



Technical Report

RobPer: An R Package to Calculate Periodograms for Light Curves Based On Robust Regression

Anita Monika Thieler, Roland
Fried, Jonathan Rathjens

02/2013



Part of the work on this technical report has been supported by Deutsche Forschungsgemeinschaft (DFG) within the Collaborative Research Center SFB 876 "Providing Information by Resource-Constrained Analysis", project C3 and within the Graduate School GK 1032 "Statistische Modellbildung". We thank the ITMC at TU Dortmund University for providing computer resources on LiDO.

Speaker: Prof. Dr. Katharina Morik
Address: TU Dortmund University
Joseph-von-Fraunhofer-Str. 23
D-44227 Dortmund
Web: <http://sfb876.tu-dortmund.de>

Abstract

An important task in astroparticle physics is the detection of periodicities in irregularly sampled time series, called light curves. The classic Fourier periodogram cannot deal with irregular sampling and with the measurement accuracies that are typically given for each observation of a light curve. Hence, methods to fit periodic functions using weighted regression were developed in the past to calculate periodograms.

We present the R Package **RobPer** which allows to combine different periodic functions and regression techniques to calculate periodograms. Possible regression techniques are least squares, least absolute deviation, least trimmed, M-, S- and τ -regression. Measurement accuracies can be taken into account including weights. Our periodogram function covers most of the attempts that have been tried earlier and provides new model-regression-combinations that have not been used before.

To detect valid periods, we apply an outlier search on the periodogram instead of using fixed critical values that are theoretically only justified in case of least squares regression, independent periodogram bars and a null hypothesis allowing only normal white noise. This outlier search can be performed using RobPer as well.

Finally, the package also includes a generator to generate artificial light curves e.g., for simulation studies.

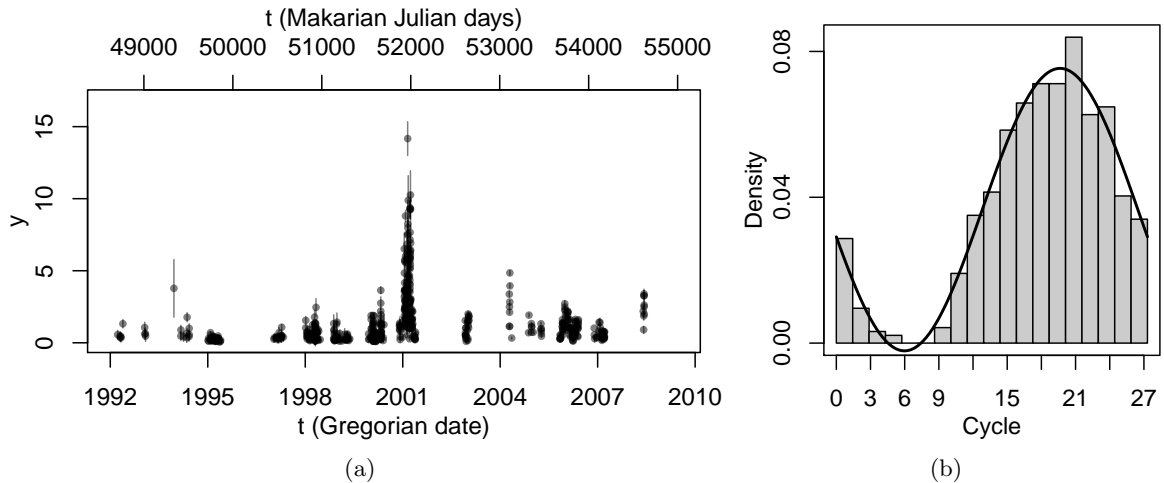


Figure 1: Light curve with γ -ray emissions for the very high energy gamma-ray source Mrk 421 (see [Tluczykont *et al.* 2010](#), and references therein). Panel a shows the light curve, vertical lines at each point show the reported measurement accuracies. Panel b shows a histogram of the observation times t_i modulo $p_s = 27.31$. A sine represents the shape rather well.

1. Introduction

We introduce the R package **RobPer**, which can be used to calculate periodograms and detect periodicities in irregularly sampled time series. Our special objective are light curves, which occur in astroparticle physics and are irregularly sampled times series $(t_i, y_i, s_i)_{i=1, \dots, n}$ consisting of unequally spaced observation times t_1, \dots, t_n , observed values y_1, \dots, y_n and measurement accuracies s_1, \dots, s_n . The measurement accuracies s_i give information about how precise the y_i were measured. They can be interpreted as estimates for the standard deviations of the observed values. The observed values possibly contain a periodic fluctuation y_f with fluctuation period p_f and the observation times t_i are realisations of irregularly distributed random variables with a periodically shaped density.

Such periodicity in the pattern of the observation times is a typical phenomenon, as the sampling of astroparticle physics' time series is influenced among others by astronomical constellations. For example, plotting a histogram of the observation times for the gamma particle source Mrk 421 modulo $p_s = 27.31$ shows an unequal distribution over a cycle of this length (see Figure 1). This is due to the fact that observations cannot be sampled in presence of full moon and the moon period is similar to p_s .

So for $i = 1, \dots, n$, we assume the following model:

$$T_i = T_{(i)}^*, \quad T_1^*, \dots, T_n^* \sim \mathcal{D}(p_s), \quad (1)$$

$$Y_i = Y_{f;i} + Y_{w;i}, \quad (2)$$

$$Y_{f;i} = f\left(\frac{T_i}{p_f}\right), \quad f(\xi) = f(\xi + 1) \forall \xi \in \mathbb{R} \quad (3)$$

$$Y_{w;i} \sim \mathcal{N}(0, \sigma_i^2), \quad (4)$$

s_i : given estimate for σ_i independent from Y_1, \dots, Y_n ,

where $T_{(i)}^*$ denotes the i -th ordered observation time in T_1^*, \dots, T_n^* and $\mathcal{D}(p_s)$ is a periodic

sampling density with period p_s . The observation times t_1, \dots, t_n and the observed values y_1, \dots, y_n are realisations of T_1, \dots, T_n and Y_1, \dots, Y_n , respectively. We assume the observation times to be measured error-free.

To detect a periodic fluctuation with period p_f in the observed values y_i , it is not possible to use the standard periodogram of Fourier analysis. This method can only be applied to time series with equidistant observation times, while light curves are typically irregularly sampled. A setting-adapted procedure, the Deeming periodogram (Deeming 1975), is not recommendable either in this case, because it is known to react to periodicity p_s in the sampling (see Hall and Li 2006).

In order to determine periodicity in light curves, hence other methods than the classical Fourier periodogram or the Deeming periodogram should be used. Popular periodogram methods in astroparticle physics are for example the Lomb-Scargle periodogram (Scargle 1982) or the Phase Dispersion Minimization periodogram (Stellingwerf 1978). Both of them and many other approaches can be generalized to fitting periodic functions to the light curve using least squares regression and calculating periodogram bars based on SE and SY, where SE is the remaining variance in the residuals of the fit and SY is the overall variance in the observed values y_i . An even broader class of periodogram methods allows also robust regression instead of least squares regression and weighted regression in order to take into account the measurement accuracies s_i .

The function `RobPer` in our homonymous R package calculates a periodogram of a light curve based on fitting periodic functions to $(t_i, y_i)_{i=1, \dots, n}$ using least squares or a robust regression technique, optionally taking measurement accuracies s_i into account using weighted regression. The coefficient of determination corresponding to the objective function of the regression technique is used as periodogram bar. This proceeding incorporates analogues to most of the existing periodograms and introduces several new techniques. Preliminary implementations of most of these periodogram methods have been compared by Thieler, Backes, Fried, and Rhode (2013). Here, we explain the usage of `RobPer`, which makes improved and extended methods for period detection publicly available.

This article is organized as follows: In Section 2, the usage and the structure of the function `RobPer` are explained. Especially, the different periodic functions and regression techniques are discussed and related to the existing periodogram methods. Diagrams which show how this R function is implemented in detail are displayed in Appendix A. Section 3 is devoted to the question how to find valid periods using a periodogram. Thieler *et al.* (2013) propose robustly fitting a Beta distribution combined with outlier detection. The function `betaCvMfit` in the package `RobPer` performs this. In Section 4, the function `tsgen` is presented to generate artificial light curves. Some examples for how to use the package are given in Section 5. Section 6 concludes with a summary. The `RobPer` software package is available from the Comprehensive R (R Core Team 2013) Archive Network at <http://CRAN.R-project.org/package=RobPer>.

2. Calculate periodograms with RobPer

The R function `RobPer` calculates a periodogram of a given light curve $(t_i, y_i, s_i)_{i=1, \dots, n}$. This is done by fitting a function $g\left(\frac{t}{p_j}\right)$ to $(t_i, y_i)_{i=1, \dots, n}$, with g being periodic with period 1, for each given trial period $(p_j)_{j=1, \dots, q}$. The periodogram bars are defined as the coefficients of determination of the respective fits. Using weighted regression with weights $1/s_i$ makes it possible to take into account the measurement accuracies. As the shape of the true fluctuation f (Equation 3) is usually unknown, we will typically have $g \neq f$.

Table 1 gives an overview over all input variables for `RobPer`. The possible shapes of the function g that may be fitted using `RobPer` are presented in Section 2.1. Fitting them using least squares regression leads in many cases to equivalents of already existing periodogram methods (see Table 2 or Thieler *et al.* 2013 for a more detailed discussion).

In addition to least squares regression, `RobPer` offers a selection of robust regression techniques to fit $g\left(\frac{t}{p_j}\right)$. They are presented in Section 2.2. All regression techniques implemented in `RobPer` base on minimizing an objective value

$$\text{SE} = \zeta(y - X\beta) \quad (5)$$

subject to β , where X presents the design matrix of a linear presentation of $g\left(\frac{t}{p}\right)$ with p being a trial period. Using the same regression technique, the location μ of the observations y_1, \dots, y_n can be estimated minimizing

$$\text{SY} = \zeta(y - \mu\mathbf{i}) \quad (6)$$

with $\mathbf{i} = \mathbf{1}_n$ being an n -variate vector of ones in case of unweighted regression. The periodogram bar may be calculated as $R^2 = 1 - \frac{\text{SE}}{\text{SY}}$. This definition for the coefficient of determination does not only apply for least squares regression, but also for L_1 - and M-regression in general (see Maronna, Martin, and Yohai 2006, p. 171) as well as for S-, LTS- and τ -regression (see Croux and Dehon 2003).

Table 2 displays periodogram methods following the principle of fitting periodic functions. Up to now, weighted regression or robust regression in affiliation with periodic step functions has only been performed by Thieler *et al.* (2013), though the unweighted least squares versions belong to the most popular periodogram methods in this area of research. S- or τ -regression, which are also available in `RobPer`, have not been proposed at all up to now in this context.

2.1. Periodic function fitted: Input variable model

For each trial period p_i , $i \in \{1, \dots, q\}$ (given by the argument `periods`, see Table 1), a periodic function (defined by `model`) is fitted to the light curve (using regression technique `regression`). Implemented periodic functions include step functions, sine functions, Fourier series and spline functions.

Step functions

Many periodogram methods from astroparticle physics such as Epoch Folding Periodogram (Leahy *et al.* 1983) or the Analysis of Variance Periodogram (Schwarzenberg-Czerny 1989) can be interpreted as fitting a step function to a light curve (see Schwarzenberg-Czerny 1998

Input	Comment
$\mathbf{ts} \in \mathbb{R}^{n \times 3}$ or $\mathbb{R}^{n \times 2}$	Light curve (t_i, y_i, s_i) or (t_i, y_i) , $i = 1, \dots, n$; If weighting = FALSE the measurement accuracies s_i column may be omitted.
weighting $\in \{\mathbf{T}, \mathbf{F}\}$	If s_i should be taken into account performing weighted regression
periods $\in \mathbb{R}_{>0}^q$	Trial periods p_1, \dots, p_q
regression	Regression technique (see Section 2.2), possible choices: "L2", "L1", "LTS", "S", "huber", "bisquare", "tau"
model	Periodic fluctuation to be fitted (see Section 2.1), possible choices: "step", "2step", "sine", "fourier(2)", "fourier(3)", "splines"
steps $\in \mathbb{N}$	number of steps per cycle for periodic step functions. Default: 10
var1 $\in \{\mathbf{T}, \mathbf{F}\}$	TRUE sets variance estimate to one for weighted M-regression. Default: weighting
tol $\in \mathbb{R}_{>0}$	Precision for convergence criteria Used in case of M-regression and in case of LTS-regression if LTSopt=TRUE . Default: 10^{-3}
genoudcontrol $\in \mathbb{N}^3$	Settings for genoud (see paragraph about LTS-regression in Section 2.2): max.generations , wait.generations , pop.size Used if regression = "bisquare" or LTSopt = TRUE & regression = "LTS". Default: {50,5,50}
LTSopt $\in \{\mathbf{T}, \mathbf{F}\}$	If regression result of ltsReg should be optimized. Default: TRUE if regression = "LTS"
taucontrol $\in \mathbb{N}^4 \times \{\mathbf{T}, \mathbf{F}\}$	Settings for τ -regression: N , kk , tt , rr , approximate Used if regression = "tau", rr only necessary for approximate = TRUE . Default: {100, 2, 5, 2, FALSE }
Scontrol $\in \mathbb{N}^3 \times \mathbb{R}_{>0}^2 \times \mathbb{N}$	Settings for S-regression: N , kk , tt , b , cc , seed Used in case of regression = "S". seed can be fixed in order to get reproducible results or can be left empty. Default: { N , 2, 5, 0.5, 1.547, NULL } with N=50 if weighting=FALSE and N=200 if weighting=TRUE
<hr/>	
Output	
periodogram $\in \mathbb{R}^q$	Vector of periodogram bars belonging to the trial periods
<hr/>	
Possibly warnings	

Table 1: In- and output of the function **RobPer**. {T, F} means {TRUE, FALSE}.

