

## Appendix B: Parameter Distributions

Distribution of the  $\mathbf{S}_{slp}^2$  estimate, Eq. 2

Consider for the moment, an idealized  $\hat{\mathbf{S}}_{slp}^2$  estimate using subsampled data to eliminate overlap in the  $N_{t+\mathbf{t}}/N_t$  ratios and  $L=1$ . Let's call it  $\hat{\mathbf{S}}_{isl p}^2$ . We can derive the distribution of

$\hat{\mathbf{S}}_{isl p}^2 / \mathbf{S}_p^2$  by observing that the slope of  $\text{var}(\ln(N_{t+\mathbf{t}}) - \ln(N_t))$  versus  $\mathbf{t} (\mathbf{t} = 1, 2, \dots, \mathbf{t}')$  is basically

$$\hat{\mathbf{S}}_{isl p}^2 = (\hat{\mathbf{S}}_{\mathbf{t}'}^2 - \hat{\mathbf{S}}_1^2) / (\mathbf{t}' - 1)$$

$$\text{where } \hat{\mathbf{S}}_i^2 = \text{var}(\ln(N_{t+i}) - \ln(N_t)) \text{ for } t = 1, i+2, 2i+3, \dots, ni+n+1$$

since the  $\text{var}(\ln(N_{t+\mathbf{t}}) - \ln(N_t))$  versus  $\mathbf{t}$  line is generally a straight. Using

$$\frac{\hat{\mathbf{S}}_{isl p}^2}{\mathbf{S}_p^2} \sim \frac{\mathbf{S}_t^2}{df_t \mathbf{S}_p^2} \mathbf{C}_{df_t}^2 \sim \mathbf{g}\left(\frac{df_t}{2}, \frac{\mathbf{S}_t^2}{\mathbf{S}_p^2} \frac{2}{df_{t'}}\right) \text{ where } \mathbf{g}(\mathbf{a}, \mathbf{b}) \text{ is a gamma distribution with shape } \mathbf{a} \text{ and}$$

scale  $\mathbf{b}$ , the distribution of  $\hat{\mathbf{S}}_{isl p}^2 / \mathbf{S}_p^2$  is straight-forward to derive.

$$\frac{\hat{\mathbf{S}}_{isl p}^2}{\mathbf{S}_p^2} = \frac{(\hat{\mathbf{S}}_{\mathbf{t}'}^2 - \hat{\mathbf{S}}_1^2) / (\mathbf{t}' - 1)}{\mathbf{S}_p^2} \sim \mathbf{g}\left(\frac{df_{t'}}{2}, \frac{2\mathbf{S}_{np}^2 + \mathbf{t}'\mathbf{S}_p^2}{(\mathbf{t}' - 1)\mathbf{S}_p^2} \frac{2}{df_{t'}}\right) - \mathbf{g}\left(\frac{df_1}{2}, \frac{2\mathbf{S}_{np}^2 + \mathbf{S}_p^2}{(\mathbf{t}' - 1)\mathbf{S}_p^2} \frac{2}{df_1}\right)$$

$$\text{where } \hat{\mathbf{S}}_i^2 = \text{var}(\ln(N_{t+i} / N_t)) \text{ for } t = 1, i+2, 2i+3, \dots, ni+n+1$$

$$\text{and } E(\hat{\mathbf{S}}_i^2) = 2\mathbf{S}_{np}^2 + i\mathbf{S}_p^2$$

Eq. A1

$$df_{t'} = \text{number of } \ln(N_{t+\mathbf{t}'} / N_t) \text{ ratios minus } 1$$

$$df_1 = \text{number of } \ln(N_{t+1} / N_t) \text{ ratios minus } 1$$

Note that the sequential  $N_{t+\mathbf{t}'} / N_t$  ratios are chosen so that there is no overlap thus each ratio is

independent. Assume for the moment, that the two gamma distributions are independent -- which

they are not. In this case, we can show that the limiting distribution of Eq. A1 as  $df_{t'}$  and  $df_1$

become large is  $\chi^2$  with  $\mathbf{y} df_{t'}$  degrees of freedom:

