

## PERIODOGRAMS FOR MULTIBAND ASTRONOMICAL TIME SERIES

JACOB T. VANDERPLAS<sup>1</sup> AND ŽELJKO IVEZIĆ<sup>2</sup>

*Draft version August 26, 2015*

### ABSTRACT

This paper introduces the *multiband periodogram*, a general extension of the well-known Lomb-Scargle approach for detecting periodic signals in time-domain data. In addition to advantages of the Lomb-Scargle method such as treatment of non-uniform sampling and heteroscedastic errors, the multiband periodogram significantly improves period finding for randomly sampled multiband light curves (e.g., Pan-STARRS, DES and LSST). The light curves in each band are modeled as arbitrary truncated Fourier series, with the period and phase shared across all bands. The key aspect is the use of Tikhonov regularization which drives most of the variability into the so-called base model common to all bands, while fits for individual bands describe residuals relative to the base model and typically require lower-order Fourier series. This decrease in the effective model complexity is the main reason for improved performance. After a pedagogical development of the formalism of least-squares spectral analysis which motivates the essential features of the multiband model, we use simulated light curves and randomly subsampled SDSS Stripe 82 data to demonstrate the superiority of this method compared to other methods from the literature, and find that this method will be able to efficiently determine the correct period in the majority of LSST’s bright RR Lyrae stars with as little as six months of LSST data, a vast improvement over the years of data reported to be required by previous studies. A Python implementation of this method, along with code to fully reproduce the results reported here, is available on GitHub.

*Subject headings:* methods: data analysis — methods: statistical

### 1. INTRODUCTION

Many types of variable stars show periodic flux variability (Eyer & Mowlavi 2008). Periodic variable stars are important both for testing models of stellar evolution and for using such stars as distance indicators (e.g., Cepheids and RR Lyrae stars). One of the first and main goals of the analysis is to detect variability and to estimate the period and its uncertainty. A number of parametric and non-parametric methods have been proposed to estimate the period of an astronomical time series (e.g., Graham et al. 2013, and references therein).

The most popular non-parametric method is the phase dispersion minimization (PDM) introduced by Stellingwerf (1978). Dispersion per bin is computed for binned phased light curves evaluated for a grid of trial periods. The best period minimizes the dispersion per bin. A similar and related non-parametric method that has been recently gaining popularity is the Supersmoother routine (Reimann 1994). It uses a running mean or running linear regression on the data to fit the observations as a function of phase to a range of periods. The best period minimizes a figure-of-merit, adopted as weighted sum of absolute residuals around the running mean. Neither the Supersmoother algorithm nor the PDM method require *a priori* knowledge of the light curve shape.

The most popular parametric method is the Lomb-Scargle periodogram, which is discussed in detail in Section 2. The Lomb-Scargle periodogram is related to the  $\chi^2$  for a least-square fit of a single sinusoid to data and can treat non-uniformly sampled time series with heteroscedastic measurement uncertainties. The underly-

ing model of the Lomb-Scargle periodogram is nonlinear in frequency and so the likelihood surface in frequency is non-convex. This non-convexity is readily apparent in the many local maxima of the typical periodogram, which makes it difficult to find the maximum via standard numerical optimization routines. Thus in practice the global maximum of the periodogram is often found by a brute-force grid search (for details see, e.g. Ivezić et al. 2014).

A more general parametric method based on the use of continuous-time autoregressive moving average (CARMA) model was recently introduced by Kelly et al. (2014). CARMA models can also treat non-uniformly sampled time series with heteroscedastic measurement uncertainties, and can handle complex variability patterns.

A weakness of all these standard methods is that they require homogeneous measurements – for astronomy data, this means that successive measurements must be taken through a single photometric bandpass (filter). This has not been a major problem for past surveys because measurements are generally taken through a single photometric filter (e.g. LINEAR, Sesar et al. 2011), or nearly-simultaneously in all bands at each observation (e.g. SDSS, Sesar et al. 2010). For the case of simultaneously taken multiband measurements, Süveges et al. (2012) utilized the principal component method to optimally extract the best period. Their method is essentially a multiband generalization of the well-known two-band Welch-Stetson variability index (Welch & Stetson 1993). Unfortunately, when data in each band are taken at different times, such an approach is not applicable. In such cases, past studies have generally relied on *ad hoc* methods such as a majority vote among multiple single-

<sup>1</sup> eScience Institute, University of Washington

<sup>2</sup> Department of Astronomy, University of Washington

band estimates of the periodogram (e.g., Oluseyi et al. 2012).

For surveys that obtain multiband data one band at a time, such as Pan-STARRS (Kaiser et al. 2010) and DES (Flaugher 2008), and for future multicolor surveys such as LSST (Ivezić et al. 2008), this *ad hoc* approach is not optimal. In order to take advantage of the full information content in available data, it would be desirable to have a single estimate of the periodogram which accounts for all observed data in a manner independent of assumptions about the underlying spectrum of the object. We propose such a method in this paper.

The proposed method is essentially a generalization of the Lomb-Scargle method to multiband case. The light curves in each band are modeled as arbitrary truncated Fourier series, with the period, and optionally the phase, shared across all bands. The key aspect enabling this approach is the use of Tikhonov regularization (discussed in detail in Section 4.3) which drives most of the variability into the so-called *base model* common to all bands, while fits for individual bands describe residuals relative to the base model and typically require lower-order Fourier series. This regularization-driven decrease in effective model complexity is the main reason for improved performance.

The remainder of the paper is organized as follows. Sections 2-4 offer a review of essential concepts in least squares modeling and least squares spectral analysis, as well as their relationship to common periodogram estimates: in Section 2 we provide a brief review of least-squares periodic fitting, and in Section 3 derive the matrix-based formalism for single-band least squares spectral analysis used through the rest of this work. Section 4 introduces several extensions and generalizations to the single-band model that the matrix formalism makes possible, including floating mean models, truncated Fourier models, and regularized models. Sections 5-7 present our new developments: in Section 5, we use the ideas and formalism of Sections 2-4 to motivate the *multiband periodogram*, and show some examples of its use on simulated data. In Section 6 we apply this method to measurements of 483 RR Lyrae stars first explored by Sesar et al. (2010, hereafter S10), and in Section 7 explore the performance of the method for simulated observations from the LSST survey. We conclude in Section 8.

## 2. BRIEF OVERVIEW OF PERIODIC ANALYSIS

The detection and quantification of periodicity in time-varying signals is an important area of data analysis within modern time-domain astronomical surveys. For evenly-spaced data, the *periodogram*, a term coined by Schuster (1898), gives a quantitative measure of the periodicity of data as a function of the angular frequency  $\omega$ . For data  $\{y_k\}_{k=1}^N$  measured at equal intervals  $t_k = t_0 + k\Delta t$ , Schuster's periodogram, which measures the spectral power as a function of the angular frequency, is given by

$$C(\omega) = \frac{1}{N} \left| \sum_{k=1}^N y_k e^{i\omega t_k} \right|^2, \quad (1)$$

and can be computed very efficiently using the Fast Fourier Transform.

Because astronomical observing cadences are rarely so uniform, many have looked at extending the ideas behind the periodogram to work with unevenly-sampled data. Most famously, Lomb (1976) and Scargle (1982) extended earlier work to define the *normalized periodogram*:

$$P_N(\omega) = \frac{1}{2V_y} \left[ \frac{[\sum_k (y_k - \bar{y}) \cos \omega(t_k - \tau)]^2}{\sum_k \cos^2 \omega(t_k - \tau)} + \frac{[\sum_k (y_k - \bar{y}) \sin \omega(t_k - \tau)]^2}{\sum_k \sin^2 \omega(t_k - \tau)} \right], \quad (2)$$

where  $\bar{y}$  is the mean and  $V_y$  is the variance of the data  $\{y_k\}$ , and  $\tau$  is the time-offset which orthogonalizes the model and makes  $P_N(\omega)$  independent of a translation in  $t$  (see Press et al. 2007, for an in-depth discussion). Lomb (1976) showed that this time-offset has a deeper effect: namely, it gives  $P_N$  a similar form to previous extensions of  $C(\omega)$ , while leaving  $P_N$  identical to the estimate of harmonic content given a least-squares fit to a single-component sinusoidal model,

$$d(t) = A \sin(\omega t + \phi). \quad (3)$$

This long-recognized connection between spectral power and least squares fitting methods was solidified by Jaynes (1987), who demonstrated that the least-squares periodogram method is a sufficient statistic for inferences about a stationary-frequency signal in the presence of Gaussian noise. Building on this result, Brethorst (1988) explored the extension of these methods to more complicated models with multiple frequency terms, non-stationary frequencies, and other more sophisticated models within a Bayesian framework.

While the important features of least squares frequency estimation via Lomb-Scargle periodograms have been discussed elsewhere, we will present a brief introduction to the subject in the following section. In particular, we re-express the problem in a matrix-based formalism that makes clear how the basic approach motivated by Lomb (1976), Scargle (1982), and others can be extended to more sophisticated models, including the multiband periodogram proposed in this work.

## 3. STANDARD LEAST SQUARES SPECTRAL FITTING

In this section we present a brief quantitative introduction to the least squares fitting formulation of the normalized periodogram of Equation (2). We denote  $N$  observed data points as

$$D = \{t_k, y_k, \sigma_k\}_{k=1}^N \quad (4)$$

where  $t_k$  is the time of observation,  $y_k$  is the observed value (typically a magnitude), and  $\sigma_k$  describes the Gaussian errors on each value. For notational simplicity we will assume without loss of generality that the data  $y_k$  are centered such that the measurements within each band satisfy

$$\frac{\sum_k w_k y_k}{\sum_k w_k} = 0 \quad (5)$$

where the weights are  $w_k = \sigma_k^{-2}$ . Though this assumption is essential to the simpler models presented in this

section, it will become superfluous with the floating-mean models described in later sections.

### 3.1. Stationary Sinusoid Model

The normalized periodogram of Equation (2) can be derived from the normalized  $\chi^2$  of the best-fit single-term stationary sinusoidal model given in Equation (3). To make the problem linear, we can re-express the model in terms of the parameter vector  $\theta = [A \cos \phi, A \sin \phi]$  so that our model is

$$y(t|\omega, \theta) = \theta_1 \sin(\omega t) + \theta_2 \cos(\omega t). \quad (6)$$

For a given  $\omega$ , the maximum likelihood estimate of the parameters  $\theta$  can be found by minimizing the  $\chi^2$  of the model, which is given by

$$\chi^2(\omega) = \sum_k \frac{[y_k - y(t_k|\omega, \theta)]^2}{\sigma_k^2}. \quad (7)$$

For the single-term Fourier model, it can be shown (see, e.g. Ivezić et al. 2014) that

$$\chi_{min}^2(\omega) = \chi_0^2[1 - P_N(\omega)] \quad (8)$$

where  $P_N(\omega)$  is the normalized periodogram given in Equation (2)<sup>3</sup> and  $\chi_0^2$  is the reference  $\chi^2$  for a constant model, which due to the assumption in Equation (5) is simply  $\chi_0^2 = \sum_k (y_k/\sigma_k)^2$ .

### 3.2. Matrix Formalism

A standard way of compactly expressing least squares models is via matrix expressions (See e.g. Brandt 1970). Likewise, the expressions related to the stationary sinusoid model can be expressed more compactly by defining the following matrices:

$$X_\omega = \begin{bmatrix} \sin(\omega t_1) & \cos(\omega t_1) \\ \sin(\omega t_2) & \cos(\omega t_2) \\ \vdots & \vdots \\ \sin(\omega t_N) & \cos(\omega t_N) \end{bmatrix};$$

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}; \quad \Sigma = \begin{bmatrix} \sigma_1^2 & 0 & \cdots & 0 \\ 0 & \sigma_2^2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_N^2 \end{bmatrix} \quad (9)$$

With these definitions, the model in Equation (6) can be expressed as a simple linear product,  $y(t|\omega, \theta) = X_\omega \theta$ , and the model and reference  $\chi^2$  can be written

$$\chi^2(\omega) = (y - X_\omega \theta)^T \Sigma^{-1} (y - X_\omega \theta) \quad (10)$$

$$\chi_0^2 = y^T \Sigma^{-1} y \quad (11)$$

The expression for the normalized periodogram can be computed by finding via standard methods the value of

<sup>3</sup> An important feature of the Lomb-Scargle approach is the modification of the model with the time-offset  $\tau$  tuned to orthogonalize the harmonic basis across the irregular times  $\{t_i\}$ . This orthogonalization cancels cross-terms in the expression of  $\chi^2$ , greatly reducing the complexity of computing  $P_N$ . As discussed in footnote 4, however, this orthogonalization does not change the resulting periodogram and so it can safely be ignored for the purposes of this work.

$\theta$  which minimizes  $\chi^2(\omega)$ , and plugging the result into Equation (8). This yields

$$P_N(\omega) = \frac{y^T \Sigma^{-1} X_\omega [X_\omega^T \Sigma^{-1} X_\omega]^{-1} X_\omega^T \Sigma^{-1} y}{y^T \Sigma^{-1} y}. \quad (12)$$

We note that Equation (12) is equivalent to Equation (2) in the homoscedastic case with  $\Sigma \propto V_y I$ . <sup>4</sup>

### 3.3. Simple Single-band Period Finding

As an example of the standard periodogram in action, we perform a simple single-band harmonic analysis of simulated  $r$ -band observations of an RR Lyrae light curve, based on empirical templates derived in S10 (Figure 1). The observations are of a star with a period of 0.622 days, and take place on 60 random nights over a 6-month period, as seen in the left panel.

The upper-right panel shows the normalized periodogram for this source as a function of period. While the power does peak at the true period of 0.622 days, an aliasing effect is readily apparent near  $P = 0.38$ . This additional peak is due to beat frequency between the true period  $P$  and the observing cadence of  $\sim 1$  day. This beat frequency is the first in a large sequence: for nightly observations, we'd expect to find excess power at periods  $P_n = P/(1 + nP)$  days, for any integer  $n$ . The strong alias in Figure 1 corresponds to the  $n = 1$  beat period  $P_n = 0.383$ . Though it is possible to carefully correct for such aliasing by iteratively removing contributions from the estimated window function (e.g. Roberts et al. 1987), we'll ignore this detail in the current work.

The lower-right panel of Figure 1 shows the maximum likelihood interpretation of this periodogram: it is a measure of the normalized  $\chi^2$  for a single-term sinusoidal model. Here we visualize the data from the left panel, but folded as a function of phase, and overplotted with the best-fit single-term model. This visualization makes it apparent that the single-term model is highly biased: RR Lyrae light curves are, in general, much more complicated than a simple sinusoid. Nevertheless, the simplistic sinusoidal model is able to recover the correct frequency to a high degree of accuracy (roughly related to the width of the peak) and significance (roughly related to the height of the peak; see Scargle (1982) for details). For a more complete introduction to and discussion of the single-term normalized periodogram, refer to, e.g. Bretthorst (1988) or Ivezić et al. (2014).

## 4. GENERALIZING THE PERIODOGRAM MODEL

We have shown two forms of the classic normalized periodogram: Equation (2) and Equation (12). Though the two expressions are equivalent, they differ in their utility. Because the expression in Equation (2) avoids the explicit construction of a matrix, it can be computed very efficiently. Furthermore, through clever use of the

<sup>4</sup> For direct comparison to the Lomb-Scargle approach, we need the equivalent of the  $\tau$  parameter which orthogonalizes the basis across the observed times  $\{t_i\}$ . Such an orthogonalization is accomplished via the transformations  $X_\omega \rightarrow X_\omega V_\omega$  and  $\theta \rightarrow V_\omega^T \theta$ , where  $V_\omega$  is the orthogonal matrix of eigenvectors of the covariance  $X_\omega^T \Sigma^{-1} X_\omega$ . The  $V_\omega$  terms straightforwardly cancel out of Equations (10)-(12) and the results of this section are unchanged. The general matrix formalism used here makes clear that this result applies to all the periodogram extensions mentioned in this work.

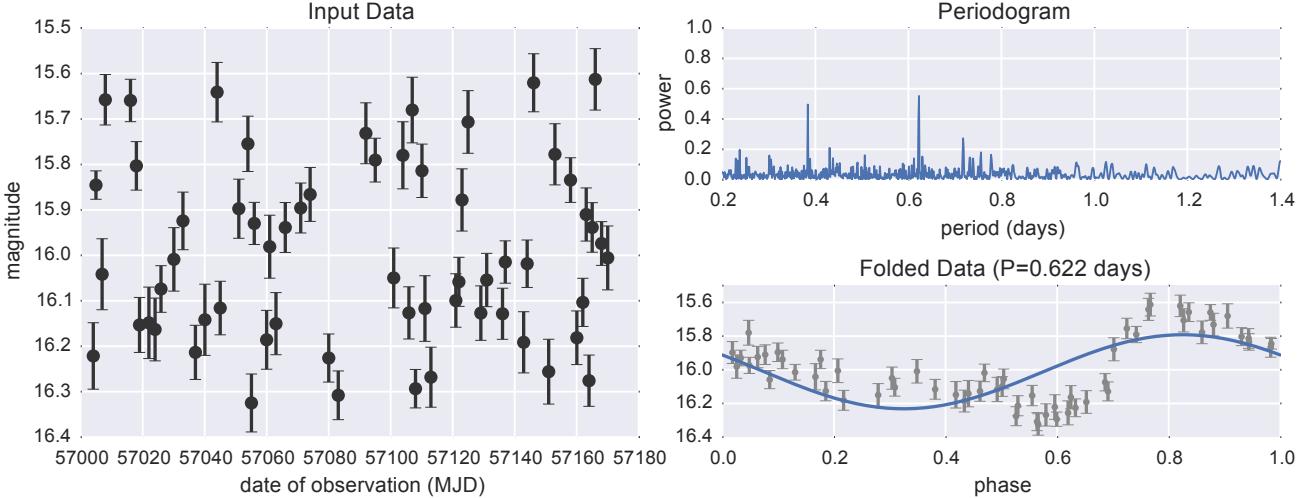


FIG. 1.— An illustration of the basic periodogram and its relationship to the single-term sinusoid model. The left panel shows the input data, while the right panels show the fit derived from the data. The upper-right panel shows the periodogram with a clear peak at the true period of 0.622 days, and the bottom-right panel shows the data as a function of the phase associated with this period. Note in the periodogram the presence of the typical aliasing effect, with power located at beat frequencies between the true period and the 1-day observing cadence (see Section 3.3 for further discussion).

Fast Fourier Transform, expressions of the form of Equation (2) can be evaluated exactly for  $N$  frequencies in  $\mathcal{O}[\log N]$  time (Press & Rybicki 1989).

The matrix-based formulation of Equation (12), though slower than the Fourier-derived formulation, is a more general expression and allows several advantages:

1. It is straightforwardly extended to heteroscedastic and/or correlated measurement noise in the data  $y_k$  through appropriate modification of the *noise covariance matrix*  $\Sigma$ .
2. It is straightforwardly extended to more sophisticated linear models by appropriately modifying the *design matrix*  $X_\omega$ .
3. It is straightforwardly extended to include Tikhonov/L2-regularization terms (see Section 4.3 for more details) by adding an appropriate diagonal term to the *normal matrix*  $X_\omega^T \Sigma^{-1} X_\omega$ .

In the remainder of this section, we will explore a few of these modifications and how they affect the periodogram and resulting model fits.

#### 4.1. Stationary Sinusoid with Floating Mean

As an example of one of these generalizations, we'll consider what has variously been called the *Date-compensated Discrete Fourier Transform* (Ferraz-Mello 1981), the *floating-mean periodogram* (Cumming et al. 1999), and the *generalized Lomb-Scargle method* (Zechmeister & Kürster 2009). Here we use the term *floating-mean periodogram*. This method adjusts the classic normalized periodogram by fitting the mean of the model alongside the amplitudes:

$$y(t | \omega, \theta) = \theta_0 + \theta_1 \sin \omega t + \theta_2 \cos \omega t \quad (13)$$

The periodogram derived from this model can be more accurate than the standard Phys. Rev. E-centered periodogram for certain observing cadences and selection

functions, and especially when searching for long-period variability or working with very few samples (Cumming et al. 1999). Zechmeister & Kürster (2009) detail the required modifications to the orthogonalized harmonic formalism of Equation (2) to allow the mean to float in the model. In the matrix formalism, the modification is much more straightforward: all that is required is to add a column of ones to the  $X_\omega$  matrix before computing the power via Equation (12). This column of ones corresponds to a third entry in the parameter vector  $\theta$ , and acts as a uniform constant offset for all data points.

For well-sampled data, there is usually very little difference between a standard periodogram on pre-centered data and a floating-mean periodogram. Where this difference becomes important is if selection effects or observing cadences cause there to be preferentially more observations at certain phases of the light curve: a toy example demonstrating this situation is shown in Figure 2. The data are drawn from a sinusoid with Gaussian errors, and data with a magnitude fainter than 16 are removed to simulate an observational bias (left panel). Because of this observational bias, the mean of the observed data is a poor predictor of the true mean, causing the standard method to poorly fit the data and miss the input period (upper-right panel). The floating-mean approach is able to automatically adjust for this bias, resulting in a periodogram which readily detects the input period of 0.622 days (lower-right panel).

#### 4.2. Truncated Fourier Models

As mentioned above, the standard periodogram is equivalent to fitting a single-term stationary sinusoidal model to the data. A natural extension is to instead use a multiple-term sinusoidal model, with frequencies at integer multiples of the fundamental frequency (See e.g. Brethorst 1988). With  $N$  Fourier terms, there are

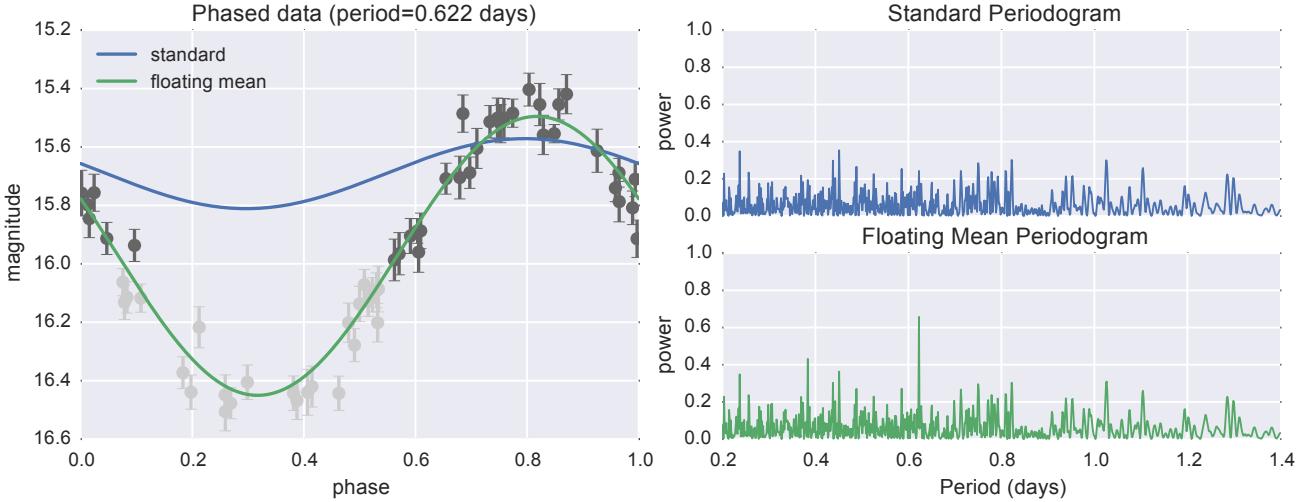


FIG. 2.— An illustration of the effect of the floating mean model for censored data. The data consist of 80 observations drawn from a sinusoidal model. To mimic a potentially damaging selection effect, all observations with magnitude fainter than 16 are removed (indicated by the light-gray points). The standard and floating-mean periodograms are computed from the remaining data; these fits are shown over the data in the left panel. Because of this biased observing pattern, the mean of the observed data is a biased estimator of the true mean. The standard fixed-mean model in this case fails to recover the true period of 0.622 days, while the floating mean model still finds the correct period.

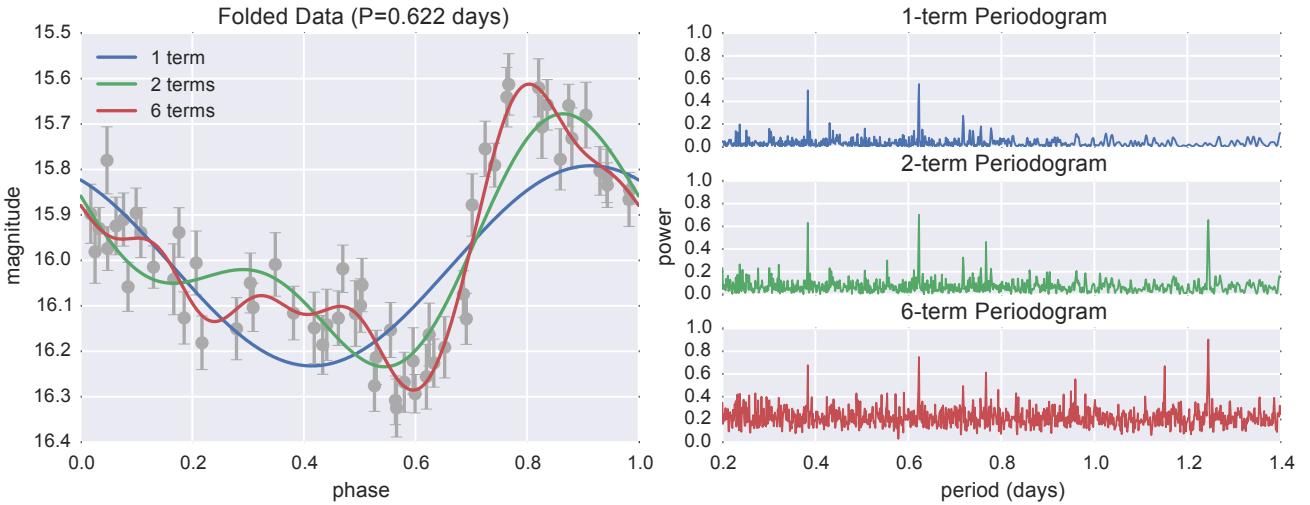


FIG. 3.— The model fits and periodograms for several truncated Fourier models. The data are the same as those in Figure 1. Note that in addition to the previously-seen 0.38-day alias, the higher-order models will generally show periodogram peaks at multiples of the true fundamental frequency  $P_0$ : this is because for integer  $n$  less than the number of Fourier terms in the model,  $P_0$  is a higher harmonic of the model at  $P = nP_0$ . Additionally, the increased degrees of freedom in the higher-order models let them fit better at any frequency, which drives up the “background” level in the periodogram.

$2N + 1$  free parameters, and the model is given by

$$y(t|\omega, \theta) = \theta_0 + \sum_{n=1}^N [\theta_{2n-1} \sin(n\omega t) + \theta_{2n} \cos(n\omega t)]. \quad (14)$$

Because this model remains linear in the parameters  $\theta$ , it can be easily accommodated into the matrix formalism of Section 3.2. For example, an  $N = 2$ -term floating-mean model can be constructed by building a design matrix

$X_\omega$  with  $2N + 1 = 5$  columns:

$$X_\omega^{(2)} = \begin{bmatrix} 1 & \sin(\omega t_1) & \cos(\omega t_1) & \sin(2\omega t_1) & \cos(2\omega t_1) \\ 1 & \sin(\omega t_2) & \cos(\omega t_2) & \sin(2\omega t_2) & \cos(2\omega t_2) \\ 1 & \sin(\omega t_3) & \cos(\omega t_3) & \sin(2\omega t_3) & \cos(2\omega t_3) \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & \sin(\omega t_N) & \cos(\omega t_N) & \sin(2\omega t_N) & \cos(2\omega t_N) \end{bmatrix} \quad (15)$$

Computing the power via Equation (12) using  $X_\omega^{(2)}$  will give the two-term periodogram. For larger  $N$ , more columns are added, but the periodogram can be computed in the same manner. Figure 3 shows a few examples of this multiterm Fourier approach as applied to

the simulated RR Lyrae light curve from Figure 1, and illustrates several important insights into the subtleties of this type of multiterm fit.

First, we see in the right panel that all three models show a clear signal at the true period of  $P_0 = 0.622$  days. The higher-order models, however, also show a spike in power at  $P_1 = 2P_0$ : the reason for this is that for an  $N > 1$ -term model, the period  $P_0$  is the first harmonic of a model with fundamental frequency  $2P_0$ , and the higher-order models contain the single-period result.

Second, notice that as the number of terms is increased, the general “background” level of the periodogram increases. This is due to the fact that the periodogram power is inversely related to the  $\chi^2$  of the fit at each frequency. A more flexible higher-order model can better fit the data at all periods, not just the true period. Thus in general the observed power of a higher-order Fourier model will be everywhere higher than the power of a lower-order Fourier model.

One might hope that when adding terms, the correct-period model would show more of an improvement than the incorrect-period model (and thus the periodogram maximum would become more pronounced in comparison to the background), but this does not generally hold. Consider that in the extreme limit in which the number of model parameters is equal to the number of data points, the model has enough flexibility to fit the data perfectly at *every* frequency, and the resulting periodogram would be everywhere unity! This can only be the case if, on average, addition of terms preferentially boosts the background level.

### 4.3. Regularized Models

The previous sections raise the question: how complicated a model should we use? We have seen that as we add more terms to the fit, the model will more closely describe the observed data. For very high-order models, however, such a close fit *over-fits* the data: that is, the fit is more responsive to statistical noise in the observations than to the underlying signal. This can be addressed by explicitly truncating the series at some number of terms, but we can also use a *regularization* term to mathematically enforce model simplicity.

A regularization term is an explicit penalty on the magnitude of the model parameters  $\theta$ , and can take a number of forms. For computational simplicity here we’ll use an *L2 regularization* – also known as Tikhonov Regularization (Tikhonov 1963) or Ridge Regression (Hoerl & Kennard 1970) – which is a quadratic penalty term in the model parameters added to the  $\chi^2$ . Mathematically, this is equivalent in the Bayesian framework to using a zero-mean Gaussian prior on the model parameters.

We encode our regularization in the matrix  $\Lambda = \text{diag}([\lambda_1, \lambda_2 \dots \lambda_M])$  for a model with  $M$  parameters, and construct a “regularized”  $\chi^2$ :

$$\chi_{\Lambda}^2(\omega) = (y - X_{\omega}\theta)^T \Sigma^{-1} (y - X_{\omega}\theta) + \theta^T \Lambda \theta \quad (16)$$

Minimizing this regularized  $\chi^2$ , solving for  $\theta$ , and plugging into the expression for  $P_N$  gives us the regularized

counterpart of Equation (12):

$$P_{N,\Lambda}(\omega) = \frac{y^T \Sigma^{-1} X_{\omega} [X_{\omega}^T \Sigma^{-1} X_{\omega} + \Lambda]^{-1} X_{\omega}^T \Sigma^{-1} y}{y^T \Sigma^{-1} y}. \quad (17)$$

Notice that the effect of this regularization term is to add a diagonal penalty to the normal matrix  $X_{\omega}^T \Sigma^{-1} X_{\omega}$ , which has the additional feature that it can correct ill-posed models where the normal matrix is non-invertible. This feature of the regularization will become important for the multiband models discussed below.

In Figure 4, we compare a regularized and unregularized 20-term truncated Fourier model on our simulated RR Lyrae light curve. We use  $\lambda = 0$  on the offset term, and make the penalty  $\lambda_j$  progressively larger for each harmonic component. The regularization prevents overfitting (left panel), and results in more prominent periodogram peaks (right panel).

## 5. A MULTIPLE-BAND MODEL

In this section we will combine the ideas of the previous sections to construct the *multiband periodogram* which flexibly accounts for heterogeneous sources of data for a single object. To start with, we might consider one of two naïve approaches to the multi-band problem:

First, we might ignore band labels entirely and simply compute a single standard Lomb-Scargle periodogram over the full dataset. This amounts to the assumption that one global model suitably fits each band, and in practice will perform poorly due to the astrophysical variability between bands: in other words, the model is too simple and under-fits the data.

