

Review of degrees of freedom / time indep.

Sample size impacts the variance of sampling distribution means. Persistence impacts the variance of the sampling distribution since it leads to an overestimate of the sample size

Effective sample size (t-stat)

$$N^* = N \left(\frac{1 - \rho_1}{1 + \rho_1} \right) \quad \begin{aligned} \rho_1 &= \text{lag-1 autocorrelation} \\ N &= \text{orig. # samples} \\ N^* &= \text{effective samples.} \end{aligned}$$

Assumes a first order Markov process

For white noise $N^* = N$, then as τ increases N^* decreases.

Leith (1973) proposed:

$$N^* = N \left(\frac{\Delta t}{2T} \right) \quad \begin{aligned} T &= \text{e-folding time} \\ &\text{scale of autocorrelation function.} \end{aligned}$$

Factor of 2 included because any given point in a red noise time series can be predicted by points before and after that point.

As before, as the the redness of time series goes up, N^* decreases.

For variance and covariance, Bretherton et al. (1999) suggest

$$N^* = N \left(\frac{1 - p_1^2}{1 + p_1^2} \right)$$

Recall we encountered something along these lines when estimating N^* for EOF significance testing.

Whatever you may use, should have an idea of how it affects your dof.

Autocorrelation | 0.15 0.3 0.5 0.7 0.9

Leith $\frac{N^*}{N}$ | 1 0.6 0.35 0.18 0.05

Bretherton $\frac{N^*}{N}$ | 1 0.83 0.6 0.34 0.1

The latter one can yield 2x the dof!

Harmonic analysis

Basic idea: Interpret a time or space series as a summation of contributions from harmonic functions - each with a unique temporal or spatial scale.

Recall the normal equations for the least squares best fit of ' y ' to predictors ' x '

$$y = a_0 + a_1 x_1 + a_2 x_2 \dots a_N x_N$$

$$\text{where } \bar{y} = a_0 + a_1 \bar{x}_1 + a_2 \bar{x}_2 \dots a_N \bar{x}_N$$

and:

$$\bar{x}_1' y' = a_1 \bar{x}_1' + a_2 \bar{x}_2' + a_3 \bar{x}_3' \dots$$

$$\bar{x}_2' y' = a_1 \bar{x}_1' + a_2 \bar{x}_2' + a_3 \bar{x}_3' \dots$$

In the case of harmonic analysis, the functions to be fit (i.e. the x 's), or basis functions, are of the form:

$$\cos\left(2\pi k \frac{t}{T}\right) \quad \text{and} \quad \sin\left(2\pi k \frac{t}{T}\right)$$

$$k=1, 2, 3, \dots$$

Hence y can be written as:

$$y = a_0 + \sum_{k=1}^{N/2} A_k \cos\left(2\pi k \frac{t}{T}\right) + \sum_{k=1}^{N/2} B_k \sin\left(2\pi k \frac{t}{T}\right)$$

$k=1, 2, 3 \dots$

T = total time length of the record

N = # of grid points or time steps

A_k, B_k = regression coefficients for each predictor.

Notes

- Each predictor is a harmonic function with frequency k/T , and hence fits into the interval $[0, T]$ 'k' times

- If $k=1$, the frequency is $1/T$ and hence predictor completes 2π radians in $[0, T]$

* - If $k=N/2$, the frequency is the highest that can be resolved, the predictor completes π radians in one timestep
→ Nyquist frequency

- for a given K , fitting K cosine and K sine waves with variable amplitudes $A \in B$ into $[0, T]$: Equivalent to fitting K cosine waves (w/ no sine waves) with variable amp & phase

The amplitude of each predictor is found by solving the normal equations (like in regression)

Discrete Fourier transform

In the special case of evenly spaced data points (e.g. a constant time interval)

Consider domain $[0, T]$ where 0 and T coincide with data points $i=1$ and $i=N+1$ where N is an even number.

1) The functions assume the form

$$\cos\left(2\pi k \frac{i\Delta t}{T}\right) \text{ and } \sin\left(2\pi k \frac{i\Delta t}{T}\right)$$

where Δt is the spacing between grid points.

2) the average of each sine / cosine function on the interval $[0, T]$ is zero.

$$\text{Hence } a_0 = \bar{y}$$

3) the harmonic functions are mutually orthogonal on $[0, T]$ because the correlation between sine and cosine functions is zero. Therefore, the off-diagonal elements in covariance matrix = 0.

Hence, the normal equations reduce to a set of equations of the form:

$$a_k = \frac{\overline{x'_k y'}}{\overline{x'^2}}$$

- 4) the functions (the k 's) each have variance $\overline{x'^2} = \frac{1}{2}$ except for $A_{N/2}$ and $B_{N/2}$, of which the former takes on a value of \pm unity at each data point, and the latter a value of zero at each grid point.

Hence, the variance of the predictor with amplitude $A_{N/2}, B_{N/2}$ is one and zero, respectively.