Model	Regression technique	Publication (Name of the method)
<code>step</code>	L2	Leahy <i>et al.</i> (1983) (Epoch Folding)
	L2	Schwarzenberg-Czerny (1989) (Analysis of Variance)
	L2, L1, huber,bisquare	Thieler <i>et al.</i> (2013)
<code>2step</code>	L2	Stellingwerf (1978) (Phase Dispersion Minimization)
	L2, L1, huber,bisquare	Thieler <i>et al.</i> (2013)
<code>sine</code>	L2	Scargle (1982) (Lomb-Scargle)
	L2	Zechmeister and Kürster (2009) (Generalized Lomb-Scargle)
	L2	Cumming <i>et al.</i> (1999) (Floating Mean)
	L2	Ferraz-Mello (1981) (Date Compensated Fourier Transform*)
	L2	Reegen (2007) (SigSpec*)
	L1	Li (2009)* , Li (2010)*
	LTS	Ahdesmäki <i>et al.</i> (2007)*
	bisquare	Ahdesmäki <i>et al.</i> (2007)*
	huber	Zhang and Chan (2005)*
<code>fourier(2), fourier(3)</code>	L2	Hall <i>et al.</i> (2000)
	L2	Palmer (2009) (Fast- χ^2)
	L2, L1, huber,bisquare	Thieler <i>et al.</i> (2013)
<code>splines</code>	L2	Akerlof <i>et al.</i> (1994)
	L2	Hall <i>et al.</i> (2000)
	L2	Oh <i>et al.</i> (2004) (Generalized Cross Validation)
	huber	Oh <i>et al.</i> (2004) (Robust Cross Validation)
	L2, L1, huber,bisquare	Thieler <i>et al.</i> (2013)

Table 2: Published periodogram methods that rely on fitting a periodic model g to a light curve using a regression technique. Models (see Section 2.1): periodic step functions and pairwise overlapping step functions (`step` and `2step`), the sine function (`sine`), Fourier series of second and third degree and periodic spline functions (`fourier(2)`, `fourier(3)` and `splines`). Regression techniques: See Table 3 for labels. The underlined methods can take into account measurement accuracies using weighted regression. The periodogram bars of methods marked by * do not base on SE or SY, but on the parameter vector of the function fitted (e.g., squared amplitude).

or [Thieler et al. 2013](#)). They use periodogram bars related to R^2 , n and the numbers of steps per cycle.

Another typical periodogram method in astroparticle physics is the Phase Dispersion Minimization Periodogram (PDM, [Stellingwerf 1978](#)). Depending on the particular setting the periodogram bar is in many cases equal to the mean of the coefficients of determination of two fits with different step functions with opposed jumps (see [Thieler et al. 2013](#) or [Thieler 2013](#) for more details).

`RobPer` provides two options to fit periodic step functions. The number of steps per cycle is controlled by the input parameter `steps`. Using `model = "step"`, a single periodic step function with steps of equal width is fitted for each trial period. Performing `regression = "L2"`, `model = "step"` is equivalent to calculating an Epoch Folding or Analysis of Variance periodogram. Using `model="2step"`, two different step functions with opposed jump times and steps of equal width are fitted separately and the periodogram bar is the mean of both coefficients of determination. This is the only option where two periodic functions are fitted for one trial period. It is included to provide the PDM periodogram with overlapping bins.

Sine functions

Sine functions are periodic and it is quite popular to use them to investigate periodicity. The classic periodogram of Fourier analysis for equally sampled time series represents the explained variance SE of a least squares fit of a sine model to the zero-centered time series. The Lomb-Scargle periodogram ([Scargle 1982](#)) works equivalently for unequally sampled time series.

As the mean of an irregularly sampled time series is not identical to the least squares fit of an intercept in a sine model, more recent methods use the uncentered data and fit a model with intercept, e.g., the Floating Mean Periodogram by [Cumming et al. \(1999\)](#) and the Generalized Lomb Scargle Periodogram by [Zechmeister and Kürster \(2009\)](#). Performing `regression = "L2"`, `model = "sine"` is equivalent to calculating those periodograms and in case of equidistant observation times also equivalent to the Fourier periodogram.

Some other methods as the Date Compensated Fourier Transform by [Ferraz-Mello \(1981\)](#), the SigSpec periodogram by [Reegen \(2007\)](#) or robust approaches by [Ahdesmäki et al. \(2007\)](#) and [Zhang and Chan \(2005\)](#) apply the same regression step as the Floating Mean and the Generalized Lomb Scargle Periodogram, but use the squared amplitude of the fitted sinusoid as the periodogram bar. In case of regular sampling, this is another representation of the classical periodogram of Fourier analysis. As the amplitude is a concept closely related to trigonometric functions, `RobPer` uses the coefficient of determination only, to obtain a general method independent of the periodic function chosen.

Further periodic functions

Recently, fitting more complex periodic functions has been proposed for periodograms. Fourier series (see [Hall et al. 2000](#) and [Palmer 2009](#)) and periodic splines (see [Akerlof et al. 1994](#), [Hall et al. 2000](#) and [Oh et al. 2004](#)) may provide better adaptivity compared to sine functions, but still present a continuous function, unlike the step function. `RobPer` offers the possibility to fit Fourier series of second (`model = "fourier(2)"`) or third (`model = "fourier(3)"`) degree or a periodic spline function with four knots per cycle (`model = "splines"`). For the latter option, B-splines are generated using the function `spline.des` from the package `splines`.

Regression technique	<code>regression</code>	R function (package)
Least squares	"L2"	<code>lm</code> (<code>stats</code>)
Least absolute deviations	"L1"	<code>rq</code> (<code>quantreg</code> , Koenker 2012)
Least trimmed squares	"LTS"	<code>ltsReg</code> (<code>robustbase</code> , Rousseeuw <i>et al.</i> 2012)
M-regression		
... with Huber function	"huber"	own implementation
... with Bisquare function	"bisquare"	own implementation
S-regression	"S"	slightly modified code from Salibian-Barrera and Yohai (2006)
τ -regression	"tau"	slightly modified code from Salibian-Barrera <i>et al.</i> (2008)

Table 3: Regression techniques implemented in `RobPer` and R functions used to perform the regression technique. For more details see Section 2.2.

2.2. Regression techniques: Input variable regression

Instead of fitting the models mentioned above by least squares regression, one can also apply robust regression techniques. In `RobPer`, the user can choose between seven regression techniques (see Table 3): Least squares regression is the most popular approach (see Table 2). Robust regression techniques like least absolute deviations, least trimmed squares (Rousseeuw and Yohai 1984) and M-regression (Huber and Ronchetti 1981) have already been used to fit sines (evaluating the squared amplitude) by Zhang and Chan (2005), Ahdesmäki *et al.* (2007), Li (2009) and Li (2010). M-regression with the Huber function was additionally applied to fit periodic splines by Oh *et al.* (2004). Thieler *et al.* (2013) use least absolute deviations and M-regression and all models described in this article to calculate periodograms based on the coefficient of determination.

To the best of our knowledge, S- (Rousseeuw and Yohai 1984) and τ -regression (Yohai and Zamar 1988) have not been used before in periodogram calculation. For the latter, we use the Algorithms Fast-S from Salibian-Barrera and Yohai (2006) and Fast- τ from Salibian-Barrera, Willems, and Zamar (2008) and slightly modified versions of the code distributed with the respective publication (see the respective paragraphs entitled in Section 2.2).

LTS-regression

The R function `ltsReg` from package `robustbase` (Rousseeuw *et al.* 2012) is used to perform LTS-regression in `RobPer`. In preliminary studies we observed that the function sometimes has problems finding a good solution. This results in a coefficient of determination which is too small or sometimes even negative. By setting `LTSopt = TRUE`, it is possible to let `RobPer` further optimize the solution of `ltsReg` by using the R function `genoud`, package `rgenoud` (Mebane Jr. and Sekhon 2011). This function uses an evolutionary approach to improve the given solution, locally optimizing the temporarily best solutions in a gradient descent algorithm. Further input parameters `pop.size` (size of one generation), `max.generations` (maximum of generations before stopping the algorithm) and `wait.generations` (maximum number of generations to wait for an improvement of the optimization criterion) control the behavior of the algorithm and can be set in `RobPer` by the input variable `genoudcontrol` (see Table 1). The input variable `tol` controls the precision for convergence criteria.

A further problem observed is that `ltsReg` sometimes aborts the fit. However it is typically able to perform the fit if it is run again. In case of a crash, `RobPer` calls `ltsReg` up to three times. After the third failed attempt, the respective periodogram bar is set to NA, or a least absolute deviation regression is performed. The latter is done, if the `ltsReg`-regression result should be further processed, using the `genoud` algorithm or using the LTS result as initial estimate for an M-regression fit (see next paragraph).

M-regression

In case of M-regression, a periodogram bar, i.e., the coefficient of determination $R^2 = 1 - \frac{SE}{SY}$ is calculated from the values

$$SE = \min_{\beta} \sum_{i=1}^n \rho \left(\frac{y_i - x_i^{\top} \beta}{\hat{\sigma}_{\beta}} \right) \quad (7)$$

and

$$SY = \min_{\mu} \sum_{i=1}^n \rho \left(\frac{y_i - \mathbf{i}_i \mu}{\hat{\sigma}_{\mu}} \right). \quad (8)$$

As explained above, they represent the minimization criteria of the fits of the chosen periodic fluctuation (SE, Equation 5) and of a location estimate (SY, Equation 6), respectively. The function ρ is a distance measure. The vector \mathbf{i} consists of ones in case of unweighted regression. In case of weighted regression, y_i , \mathbf{i}_i and the rows x_i of the design matrix are standardized by the measurement accuracy s_i (see Figure 11 in the Appendix).

The value $\hat{\sigma}_{\beta}$ is obtained in an initial estimation of the periodic fluctuation, calculating a scale estimate of the fitted residuals. [Maronna *et al.* \(2006, S. 171\)](#) recommend setting $\hat{\sigma}_{\mu}$ to $\hat{\sigma}_{\beta}$. This means that SY depends on the trial period and cannot be calculated globally. On the other hand this ensures that the full model $Y = X\beta + \epsilon$, $\epsilon \sim \mathcal{N}(0, \sigma_{\beta})$ is a generalisation of the intercept model $Y = \mathbf{i}\mu + \epsilon$, $\epsilon \sim \mathcal{N}(0, \sigma_{\mu})$ and thus $SE \leq SY$ and $R^2 \geq 0$.

So for this regression technique, an implementation is needed where the scale estimate can be fixed in advance. The R functions known to us for M-regression (`r1m`, package **MASS** from [Venables and Ripley 2002](#), `lmrob.M.fit`, package **robustbase** from [Rousseeuw *et al.* 2012](#), `iwlsm`, package **RSiena** from [Ripley, Boitmanis, and Snijders 2013](#) and `robustregBS` and `robustRegH`, package **robustreg** from [Johnson 2012](#)) do not fulfill this requirement. Hence M-regression is newly implemented for **RobPer**. Like the functions specified before, this implementation is based on an Iteratively Reweighted Least Squares (IRWLS) approach (see [Maronna *et al.* 2006](#), pp. 104–105), but meets our special requirements. For M-regression using the biweight function, the implementation makes also use of the function `genoud`, package **rgenoud** (see previous paragraph) because IRWLS does not deal well with objective functions with local optima.

As noted above, observed values and design matrices are standardized by the measurement accuracies in case of weighted regression. The variance of the error is expected to be about one then. Hence it can be reasonable to set $\hat{\sigma}_{\beta}$ to one. This can be done in `RobPer` setting the input parameter `var1` to `TRUE`. Our experience is that this is recommendable (in case of weighted M-regression).

To calculate a periodogram bar using M-regression with IRWLS, three initial estimates are needed: A scale estimate $\hat{\sigma}_{\beta}$ (if not set to one) and initial location estimates $\hat{\beta}^{(0)}$ and $\hat{\mu}^{(0)}$

for β and μ . The initial estimates should be obtained using robust techniques. As proposed by Maronna *et al.* (2006, p. 105) we use the median (weighted if the s_i shall be taken into account) to initially estimate μ . For β , LTS-regression (see previous paragraph) is used. It has a high breakdown point and is thus appropriate in situations with many observations not agreeing with the best fit. This situation will often occur in periodogram calculation, as many trial periods and thus many wrong models are fitted to the light curve. The scale estimate $\hat{\sigma}_\beta$ is calculated as the (weighted) median of the residuals of the LTS-fit.

S-regression

In case of `regression = "S"`, `RobPer` uses the Fast-S-Algorithm by Salibian-Barrera and Yohai (2006) to efficiently perform S-regression for fitting the periodic function. The algorithm starts with a set of N parameter candidates, locally optimizing them using `kk` iterations, then optimizing the `tt` best of these candidates until convergence and finally choosing the best parameter candidate.