Second, we might treat each band entirely independently and compute a standard Lomb-Scargle periodogram on each, and use the additivity of  $\chi^2$  along with Equation (8) to construct a multiband periodogram. This amounts to the assumption that the bands have completely independent phases and amplitudes, and has too many free parameters to be useful in most cases of interest. In other words, the model is too complex and over-fits the data (see Section 5.1 for further discussion).

To compute a periodogram which strikes a balance between these two extremes, we will take advantage of the easy extensibility of the matrix formalism which led to our generalizations above. The multiband model presented here contains the following features:

1. An  $N_{base}$ -term truncated Fourier “base model” which models the shared variability among all  $K$  observed bands.
2. A set of  $N_{band}$ -term truncated Fourier fits, each of which models the residual within a single band from the shared variability accounted for in the base model.

The total number of parameters for  $K$  bands is then  $M_K = (2N_{base} + 1) + K(2N_{band} + 1)$ . As a result, for each band  $k$  we have the following model of the observed magnitudes:

$$y_k(t|\omega, \theta) = \theta_0 + \sum_{n=1}^{N_{base}} [\theta_{2n-1} \sin(n\omega t) + \theta_{2n} \cos(n\omega t)] + \theta_0^{(k)} + \sum_{n=1}^{N_{band}} [\theta_{2n-1}^{(k)} \sin(n\omega t) + \theta_{2n}^{(k)} \cos(n\omega t)]. \quad (18)$$

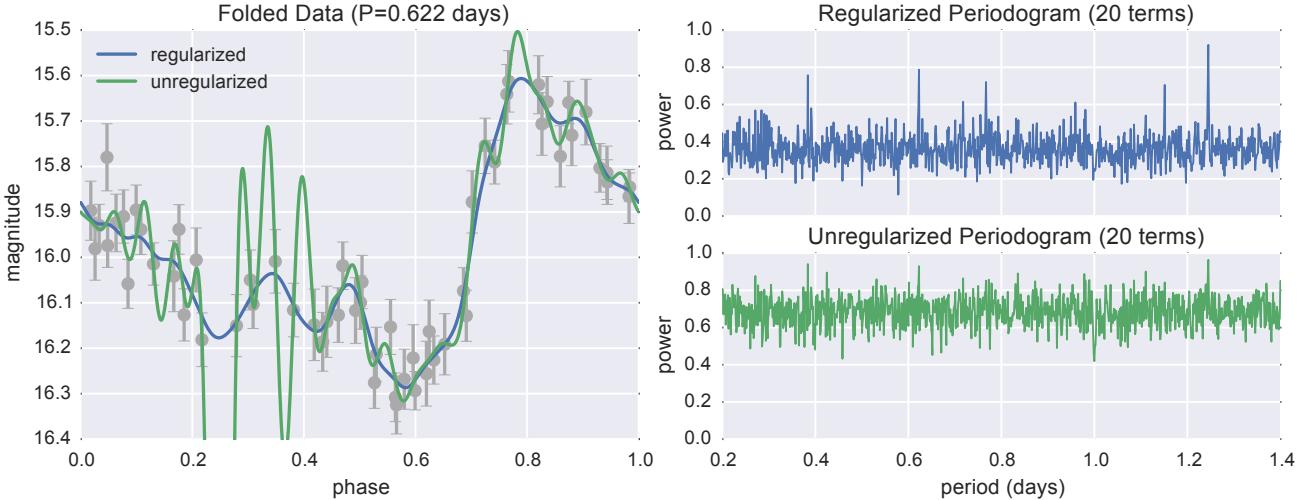


FIG. 4.— The effect of regularization on a high-order model. The data is the same as those in Figure 1. We fit a 20-term truncated Fourier model to the data, with and without a regularization term. Without regularization, the model oscillates widely to fit the noise in the data. The regularization term effectively damps the higher-order Fourier modes and removes this oscillating behavior, leading to a more robust model with stronger periodogram peaks.

The important feature of this model is that *all bands* share the same base parameters  $\theta$ , while their offsets  $\theta^{(k)}$  are determined individually. Note the potential for confusion:  $N_{band}$  here is not the number of observed bands, but the number of Fourier components fit to the residuals in each of the  $K$  observed bands.

We can construct the normalized periodogram for this model by building a sparse design matrix with  $M_K$  columns. Each row corresponds to a single observation through a single band. Columns corresponding to the base model and the matching observation band will have nonzero entries; all other columns will be filled with zeros. For example, the  $(N_{base}, N_{band}) = (1, 0)$  model corresponds to one with a simple single-term periodic base frequency, and an independent constant offset term in each band. The associated design matrix depends on the particular data, but will look similar to this:

$$X_\omega^{(1,0)} = \begin{bmatrix} 1 & \sin(\omega t_1) & \cos(\omega t_1) & 1 & 0 & 0 & 0 & 0 \\ 1 & \sin(\omega t_2) & \cos(\omega t_2) & 0 & 0 & 0 & 0 & 1 \\ 1 & \sin(\omega t_3) & \cos(\omega t_3) & 0 & 0 & 0 & 1 & 0 \\ \vdots & \vdots \\ 1 & \sin(\omega t_N) & \cos(\omega t_N) & 0 & 0 & 1 & 0 & 0 \end{bmatrix} \quad (19)$$

Here the nonzero entries of the final five columns are binary flags indicating the  $(u, g, r, i, z)$ -band of the given observation: for this example, the first row is a  $u$ -band measurement, the second is a  $z$ -band, the third is a  $i$ -band, etc., as indicated by the position of the nonzero matrix element within the row.

On examination of the above matrix, it's clear that the columns are not linearly independent (i.e.  $X_\omega$  is low-rank), and thus the parameters of the best-fit model will be degenerate. Intuitively, this is due to the fact that if we add an overall offset to the base model, this can be perfectly accounted for by subtracting that same offset from each residual model. Mathematically, the result of this is that the normal matrix  $X_\omega^T \Sigma^{-1} X_\omega$  will be non-invertible, and thus the periodogram is ill-defined. In order to proceed, then, we'll either have to use a different

model, or use a cleverly-constructed regularization term on one of the offending parameters.

We'll choose the latter here, and regularize all the band columns while leaving the base columns un-regularized: for the above  $X_\omega$  matrix, this regularization will look like

$$\Lambda^{(1,0)} = \text{diag}([0, 0, 0, \lambda, \lambda, \lambda, \lambda, \lambda]) \quad (20)$$

where  $\lambda$  controls the degree of regularization. As  $\lambda$  grows large, the model will preferentially push power into the base terms, while minimizing the deviations of the model for each individual band.

Here we will choose  $\lambda$  to be some small fraction of the trace of the normal matrix  $[X_\omega^T \Sigma^{-1} X_\omega]$ . This choice ensures the multiband periodogram is well-defined, while maintaining the flexibility of the model in accounting for independent band-to-band variation. With this regularization in place, the model is well-posed and Equation (17) can be used to straightforwardly compute the power. The effective number of free parameters for such a regularized  $(N_{base}, N_{band})$  model with  $K$  filters is  $M_K^{eff} = 2N_{base}^{eff} + K(2N_{band} + 1)$  where  $N_{base}^{eff} = \max(0, N_{base} - N_{band})$  is the effective number of base terms.

The final remaining piece to mention is our assumption in Equation (5) that the data are centered. This is required so that the simple form of the reference  $\chi_0^2$  remains valid. For the multiband model, this assumption requires that the data satisfy Equation (5) *within each band*: equivalently, we could lift this assumption and compute the reference  $\chi_0^2$  of the multiband model with an independent floating mean within each band; the results will be identical.

This multiband approach, then, actually comprises a set of models indexed by their value of  $N_{base}$  and  $N_{band}$ . The most fundamental models have  $(N_{base}, N_{band}) = (1, 0)$  and  $(0, 1)$ , which we'll call the *shared-phase* and *multi-phase* models respectively. In the shared-phase model, all variability is assumed to be shared between the bands, with only the fixed offset between them allowed to float. In the multi-phase model, each band has

independent variability around a shared fixed offset.

### 5.1. Relationship of Multiband and Single-band approaches

With this formalism in place, we can return briefly to the naïve models discussed at the beginning of Section 5. The first, which ignores band information, is simply a standard Lomb-Scargle over the heterogeneous data. The second, in which each band is fit independently, turns out to be equivalent to the  $(N_{base}, N_{band}) = (0, 1)$  model defined above. Here the base model is a simple global offset which is degenerate with the offsets in each band, so that the design matrix  $X_\omega$  can be straightforwardly rearranged as block-diagonal. A block-diagonal design matrix in a linear model indicates that components of the model are being solved independently: here these independent components amount to the single-band floating-mean model from Section 4.1, fit independently for each of the  $K$  bands.

For band  $k$ , we'll denote the single-band floating-mean periodogram as

$$P_N^{(k)}(\omega) = 1 - \frac{\chi_{min,k}^2(\omega)}{\chi_{0,k}^2} \quad (21)$$

The full multiband periodogram is given by

$$P_N^{(0,1)}(\omega) = 1 - \frac{\sum_{k=1}^K \chi_{min,k}^2(\omega)}{\sum_{k=1}^K \chi_{0,k}^2} \quad (22)$$

and it can be shown straightforwardly that  $P_N^{(0,1)}$  can be constructed as a weighted sum of  $P_N^{(k)}$ :

$$P_N^{(0,1)}(\omega) = \frac{\sum_{k=1}^K \chi_{0,k}^2 P_N^{(k)}}{\sum_{k=1}^K \chi_{0,k}^2}. \quad (23)$$

Thus the  $(N_{base}, N_{band}) = (0, 1)$  multiband periodogram is identical to a weighted sum of standard periodograms in each band, where the weights  $\chi_{0,k}^2$  are a reflection of both the number of measurements in each band and how much those measurements deviate from a simple constant reference model.

### 5.2. Multiband Periodogram for Simulated Data

Before applying the multiband method to real data, we will here explore its effectiveness on a simulated RR Lyrae lightcurve. The upper panels of Figure 5 show a multiband version of the simulated RR Lyrae light curve from Figure 1. The upper-left panel shows 60 nights of observations spread over a 6-month period, and for each night all five bands ( $u, g, r, i, z$ ) are recorded. Using the typical approach from the literature, we individually compute the standard normalized periodogram within each band: the results are shown in the upper-right panel. The data are well-enough sampled that a distinct period of 0.622 days can be recognized within each individual band, up to the aliasing effect discussed in Section 3.3. Previous studies have made use of the information in multiple bands to choose between aliases and estimate uncertainties in determined periods (e.g. Sesar et al. 2010; Oluseyi et al. 2012). While this approach is sufficient for well-sampled data, it becomes problematic when the multiband data are sparsely sampled.