- 5) From #3, the solutions for each predictor's amplitude:

$$a_k = \frac{\overline{x'_k y'}}{\overline{x'^2}} \text{ reduce to .}$$

$$A_k = \frac{2}{N} \sum_{i=1}^N y_i \cos\left(2\pi k \frac{i\Delta t}{T}\right)$$

$$B_k = \frac{2}{N} \sum_{i=1}^N y_i \sin\left(2\pi k \frac{i\Delta t}{T}\right)$$

$$k = 1, N/2 - 1$$

$$A_{N/2} = \frac{1}{N} \sum_{i=1}^N y_i \cos\left(\pi N \frac{\Delta t}{T}\right)$$

$$B_{N/2} = 0.$$

The function y can now be written as a sum of sines and cosines:

$$y(t) = \bar{y} + \sum_{k=1}^{N/2-1} A_k \cos\left(2\pi k \frac{t}{T}\right)$$

$$+ \sum_{k=1}^{N/2-1} B_k \sin\left(2\pi k \frac{t}{T}\right)$$

$$+ A_{N/2} \cos\left(\pi N \frac{t}{T}\right)$$

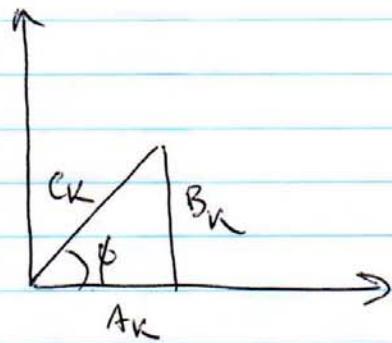
Can also write this ~~in~~ in amplitude (phase form):

$$y(t) = \bar{y} + \sum_{k=1}^{N/2-1} c_k \cos\left(2\pi k \frac{(t-t_k)}{T}\right) + A_{N/2} \cos\left(\pi N \frac{t}{T}\right)$$

$$\text{where } c_k^2 = A_k^2 + B_k^2$$

$$t_k = \frac{T}{2\pi k} \tan^{-1}\left(\frac{B_k}{A_k}\right)$$

In geometric terms:



$$\phi = \tan^{-1} \left(\frac{B_K}{A_K} \right)$$

Fraction of variance explained for each k :

$$r^2(y, x_k) = \frac{(\overline{x_k y'})^2}{\overline{x_k'^2} \overline{y_k'^2}} \quad \cancel{\text{---}}$$

Using the expressions for A_k and B_k ,
and the fact that $\overline{x_k'^2} = \frac{1}{2}$

$r^2(y, x_k)$ is $\frac{A_k^2}{2 \overline{y_k'^2}}$ for the cosine functions

$r^2(y, x_k)$ is $\frac{B_k^2}{2 \overline{y_k'^2}}$ for the sine functions:

... $\frac{A_{N/2}}{\overline{y_k'^2}}$ for $N/2$

Notes on expressing power spectrum coefficients.

Can also express $y(t)$ with complex notation:

$$y(t) = \bar{y} + \sum_{k=1}^{N/2-1} h_k e^{i(2\pi k \frac{t}{T})}$$

Recall $e^{i\theta} = \cos \theta + i \sin \theta$

Real part of $h_k = A_k$

Imaginary part $h_k = B_k$

The term with the "k" is more compactly expressed in terms of angular frequency (ω)

$$\boxed{\omega = \frac{2\pi k}{T} \quad \text{or} \quad \frac{2\pi k}{N}}$$

* Be aware of these notational conventions in the Hartmann and Wilks references — or you may find yourself confused!

The Power Spectrum

The plot of $C_k^2/2$ vs. k is called the power spectrum of $y(t)$

Frequency spectrum if t is time

Wavenumber spectrum if t is space

The 'line spectrum' (i.e. computing the Fourier transform using the entire length of record) is fine if you have infinite record.

For a finite spectrum, the line spectrum has 3 main drawbacks.

1) Integral values of k do not have any special relationship to population being sampled. Simply chosen by the length of data record; which is usually just what is available.

2) Individual spectral lines have ~~are~~ very few degrees of freedom (≈ 2 dof)

N data points determine: mean, $N/2$ amplitudes, $N/2 - 1$ phases
~~(^{Recall} dof)~~ dof relates to samples - # of estimated parameters)

3) Most geophysical data are not truly periodic, but only quasi-periodic. So they're better represented by spectral bands rather than individual spectral lines.

Just about any thing in climate that has the word 'oscillation' in it behaves this way!

e.g. ENSO \rightarrow 3-7 years

PDO \rightarrow 20-30 years

Given these drawbacks, most commonly used is the continuous power spectrum, where variance of y is given per unit frequency!

$$\overline{y^2} = \int_0^{k^*} \Phi(k) dk$$

$\Phi(k)$ = continuous power spectrum

k^* = Nyquist frequency

One cycle per $2\Delta t$

Highest frequency that can be resolved.

Can also write as $\Phi(\omega)$, which is commonly done.

Then for a spectral band:

