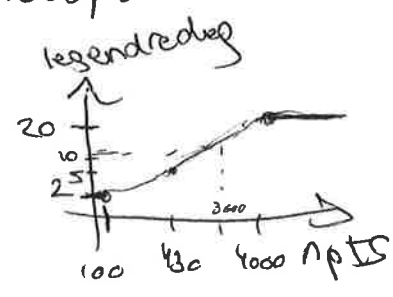


Refactor notes

- leave ~~whiten~~-allquarters as-is ... write in charges as ^{agree} ~~iter~~-whiten-allquarters, which should just make it redundant.

legendre deg selection: somewhere btwn ~~4~~⁶ & 20 depending on npts.

today's, 48pts/day \approx 4300pts.
 6⁹ days \approx 420 pts.
 2 days \approx 100 pts

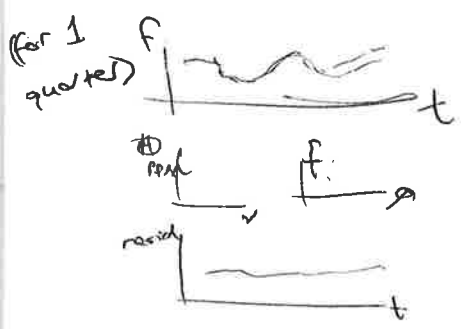


note residual becomes data for next iteration

(factor of 2 in storage)

Visualizer:

- Basically same as 3row plots but w/ item in title so



ToDO:

- add a "redetrend" step (on $f(t)$, $n \approx 29$ legendre)

rms. floor:
 set at 0.25%
 = 0.0005

Algorithm for removing all periodic components of EB variability,
while keeping planet dips:

17/03/16

- For each quarter, remove ($n=20$ finite Legendre series - a medium-order polynomial) slow-varying $f(t)$ variability, (Δ normalize to relative flux units).
- while $RMS > 0.1\%$ (SNR per transit of a $\sim 4R_{\oplus}$ CBP being ~ 1):
 - Stellingwerf phase dispersion minimizer (coarse freq bins) & get coarse period.
 - ↳ repeat over narrow bins centered on peak signal.
 - Select ~~the~~ period, which we will "whiten" at.
 - Phase-fold ~~the period~~ at period and fit (finite-order Legendre series again. Order can be determined by cross-validation, or AIC/BIC, ~~with~~ to avoid overfitting & underfitting. It should be low enough that it will gloss over transits [too sharp]).
 - Subtract fit. Compute new RMS.

(n.b. astrobase/varbase/signals.gls - whiten does a similar thing, but with GLS $\times 1/2$ then fit out n best periods) (1103 ppm)

The main assumption this makes is that removing the periodic components of the EB will be a sufficient way to drive everything below RMS of 0.1%. (If we're optimistic & write the program well, possibly even to $< 0.05\%$, which would really put our completeness at "beyond a doubt" levels.)

While aperiodic variability (attributed to starspots / magnetic activity) could be important, fitting models to it is generally trickier.

(See: ARMA, ARIMA, ARFIMA, & think further about how autoregression - "does the ^{pulsative} EB signal here match whatever is in the data there? - could be a path forward there... note similarly that 'george' might be faster for GP regression than the sklearn implementation I tried).

Most importantly, all of the above fancy statistical crap MAY NOT BE NECESSARY (ideally, it WILL NOT be necessary) for the main result - completeness for at least $4R_{\oplus}$ CBPs about big pulsators - to hold.

Stellingwerf (1978) Phase Dispersion Minimization

- Nonparametric whitening: for a proposed period, phase-fold & bin over N bins. Compute mean & std devⁿ for each bin. These are related to Θ , a statistic which becomes small (negative?) for the lowest std devⁿ of all proposed periods, (n.b. the binned std devⁿ / variance is compared to the non-phased variance - the ratio will be ~ 1 for a false period, and $\ll 1$ for a true period).

Mags $\rightarrow \vec{x}$, times \vec{t} (x_i, t_i) for $i = \{1, \dots, N\}$.

$$\sigma^2 = \frac{\sum (x_i - \bar{x})^2}{N-1} \quad (1)$$

\uparrow
variance of mags

Choose M distinct samples, with variances s_j^2 ($j = \{1, \dots, M\}$). Say each contains n_j data points. The variance for all the samples is

$$s^2 = \frac{\sum (n_j - 1) s_j^2}{\sum n_j - M} \quad (2)$$

• Total period Π , has phase vector $\vec{\Phi}$: $\Phi_i = t_i / \Pi - \lfloor t_i / \Pi \rfloor$
Divide $\vec{\Phi}$ into bins (e.g., 10 or 100). Compute
 \downarrow
floor, or "take integer part"

$$\Theta = \frac{s^2}{\sigma^2}$$

Note $\Theta \approx 1$ for false periods, but is a local minimum for true periods.