The lower panels of Figure 5 show the same 60 nights of data, except with only a *single* band observation recorded each night. The lower-left panel shows the observations as a function of phase, and the lower-right panels show the periodograms derived from the data. With only 12 observations for each individual band, it is clear that there is not enough data to accurately determine the period within each single band. The shared-phase ( $N_{base}, N_{band}$ ) =  $(1, 0)$  multiband approach, shown in the lower-right panel, fits a single model to the full data and clearly recovers the true frequency of 0.622 days. The key result is that while methods based on the standard periodogram are suitable for densely-sampled data, the multiband periodogram is superior for sparsely-sampled multiband observations.

This shared-phase  $(1, 0)$  model is only one of the possible multiband options, however: Figure 6 compares multiband fits to this data for models with various choices of  $(N_{base}, N_{band})$ . We see here many of the characteristics noted above for single-band models: as discussed in Section 4.2, increasing the number of Fourier terms leads to power at multiples of the fundamental period, and increased model complexity (roughly indexed by the effective number of free parameters  $M^{eff}$ ) tends to increase the background level of the periodogram, obscuring significant peaks. For this reason, models with  $N_{base} > N_{band}$  are the most promising: they allow a flexible fit with minimal model complexity. Motivated by this, in the next section we'll apply the simplest of this class of models, the  $(1, 0)$  shared-phase model, to data from the Stripe 82 of the Sloan Digital Sky Survey.

## 6. APPLICATION TO STRIPE 82 RR LYRAE

Stripe 82 is a three hundred square degree equatorial region of the sky which was repeatedly imaged through multiple band-passes during phase II of the Sloan Digital Sky Survey (SDSS II, see Sesar et al. 2007). Here we consider the SDSS II observations of 483 RR Lyrae stars compiled and studied by S10, in which periods for these stars were determined based on empirically-derived light curve templates. Because the template-fitting method is extremely computationally intensive, S10 first determined candidate periods by taking the top 5 results of the Supersmooth (Reimann 1994) algorithm applied to the  $g$ -band; template fits were then performed at each candidate period and the period with the best template fit was reported as the true period. In this section, we make use of this dataset to quantitatively evaluate the effectiveness of the multiband periodogram approach.

### 6.1. Densely-sampled Multiband Data

The full S10 RR Lyrae dataset consists of 483 objects with an average of 55 observations in each of the five SDSS  $ugriz$  bands spread over just under ten years. In the upper panels of Figure 7 we show the observed data for one of these objects, along with the periodogram derived with the single-band supersmooth model<sup>5</sup> and the shared-phase  $(0, 1)$ -multiband model. Here we have a case which is analogous to that shown for simulated

<sup>5</sup> The supersmooth “periodogram”  $P_{SS}$  is constructed from the minimum sum of weighted model residuals  $\bar{r}_{min}$  in analogy with Equation (8):  $P_{SS}(\omega) = 1 - \bar{r}_{min}(\omega)/\bar{r}_0$ , where  $\bar{r}_0$  is the mean absolute residual around a constant model.

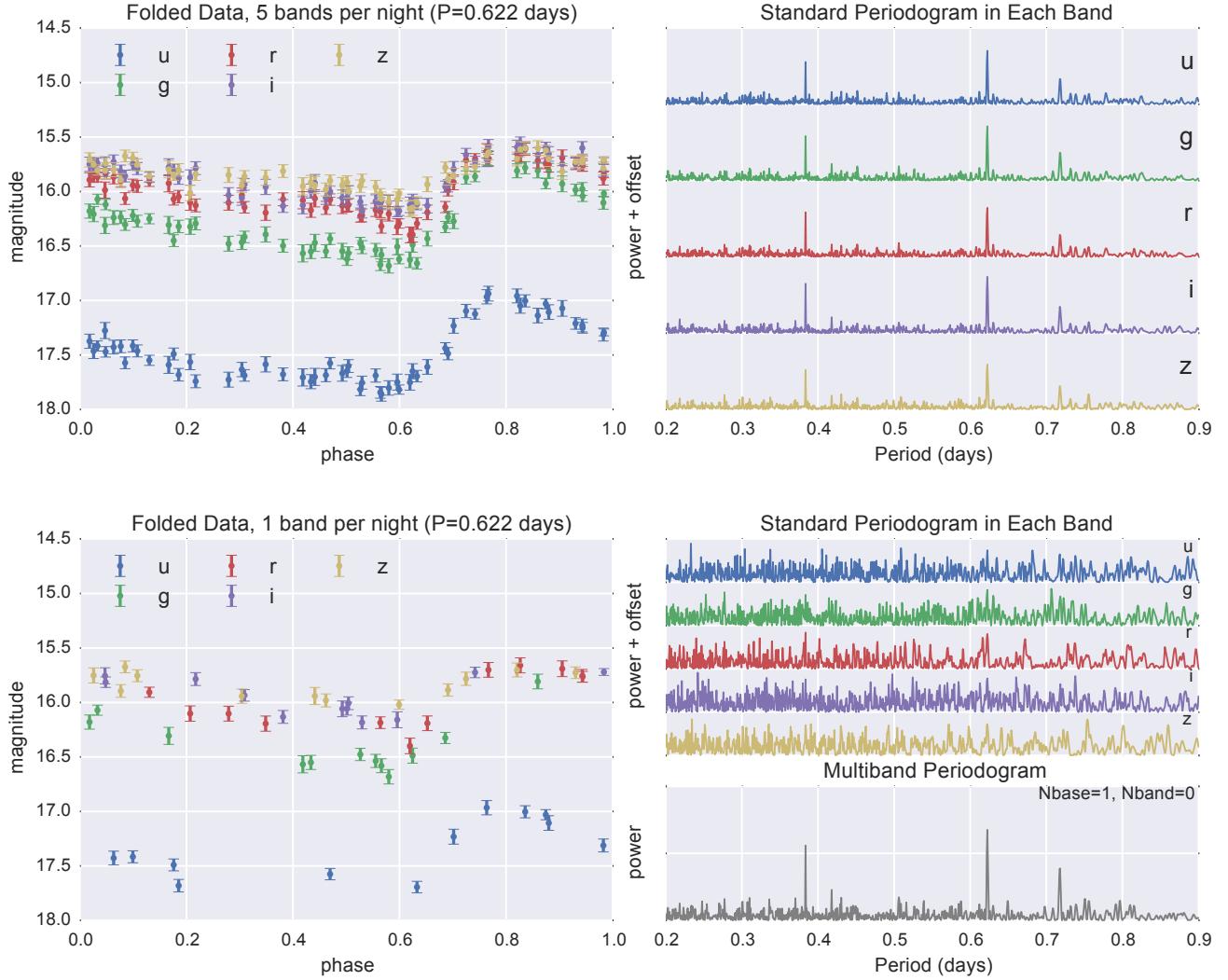


FIG. 5.— An illustration of the performance of the multiband periodogram. The upper panels show simulated *ugriz* observations of an RR Lyrae light curve in which all 5 bands are observed each night. With 60 observations in each band, a periodogram computed from any single band is sufficient to determine the true period of 0.622 days. The lower panels show the same data, except with only a single *ugriz* band observed each night (i.e. 12 observations per band). In this case, no single band has enough information to detect the period. The shared-phase multiband approach of Section 5 (lower-right panel) combines the information from all five bands, and results in a significant detection of the true period. This indicates that while methods based on the standard periodogram are suitable for densely-sampled multiband data, the multiband periodogram is superior for sparsely-sampled multiband observations.

data in the top panels of Figure 5: each band has enough data to easily locate candidate peaks, the best of which is selected via the S10 template-fitting procedure.

The lower panels of Figure 7 compare the S10 period with the best periods obtained from the 1-band supersmoother (lower-left) and from the shared-phase multiband model (lower-right). To guide the eye, the figure includes indicators of the locations of beat aliases (dotted lines) and first harmonic aliases (dashed lines) of the S10 period. Numerical results are summarized in the upper rows of Table 1.

The best-fit supersmoother period matches the S10 period in 87% of cases (421/483), while the best-fit multiband period matches the S10 period in 79% of cases (382/483). The modes of failure are instructive: when the supersmoother model misses the S10 period, it tends to land on a harmonic alias (i.e. the dashed line). This is due to the flexibility of supersmoother: a doubled pe-

riod spreads the points out, leading to fewer constraints in each neighborhood and thus a smaller average residual around model. In other words, the SuperSmoothen tends to over-fit data which is sparsely-sampled. On the other hand, when the multiband model misses the S10 period, it tends to land on a beat alias between the S10 period and the 1-day observing cadence (i.e. the dotted line). This is due to the fact that the single-frequency periodic model is biased, and significantly under-fits the data: it cannot distinguish residuals due to underfitting from residuals due to window function effects.

In both models, the S10 period appears among the top 5 periods 99% of the time: 477/483 for supersmoother, and 480/483 for multiband.<sup>6</sup> This suggests that had S10

<sup>6</sup> We might expect this correspondence to be 100% in the case of the *g*-band supersmoother, which was the model used in the first pass of the S10 computation. This discrepancy here is likely due to the slightly different supersmoother implementations used in S10

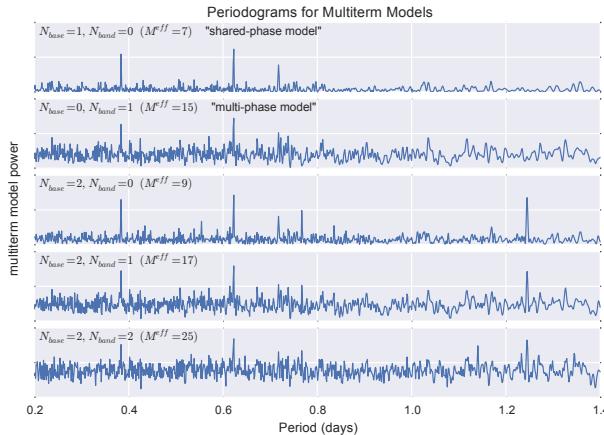


FIG. 6.— Comparison of the periodograms produced by various multiband models. The data is the same as that used in Figure 5.  $N_{\text{base}}$  gives the number of Fourier terms in the base model, and  $N_{\text{band}}$  gives the number of Fourier terms used to fit the residuals around this model within each band. The characteristics discussed with previous figures are also seen here: in particular, the level of “background noise” in the periodogram grows with the model complexity  $M$ ,

used the multiband Lomb-Scargle rather than the supersmoother in the first pass for that study, the final results presented there would be for the most part unchanged.

The results of this subsection show that the shared-phase multiband approach is comparable to the single-band supersmoother approach for densely-sampled multiband data, although it has a tendency to get fooled by structure in the survey window. Correction for this based on the estimated window power may alleviate this (see Roberts et al. (1987) for an example of such an approach) though in practice selecting from among the top 5 peaks appears to be sufficient.

## 6.2. Sparsely-sampled Multiband Data

Above we saw that the multiband model is comparable to methods from the literature for densely-sampled data. Where we expect the multiband approach to gain an advantage is when the data are sparsely sampled, with data through only a single band at each observation time. To simulate this, we reduce the size of the Stripe 82 RR Lyrae dataset by a factor of 5, keeping only a single band of imaging each night: an average of 11 observations of each object per band. This is much closer to the type of data which will be available in future multiband time-domain surveys.