The R function `FastS` used in `RobPer` is a slightly modified version of the function `fast.s` published by Salibian-Barrera and Yohai (2006). It was changed in order to work more efficiently in the context given here, especially when fitting step functions, and to specify one parameter candidate in advance. This candidate is set to

$$\hat{\beta}_\mu = \begin{cases} (\hat{\mu}, \dots, \hat{\mu})^\top \in \mathbb{R}^m & \text{model} \in \{\text{"step"}, \text{"2step"}, \text{"splines"}\} \\ (\hat{\mu}, 0, \dots, 0)^\top \in \mathbb{R}^m & \text{model} \in \{\text{"sine"}, \text{"fourier(2)"}, \text{"fourier(3)"}\} \end{cases} \quad (9)$$

where m denotes the dimension of the linear model of the periodic function and $\hat{\mu}$ denotes the location estimate. $\hat{\beta}_\mu$ corresponds to the fit obtained from the location model in the parametrization of the full model. This ensures that fitting the full periodic function will not give a worse fit than fitting only a location parameter. Otherwise it could happen that $SY < SE$ and the coefficient of determination (which has to lie in $[0, 1]$) would be negative.

Further changes in `FastS` are:

1. The input variables `k` and `best.r` are renamed to `kk` and `tt` to unify notation as in `FastTau`. The input variables `int`, `N`, `kk`, `tt`, `b`, `cc` and `seed` are merged to a list `Scontrol`, which is also input to `RobPer` (except for `int`, which is fixed in `RobPer`).
2. If an intercept column is added to the designmatrix (using `Scontrol$int = TRUE`), this is done before the dimension of the designmatrix is determined (instead of doing this first and redoing it in case of `Scontrol$int = TRUE`).
3. To find a subsample in general position, regressors $x_{i^*}^\top$ are sampled from the set of rows of the designmatrix X ignoring the frequency of occurrence in X . For each regressor $x_{i^*}^\top$, one value y_i is then sampled from the entries of y belonging to this regressor. In case of a step function to be fitted, one observation per step is drawn to get a subsample.
4. If no subsample can be found in 100 trials, `FastS` returns `NA`. `RobPer` then releases a warning, but can calculate further periodogram bars for other trial periods.
5. The internal functions `loss.S`, `re.s`, `f.w`, `scale1`, `our.solve` and `rho` are now defined outside `FastS`. Otherwise R would have to redefine them for each periodogram bar.

6. The subfunction `norm` is replaced by the function `norm(..., "2")` from the **base** package.
7. The labels of the output are changed for better interpretation.

τ -regression

In case of `regression = "tau"`, τ -regression is used to fit the periodic function. `RobPer` uses the Fast- τ -Algorithm of [Salibian-Barrera *et al.* \(2008\)](#) which works according to the same optimizing principle as `FastS` for S-regression (see previous paragraph), i.e., optimizing `N` candidates in `kk` iterations and further optimizing the `tt` best of these until convergence. Since computation of the objective value is expensive, it is possible to approximate it with `rr` iteration steps when choosing `approximate = TRUE`. For more details see [Salibian-Barrera *et al.* \(2008\)](#).

The R function `FastTau` used in **RobPer** is a slightly modified version of the R-code published in [Salibian-Barrera *et al.* \(2008\)](#) with similar changes as in `FastS` compared to `fast.s` (see previous paragraph). The changes are:

1. A candidate β_μ (see Equation 9) is allowed.
2. Input variables `N`, `kk`, `tt`, `rr` and `approximate` are combined to a list `taucontrol`, which is also input to `RobPer`.
3. Subsamples in general position are found as in `FastS` (change 3 in the previous paragraph).
4. If no subsample can be found, `FastTau` returns `NA` instead of a break using the `stop` function. This behavior allows `RobPer` to release a warning, but calculate further periodogram bars for other trial periods.
5. A multiple times used block of code to check new regression parameter candidates for providing the best optimization value so far has been outsourced to the subfunction `checkbest`.
6. Due to rounding errors, it may happen in the IRWLS algorithm that negative values close to zero occur, although they have to be non-negative by theory. This is avoided by rounding such values to eight digits.
7. The subfunction `randomset` is replaced by the R function `sample` from the **base** package as both functions fulfill the same task and `sample` is faster.
8. The labels of the output are changed for better interpretation.

3. Fit Beta distributions with betaCvMfit

In this section we present the function `betaCvMfit` to robustly fit a Beta distribution to a sample using Cramér-von-Mises (CvM) distance minimization. The function is adapted from R code by Brenton R. Clarke for fitting a Gamma distribution (see [Clarke, McKinnon, and Riley 2012](#)) using CvM distance minimization. Section 3.1 motivates the application of this function, while its usage is explained in more detail in Section 3.2.

3.1. Motivation

After a periodogram is calculated, one might be interested in the automatic detection of significant periods. A period shall be called significant, if the respective periodogram bar is atypical from the distribution of the applied criterion under the null hypothesis of no periodic fluctuation. To determine significance, this distribution needs to be known or estimated. Let β_α be the α -quantile of this distribution. Assuming independent identically distributed periodogram bars $\text{Per}(p_1), \dots, \text{Per}(p_q)$ we get

$$P\left(\max(\text{Per}(p_1), \dots, \text{Per}(p_q)) \geq \beta_{\sqrt[q]{1-\alpha}}\right) = \alpha. \quad (10)$$

A single periodogram bar calculated as described in Section 2 using unweighted least squares regression is $\mathcal{B}(\frac{m-1}{2}, \frac{n-m}{2})$ -distributed, where \mathcal{B} denotes the Beta distribution and m is the dimension of the model. This result can be found in [Schwarzenberg-Czerny \(1998\)](#) or easily be deduced from [Seber and Lee \(2003, p. 110\)](#) and [Gupta and Nadarajah \(2004, p. 51\)](#). Already small violations of the assumptions made about the method or the light curve disturb this proceeding. In this work, we consider weighted and robust regression in addition to ordinary least squares. Beside, we have to take into account small deviations from our model assumptions like bad estimates s_i . An example is shown in Figure 2. Panel a shows the weighted least squares periodogram (using a sine model) of a light curve only consisting of white noise. The observed values were generated as

$$y_i = y_{w;i} + c \cdot y_{r;i}, i = 1, \dots, n \quad (11)$$

with $y_{w;i}$ and $y_{r;i}$ being realisations from

$$Y_{w;i} \sim \mathcal{N}(0, s_i^2), \quad (12)$$

$$Y_{r;i} \sim \mathcal{N}(0, 1). \quad (13)$$

The value of s_i is given for all i , and c is chosen to fulfill

$$\frac{\text{var}(c \cdot y_r)}{\text{var}(y_w) + \text{var}(c \cdot y_r)} = 0.2, \quad (14)$$

where var denotes the empirical variance. This means, there is roughly an extra 20 percent noise which is not explained by the measurement accuracies. Evidently, no periodogram bar is outstanding, but using the $\sqrt[q]{0.95}$ quantile of a $\mathcal{B}(\frac{m-1}{2}, \frac{n-m}{2})$ distribution (dashed line), several periods are found automatically.

To circumvent these problems, [Thieler et al. \(2013\)](#) propose to relax the assumption of a predefined $\mathcal{B}(\frac{m-1}{2}, \frac{n-m}{2})$ -distribution and only assume that the periodogram values can be

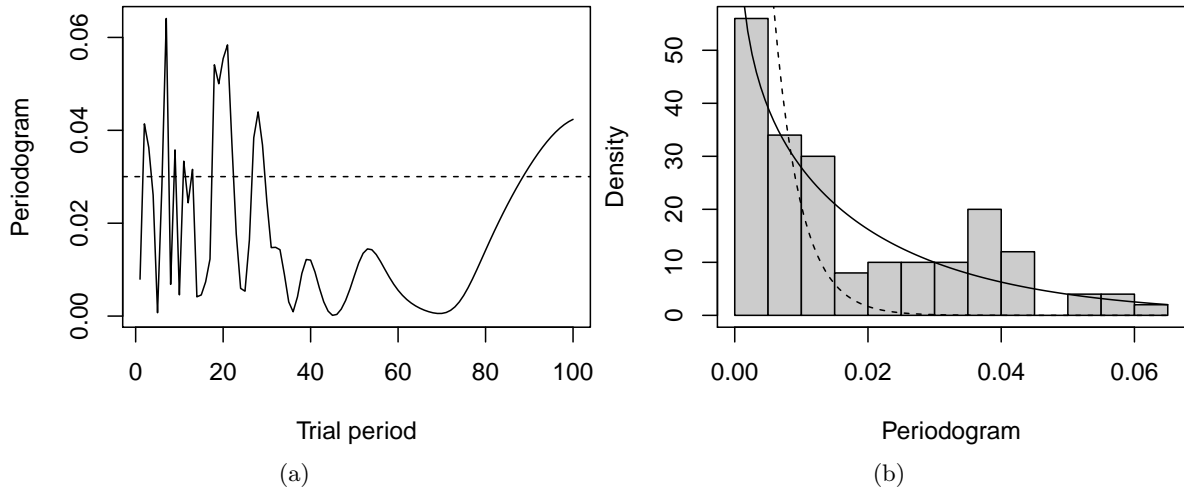


Figure 2: Example illustrating that a predefined $\mathcal{B}\left(\frac{m-1}{2}, \frac{n-m}{2}\right)$ distribution is sometimes not flexible enough if the model restrictions are slightly violated (see text for details). Panel a shows the periodogram of a light curve not completely following the assumed data model with the $\sqrt[3]{0.95}$ quantile of a $\mathcal{B}\left(\frac{m-1}{2}, \frac{n-m}{2}\right)$ distribution (dashed line). Panel b shows a histogram of the periodogram bars, with the density of the $\mathcal{B}\left(\frac{m-1}{2}, \frac{n-m}{2}\right)$ (dashed) and the CvM-fitted Beta distribution with parameters $0.8 < 1 = \frac{m-1}{2}$ and $40.18 < 248.5 = \frac{n-m}{2}$ (solid).

approximated by any Beta distribution. As peculiar periods are expected to show up as outliers, robustly fitting a $\mathcal{B}(\theta_1, \theta_2)$ -distribution to $\text{Per}(p_1), \dots, \text{Per}(p_q)$ is proposed. The authors use CvM distance minimization for this, which has been a recommendable technique in the past to fit Gamma distributions in the presence of outliers (see [Clarke et al. 2012](#)). The CvM is defined as

$$\int_0^\infty (F_n(u) - F_\theta(u))^2 dF_\theta(u) = \frac{1}{n} \sum_{i=1}^n \left(F_\theta(u_{(i)}) - \frac{i-0.5}{n} \right)^2 + \frac{1}{12n^2}, \quad (15)$$

where $u_{(1)}, \dots, u_{(n)}$ is the ordered sample, F_n is the empirical distribution function and F_θ is the distribution function of $\mathcal{B}(\theta_1, \theta_2)$.

Figure 2(b) shows the predefined (solid) and the CvM-fitted (dashed) Beta density for a periodogram calculated from the only-noise-data described above. While the $\sqrt[3]{0.95}$ quantile of the predefined distribution is about 0.03, the related quantile of the fitted distribution is 0.16 and no period is detected automatically.

The above approach falls within the framework of outlier detection described by [Davies and Gather \(1993\)](#) and is successfully used by [Thieler et al. \(2013\)](#) in the context discussed here. However, it assumes independent periodogram bars. This may cause problems when the periodogram peaks are broad (because the assumption of independency of the periodogram bars is violated): Then it can be hard for the automatism to find any outlying periodogram value, as there are many high values. One might try to ease this problem choosing a selection of trial periods with large distances or considering only the periods referring to local maxima in the periodogram as (roughly) independent trial periods (modifying and expanding an approach of [Zechmeister and Kürster 2009](#)) and fit the Beta distribution to them using a CvM fit.

Simulations indicate that the Beta distribution describes the distribution under the null hypothesis rather well for the different periodograms. Nevertheless, in the following we will call detected periods 'valid' and not 'significant' to stress that our approach to detect periods lacks a theoretical justification.

3.2. The R function `betaCvMfit`

The function `betaCvMfit` fits a $\mathcal{B}(\theta_1, \theta_2)$ -distribution with mean $\theta_1/(\theta_1 + \theta_2)$ to a sample vector `data` using CvM distance minimization and has been applied in [Thieler *et al.* \(2013\)](#) for fitting Beta distributions to periodograms to detect valid periods.

As it may happen that the periodogram bars become negative due to fitting problems, the function sets all negative entries of `data` to zero. If the logical input variable `CvM` is set to `TRUE`, a CvM fit is performed. As initial values for the optimization, the moment estimates of the Beta distribution

$$\hat{\theta}_1 = -\frac{\bar{x} \cdot (-\bar{x} + \bar{x}^2 + \hat{s}^2)}{\hat{s}^2}, \quad \hat{\theta}_2 = \frac{\hat{\theta}_1 - \hat{\theta}_1 \cdot \bar{x}}{\bar{x}} \quad (16)$$

are used. If the input variable `rob` is set to `TRUE`, the median and the median absolute deviation from the median (MAD) are used instead of the arithmetic mean for \bar{x} and the standard deviation for \hat{s} , respectively. In case of a very small estimate for \hat{s} (which in our experience mostly happens if \hat{s} is the MAD), the function stops as it is not possible to calculate the estimates $\hat{\theta}_1$ and $\hat{\theta}_2$ shown above. The parameters of a Beta distribution are strictly positive. Since it can happen that $\hat{\theta}_1$ or $\hat{\theta}_2$ are negative, the initial estimates are clipped to be at least 0.00001. If `CvM` is set to `FALSE`, the CvM distance is not optimized, but the initial estimates $\hat{\theta}_1$ and $\hat{\theta}_2$ are returned.