The moment generating function of  $\mathbf{g}(\mathbf{a}_1, \mathbf{b}_1) - \mathbf{g}(\mathbf{a}_2, \mathbf{b}_2)$  is  $(1 - \mathbf{b}_1 t)^{-a_1} (1 + \mathbf{b}_2 t)^{-a_2}$ . Thus

the moment generating function for the distribution in Eq. A1 is

$$\text{mgf} = \left(1 - \frac{4\mathbf{V} + 2\mathbf{t}'}{(\mathbf{t}' - 1)df_{t'}}t\right)^{-df_{t'}/2} \left(1 + \frac{4\mathbf{V} + 2}{(\mathbf{t}' - 1)df_1}t\right)^{-df_1/2} \quad \text{where } \mathbf{V} = \mathbf{s}_{np}^2 / \mathbf{s}_p^2.$$

Take the natural log of this to get,

$$\ln(\text{mgf}) = -\frac{df_{t'}}{2} \ln\left(1 - \frac{4\mathbf{V} + 2\mathbf{t}'}{(\mathbf{t}' - 1)df_{t'}}t\right) - \frac{df_1}{2} \ln\left(1 + \frac{4\mathbf{V} + 2}{(\mathbf{t}' - 1)df_1}t\right).$$

Using the Taylor expansion for  $\ln(1+x)$  and multiplying the second element by  $df_{t'}/df_{t'}$ ,

$$\begin{aligned} \ln(\text{mgf}) &\approx -\frac{df_{t'}}{2} \left( -\frac{4\mathbf{V} + 2\mathbf{t}'}{(\mathbf{t}' - 1)df_{t'}}t - \frac{(4\mathbf{V} + 2\mathbf{t}')^2}{2(\mathbf{t}' - 1)^2 df_{t'}^2}t^2 + O(1/df_{t'}^{3+}) \right) \\ &\quad - \frac{df_{t'}}{2} \left( \frac{4\mathbf{V} + 2}{(\mathbf{t}' - 1)df_{t'}}t - \frac{(4\mathbf{V} + 2)^2 (df_{t'}/df_1)}{2(\mathbf{t}' - 1)^2 df_{t'}^2}t^2 + O(1/df_{t'}^{3+}) \right) \\ &= -\frac{df_{t'}}{2} \left( -\frac{2}{df_{t'}}t - \frac{(4\mathbf{V} + 2\mathbf{t}')^2 + (4\mathbf{V} + 2)^2 (df_{t'}/df_1)}{2(\mathbf{t}' - 1)^2 df_{t'}^2}t^2 + O(1/df_{t'}^{3+}) \right) \end{aligned}$$

Ignoring higher order terms, the  $\ln(\text{mgf})$  has the form:

$$-\frac{\mathbf{y}df_{t'}}{2} \left( -\frac{2}{\mathbf{y}df_{t'}}t - \frac{1}{2} \frac{4}{\mathbf{y}^2 df_{t'}^2}t^2 \right) \approx -\frac{\mathbf{y}df_{t'}}{2} \ln\left(1 - \frac{2}{\mathbf{y}df_{t'}}t\right) \text{ for } df_{t'} \text{ big.}$$

which is the  $\ln(\text{mgf})$  of the following  $\chi^2$  distribution:

$$\frac{1}{\mathbf{y}df_{t'}} \mathbf{c}_{ydf_{t'}}^2 \quad \text{where } \mathbf{y} = \frac{4(\mathbf{t}' - 1)^2}{(4\mathbf{V} + 2\mathbf{t}')^2 + (4\mathbf{V} + 2)^2 (df_{t'}/df_1)}.$$

As noted, the gamma distributions for the variances of  $\ln(N_{t+t'}/N_t)$  and  $\text{var}(\ln(N_{t+\mathbf{1}}/N_t))$  are actually correlated. The effect of the correlation, as seen from numerical experiments, is to cause the distribution in Eq. A1 to approach the limiting distribution faster (i.e. when the  $df$ s in the gamma distributions are smaller).