The upper panels of Figure 8 show an example light curve from this reduced dataset, along with the supersmoother and multiband periodograms derived from this data. Analogously to the lower panels of Figure 5, the single-band supersmoother model loses the true period within the noise, while the shared-phase multiband model still shows prominent signal near the S10 period.

The lower panels of Figure 8 show the relationship between the S10 periods (based on the full dataset) and the periods derived with each model from this reduced dataset, and these results are summarized in the lower rows of Table 1. It is clear that the supersmoother model

and in this work. Objects showing this discrepancy are those with very low signal-to-noise.

is simply over-fitting noise with this few data points: the top period matches S10 in only 23% of cases (compared to 87% with the full dataset), and the top 5 periods contain the S10 period only 45% of the time. The failure mode is much less predictable as well: rather than being clustered near aliases, most of the period determinations are scattered seemingly randomly around the parameter space.

While the multiband method performed comparably to the S10 method on dense data, it far outperforms S10 on the sparse dataset. Even with an 80% reduction in the number of observations, the multiband method matches the S10 period 64% of the time (compared to 79% with the full dataset), and the top 5 peaks contain the S10 period 94% of the time (compared to 99% with the full dataset). This performance is due to the fact that the multiband algorithm has relatively few parameters, but is yet able to flexibly accommodate noisy data from multiple observing bands. In particular, this suggests that with the multiterm periodogram, the S10 analysis could have been done effectively with only a small fraction of the available data. This bodes well for future surveys, where data on variable stars will be much more sparsely sampled.

## 6.3. Potential Improvements to the Multiband Method

A well-known (though often unrecognized) difficulty of Lomb-Scargle-type periodograms on unevenly-sampled data is that they do not measure the power of the signal in question, but the power of the signal *convolved with the observing with the survey window function*. For regularly-sampled timeseries, this convolution is the source of the perfect aliasing beyond the Nyquist sampling limit; for non-regular sampling, this aliasing generally happens to some degree at *all* frequencies! Because of this, even a signal with a single well-defined period will result in a Lomb-Scargle periodogram with multiple maxima at locations which depend on both the underlying signal and the precise observing window.

The multiband periodogram, as a generalization of Lomb-Scargle, shares this difficulty: it tends to respond to frequency structure in the window function as well as frequency structure in the data. This can be viewed as a result of the very model simplicity which causes its success in the case of sparse multiband data: it cannot disentangle bias in the model from bias due to features in the survey window.

This could potentially be accounted for by correcting for the effect of the estimated window function; one potential method for this involves estimating the deconvolution of the window power and the observed power (Roberts et al. 1987). It may also be possible to propose a multiband extension of, e.g., CARMA (Kelly et al. 2014) or another forward-modeling approach to detecting periodicity.

Another potentially fruitful avenue of research which we do not study here is the adjustment of the regularization terms in the model, and the application of other types of regularization to the higher-order periodogram. In particular, L1 regularization (also known as Lasso regression) could lead to interesting results: L1 regularization is similar in spirit to the Tikhonov regularization discussed in Section 4.3, but tends toward sparsity in the model parameters (see, e.g. Ivezić et al. 2014, for

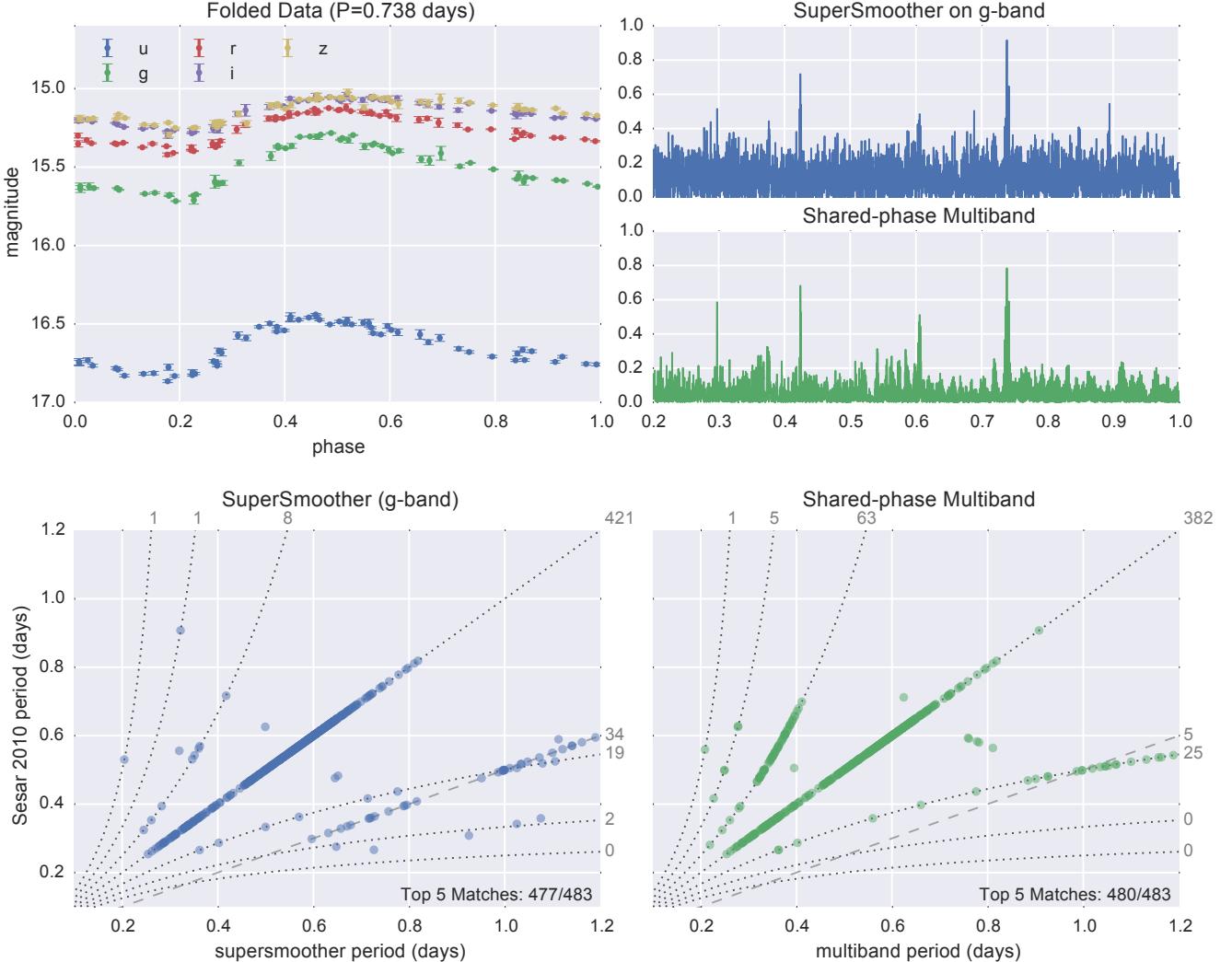


FIG. 7.— Comparison of the Multiband algorithm and single-band supersmoother algorithm on 483 well-sampled RR Lyrae light curves from Stripe 82. The upper panels show a representative lightcurve and periodogram fits, while the bottom panels compare the derived periods to the template-based periods reported in S10. Shown for reference are the beat aliases (dotted lines) and the first harmonic alias (dashed lines); numbers along the top and right edges of the panels indicate the number of points aligned with each trend. The single-band supersmoother model tends to err toward harmonic aliases, while the multiband model tends to err toward beat frequency aliases. Both methods find the correct period among the top 5 significant peaks around 99% of the time. This suggests that for densely-sampled multiband surveys, the multiband periodogram will match the results of standard methods (but see Figure 8).

TABLE 1  
PERIOD DETERMINATION FROM DENSE AND SPARSE DATA (483 TOTAL)

Data	Method	Match among top 5	Top peak matches	Beat Aliases	Harmonic Aliases
Dense data (Figure 7)	g-band Supersmoother	477 (98.8%)	421 (87.2%)	31	34
	Multi-band Periodogram	480 (99.4%)	382 (79.1%)	94	5
Sparse data (Figure 8)	g-band Supersmoother	219 (45.3%)	113 (23.4%)	101	4
	Multi-band Periodogram	449 (93.0%)	308 (63.8%)	136	7

a discussion). Such an approach could provide a useful tradeoff between model complexity and bias in the case of higher-order truncated Fourier models, though comes at a higher computational cost.

Another potentially interesting extension of the multi-band case would be to define and make use of physically-motivated priors in the light-curve shape. This approach could allow the model bias to be decreased without a commensurate increase in model complexity, which is

what causes poor performance in the case of sparsely-sampled noisy data. As an example of such a physically-motivated prior, consider that the paths of RR Lyrae stars through color-color and color-magnitude space are constrained by known astrophysical processes in the structure of the stars (e.g., see Fig. 5 in Szabó et al. 2014). Making use of this information could help break degeneracies in period determination with higher-order models.