Figure 3 shows the different fits varying the input variables `CvM` and `rob` for 50 $\mathcal{B}(4, 15)$ -distributed observations containing 10 percent outliers between 0.8 and 1.

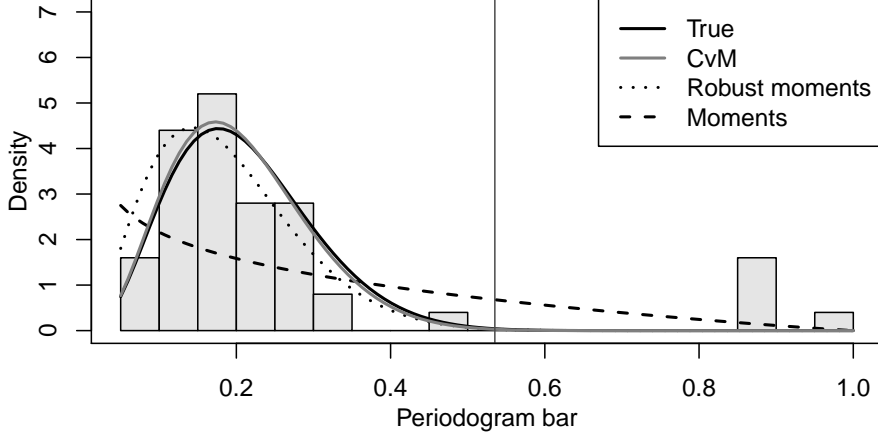


Figure 3: Grey-scale-version of the example for `betaCvMfit` given in the **RobPer**-manual: Histogram of 45 $\mathcal{B}(4, 15)$ -distributed observations and 5 outliers uniformly distributed between 0.8 and 1. The black solid line shows the $\mathcal{B}(4, 15)$ -distribution, the other curves show different fits using `betaCvMfit` (in case of `CvM = TRUE`, the different settings for `rob` lead to the same result).

4. Generate light curves with `tsgen`

To investigate our periodogram methods in simulations, we implemented the R function `tsgen` to generate artificial light curves. A preliminary version of this function is used in [Thieler et al. \(2013\)](#). The light curves $(t_i, y_i, s_i)_{i=1, \dots, n}$ are generated as realisations of the model

$$T_i = T_{(i)}^*, \quad T_1^*, \dots, T_n^* \sim \mathcal{D}(p_s), \quad (17)$$

$$Y_i = \begin{cases} Y_{f;i} + Y_{w;i} + Y_{r;i}, & Y_i \text{ 'behaves regularly' } \\ Y_i^*, & Y_i \text{ is an outlier} \end{cases}, \quad (18)$$

$$Y_{f;i} = f\left(\frac{T_i}{p_f}\right), \quad f(\xi) = f(\xi + 1) \forall \xi \in \mathbb{R} \quad (19)$$

$$Y_{w;i} \sim \mathcal{N}(0, \sigma_i^2), \quad (20)$$

$$s_i = \begin{cases} \text{given estimate for } \sigma_i \text{ independent from } Y_1, \dots, Y_n, & s_i \text{ 'behaves regularly' } \\ s_i^*, & s_i \text{ is an outlier} \end{cases}, \quad (21)$$

where $T_{(i)}^*$ denotes the i -th ordered observation time in T_1^*, \dots, T_n^* and $\mathcal{D}(p_s)$ is a periodic sampling density with period p_s . The noise component Y_r is a power law noise (see [Timmer and König 1995](#)) with power exponent α and is white noise in case of $\alpha = 0$. Inserting another noise component and two types of outliers, this extended model allows to generate data violating the model introduced in Section 1.

The function calls several autonomous subfunctions one by one which perform individual simulation steps. These are:

1. Generate a sampling t_1, \dots, t_n (using `sampler`, see Section 4.1)
2. Generate a periodic signal $y_{f;1}, \dots, y_{f;n}$ (using `signalgen`, see Section 4.2)

3. Add noise $y_{w;1}, \dots, y_{w;n}$ with related measurement accuracies s_1, \dots, s_n and a noise component $y_{r;1}, \dots, y_{r;n}$ unrelated to the s_i (using `lc_noise`, see Section 4.3)
4. Disturb the light curve replacing measurement accuracies s_i by outliers, or replacing observations $y_i = y_{f;i} + y_{w;i} + y_{r;i}$ by aperiodic features (using `disturber`, see Section 4.4)

Table 4 lists all input variables for the subfunctions. The grey-shaded arguments are also input arguments to `tsgen`, which passes them to the respective subfunction.

4.1. Generate sampling using sampler

The R function `sampler` is used to sample observation times t_1, \dots, t_n in the interval $[0, n_s \cdot p_s]$ with a possibly periodic sampling of period p_s . The sampling pattern depends on the argument `ttype` (see Table 4). If a periodic pattern is chosen, the observed time interval covers n_s cycles of it.

In case of `ttype = "equi"`, the observation times are equidistantly sampled with $t_i = i \frac{p_s \cdot n_s}{n}$. For `ttype = "unif"`, the observation times are independently drawn from a uniform distribution on $[0, n_s \cdot p_s]$. Both these sampling schemes are aperiodic, the sampling period p_s only influences the duration $t_n - t_1$ of the sampling.

For `ttype = "sine"` and `ttype = "trian"`, the observation times are sampled from a periodic density with sampling period p_s . First, observation cycles z_i^* are drawn from a discrete uniform distribution on $\{1, \dots, n_s\}$ to determine the cycle the i -th observation is part of. Second, observation phases φ_i^* are sampled with density

$$d_{sine}(x) = \sin(2\pi x) + 1 \quad (\text{for } \text{ttype} = \text{"sine"}) \quad (22)$$

$$\text{or } d_{trian}(x) = \begin{cases} 3x, & 0 \leq x \leq \frac{2}{3}, \\ 6 - 6x, & \frac{2}{3} < x \leq 1 \end{cases} \quad (\text{for } \text{ttype} = \text{"trian"}). \quad (23)$$

To sample from d_{sine} , the function `BBsolve`, package `BB` by [Varadhan and Gilbert \(2009\)](#), is used.

The unsorted observation times t_i^* are then generated using

$$t_i^* = \varphi_i^* + (z_i^* - 1)p_s. \quad (24)$$

The sine-shaped density is motivated by sampling patterns observed in real data (see Figure 1(b)), the triangular shaped density offers an alternative periodic sampling. Separately sampling observation cycle and phase was proposed by [Hall and Yin \(2003\)](#).

As the result, `sampler` returns the ordered observation times t_1, \dots, t_n .

4.2. Generate periodic signal using signalgen

To generate the periodic component in the observed values, the R function `signalgen` is used. The values $y_{f;1}, \dots, y_{f;n}$ with fluctuation period p_f at observation times t_1, \dots, t_n are generated using

$$y_{f;i} = f\left(\frac{t_i}{p_f}\right), \quad i = 1, \dots, n. \quad (25)$$

Input	Subfunction	Comment
<code>ps</code> $\in \mathbb{R}_{>0}$	<code>sampler</code> , <code>disturber</code>	Sampling period p_s . Default 1.
<code>ncycles</code> $\in \mathbb{N}$	<code>sampler</code>	Number n_s of sampling cycles.
<code>npoints</code> $\in \mathbb{N}$	<code>sampler</code>	Sample size n .
<code>ttype</code>	<code>sampler</code>	Distribution $\mathcal{D}(p_s)$ of the unsorted observation times. Options are: " <code>equi</code> " (equidistant sampling), " <code>unif</code> " (uniform sampling), " <code>sine</code> " (sine-shaped density, see Section 4.1) and " <code>trian</code> " (triangular density, see Section 4.1).
<code>tt</code> $\in \mathbb{R}^n$	<code>signalgen</code> , <code>lc_noise</code> , <code>disturber</code>	Observation times t_1, \dots, t_n , e.g., output from <code>sampler</code> .
<code>pf</code> $\in \mathbb{R}_{>0}$	<code>signalgen</code>	Fluctuation period p_f . Default 1.
<code>ytype</code>	<code>signalgen</code>	Type of periodic fluctuation f . Options: " <code>const</code> " (constant), " <code>sine</code> " (sine), " <code>trian</code> " (triangular function) and " <code>peak</code> " (peak function).
<code>sig</code> $\in \mathbb{R}^n$	<code>lc_noise</code>	Values $y_{f;1}, \dots, y_{f;n}$ of the periodic fluctuation, e.g., output from <code>signalgen</code> .
<code>SNR</code> $\in \mathbb{R}_{>0}$	<code>lc_noise</code>	Relation $\text{var}(y_f) / \text{var}(y_w + y_r)$.
<code>redpart</code> $\in [0, 1]$	<code>lc_noise</code>	Fraction $\text{var}(y_r) / (\text{var}(y_w) + \text{var}(y_r))$ of noise not related to measurement accuracies.
<code>alpha</code>	<code>lc_noise</code>	Power law coefficient of the noise component y_r . Set to zero for $y_{r;1}, \dots, y_{r;n} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$.
<code>y</code> $\in \mathbb{R}^n$	<code>disturber</code>	Observed values y_1, \dots, y_n , e.g., output from <code>lc_noise</code> .
<code>s</code> $\in \mathbb{R}_{>0}^n$	<code>disturber</code>	Measurement accuracies s_1, \dots, s_n , e.g., output from <code>lc_noise</code> .
<code>s.outlier.fraction</code> $\in [0, 1]$	<code>disturber</code>	Fraction of measurement accuracies to be replaced by outliers.
<code>interval</code> $\in \{\text{TRUE}, \text{FALSE}\}$	<code>disturber</code>	If TRUE, the y_i belonging to a random time interval are disturbed.

Table 4: Input values for the subfunctions of `tsgen`. See the respective section for more details. Grey-shaded values are also input for `tsgen`, which passes the values to the respective subfunction. `var` denotes the empirical variance.

The observation times, the fluctuation period and the shape of f are arguments of the function (see Table 4). In case of `ytype = "const"`, f is defined as

$$f(t) = 0, \quad (26)$$

so there is no (periodic) fluctuation. This setting can be used to investigate the false alarm probability of a period detection method. In case of `ytype = "sine"`, f is defined as

$$f(t) = \sin\left(\frac{2\pi t}{p_f}\right). \quad (27)$$

This is a typical assumption in the literature. For `ytype = "trian"`,

$$f(t) = \begin{cases} 3\varphi_1(t), & 0 \leq \varphi_1(t) \leq \frac{2}{3} \\ 6 - 6\varphi_1(t), & \frac{2}{3} < \varphi_1(t) \leq 1 \end{cases} \quad (28)$$

with $\varphi_1(t) = t \bmod 1 = (t - \lfloor t \rfloor)$ is used. This triangular shaped function was originally implemented in order to be able to choose between different periodic shapes. The light curve observed for CoRoT ID 0105288363 (Chadid *et al.* 2011) shows that functions with a similar shape are quite realistic. When choosing `ytype = "peak"`, y_f is generated using

$$f(t) = \begin{cases} 9 \exp\left(-3p_f^2 \left(\varphi_1(t) - \frac{2}{3}\right)^2\right), & 0 \leq \varphi_1(t) \leq \frac{2}{3}, \\ 9 \exp\left(-12p_f^2 \left(\varphi_1(t) - \frac{2}{3}\right)^2\right), & \frac{2}{3} < \varphi_1(t) \leq 1 \end{cases}. \quad (29)$$

This function mostly shows values close to zero and large values for only one time unit per cycle. This 'peak' occurring in each cycle has an asymmetric shape.

As the result, `signalgen` returns the periodic component $y_{f;1}, \dots, y_{f;n}$ of the observed values.

4.3. Add noise and measurement accuracies using `lc_noise`

The R function `lc_noise` is used to generate measurement accuracies s_1, \dots, s_n and add noise to a periodic fluctuation (see Table 4). The measurement accuracies are sampled from a $\text{Gamma}(3,10)$ distribution. This choice is motivated by real data from Tluczykont *et al.* (2010). As shown in Equation 4, the noise component $y_w = (y_{w;1}, \dots, y_{w;n})^\top$ is a realisation of Y_w with $Y_{w;i} \sim \mathcal{N}(0, s_i^2)$.

A second noise component y_r does not depend on the s_i . It is generated as red noise, i.e., following a power law with power law index α . For $\alpha = 0$ we get white noise, flicker noise (pink noise) is generated using $\alpha = 1$ and brown noise using $\alpha = 2$. The power law noise is generated using subfunctions `TK95_uneq` and `TK95`. The latter generates an equidistant time series of power law noise according to Timmer and König (1995). For irregular observation times, a noise series resulting from `TK95` is used and an unequally sampled noise series is generated following Uttley, McHardy, and Papadakis (2002).

The noise components are scaled so that the variance of the $y_{r;i}$ has approximately the proportion `redpart` in the overall noise variance and that `SNR` is the ratio $\text{var}(y_f)/\text{var}(y_w + y_r)$, where $\text{var}(x)$ is the empirical variance of vector x .