The  $\hat{\mathbf{S}}_{slp}^2$  used in the Dennis-Holmes method is somewhat different than the idealized  $\hat{\mathbf{S}}_{slp}^2$  used in this derivation. First, the  $N_{t+t'}/N_t$  ratios cannot generally be subsampled due to short

time series. This means the ratios are correlated and  $df_{t'}$  is substantially less than the number of ratios minus one; additionally the lack of subsampling makes  $\hat{\mathbf{S}}_{slp}^2$  biased. The data are running sum transformed ( $L>1$ ); this leads to further bias. These are trade-offs that improve estimation for short corrupted time series by reducing the number of negative variance estimates (percent errors column in Table B1). Despite the differences, understanding the limiting distribution for the idealized  $\hat{\mathbf{S}}_{islp}^2$  helps us understand why when we estimated a non-idealized  $\hat{\mathbf{S}}_{slp}^2$  ( $L > 1$  and data not subsampled) from simulated data, we observed that  $\hat{\mathbf{S}}_{slp}^2 / \text{mean}(\hat{\mathbf{S}}_{slp}^2)$  showed a distribution of the form  $(1/df) \mathbf{C}_{df}^2$  for a wide range of time series lengths, non-process to process error ratios, and filter lengths (Table B1).

Monte Carlo estimation was used to numerically estimate the  $\mathbf{C}^2$  distributions for  $\hat{\mathbf{S}}_{slp}^2$  estimates used in the Dennis-Holmes method (= the slope of  $\ln(R_{t+1}/R_t)$  versus  $t$  for  $t = 1, 2, 3, 4$ ). Monte Carlo estimation uses parameter estimates from samples of data generated with simulations to calculate the distribution of the parameter estimate (this is akin to parametric bootstrapping). We generated 5000 time series of length  $n$  using the model,  $N_{t+1} = N_t \exp(\mathbf{m} + \mathbf{e}_p)$ ,  $O_t = N_t \exp(\mathbf{e}_{np})$  where the process error,  $\mathbf{e}_p \sim \text{Normal}(0, \mathbf{S}_p)$ , and the non-process error,  $\mathbf{e}_{np} \sim \text{Normal}(0, \mathbf{S}_{np})$ . Let  $\text{mean}(\hat{\mathbf{S}}_{slp}^2)$  denote the mean of all 5000  $\hat{\mathbf{S}}_{slp}^2$  estimates. For each simulation, we calculated the statistic  $\hat{\mathbf{S}}_{slp}^2 / \text{mean}(\hat{\mathbf{S}}_{slp}^2) = \eta$ . We then found the best fitting  $df_{slp}$  parameter such that

$$\mathbf{h} \sim \frac{1}{df_{slp}} \mathbf{C}_{df_{slp}}^2$$

We did this by finding the  $df_{slp}$  that maximized the p-value from a Kolmogorov-Smirnov goodness of fit test. We repeated the fitting process for different time series lengths ( $n$ ), filter lengths ( $L$ ), ratios of process to non-process error ratios ( $\mathbf{S}_p/\mathbf{S}_{np}$ ) and  $\mathbf{m}$  and  $\mathbf{S}_p$ . The best fitting  $df_{slp}$  values for

different  $n$ ,  $L$  and  $(\mathbf{S}_p/\mathbf{S}_{np})$  are given in Table B1 with the p-values for the fitted distribution. The observed bias and  $\gamma$  parameters from the simulations are given in Table B2. The degrees of freedom depended mainly on the length of the time series,  $n$ , and the length of the filter,  $L$ . There was an approximately linear relationship between  $n$ ,  $L$  and the  $df_{slp}$  values in Table B1. The following formula gives a close approximation of the numerically calculated  $df_{slp}$ :

$$df_{slp} = 0.333 + 0.212 n - 0.387 L \quad \text{for } n > 15.$$

*Variance of  $\hat{\mathbf{m}}_k$ , Eq. 3*

Given  $n$  observations,  $O_1, O_2, O_3 \dots O_n$ , of the true population size,  $N_1, N_2, N_3 \dots N_n$ , the  $O_t$  series is transformed into a running sum,  $R_1, R_2, R_3 \dots R_r$  where  $r = n-L+1$  and  $R_t = \sum_{i=t}^{t+L-1} O_i$ .