But how significant are these peaks? Schwarzenberg-Czerny (1997) notes that comparing Θ with the "Fischer-Snedecor F statistic" isn't quite right - in fact Θ follows a Beta distribution.

(Schwarzenberg-Czerny also notes that "the high-performance Fourier series method based on orthogonal projectors" weakens the use of PDM).

(He also notes that the step-function binning has some undesirable properties).

\uparrow which can be circumvented w/ linear fits)

we want a low pass filter.

17/03/16

Butterworth filter has gain $G(\omega)$ at n^{th} order of

$$G^2(\omega) = |H(j\omega)|^2 = \frac{G_0^2}{1 + \left(\frac{j\omega}{j\omega_c}\right)^{2n}}$$

n = filter order
 ω_c = cutoff freq ($\sim -3\text{dB}$ freq)
 G_0 = DC gain.

as $n \rightarrow \infty$, the gain becomes a rectangle funcⁿ.

Samp's implementation...

N.b. digital filters operate on discrete data, analog continuous(?)

gives numerator 'b' & denominator 'a' of an IIR filter.
(or FIR)

"IIR" = infinite impulse response.

$$y[n] = \frac{1}{a_0} \left(\sum_{i=0}^P b_i x[n-i] - \sum_{j=1}^Q a_j y[n-j] \right)$$

output
signal

$x[n]$
input signal

a_n feedback
filter coeff

Q : feedback
order

b_n : feedforward
filter coeff

P : feedforward order

Perhaps unsurprisingly, in practice this looks like not so good.

Maybe what we want is a finite impulse response filter
constructed from windows?

• what difference ~~does~~ does having planets make in
frequency domain?

• Another route for improvement:

better initial detrending. When there are systematics there, they
come into the fold.

target here: 0.1% RMS.

OR just do iterative POM (see
if we get any where).

good example EB subtraction codes:

• 4273411, SAP, $P = \frac{1}{2} \cdot 16 = 8$ ~~0.0000625~~ $= 0.0625\%$ dips.

• 4554004, ONLY 2 quarters!

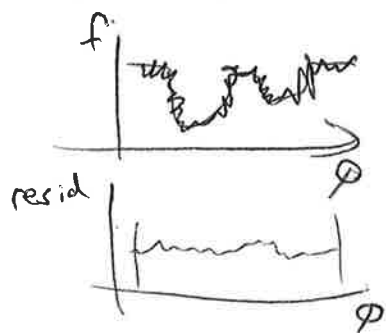
↳ short-term periodicity, with a BEAT.
(≈ 1 day)

• 4660997: more detached, w/ fit too fit to begin ✓.

• 4850874: short-term periodicity at ≈ 1 day level?

~~You want the residuals Call of them to be subtraction, not division!~~
after coming up w/ fit
which messes things up)

in phase-folded plots, the amount of "fuzz" around the fit all will slow up...



So rather than iteratively whitening by ~~rediscovers~~ the running PDM, refitting, resubtracting, ...)

use cross-validation to select preferred smoothing scale...
(w/in a quarter? or inter-quarter...)

should ideally be doing this at EVERY fitting step!

(implement it in the "EB signal subtraction testing" subroutine).

17/03/16:

• Butterworth filter

• Iterative whitening on PDM periods

Gaussian process \equiv collection of random variables in parameter space, any subset of which can be specified by a joint Gaussian distribⁿ.

17/03/15

To specify, need

- 1) mean
- 2) covariance function ("kernel").

Different kernels (covariance functions) have different "hyperparameters".

E.g. the squared exponential covariance function,

$$\text{Cor}(x_1, x_2; h) = \exp\left(-\frac{(x_1 - x_2)^2}{2h^2}\right)$$

↑
hyperparam.

if you have periodic functions, consider the Exp. Sine Squared kernel. has length parameter l and periodicity param $P > 0$

$$\text{cov}(x_i, x_j, l, P) \equiv k(x_i, x_j) = \exp\left[-\frac{2 \sin\left(\frac{\pi}{P} d(x_i, x_j)\right)}{l}\right]$$

can try out above, OR just try

$|x_i - x_j|$, or $\sqrt{(x_i - x_j)^2}$

• iterative whitening.

• remove data near gaps

meaning:
iterative spline fitting!
per F Daw's code.

↳ Implement as
iterative whitening, w/
all same options (default:
spline)

where are failures?

• 10905804

~~10905804~~

• 10965091

• 11036692

• 10032392: 1% dip, but periodicity (short P) draws out