Note that the white noise components' variances are exactly s_i^2 , so the s_i are not estimates, but true values. In this sense, the measurement accuracies of a generated light curve are more informative than for real light curves, where the measurement accuracies are estimates. Allowing for a second noise component makes it possible to lower the information of the measurement accuracies with respect to the overall noise in the observed values.

The function `lc_noise` returns the observed values $y_i = y_{f;i} + y_{w;i} + y_{r;i}$, $i = 1, \dots, n$.

4.4. Disturb light curve using disturber

The last subfunction applied in `tsген` is `disturber`, which can be used to disturb a given light curve (see Table 4). It replaces a given fraction of measurement accuracies by the smaller value $s_i^* = \frac{1}{2} \min(s_1, \dots, s_n)$, i in a subset of $\{1, \dots, n\}$. As small measurement accuracies stand for precise observations, the influence of observations with disturbed measurement accuracies s_i^* rises in case of a weighted fit. For unweighted regression, this type of disturbance does not affect the result of the fit.

Optionally, `disturber` also replaces observed values y_i by atypical values. For this, a time interval $[t_{start}, t_{start} + 3p_s]$ within the interval $[t_1, t_n]$ is randomly chosen and all observed values belonging to this time interval are replaced by a peak function:

$$y_i^* = 6 \tilde{y}_{0.9} \frac{d_{\mathcal{N}(t_{start}+1.5p_s, p_s^2)}(t_i)}{d_{\mathcal{N}(0, p_s^2)}(0)} \quad \forall i : t_i \in [t_{start}, t_{start} + 3p_s], \quad (30)$$

where $d_{\mathcal{N}(a, b^2)}(x)$ denotes the density of a normal distribution with mean a and variance b^2 at x . If the y_i are intended to be disturbed and the light curve is shorter than $3p_s$, the function will stop with an error message.

The function returns the modified vectors $y = (y_1, \dots, y_n)^\top$ and $s = (s_1, \dots, s_n)^\top$. If the option to change y -values is not used (see Table 4) and the fraction of outlying measurement accuracies is set to zero, y and s are returned unchanged.

5. Application

In this Section, we give examples how to use the **RobPer** package for light curve analysis. We start with an artificial example also given in the manual and then analyze some real data.

5.1. Artificial example

To generate an artificial light curve, `tsgen` can be used:

```
set.seed(22)
lightcurve <- tsgen(ttype = "sine", ytype = "peak", pf = 7, redpart = 0.1,
  s.outlier.fraction = 0, interval = FALSE, npoints = 200,
  ncycles = 100, ps = 5, SNR = 3, alpha = 0)
```

This light curve has a sine-shaped sampling (`ttype`) with sampling period 5 (`ps`) and covers a time interval of about 100 sampling cycles (`ncycles`), so 500 time units. It consists of 200 observations (`npoints`) and the observed values contain a peak-shaped periodic fluctuation (`ytype`) with fluctuation period 7 (`pf`). The measurement accuracies are related to about 90 percent of the noise component (`1-redpart`), the rest of the noise is white as well (`alpha`). The empirical variance of the periodic fluctuating component in the observed values is three times larger than the empirical variance in the noise component (`SNR`). The light curve contains no outliers in the measurement accuracies (`s.outlier.fraction`), or in the observed values (`interval`).

Alternatively, the functions `sampler`, `signalgen`, `lc_noise` and `disturber` can be used to generate the same light curve, see Section 4:

Sampling observation times:

```
set.seed(22)
tt <- sampler(ttype = "sine", npoints = 200, ncycles = 100, ps = 5)
```

Generate periodic fluctuation:

```
yf <- signalgen(tt, ytype = "peak", pf = 7)
```

Add noise and scale signal to the right SNR:

```
temp <- lc_noise(tt, sig = yf, SNR = 3, redpart = 0.1, alpha = 0)
y <- temp$y
s <- temp$s
```

Replace measurement accuracies by tiny outliers or include a peak:

```
temp <- disturber(tt, y, s, ps = 5, s.outlier.fraction = 0,
  interval = FALSE)
```

The result is the same:

```
all(cbind(tt, temp$y, temp$s) == lightcurve)
```

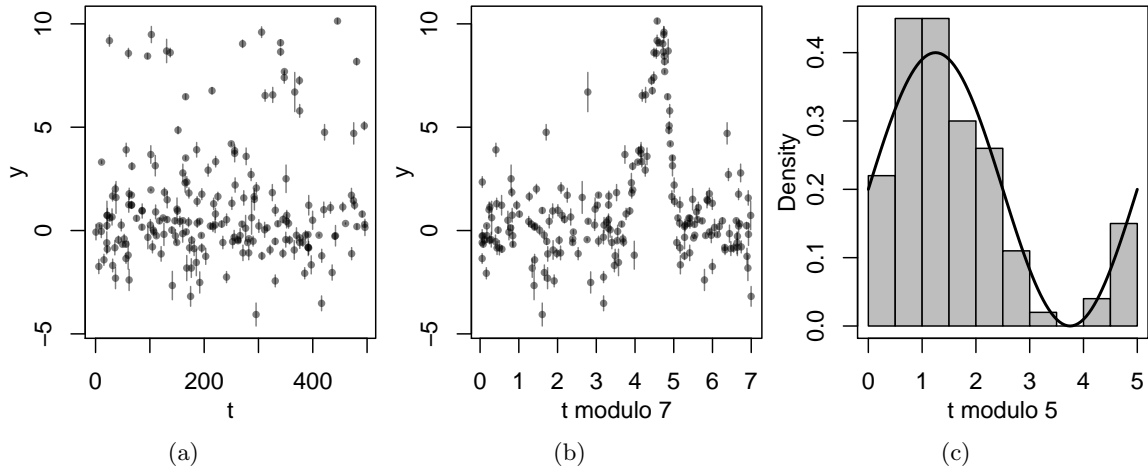


Figure 4: (a) Artificial light curve. The vertical bars mark the s_i . (b) Plotting time axis modulo 7 to show the periodic fluctuation of $p_f = 7$. (c) Histogram and sampling density of the observation times modulo 5 to show the sampling periodicity of $p_s = 5$.

Figure 4 shows plots of the generated light curve.

In the next step, we calculate a periodogram of the light curve. The periodogram is calculated fitting a spline model using unweighted M-regression with the Huber function. The light curve spans a time interval of approximately $\text{ncycles} \cdot \text{ps} = 500$ time units, so it is sensible to investigate periods up to 50 (one tenth, see Halpern, Leighly, and Marshall 2003).

```
PP <- RobPer(lightcurve, model = "splines", regression = "huber",
             weighting = FALSE, var1 = FALSE, periods = 1:50)
```

Outstanding periodogram bars are sought fitting a Beta distribution to the periodogram values using Cramér-von-Mises distance minimization (CvM) and determining the $\sqrt[5]{0.95}$ -quantile with $q = 50$ as the number of periodogram bars.

```
betavalues <- betaCvMfit(PP)
crit.val <- qbeta((0.95)^(1 / 50), shape1 = betavalues[1],
                 shape2 = betavalues[2])
```

Panel 5(a) depicts the histogram of the periodogram bars, the Beta distribution fitted (solid line) and its $\sqrt[5]{0.95}$ -quantile (solid vertical line). Further fits of a Beta distribution (method of moments, dotted and robust method of moments, dashed) and their respective $\sqrt[5]{0.95}$ -quantiles are shown as well.

```
hist(PP, breaks = 20, freq = FALSE, ylim = c(0, 250), xlim = c(0, 0.2),
     col = "grey", main = "")
betafun <- function(x) dbeta(x, shape1 = betavalues[1],
                             shape2 = betavalues[2])
curve(betafun, add = TRUE, lwd = 2)
abline(v = crit.val, lwd = 2)
```

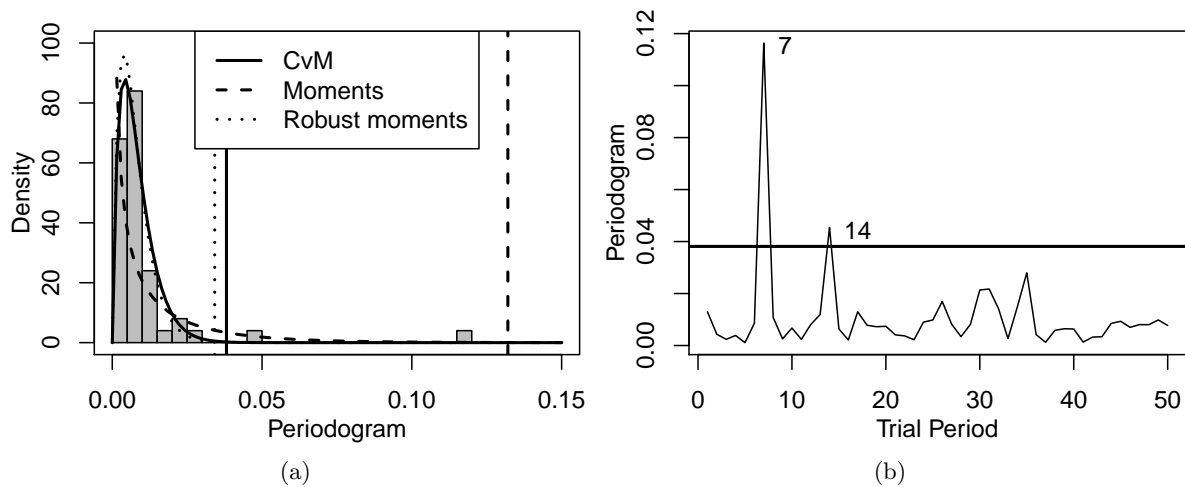


Figure 5: Periodogram bars calculated fitting a spline model using unweighted M-regression with the Huber function to the artificial example from Figure 4: Robustly fitting a Beta distribution to the periodogram bars in (a) leads to two outstanding trial periods (b).

Method of moments:

```
par.mom <- betaCvMfit(PP, rob = FALSE, CvM = FALSE)
myf.mom <- function(x) dbeta(x, shape1 = par.mom[1], shape2 = par.mom[2])
curve(myf.mom, add = TRUE, lwd = 2, lty = 2)
crit.mom <- qbeta((0.95)^(1 / 50), shape1 = par.mom[1], shape2 = par.mom[2])
abline(v = crit.mom, lwd = 2, lty = 2)
```

Robust method of moments:

```
par.rob <- betaCvMfit(PP, rob = TRUE, CvM = FALSE)
myf.rob <- function(x) dbeta(x, shape1 = par.rob[1], shape2 = par.rob[2])
curve(myf.rob, add = TRUE, lwd = 2, lty = 3)
crit.rob <- qbeta((0.95)^(1 / 50), shape1 = par.rob[1], shape2 = par.rob[2])
abline(v = crit.rob, lwd = 2, lty = 3)
legend("top", lty = 1:3, legend = c("CvM", "Moments", "Robust moments"),
      bg = "white", lwd = 2)
box()
```

Using the $\sqrt[50]{0.95}$ quantile of the CvM fit (solid line), a period of 7 time units seems to be valid (see Panel 5(b)). A multiple of the period (14) look valid, too. So the real periodic fluctuation of $p_f = 7$ is well recognized within the disturbed signal, as intended. Of course, a periodic function with period p is also periodic with period $k \cdot p$, $p \in \mathbb{N}$.

```
plot(1:50, PP, xlab = "Trial period", ylab = "Periodogram", main = "",
     type = "l")
abline(h = crit.val, lwd = 2)
text(c(7, 14), PP[c(7, 14)], c(7, 14), adj = 1, pos = 4)
```


5.2. Disturbed data from Gro J 0422-32

The first real data set we analyze is a light curve for gamma ray emission of the source Gro J 0422-32 obtained with the BATSE Telescope. The data was downloaded in June 2011 from <http://lhea-www.gsfc.nasa.gov/users/craigm/batse-lc/> and is shown in Figure 6(a). A large peak is visible starting at about 48900 Makarian Julian days (which corresponds to December 10th 1991 in the Gregorian calendar), a so called gamma ray burst. It occasionally occurs in gamma ray observations and can be considered as outlier. The light curve covers a time interval of about 3312 days, so following Halpern *et al.* (2003) we consider periods up to 330 days (about one tenth of the overall duration of the light curve). Figure 6(b) shows the periodogram obtained fitting a sine function using least squares regression, which is the classical approach in astroparticle physics. It is calculated using

```
RobPer(star_groj0422.32, periods = 1:330, model = "sine", regression = "L2",  
       weighting = FALSE)
```

Periodograms for τ -regression and M-regression using the Huber function are obtained replacing "L2" by "tau" or "huber" in the code line above. The respective periodograms are shown in Panels 6(c) and 6(d). All three periodograms do not show any outstanding peak. Apart from this, the periodograms using robust regression have a completely different shape than the least squares periodogram, which seems to have problems with the gamma ray burst. It might be questionable if the least squares periodogram would be able to find a periodic structure in the observations not influenced by the gamma ray burst. We add a sine with period 30 and amplitude 0.005 to the observed values and repeat the analysis. The results can be seen in Figure 7. In Panel a it is visible that we did not introduce a strong periodic behavior. Nevertheless, the robust periodograms (Panels c and d) easily detect it, while there is only a small local peak in the least squares periodogram (Panel b). The horizontal lines in (c) and (d) show the respective $\sqrt[330]{0.95}$ -quantiles of the CvM-fitted Beta distribution and are calculated from a periodogram PP using

```
shapes <- betaCvMfit(PP)  
Crit <- qbeta(0.95^(1 / 330), shape1 = shapes[1], shape2 = shapes[2])
```

So, as opposed to least squares regression, robust techniques are able to detect an (added) periodic fluctuation although the data are disturbed seriously by the gamma ray burst.