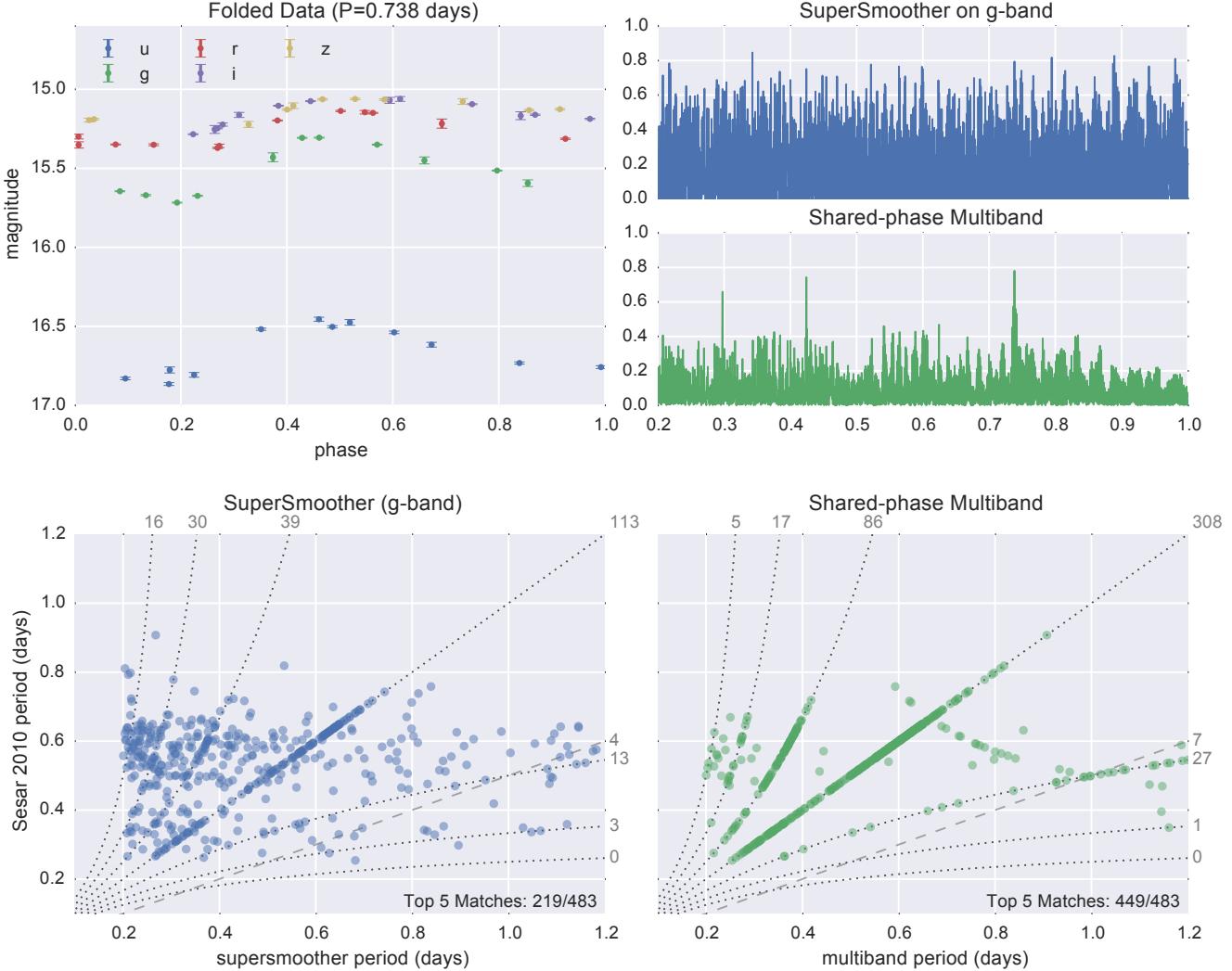


FIG. 8.— This figure repeats the experiment shown in Figure 7 (see caption there for description), but the data is artificially reduced to only a single-band observation on each evening, a situation reflective of the observing cadence of future large-scale surveys. In this case, the single-band SuperSmoker strategy used as a first pass in S10 fails: there is simply not enough data in each band to recover an accurate period estimate. The correct period is among the top 5 candidates in fewer than 50% of cases. The shared-phase multiband approach utilizes information from all five bands, and returns much more robust results: even with the greatly-reduced data, the true period is among the top 5 candidates in 93% of cases. This suggests that for sparsely-sampled multiband survey data (such as that expected from LSST) the multiband periodogram will produce superior results when compared to standard methods – see Figure 9.

## 7. PROSPECTS FOR MULTIBAND PERIODOGRAMS WITH LSST

Previously, Oluseyi et al. (2012) evaluated the prospects of period finding in early LSST data, and found results which were not encouraging. Using the conservative criterion of a 2/3 majority among the top single-band supersmoother periods in the  $g$ ,  $r$ , and  $i$  bands, they showed that, depending on spectral type, finding reliable periods for the brightest ( $g \sim 20$ ) RR Lyrae stars will require several years of LSST data, while periods for some of the faintest ( $g \sim 25$ ) stars will not be reliable with even ten years of data!

One potential remedy is to move away from general models like supersmooth and lomb-scargle to specific template-fitting methods such as those used in S10. Indeed, such methods perform well even for sparsely-sampled multiband data such as those from the PanSTARRS survey; the primary drawback is that such

blind template fits are computationally extremely expensive: they involve nonlinear optimizations over each of several hundred candidate templates at each of tens of thousands of candidate frequencies (B. Sesar, private communication). Thus the template-fitting method, though it can produce accurate periods, in practice requires several hours of CPU time for a well-sampled period grid for a single source (compared to several seconds for the multiband periodogram proposed here). Note that several hours per object is orders-of-magnitude too slow in the case of LSST; to estimate periods for a billion stars on a 1000-core machine in a year requires a compute-time budget of only 30 seconds per light curve.

Because of the computational expense of the pure template-fitting method, when working with SDSS II data S10 performed a first-pass with a single-band supersmooth to establish candidate periods, which were in turn evaluated with template-fitting approach. Here

we show that such a hybrid strategy combining the multiband periodogram and the S10 template fits will be useful for determining periodicity of variables in early LSST data releases, greatly improving on the outlook presented in Oluseyi et al. (2012).

We suggest the following procedure for determining periods in future multiband datasets:

1. As a first pass, find a set of candidate frequencies using the multiband periodogram. This is a fast linear optimization that can be straightforwardly parallelized.
2. Within these candidate frequencies, use the more costly template-fitting procedure to choose the optimal period from among the handful of candidates.
3. Compute a goodness-of-fit statistic for the best-fit template to determine whether the fit is suitable; if not, then apply the template-fitting procedure across the full period range.

Here we briefly explore simulated LSST observations of RR Lyrae stars in order to gauge the effectiveness of the first step in this strategy; the effectiveness of the template-fitting step will be explored further in future work. Rather than doing the full analysis including the final template fits, we will focus on the ability of the multiband periodogram to quickly provide suitable candidate periods under the assumption that the S10 template algorithm will then select or reject the optimal period from this set.

### 7.1. LSST Simulations

We use a simulated LSST cadence (Delgado et al. 2006; Ridgway et al. 2012; Jones et al. 2014) in 25 arbitrarily chosen fields that are representative of the anticipated main survey temporal coverage. We simulate a set of 50 RR Lyrae observations with the S10 templates, with a range of apparent magnitudes between  $g = 20$  and  $g = 24.5$ , corresponding to bright-to-faint range of LSST main-survey observations, and with expected photometric errors computed using eqs. 4–6 from Ivezić et al. (2008). Given the capability of template-fitting to choose among candidate periods, we use a more relaxed period-matching criterion than in Oluseyi et al. (2012): when evaluating the single-band supersmoother, we require that the true period is among the five periods determined independently in the  $u, g, r, i, z$  bands; in the multiband case we require that the true period is among the top five peaks in the multiband periodogram.

Figure 9 shows the fraction of stars where this period matching criterion is met as a function of  $g$ -band magnitude and subset of LSST data. The solid lines show the multiband results; the dashed lines show the single-band supersmoother results; and the shading helps guide the eye for the sake of comparison. Because of our relaxed matching criteria, even the single-band supersmoother results here are much more optimistic than the Oluseyi et al. (2012) results (compare to Figure 15 in that work): the supersmoother result here can be considered representative of a best-case scenario for *ad hoc* single-band fits. Without fail, the multiband result exceeds this best-case single-band result; the improvement is most apparent for faint stars, where the greater model flexibility of the supersmoother causes it to over-fit the noisy data.

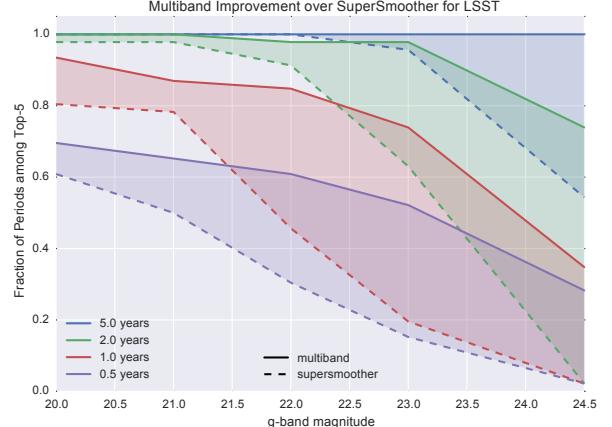


FIG. 9.— Fraction of periods correctly determined for LSST RR Lyrae as a function of the length of the observing season and the mean  $g$ -band magnitude, for the multiband periodogram approach (method of this work; solid lines) and single-band supersmoother approach (method of Oluseyi et al. 2012, dashed lines). The multiband method is superior to the single-band supersmoother approach in all cases, and especially for the faintest objects.

The performance of the multiband periodogram points to much more promising prospects for science with variable stars than previously reported. In particular, even with only six months of LSST data, we can expect to correctly identify the periods for over 60% of stars brighter than  $g = 22$ ; with the first two years of LSST observations, this increases to nearly 100%; with five years of data, the multiband method identifies the correct period for 100% of even the faintest stars. Part of this improvement is due to the performance of the shared-phase multiband model with noisy data, and part of this improvement is due to the relaxed period-matching constraints enabled by the hybrid approach of periodogram-based and template-based period determination.

## 8. DISCUSSION AND CONCLUSION

We have motivated and derived a multiband version of the classic Lomb-Scargle method for detecting periodicity in astronomical time-series. Experiments on several hundred RR Lyrae stars from the SDSS Stripe 82 dataset indicate that this method outperforms methods used previously in the literature, especially for sparsely-sampled light curves with only single bands observed each night. While there are potential areas of improvement involving corrections to window function artifacts and accounting for physically-motivated priors, the straightforward multiband model outperforms previous *ad hoc* approaches to multiband data.

Looking forward to future variable star catalogs from PanSTARRS, DES, and LSST, there are two important constraints that any analysis method must meet: the methods must be able to cope with heterogeneous and noisy observations through multiple band-passes, and the methods must be fast enough to be computable on millions or even billions of objects. The multiband method, through its combination of flexibility and model simplicity, meets the first constraint: as shown above, in the case of sparsely-sampled noisy multiband data, it outperforms previous approaches to period determination. It also meets the second constraint: it requires the solution of a simple linear model at each frequency, com-

pared to a rank-based sliding-window model in the case of supersmooth, a nonlinear optimization in the case of template-fitting, and a Markov Chain Monte Carlo analysis in the case of CARMA models. In our own benchmarks, we found the multiband method to be several times faster than the single-band supersmooth approach, and several orders of magnitude faster than the template fitting approach.

The strengths and weaknesses of the multiband method suggest a hybrid approach to finding periodicity in sparsely-sampled multiband data: a first pass with the fast multiband method, followed by a second pass using the more computationally intensive template-fitting method to select among these candidate periods. Despite pessimism in previous studies, our experiments with simulated LSST data indicate that such a hybrid approach will successfully identify periods in the majority of RR Lyrae stars brighter than  $g \sim 22.5$  in the first months of the survey, and the majority of the faintest detected stars with several years of data. This finding suggests that the

multiband periodogram could have an important role to play in the analysis of variable stars in future multiband surveys.

We have released a Python implementation of the multiband periodogram on GitHub, along with Python code to reproduce all results and figures in this work; this is described in Appendix A. As we were finalizing this manuscript, we were made aware of a preprint of an independent exploration of a similar approach to multiband light curves (Long et al. 2014); we discuss the similarities and differences between these two approaches in Appendix B.

*Acknowledgments:* JTV is supported by the University of Washington eScience institute, including grants from the Alfred P. Sloan Foundation, the Gordon and Betty Moore Foundation, and the Washington Research Foundation. The authors thank GitHub for providing free academic accounts which were essential in the development of this work.

