{T_2} w_t (y_t(j) - \overline{x_{T_2}}(j)) \\ \gamma_1(i) &= \bar{x}(i) - \overline{x_{T_1}}(i), \\ \gamma_2(i) &= \bar{x}(i) - \overline{x_{T_2}}(i). \end{aligned}$$

4 Conclusions

We have shown that mean-free correlations can be easily computed in chunks for arbitrary time series x_t , y_t , including time-dependent weights. Moreover,

symmetrized mean-free correlations can be computed for arbitrary time-series, which can also be time-lagged copies, but the addition of weights does not make sense here. Finally, we found that for time-lagged time series which are not supposed to be symmetrized, it seems to make sense to compute the means only through the x_t series, but this leads to problems when chunk-wise estimation is desired, because the updating formulas require knowledge of the final mean.

Thus, we run into problems for time-lagged correlations between when short trajectories are used and we wish to subtract the empirical mean from the data. I wonder, however, if we need to take care of this case.

5 Proofs

First, we determine an expression for the full correlation in terms of the partial sums S_{T_1} , S_{T_2} and a correction term for all cases considered here. We will see then that the correction term can be expressed in the forms given in Eqs. (1.7), (2.1) and (3.1). Let us consider the standard case:

$$C(i, j) = \sum_{t=1}^T w_t (x_t(i) - \bar{x}(i)) (y_t(j) - \bar{y}(j)) \quad (5.1)$$

$$\begin{aligned} &= \sum_{t=1}^{T_1} w_t (x_t(i) - \bar{x}(i)) (y_t(j) - \bar{y}(j)) \\ &\quad + \sum_{t=T_1+1}^{T_2} w_t (x_t(i) - \bar{x}(i)) (y_t(j) - \bar{y}(j)) \end{aligned} \quad (5.2)$$

$$\begin{aligned} &= \sum_{t=1}^{T_1} w_t ((x_t(i) - \bar{x}_{T_1}(i)) - \gamma_1^x(i)) ((y_t(j) - \bar{y}_{T_1}(j)) - \gamma_1^y(j)) \\ &\quad + \sum_{t=T_1+1}^{T_2} w_t ((x_t(i) - \bar{x}_{T_2}(i)) - \gamma_2^x(i)) ((y_t(j) - \bar{y}_{T_2}(j)) - \gamma_2^y(j)) \end{aligned} \quad (5.3)$$

where $\gamma_k^x(i) = \bar{x}(i) - \bar{x}_{T_k}(i)$ and $\gamma_k^y(i) = \bar{y}(i) - \bar{y}_{T_k}(i)$. We proceed to find

$$\begin{aligned} C(i, j) &= \sum_{t=1}^{T_1} w_t (x_t(i) - \bar{x}_{T_1}(i)) (y_t(j) - \bar{y}_{T_1}(j)) - \gamma_1^x(i) (y_t(j) - \bar{y}_{T_1}(j)) \\ &\quad - \gamma_1^y(j) (x_t(i) - \bar{x}_{T_1}(i)) + \gamma_1^x(i) \gamma_1^y(j) \end{aligned}$$

$$\begin{aligned} &+ \sum_{t=T_1+1}^{T_2} w_t (x_t(i) - \bar{x}_{T_2}(i)) (y_t(j) - \bar{y}_{T_2}(j)) - \gamma_2^x(i) (y_t(j) - \bar{y}_{T_2}(j)) \\ &\quad - \gamma_2^y(j) (x_t(i) - \bar{x}_{T_2}(i)) + \gamma_2^x(i) \gamma_2^y(j) \end{aligned} \quad (5.4)$$

$$= S_{T_1}(i, j) + S_{T_2}(i, j) + W_{T_1} \gamma_1^x(i) \gamma_1^y(j) + W_{T_2} \gamma_2^x(i) \gamma_2^y(j). \quad (5.5)$$