5.3. Data from Makarian 421 and 501

We continue the R functions' application to real data with gamma ray light curves published in [Tluczykont *et al.* \(2010\)](#). The light curve obtained for Makarian 421 (Mrk 421) is shown in Figure 1 (a) on page 3. Periodograms obtained fitting a sine are shown in Figure 8. When looking at the least squares periodogram (Panel a), one might wonder if there is a periodicity of 31 hidden in the same way as when adding a small periodic fluctuation to the Gro J 0422-32 data (see Figure 7(b)). However, the same periodograms for τ -regression (Panel b) and Huber-M-regression (Panel c) show a different behavior from Figure 7, so this does not seem to be the case. Especially, the least squares and the Huber-M- periodogram show a quite similar behavior regarding the local maxima. This could mean that there are not many observations weighted down in Huber-M-regression.

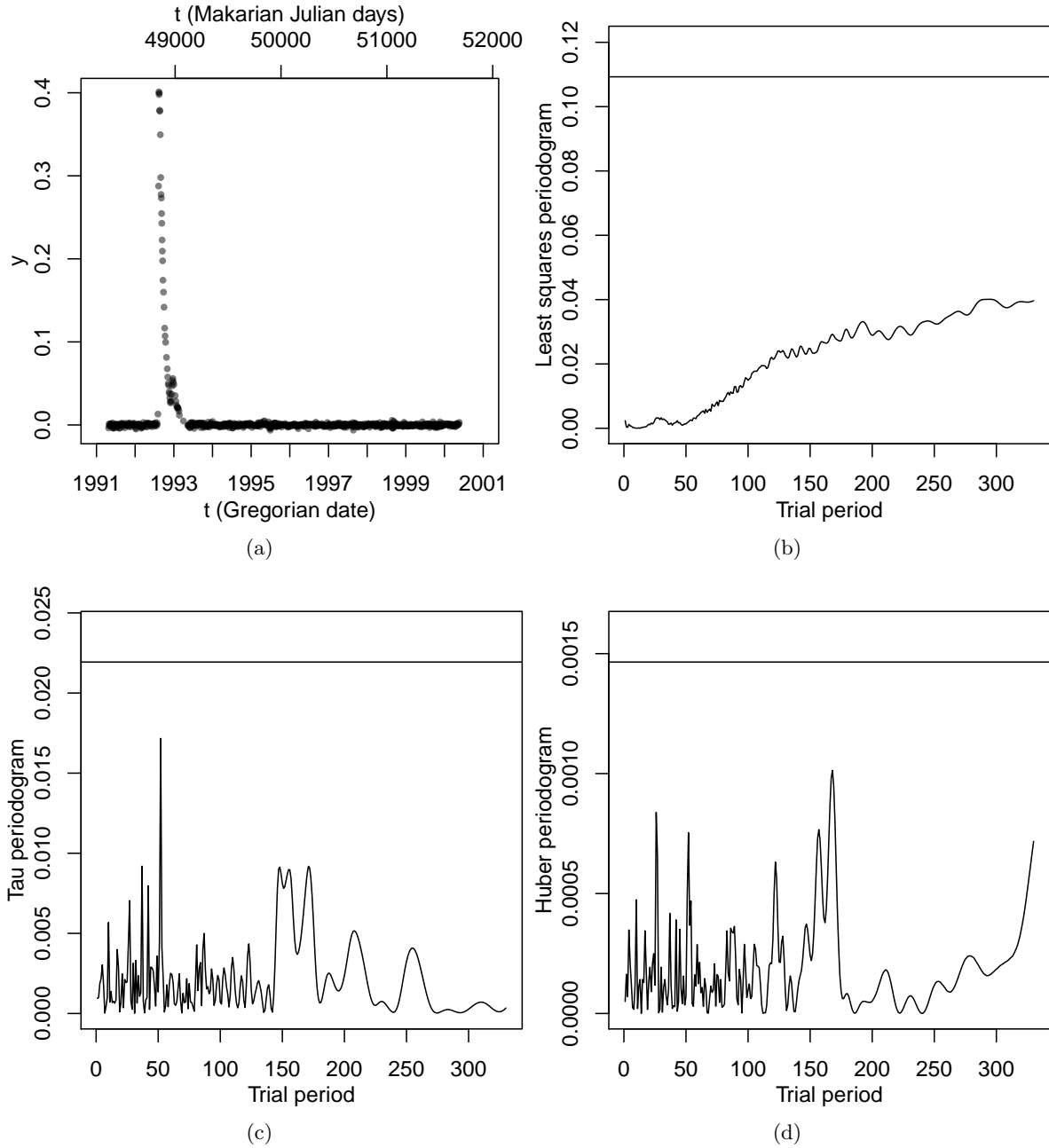


Figure 6: Analysis of Gro J 0422-32: (a) shows the light curve, while (b)–(d) show the periodograms fitting a sine using least squares (b), τ - (c), Huber-M-regression (d). No periodogram bar lies over the respective $\sqrt[330]{0.95}$ -quantile of the CvM-fitted Beta distribution (horizontal line).

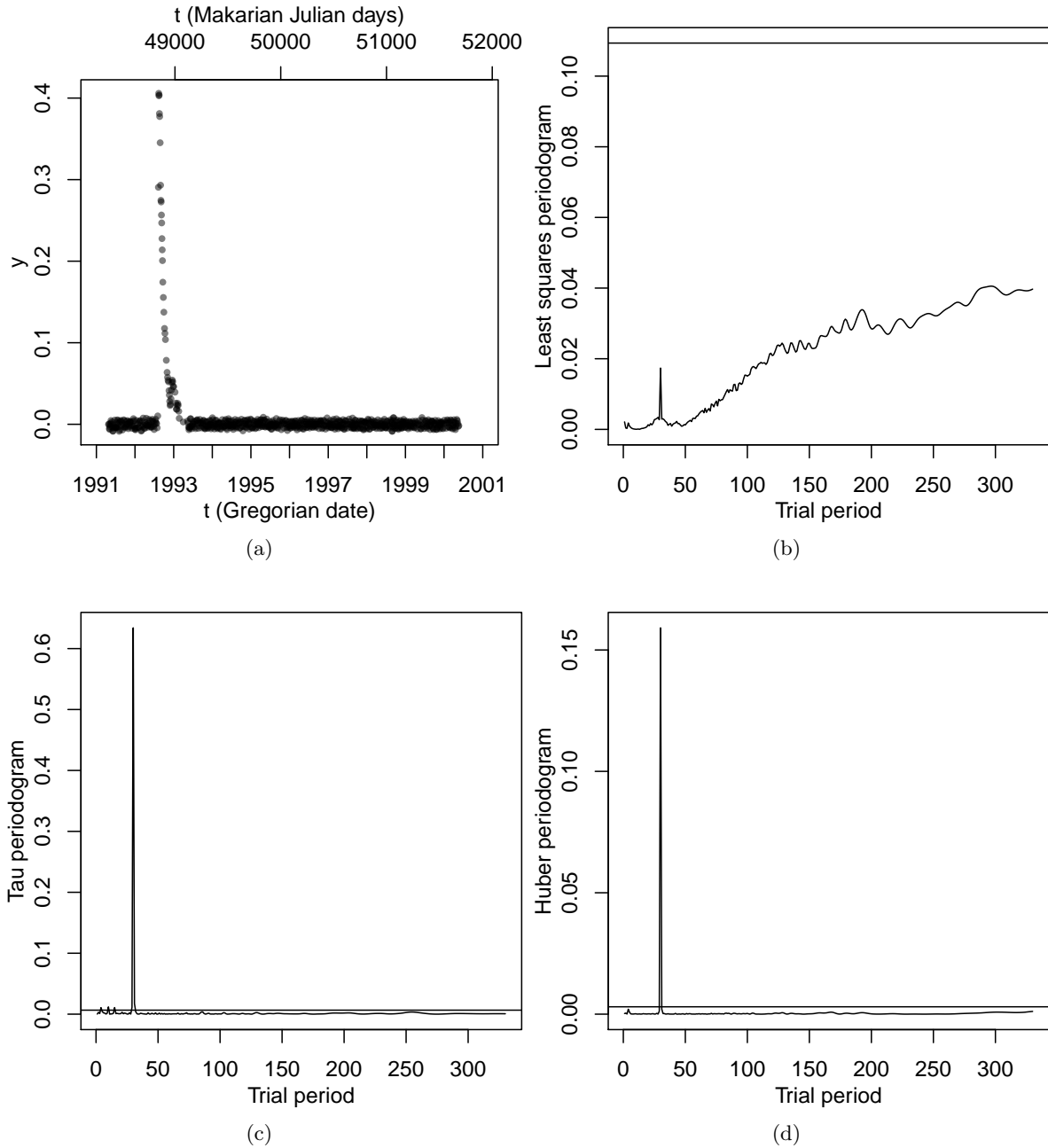


Figure 7: Adding a sine with amplitude 0.005 in Gro J 0422-32: (a) light curve, (b)–(d) periodograms fitting a sine using least squares (b), τ - (c), Huber-M-regression (d). The horizontal lines in Panels b–d show the respective $\sqrt[330]{0.95}$ - quantile of the CvM-fitted Beta distribution.

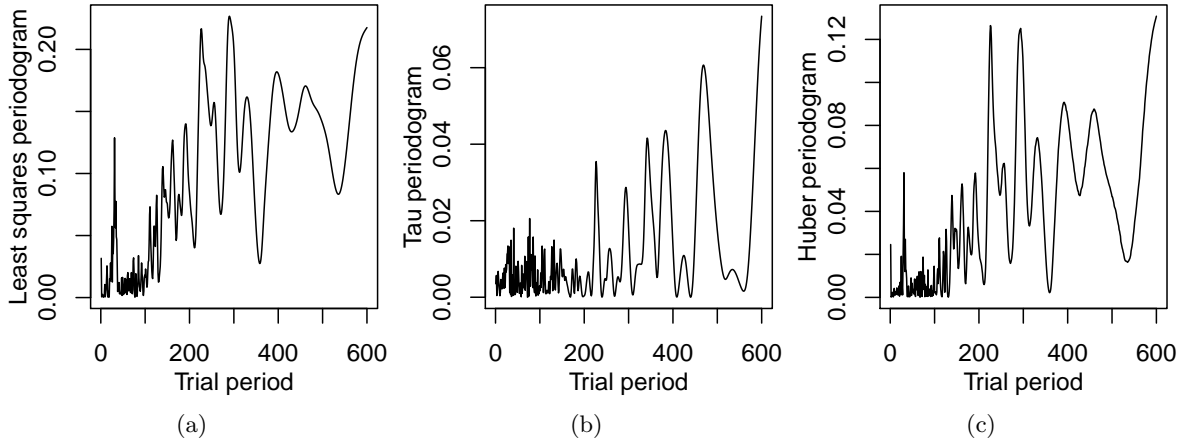


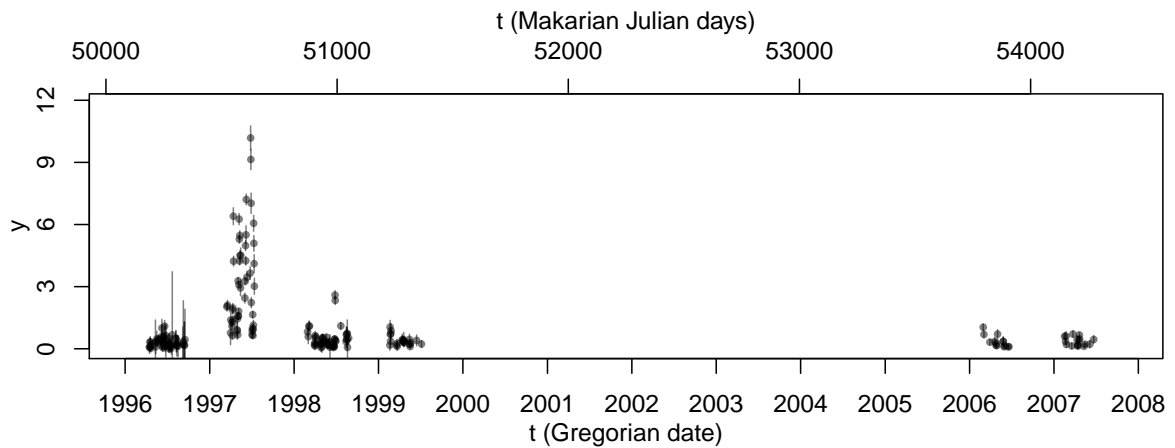
Figure 8: Periodograms for Mrk 421 (see Figure 1 (a)) obtained fitting a sine with (a) least squares regression (b) τ -regression (c) Huber-M-regression.