## REFERENCES

- Brandt, S. 1970, Statistical and computational methods in data analysis
- Bretthorst, G. 1988, Bayesian Spectrum Analysis and Parameter Estimation, Lecture Notes in Statistics (Springer)
- Buitinck, L., Louppe, G., Blondel, M., et al. 2013, CoRR, abs/1309.0238
- Cumming, A., Marcy, G. W., & Butler, R. P. 1999, ApJ, 526, 890
- Delgado, F., Cook, K., Miller, M., Allsman, R., & Pierfederici, F. 2006, in Society of Photo-Optical Instrumentation Engineers (SPIE) Conference Series, Vol. 6270, Society of Photo-Optical Instrumentation Engineers (SPIE) Conference Series, 1
- Eyer, L., & Mowlavi, N. 2008, Journal of Physics Conference Series, 118, 012010
- Ferraz-Mello, S. 1981, AJ, 86, 619
- Flaugher, B. 2008, in A Decade of Dark Energy: Spring Symposium, Proceedings of the conferences held May 5–8, 2008 in Baltimore, Maryland. (USA). Edited by Norbert Pirzkal and Henry Ferguson. <http://www.stsci.edu/institute/conference/spring2008>
- Graham, M. J., Drake, A. J., Djorgovski, S. G., et al. 2013, MNRAS, 434, 3423
- Hoerl, A. E., & Kennard, R. W. 1970, Technometrics, 12, 55
- Ivezić, Ž., Connolly, A., VanderPlas, J., & Gray, A. 2014, Statistics, Data Mining, and Machine Learning in Astronomy: A Practical Python Guide for the Analysis of Survey Data, Princeton Series in Modern Observational Astronomy (Princeton University Press)
- Ivezić, Ž., et al. 2008, arXiv:0805.2366
- Jaynes, E. 1987, in Fundamental Theories of Physics, Vol. 21, Maximum-Entropy and Bayesian Spectral Analysis and Estimation Problems, ed. C. Smith & G. Erickson (Springer Netherlands), 1–37
- Jones, R. L., Yoachim, P., Chandrasekharan, S., et al. 2014, in Society of Photo-Optical Instrumentation Engineers (SPIE) Conference Series, Vol. 9149, Society of Photo-Optical Instrumentation Engineers (SPIE) Conference Series, 0
- Kaiser, N., Burgett, W., Chambers, K., et al. 2010, in Society of Photo-Optical Instrumentation Engineers (SPIE) Conference Series, Vol. 7733, Society of Photo-Optical Instrumentation Engineers (SPIE) Conference Series
- Kelly, B. C., Becker, A. C., Sobolewska, M., Siemiginowska, A., & Uttley, P. 2014, ApJ, 788, 33
- Lomb, N. R. 1976, Ap&SS, 39, 447
- Long, J. P., Chi, E. C., & Baraniuk, R. G. 2014, arXiv:1412.6520
- Oluseyi, H. M., Becker, A. C., Culliton, C., et al. 2012, AJ, 144, 9
- Pedregosa, F., Varoquaux, G., Gramfort, A., et al. 2011, Journal of Machine Learning Research, 12, 2825
- Press, W. H., & Rybicki, G. B. 1989, ApJ, 338, 277
- Press, W. H., Teukolsky, S. A., Vetterling, W. T., & Flannery, B. P. 2007, Numerical Recipes 3rd Edition: The Art of Scientific Computing, 3rd edn. (New York, NY, USA: Cambridge University Press)
- Reimann, J. D. 1994, PhD thesis, University of California, Berkeley
- Ridgway, S. T., Chandrasekharan, S., Cook, K. H., et al. 2012, in Society of Photo-Optical Instrumentation Engineers (SPIE) Conference Series, Vol. 8448, Society of Photo-Optical Instrumentation Engineers (SPIE) Conference Series, 10
- Roberts, D. H., Lehar, J., & Dreher, J. W. 1987, AJ, 93, 968
- Scargle, J. D. 1982, ApJ, 263, 835
- Schuster, A. 1898, Terrestrial Magnetism, 3, 13
- Sesar, B., Stuart, J. S., Ivezić, Ž., et al. 2011, AJ, 142, 190
- Sesar, B., Ivezić, Ž., Lupton, R. H., et al. 2007, AJ, 134, 2236
- Sesar, B., Ivezić, Ž., Grammer, S. H., et al. 2010, ApJ, 708, 717
- Stellingwerf, R. F. 1978, ApJ, 224, 953
- Süveges, M., Sesar, B., Váradí, M., et al. 2012, MNRAS, 424, 2528
- Szabó, R., Ivezić, Ž., Kiss, L. L., et al. 2014, ApJ, 780, 92
- Tikhonov, A. 1963, in Soviet Math. Dokl., Vol. 5, 1035–1038
- van der Walt, S., Colbert, S. C., & Varoquaux, G. 2011, Computing in Science and Engg., 13, 22
- Vanderplas, J. 2015a, gatspy: General Tools for Astronomical Time Series in Python, doi:10.5281/zenodo.14833
- . 2015b, doi:10.5281/zenodo.14475
- Vanderplas, J., Connolly, A., Ivezić, Ž., & Gray, A. 2012, in Conference on Intelligent Data Understanding (CIDU), 47–54
- Welch, D. L., & Stetson, P. B. 1993, AJ, 105, 1813
- Zechmeister, M., & Kürster, M. 2009, A&A, 496, 577

## APPENDIX PYTHON IMPLEMENTATION OF MULTIBAND PERIODOGRAM

The algorithm outlined in this paper is available in `gatspy`, an open-source Python package for general astronomical time-series analysis<sup>7</sup> (Vanderplas 2015a). Along with the periodogram implementation, it also contains code to

<sup>7</sup> <http://github.com/astroml/gatspy/>

download all the data used in this work. Code to reproduce this paper, including all figures, is available in a separate repository<sup>8</sup>.

`gatspy` is a pure-Python package written to be compatible with both Python 2 and Python 3, and performs fast numerical computation through dependencies on `numpy` (van der Walt et al. 2011)<sup>9</sup> and `astroML` (Vanderplas et al. 2012)<sup>10</sup>, which offer optimized implementations of numerical methods in Python.

The API for the module is largely influenced by that of the `scikit-learn` package (Pedregosa et al. 2011; Buitinck et al. 2013)<sup>11</sup>, in which models are Python class objects which can be fit to data with the `fit()` method. Here is a basic example of how you can use `multiband_LS` to download the data used in this paper, fit a multiband model to the data, and compute the power at a few periods:

```
from gatspy.periodic import LombScargleMultiband
import numpy as np

# Fetch the Sesar 2010 RR Lyrae data
from gatspy.datasets import fetch_rrlyrae
data = fetch_rrlyrae()
t, mag, dmag, filts = data.get_lightcurve(data.ids[0])

# Construct the multiband model
model = LombScargleMultiband(Nterms_base=0, Nterms_band=1)
model.fit(t, mag, dmag, filts)

# Compute power at the following periods
periods = np.linspace(0.2, 1.4, 1000) # periods in days
power = model.periodogram(periods)
```

Other models are available as well. For example, here is how you can compute the periodogram under the supersmoother model; this implementation of the supersmoother periodogram makes use of the `supersmoother` Python package (Vanderplas 2015b).

```
from gatspy.periodic import SuperSmoother

# Construct the supersmoother model
model = SuperSmoother()
gband = (filts == 'g')
model.fit(t[gband], mag[gband], dmag[gband])

# Compute power at the given periods
power = model.periodogram(periods)
```

The models in the `gatspy` package contain many more methods, and much more functionality than what is shown here. For updates, more examples, and more information, visit <http://github.com/astroml/gatspy/>.

#### COMPARISON WITH LONG (2014)

As we were finishing this study, we learned that another group had released a preprint independently addressing the multiband periodogram case, and come up with a solution very similar to the one presented here (Long et al. 2014, hereafter LCB14). They present two methods, the “Multiband Generalized Lomb-Scargle” (MGLS) which is effectively identical to the  $(1, 0)$  multi-phase model here, and the “Penalized Generalized Lomb-Scargle” (PGLS), which is similar in spirit to our  $(0, 1)$  shared-phase model.

In the PGLS model, they start with a multi-phase model, fitting independent  $N = 1$  term fits to each band, and apply a nonlinear regularization term which penalizes differences in the amplitude and phase. In terms of the formalism used in this work, the PGLS model minimizes a regularized  $\chi^2$  of the form

$$\chi_{PGLS}^2 = \sum_{k=1}^K \left[ \chi_{GLS}^2(D^{(k)}) + J_A(A^{(k)}) + J_\phi(\phi^{(k)}) \right]. \quad (\text{B1})$$

where  $K$  is the number of bands,  $\chi_{GLS}^2(D^{(k)})$  is the  $\chi^2$  of the standard floating mean model on the single-band data  $D^{(k)}$ , and  $J_A$  and  $J_\phi$  are regularization/penalty terms which are a function of the amplitude  $A^{(k)}$  and phase  $\phi^{(k)}$  of each model. In terms of our linear model parameters  $\theta^{(k)}$ , this amplitude and phase can be expressed:

$$\begin{aligned} A^{(k)} &= \sqrt{(\theta_1^{(k)})^2 + (\theta_2^{(k)})^2} \\ \phi^{(k)} &= \arctan(\theta_2^{(k)} / \theta_1^{(k)}) \end{aligned} \quad (\text{B2})$$

The selected form of these regularization terms penalizes deviations of the amplitude and phase from a common mean between the bands; in this sense the PGLS model can be considered a conceptual mid-point between our shared-phase and multi-phase models. Within the formalism proposed in the current work, such a mid-point may be alternatively

<sup>8</sup> [http://github.com/jakevdp/multiband\\_LS/](http://github.com/jakevdp/multiband_LS/)

<sup>9</sup> <http://www.numpy.org>

<sup>10</sup> <http://www.astroml.org>

<sup>11</sup> <http://scikit-learn.org>

attained by suitably increasing the regularization parameter  $\lambda$  used in our shared-phase model, though the precise nature of the resulting regularization will differ.

Computationally, the PGGLS model requires a nonlinear optimization at each frequency  $\omega$ , and is thus much more expensive than the straightforward linear optimization of our shared-phase model. For this reason, LCB14 proposes a clever method by which nested models are used to reduce the number of nonlinear optimizations used: essentially, by showing that the (linear) MGGLS  $\chi^2$  is a lower-bound of the (non-linear) PGGLS  $\chi^2$ , it is possible to iteratively reduce the number of PGGLS computations required to minimize the  $\chi^2$  among a grid of frequencies. Such an optimization could also be applied in the case of our shared-phase model, but is not necessary here due to its already high speed. Nevertheless, when applying the method to a very large number of light curves, as in e.g. LSST, such a computational trick may prove very useful.

Given these important distinctions between the models proposed here and in LCB14, in future work we plan to do a detailed comparison of the two approaches to multiband model regularization.