$$\begin{aligned} \hat{\mathbf{m}}_k &= \text{sample mean of } \ln(R_{t+1}/R_t) = \frac{1}{n-L} (\ln(R_r) - \ln(R_1)) \\ \text{var}(\hat{\mathbf{m}}_k) &= \text{var}\left(\frac{1}{n-L} (\ln(R_r) - \ln(R_1))\right) \\ &= \frac{1}{(n-L)^2} \text{var}(\ln(R_r) - \ln(R_1)) = \frac{1}{(n-L)^2} \text{var}\left(\ln\left(\sum_{i=r}^n O_i\right) - \ln\left(\sum_{i=1}^{L-1} O_i\right)\right) \end{aligned}$$

Denote  $\bar{N}_t$  as the mean of the  $N$ 's that comprise the running sum,  $R_t$ :  $\bar{N}_t = \frac{1}{L} \sum_{i=t}^{t+L-1} N_i$ , and recall

$$O_t = \mathbf{e}_{np,t} N_t.$$

$$\text{Thus } \text{var}(\hat{\mathbf{m}}_k) = \frac{1}{(n-L)^2} \text{var}\left(\ln\left(\sum_{i=r}^n \mathbf{e}_{np,i} \mathbf{d}_{i,r} \bar{N}_r\right) - \ln\left(\sum_{i=1}^{L-1} \mathbf{e}_{np,i} \mathbf{d}_{i,1} \bar{N}_1\right)\right) \text{ where } \mathbf{d}_{i,t} = \frac{N_i}{\bar{N}_t}$$

$$\begin{aligned}
&\approx \frac{1}{(n-L)^2} \text{var} \left( \ln \left( \sum_{i=r}^n \mathbf{e}_{np,i} \right) - \ln \left( \sum_{i=1}^{L-1} \mathbf{e}_{np,i} \right) + \ln(\bar{N}_r) - \ln(\bar{N}_1) \right) \\
&\quad \text{for } L \text{ small so that } \mathbf{d}_{i,r} = \frac{N_i}{N_r} \approx 1 \\
&\text{and } \text{var}(\hat{\mathbf{m}}_k) \approx \frac{1}{(n-L)^2} \left( \frac{2\mathbf{s}_{np}^2}{L} + (n-L)\mathbf{s}_p^2 \right) \\
&\quad \text{for } \mathbf{s}_{np}^2 < 1 \text{ so that } \text{var} \left( \ln \left( \sum_{i=r}^k \hat{a}_{np,i} \right) \right) \approx \frac{\mathbf{s}_{np}^2}{L} \text{ and } L \text{ small so that } \ln \left( \frac{\bar{N}_r}{\bar{N}_1} \right) \approx (n-L)\mathbf{s}_p^2
\end{aligned}$$

Note that  $\hat{\mathbf{m}}_k$  is the mean of the  $\ln(R_{t+1}/R_t)$  ratios from the time series; however, for corrupted time series, the variance of the  $\hat{\mathbf{m}}_k$  is not  $1/(n-L)$  times the variance of the  $\ln(R_{t+1}/R_t)$  ratios, as it would be the case for uncorrupted time series:

$$\frac{1}{n-L} \text{var}(\ln(R_{t+1}/R_t)) = \frac{1}{n-L} \left( \frac{2}{L} \mathbf{s}_{np}^2 + \mathbf{s}_p^2 \right) \neq \text{var}(\hat{\mathbf{m}}_k) = \frac{1}{(n-L)^2} \left( \frac{2}{L} \mathbf{s}_{np}^2 + (n-L)\mathbf{s}_p^2 \right)$$

Using the variance of the  $\ln(R_{t+1}/R_t)$  ratios would lead to high overestimation of the variance of  $\hat{\mathbf{m}}_k$ . This overestimation is greater for smaller  $L$ .