Another light curve, obtained for Mrk 501, and periodograms using least squares regression, τ -regression and Huber-M-regression are shown in Figure 9. Here we apply step regression, which is equivalent to Epoch Folding or Phase Dispersion Minimization when using least squares regression (see Section 2). The periodogram is calculated applying

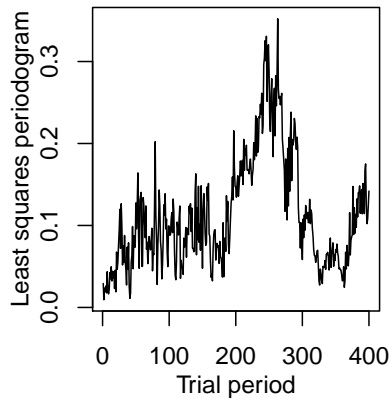
```
RobPer(Mrk501, periods = 1:400, model = "step", regression = "L2",
       weighting = FALSE)
```

in case of least squares regression and with `regression = "tau"` or `regression = "huber"` in case of τ - or Huber-M-regression, respectively. For least squares regression (Panel b) and Huber-M-regression (Panel d) we see a broad peak between the trial periods 200 and 300, much too broad to be considered as valid period (see Halpern *et al.* 2003). For τ -regression (Panel c), this behavior is not observed.

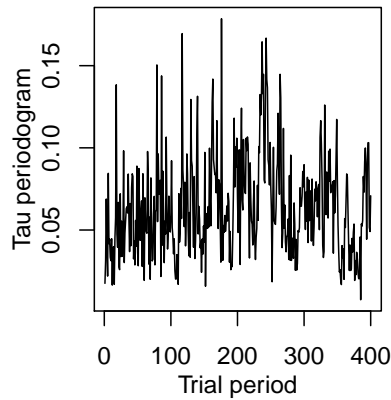
In the examples from the previous section, robust techniques recognize some periodicity in a light curve, while the least squares periodogram only provides a slightly atypical behavior for the trial period in question. Now it is the other way round: The least squares periodogram exhibits some interesting features, but robust techniques do not find evidence of a periodicity.



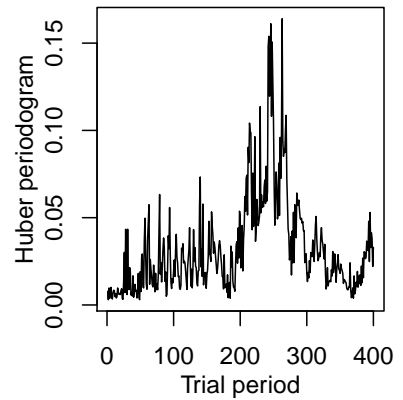
(a)



(b)



(c)



(d)

Figure 9: Light curve (a) and Periodograms for Mrk 501 obtained fitting a periodic step function with (b) least squares regression (c) τ -regression (d) Huber-M-regression.

6. Conclusions

The R package **RobPer** presented in this work is appropriate to search for periodicity in irregularly sampled time series, taking into account additional information on the precision of the measurement, if available. These are the typical characteristics of light curves, time series occurring in astroparticle physics. The periodogram is calculated fitting periodic functions to the light curve. It can be chosen between six different periodic functions and seven different regression techniques, so 42 possible combinations are offered, not taking into account further options like choosing the number of steps for the step model or using weighted regression. The function `betaCvMfit` allows to search for prominent periodogram bars as outliers in a Beta distribution robustly fitted to the periodogram. The function `tsgen` allows generation of artificial light curves for investigative use.

References

- Ahdesmäki M, Lähdesmäki H, Gracey A, Shmulevich I, Yli-Harja O (2007). “Robust Regression for Periodicity Detection in Non-Uniformly Sampled Time-Course Gene Expression Data.” *BMC Bioinformatics*, **8**(1), 233–248.
- Akerlof C, Alcock C, Allsman R, Axelrod T, Bennett DP, Cook KH, Freeman K, Griest K, Marshall S, Park HSea (1994). “Application of Cubic Splines to the Spectral Analysis of Unequally Spaced Data.” *The Astrophysical Journal*, **436**, 787–794.
- Chadid M, Perini C, Bono G, Auvergne M, Baglin A, Weiss WW, Deboscher J (2011). “CoRoT Light Curves of Blazhko RR Lyrae Stars.” *Astronomy & Astrophysics*, **527**, A146.
- Clarke BR, McKinnon PL, Riley G (2012). “A Fast Robust Method for Fitting Gamma Distributions.” *Statistical Papers*, **53**(4), 1001–1014.
- Croux C, Dehon C (2003). “Estimators of the Multiple Correlation Coefficient: Local Robustness and Confidence Intervals.” *Statistical Papers*, **44**(3), 315–334.
- Cumming A, Marcy GW, Butler RP (1999). “The Lick Planet Search: Detectability and Mass Thresholds.” *The Astrophysical Journal*, **526**(2), 890–915.
- Davies L, Gather U (1993). “The Identification of Multiple Outliers.” *Journal of the American Statistical Association*, **88**(423), 782–792.
- Deeming TJ (1975). “Fourier Analysis with Unequally-Spaced Data.” *Astrophysics and Space Science*, **36**(1), 137–158.
- Ferraz-Mello S (1981). “Estimation of Periods from Unequally Spaced Observations.” *The Astronomical Journal*, **86**(4), 619–624.
- Gupta AK, Nadarajah S (2004). *Handbook of Beta Distribution and Its Applications*. Dekker, New York, Basel.
- Hall P, Li M (2006). “Using the Periodogram to Estimate Period in Nonparametric Regression.” *Biometrika*, **93**(2), 411–424.

- Hall P, Reimann J, Rice J (2000). “Nonparametric Estimation of a Periodic Function.” *Biometrika*, **87**(3), 545–557.
- Hall P, Yin J (2003). “Nonparametric Methods for Deconvolving Multiperiodic Functions.” *Journal of the Royal Statistical Society B (Statistical Methodology)*, **65**(4), 869–886.
- Halpern JP, Leighly KM, Marshall HL (2003). “An Extreme Ultraviolet Explorer Atlas of Seyfert Galaxy Light Curves: Search for Periodicity.” *The Astrophysical Journal*, **585**, 665–676.
- Huber PJ, Ronchetti E (1981). *Robust Statistics*, volume 1. John Wiley & Sons.
- Johnson IM (2012). *robustreg: Robust Regression Functions*. R package version 0.1-3, URL <http://CRAN.R-project.org/package=robustreg>.
- Koenker R (2012). *quantreg: Quantile Regression*. R package version 4.90, URL <http://CRAN.R-project.org/package=quantreg>.
- Leahy DA, Darbro W, Elsner RF, Weisskopf MC, Kahn S, Sutherland PG, Grindlay JE (1983). “On Searches for Pulsed Emission with Application to four Globular Cluster X-Ray Sources: NGC 1851, 6441, 6624, and 6712.” *The Astrophysical Journal*, **266**(1), 160–170.
- Li TH (2009). “A Robust Spectral Analyzer for One-Dimensional and Multi-Dimensional Data Analysis.” 2009/0112954 A1. US Patent Application.
- Li TH (2010). “A Nonlinear Method for Robust Spectral Analysis.” *IEEE Transactions on Signal Processing*, **58**(5), 2466–2474.
- Maronna RA, Martin RD, Yohai VJ (2006). *Robust Statistics: Theory and Methods*. John Wiley & Sons, Chichester.
- Mebane Jr WR, Sekhon JS (2011). “Genetic Optimization Using Derivatives: The **rgenoud** Package for R.” *Journal of Statistical Software*, **42**(11), 1–26.
- Norm DIN 66261 (1985). “Sinnbilder für Struktogramme nach Nassi-Shneiderman.”
- Oh HS, Nychka D, Brown T, Charbonneau P (2004). “Period Analysis of Variable Stars by Robust Smoothing.” *Journal of the Royal Statistical Society C (Applied Statistics)*, **53**(1), 15–30.
- Palmer DM (2009). “A Fast Chi-Squared Technique for Period Search of Irregularly Sampled Data.” *The Astrophysical Journal*, **695**, 496–502.
- R Core Team (2013). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. URL <http://www.R-project.org/>.
- Reegen P (2007). “SigSpec – I. Frequency- and Phase-Resolved Significance in Fourier Space.” *Astronomy & Astrophysics*, **467**(3), 1353–1371.
- Ripley R, Boitmanis K, Snijders TAB (2013). *RSiena: Siena – Simulation Investigation for Empirical Network Analysis*. R package version 1.1-232, URL <http://CRAN.R-project.org/package=RSiena>.

- Rousseeuw PJ, Croux C, Todorov V, Ruckstuhl A, Salibian-Barrera M, Verbeke T, Koller M, Maechler M (2012). *robustbase: Basic Robust Statistics*. R package version 0.9-3.
- Rousseeuw PJ, Yohai VJ (1984). “Robust Regression by Means of S-Estimators.” In J Franke, W Härdle, D Martin (eds.), *Robust and Nonlinear Time Series Analysis*, Lecture Notes in Statistics No. 26, pp. 256–272. Springer-Verlag, Berlin, New York.
- Salibian-Barrera M, Willems G, Zamar R (2008). “The Fast- τ Estimator for Regression.” *Journal of Computational and Graphical Statistics*, **17**(3), 659–682.
- Salibian-Barrera M, Yohai VJ (2006). “A Fast Algorithm for S-Regression Estimates.” *Journal of Computational and Graphical Statistics*, **15**(2), 414–427.
- Scargle JD (1982). “Studies in Astronomical Time Series Analysis. II. Statistical Aspects of Spectral Analysis of Unevenly Spaced Data.” *The Astrophysical Journal*, **263**, 835–853.
- Schwarzenberg-Czerny A (1989). “On the Advantage of Using Analysis Of Variance for Period Search.” *Monthly Notices of the Royal Astronomical Society*, **241**, 153–165.
- Schwarzenberg-Czerny A (1998). “The Distribution of Empirical Periodograms: Lomb-Scargle and PDM Spectra.” *Monthly Notices of the Royal Astronomical Society*, **301**(3), 831–840.
- Seber GAF, Lee AJ (2003). *Linear Regression Analysis*. 2nd edition. John Wiley & Sons, Hoboken, New Jersey.
- Stellingwerf RF (1978). “Period Determination Using Phase Dispersion Minimization.” *The Astrophysical Journal*, **224**, 953–960.
- Thieler AM (2013). “Robuste Verfahren zur Periodendetektion in ungleichmäßig beobachteten Lichtkurven.” Unpublished doctoral thesis, TU Dortmund University.
- Thieler AM, Backes M, Fried R, Rhode W (2013). “Periodicity Detection in Irregularly Sampled Light Curves by Robust Regression and Outlier Detection.” *Statistical Analysis and Data Mining*, **6**(1), 73–89.
- Timmer J, König M (1995). “On Generating Power Law Noise.” *Astronomy and Astrophysics*, **300**, 707–710.
- Tluczykont M, Bernardini E, Satalecka K, Clavero R, Shayduk M, Kalekin O (2010). “Long-Term Lightcurves from Combined Unified Very High Energy Gamma-Ray Data.” *Astronomy and Astrophysics*, **524**, A48.
- Uttley P, McHardy IM, Papadakis IE (2002). “Measuring the Broad-Band Power Spectra of Active Galactic Nuclei with RXTE.” *Monthly Notices of the Royal Astronomical Society*, **332**(1), 231–250.
- Varadhan R, Gilbert P (2009). “**BB**: An R Package for Solving a Large System of Nonlinear Equations and for Optimizing a High-Dimensional Nonlinear Objective Function.” *Journal of Statistical Software*, **32**(4), 1–26.
- Venables WN, Ripley BD (2002). *Modern Applied Statistics with S*. 4th edition. Springer-Verlag, New York.

Yohai VJ, Zamar RH (1988). “High Breakdown-Point Estimates of Regression by Means of the Minimization of an Efficient Scale.” *Journal of the American Statistical Association*, **83**(402), 406–413.

Zechmeister M, Kürster M (2009). “The Generalised Lomb-Scargle Periodogram. A New Formalism for the Floating-Mean and Keplerian Periodograms.” *Astronomy and Astrophysics*, **496**(2), 577–584.

Zhang Z, Chan SC (2005). “Robust Adaptive Lomb Periodogram for Time-Frequency Analysis of Signals with Sinusoidal and Transient Components.” In *Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP’05)*, volume 4, pp. 493–496. IEEE.

A. Implementation diagrams for RobPer

In this Appendix, the structure of the RobPer function is displayed as Nassi-Shneiderman diagram (structogram after Norm DIN 66261). Figure 10 contains a reading guidance for the blocks used in the structogram. The structogram for RobPer is displayed in Figure 11, for the algorithm singleFUN in Figure 12 and for the function IRWLS in Figure 13. The input and output variables of the latter are shown in Table 5. The following definitions are used:

$$\zeta_{L_2}(r) = \sum_{i=1}^n r_i^2 \quad (31)$$

$$\zeta_{LTS}(r) = \sum_{i=1}^{h(m)} r_{(i)}, \quad h(m) = \left\lfloor \frac{n}{2} \right\rfloor + \left\lfloor \frac{m+1}{2} \right\rfloor, \quad (32)$$

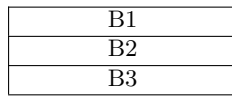
$$\zeta_{L_1}(r) = \sum_{i=1}^n |r_i|, \quad (33)$$

$$\rho_{MH}(\nu) = \begin{cases} \nu^2 & |\nu| \leq k \\ 2k|\nu| - k^2 & |\nu| > k \end{cases}, \quad \rho_{MB}(\nu) = \begin{cases} 1 - \left(1 - \left(\frac{\nu}{k}\right)^2\right)^3 & |\nu| \leq k \\ 1 & |\nu| > k \end{cases}, \quad (34)$$

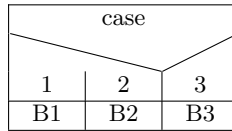
$$\zeta_{MH}(r) = \sum_{i=1}^n \rho_{MH}\left(\frac{r_i}{\hat{\sigma}}\right), \quad \zeta_{MB}(r) = \sum_{i=1}^n \rho_{MB}\left(\frac{r_i}{\hat{\sigma}}\right), \quad (35)$$

$$W_{MH}(\nu) = \begin{cases} c_{MH} & |\nu| \leq k \\ c_{MH} \cdot \frac{k}{|\nu|} & |\nu| > k \end{cases}, \quad W_{MB}(\nu) = \begin{cases} c_{MB} \cdot \left(1 - \left(\frac{\nu}{k}\right)^2\right)^2 & |\nu| \leq k \\ 0 & |\nu| > k \end{cases}. \quad (36)$$

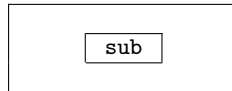
The normalization constant can be set to $c_{MH} = c_{MB} = 1$ due to the scale invariance of the least squares estimation used in the Iteratively Reweighted Least Squares (IRWLS) step.