It remains to deal with the term:

$$\begin{aligned}
W_{T_1} \gamma_1^x(i) \gamma_1^y(j) + W_{T_2} \gamma_2^x(i) \gamma_2^y(j) &= W_{T_1} (\bar{x}(i) \bar{y}(j) - \bar{x}(i) \overline{y_{T_1}}(j) - \overline{x_{T_1}}(i) \bar{y}(j) + \overline{x_{T_1}}(i) \overline{y_{T_1}}(j)) \\
&\quad + W_{T_2} (\bar{x}(i) \bar{y}(j) - \bar{x}(i) \overline{y_{T_2}}(j) - \overline{x_{T_2}}(i) \bar{y}(j) + \overline{x_{T_2}}(i) \overline{y_{T_2}}(j)) \\
&= (W_{T_1} + W_{T_2}) \bar{x}(i) \bar{y}(j) + W_{T_1} \overline{x_{T_1}}(i) \overline{x_{T_1}}(j) \\
&\quad + W_{T_1} \overline{x_{T_1}}(i) \overline{x_{T_1}}(j) - \bar{x}(i) (W_{T_1} \overline{y_{T_1}}(j) + W_{T_2} \overline{y_{T_2}}(j)) \\
&\quad - \bar{y}(j) (W_{T_1} \overline{x_{T_1}}(i) + W_{T_2} \overline{x_{T_2}}(i)). \tag{5.7}
\end{aligned}$$

Now, we use that $W_{T_1} \overline{x_{T_1}}(i) + W_{T_2} \overline{x_{T_2}}(i) = W_T \bar{x}(i)$ to find:

$$\begin{aligned}
&= W_{T_1} \overline{x_{T_1}}(i) \overline{x_{T_1}}(j) + W_{T_1} \overline{x_{T_1}}(i) \overline{x_{T_1}}(j) \\
&\quad - \bar{x}(i) (W_{T_1} \overline{y_{T_1}}(j) + W_{T_2} \overline{y_{T_2}}(j)) \tag{5.8}
\end{aligned}$$

$$\begin{aligned}
&= \frac{1}{W_T} [W_T (W_{T_1} \overline{x_{T_1}}(i) \overline{x_{T_1}}(j) + W_{T_1} \overline{x_{T_1}}(i) \overline{x_{T_1}}(j))] \\
&\quad - \frac{1}{W_T} [W_T \bar{x}(i) (W_{T_1} \overline{x_{T_1}}(i) + W_{T_2} \overline{x_{T_2}}(i))] \tag{5.9}
\end{aligned}$$

$$= \frac{W_{T_1} W_{T_2}}{W_T} [\overline{x_{T_1}}(i) \overline{y_{T_1}}(j) + \overline{x_{T_2}}(i) \overline{y_{T_2}}(j) - \overline{x_{T_1}}(i) \overline{y_{T_2}}(j) - \overline{x_{T_2}}(i) \overline{y_{T_1}}(j)] \tag{5.10}$$

This completes the proof of Eq. (1.7). For the symmetric case, the procedure from Eqs. (5.1)-(5.5) can be repeated to come up with the expression

$$\begin{aligned}
C_s(i, j) &= S_{T_1}(i, j) + S_{T_2}(i, j) + T_1 (\gamma_1(i) \gamma_1(j) + \gamma_1(j) \gamma_1(i)) \\
&\quad + T_2 (\gamma_2(i) \gamma_2(j) + \gamma_2(j) \gamma_2(i)),
\end{aligned}$$

where $\gamma_k(i) = \bar{x}(i) - \overline{x_{T_k}}(i)$. Please note that for time-dependent weights, the cross-terms in Eq. (5.4) would not cancel. Then, the steps of Eqs. (5.6)-(5.9) can be repeated in the same way. For the asymmetric case, Eqs. (5.1)-(5.5) yield the expression

$$\begin{aligned}
C_a(i, j) &= S_{T_1}(i, j) + S_{T_2}(i, j) + W_{T_1} \gamma_1(i) \gamma_1(j) + W_{T_2} \gamma_2(i) \gamma_2(j) \\
&\quad + \gamma_1(i) \sum_{t=1}^{T_1} w_t (y_t(j) - \overline{x_{T_1}}(j)) + \gamma_2(i) \sum_{t=T_1+1}^{T_2} w_t (y_t(j) - \overline{x_{T_2}}(j)).
\end{aligned}$$

Here, we have used $\gamma_k(i) = \bar{x}(i) - \overline{x_{T_k}}(i)$. The expression $W_{T_1} \gamma_1(i) \gamma_1(j) + W_{T_2} \gamma_2(i) \gamma_2(j)$ can be reformulated through Eqs. (5.6)-(5.10) to end up with Eq. (3.1).