*Estimate of the distribution  $\hat{\mathbf{m}}_k$  from data, Eq. 4*

If  $\mathbf{s}_p^2$  and  $\mathbf{s}_{np}^2$  were known, it would be straight-forward to specify the distribution of  $\hat{\mathbf{m}}_k$ , (i.e,  $\text{normal}(\hat{\mathbf{m}}_k, \mathbf{s}_{m,R}^2)$ ) however, instead we have to use estimates of  $\mathbf{s}_p^2$  and  $\mathbf{s}_{np}^2$  which themselves have some distribution. Below is outlined an estimate of the distribution of  $\hat{\mathbf{m}}_k$  which uses only  $\hat{\mathbf{s}}_{slp}^2$ . Deriving a distribution based on both  $\hat{\mathbf{s}}_{slp}^2$  and  $\hat{\mathbf{s}}_{np}^2$  appears problematic given the nature of the distribution of  $\hat{\mathbf{s}}_{np}^2$  (see below) and given that the  $\hat{\mathbf{s}}_{np}^2$  estimate is not independent of  $\hat{\mathbf{s}}_{slp}^2$ .

By simple algebra, we can rewrite  $\frac{\hat{\mathbf{m}}_k - \mathbf{m}}{\sqrt{\hat{\mathbf{S}}_{slp}^2 / (n-L)}}$  as

$$\begin{aligned} \frac{\hat{\mathbf{m}}_k - \mathbf{m}}{\sqrt{\hat{\mathbf{S}}_{slp}^2 / (n-L)}} &= \frac{\hat{\mathbf{m}}_k - \mathbf{m}}{\sqrt{(n-L)\mathbf{S}_{mR}^2 / (n-L)}} \frac{1}{\sqrt{\mathbf{S}_{slp}^2 / (n-L)\mathbf{S}_{mR}^2}} \frac{1}{\sqrt{\hat{\mathbf{S}}_{slp}^2 / \mathbf{S}_{slp}^2}} \\ &= \frac{1}{\sqrt{\hat{\mathbf{S}}_{slp}^2 / (n-L)\mathbf{S}_{mR}^2}} \frac{\hat{\mathbf{m}}_k - \mathbf{m}}{\mathbf{S}_{mR}} \frac{1}{\sqrt{\hat{\mathbf{S}}_{slp}^2 / \mathbf{S}_{slp}^2}} \end{aligned}$$

$$\text{Thus, } \frac{\hat{\mathbf{m}}_k - \mathbf{m}}{\sqrt{\hat{\mathbf{S}}_{slp}^2 / (n-L)}} \sim \frac{1}{\sqrt{\mathbf{g}}} \text{normal}(0,1) \frac{1}{\sqrt{\mathbf{C}_{dfslp}^2 / df_{slp}}} \sim \frac{1}{\sqrt{\mathbf{g}}} t_{dfslp}$$

$$\text{where } \mathbf{g} = \frac{\mathbf{S}_{slp}^2}{(n-L)\mathbf{S}_{mR}^2} = \frac{\mathbf{S}_{slp}^2}{\left( \frac{2}{L(n-L)} \mathbf{S}_{np}^2 + \mathbf{S}_p^2 \right)} \text{ for } \mathbf{S}_{np}^2 \text{ small}$$

*Point estimate of  $\mathbf{S}_{np}^2$*

A point estimate of  $\mathbf{S}_{np}^2$  can be calculated by noting that  $\text{var}(\ln(O_{t+1} / O_t)) = 2\mathbf{S}_{np}^2 + \mathbf{S}_p^2$ ,

thus  $\hat{\mathbf{S}}_{np}^2 = \frac{1}{2}(\hat{\mathbf{S}}_O^2 - \hat{\mathbf{S}}_p^2) \approx \frac{1}{2}(\hat{\mathbf{S}}_O^2 - \hat{\mathbf{S}}_{slp}^2)$  and  $\hat{\mathbf{S}}_{np}^2 \sim \frac{1}{2}\mathbf{S}_O^2 \mathbf{C}_{n-1}^2 - \frac{1}{2}\mathbf{S}_{slp}^2 \mathbf{C}_{dfslp}^2$  where

$\hat{\mathbf{S}}_O^2 = \text{var}(\ln(O_{t+1} / O_t))$  from the data.