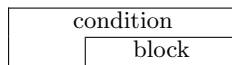
First run B1, afterwards run B2, at last run B3.
Horizontal lines between subsequent blocks are sometimes omitted for better readability.



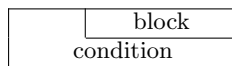
If case 1, run B1; if case 2, run B2; if case 3, run B3.



Run `sub` (some algorithm, code or function outsourced).

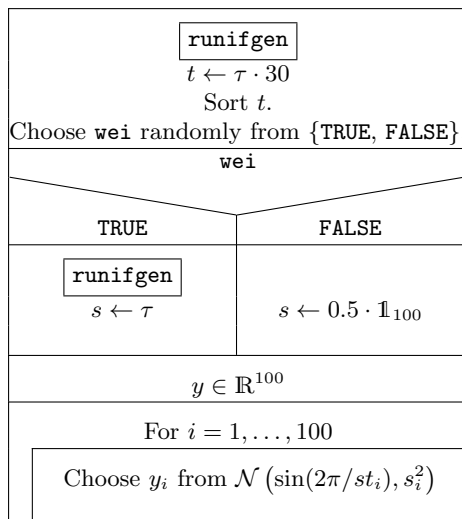


Reiteration of a block with a check in advance, whether a condition is fulfilled (e.g., a `for`-loop)



Reiteration of a block with a check afterwards, whether a condition is fulfilled (e.g., by `if(!...)...break`)

(a)



```
eval(parse(text = runifgen))
t <- tau * 30
t <- sort(t)
wei <- sample(c(TRUE, FALSE), 1)
if(wei) {
  eval(parse(text = runifgen))
  s <- tau }

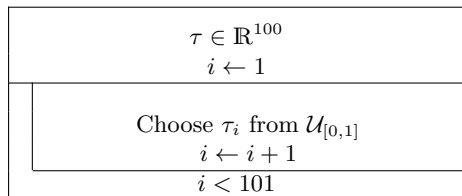
if(!wei) {
  s <- rep(0.5, 100)
}

y <- numeric(100)

for(i in 1:100) {

y[i] <- rnorm(1, mean = sin(2 * pi / 5 * t[i]),
sd = s) }
```

Block `runifgen`:



```
runifgen <- paste("
tau <- numeric(100)
i <- 1
repeat {
tau[i] <- runif(1)
i <- i+1
if(!i < 101) break }
")
```

(b)

Figure 10: Reading guidance for the structograms: (a) Blocks used for the representation of an algorithm. (b) Structogram (left) for a simple R code (right), which generates the observations $(t_i, y_i, s_i)_{i=1, \dots, 100}$ of a simple light curve with fluctuation period 5. This R code is for demonstration only and not programmed efficiently.

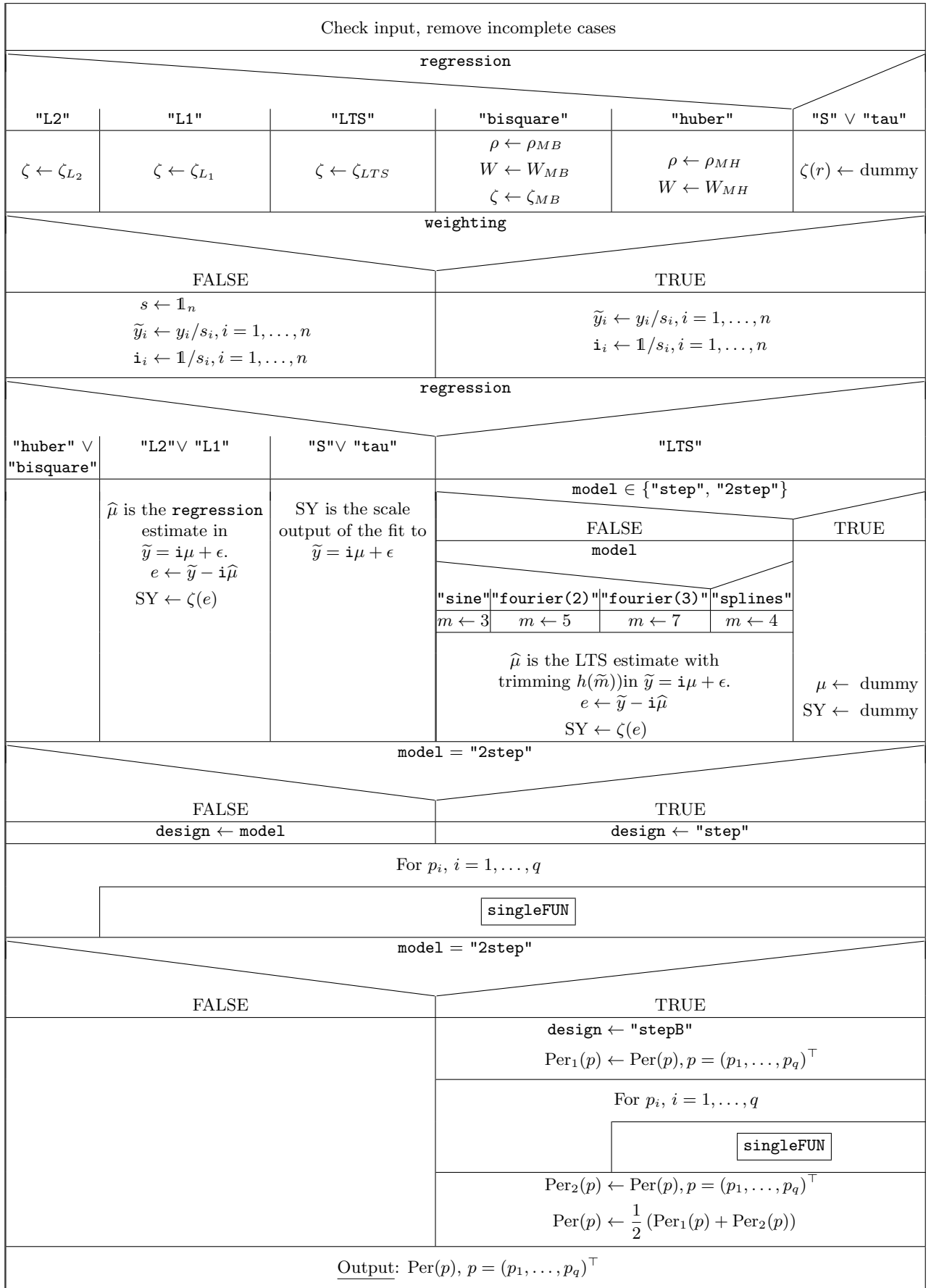


Figure 11: Structogram of RobPer. The block singleFUN is displayed in detail in Figure 12.

$X \leftarrow \boxed{\text{Xgen}(\text{model}, p_i)}$			
$\tilde{X} \leftarrow X/s$			
Enough independent rows in \tilde{X}			
FALSE	TRUE		
$\text{regression} \in \{\text{"huber"}, \text{"bisquare"}\}$			
TRUE		FALSE	
$\tilde{m} \leftarrow \text{number of columns of } \tilde{X}$ $\boxed{\text{ltsReg}(\dots, \text{nsamp}=50)}$ $\hat{\beta}$ is the LTS estimate with trimming $h(\tilde{m})$ in $\tilde{y} = \tilde{X}\beta + \epsilon$ $r \leftarrow \tilde{y} - \tilde{X}\hat{\beta}$		$\hat{\beta}$ is the regression estimate in $\tilde{y} = \tilde{X}\beta + \epsilon$ $r \leftarrow \tilde{y} - \tilde{X}\hat{\beta}$ $\text{SE} \leftarrow \zeta(r)$	
var1		$\text{regression} = \text{"LTS"} \wedge \text{design} \in \{\text{"step"}, \text{"stepB"}\}$	
FALSE	TRUE	FALSE	TRUE
$\hat{\sigma} \leftarrow \frac{\text{med}(r_j , r_j \neq 0)}{0.675}$	$\hat{\sigma} \leftarrow 1$	$\tilde{m} \leftarrow \text{number of columns of } X$ $\hat{\mu}$ is the LTS estimate with trimming $h(\tilde{m})$ in $\tilde{y} = i\mu + \epsilon$ $e \leftarrow \tilde{y} - i\hat{\mu}$ $\text{SY} \leftarrow \zeta(e)$	
$\hat{\mu}$ is the L_1 estimate in $\tilde{y} = i\mu + \epsilon$ $e \leftarrow \tilde{y} - i\hat{\mu}$		$\text{regression} = \text{"LTS"} \wedge \text{LTSopt} = \text{TRUE}$	
$\hat{\mu} \leftarrow \boxed{\text{IRWLS}(\text{yy} = \tilde{y}, \tilde{X} = i, W = W, \epsilon = e, \sigma = \hat{\sigma}, \text{tol})}$		FALSE	
$\text{regression} = \text{"bisquare"}$			
TRUE	FALSE	TRUE	
$\boxed{\text{genoud}}$ $r \leftarrow \tilde{y} - \tilde{X}\hat{\beta}$		$\boxed{\text{genoud}}$ $r \leftarrow \tilde{y} - \tilde{X}\hat{\beta}$ $\text{SE} \leftarrow \zeta(r)$	
$\hat{\beta} \leftarrow \boxed{\text{IRWLS}(\text{yy} = \tilde{y}, \tilde{X} = \tilde{X}, W = W, \epsilon = r, \sigma = \hat{\sigma}, \text{tol})}$		TRUE	
$\text{SE} \leftarrow \sum_{j=1}^n \rho \left(\frac{\tilde{y}_j - \tilde{x}_j \hat{\beta}}{\hat{\sigma}} \right)$ $\text{SY} \leftarrow \sum_{j=1}^n \rho \left(\frac{\tilde{y}_j - i\hat{\mu}}{\hat{\sigma}} \right)$			
$\text{Per}(p_i) \leftarrow \text{NA}$	$\text{Per}(p_i) \leftarrow 1 - \text{SE} / \text{SY}$		

Figure 12: Structogram of `singleFUN`. NA indicates a missing value. The block `IRWLS` is displayed in detail in Figure 13.

Input	Symbol	Explanations
$\mathbf{yy} \in \mathbb{R}^n$	\mathbf{yy}	Observed values
$\mathbf{matrix_} \in \mathbb{R}^{n \times m}$	\mathfrak{X}	Designmatrix
$\mathbf{w}: \mathbb{R} \rightarrow \mathbb{R}_{\geq 0}$	\mathbf{w}	Weight function
$\mathbf{residuals_} \in \mathbb{R}^n$	\mathbf{e}	Vector of residuals
$\mathbf{scale_} \in \mathbb{R}_{>0}$	σ	(Estimate of) Standard deviation
$\mathbf{tol} \in \mathbb{R}_{>0}$	\mathbf{tol}	Precision for convergence
Output		
$\mathbf{tempIRWLS\$coeff}$	$\hat{\mathbf{b}}$	Fitted vector of parameters

Table 5: Input and output variables of the function IRWLS.

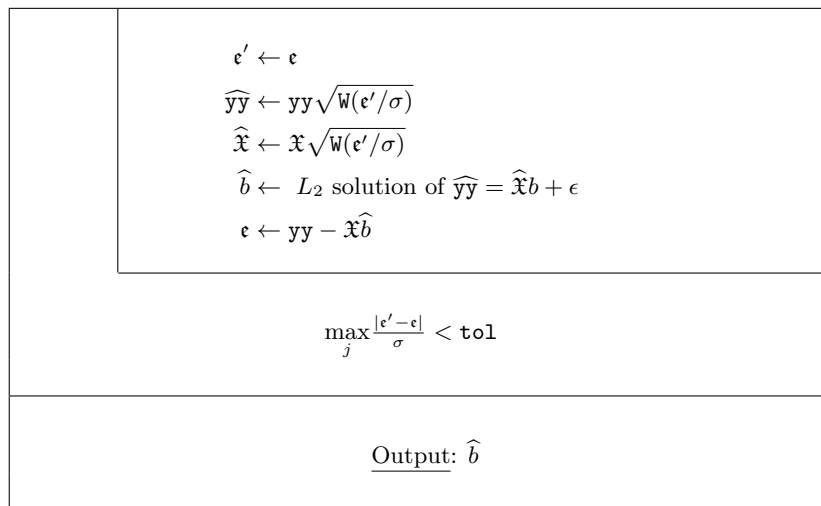


Figure 13: Function IRWLS in RobPer