

Beyond DC and MCMC: alternative algorithms and approaches to fitting light curves

A. Kochoska, K. Conroy, K. Hambleton and A. Prša

*Department of Astrophysics and Planetary Science, Villanova University,
800 Lancaster Avenue, Villanova PA 19085, USA*

Received: October 31, 2019; Accepted: December 16, 2019

Abstract. The parameter space of binary star light curve models is highly complex and degenerate, thus basic fitting approaches often fail to yield a good (and correct) estimate of the parameter values and their uncertainties. On the other hand, we have an increasingly large number of fitting and sampling algorithms available that can be relatively easily interfaced with open-source eclipsing binary packages, like PHOEBE 2. We showcase several fitting methods, including local and global minimizers, nested sampling and machine learning methods, and evaluate their performance on fitting a light curve model with PHOEBE 2.

Key words: binaries: eclipsing – methods: numerical – methods: statistical

1. Introduction

Robust fitting of the light- and radial velocity curves has been an outstanding issue in the field of binary stars for several decades now. As Prša & Zwitter (2005) showed, the parameter space of the binary star models is highly complex and degenerate, which leads to correlations between certain parameter and poses difficulties to finding the global optimum, as well as estimating the parameter uncertainties. With the advancements made in both technology and computing, we now have much better data and more precise models. The development of PHOEBE 2 (Prša et al., 2016) as a `python` package has opened up a new world of possibilities when it comes to fitting, since it can be easily interfaced with the many open-source optimizing packages that `python` offers.

In order to showcase the performance of several different approaches to fitting light curve data, we attempt to retrieve the true values of the parameters used to generate a synthetic light curve with PHOEBE 2 with added non-white noise. We fit for the mass ratio (q), inclination (incl), eccentricity (e), argument of periastron (ω), fractional radii (r_1, r_2), effective temperature of the primary ($T_{\text{eff},1}$), temperature ratio ($T_{\text{eff},2}/T_{\text{eff},1}$), passband luminosity (pblum) and third light ($l3$).

2. Minimizers

The simplest approach to fitting is through using an optimizing algorithm. Optimizing algorithms can be local, meaning they would typically find a local minimum close to the initial point, or global, which explore the parameter space stochastically and more robustly in search for a global minimum. We have applied four optimizing algorithms from the `scipy.optimize` library: Nelder-Mead Simplex (NMS), Powell's and L-BFGS-B (local) and differential evolution (DE, global). The minimization results are given in Fig. 1 and show that differential evolution outperforms the local minimizers. However, the price in accuracy is being paid by the computation time, which is significantly longer for global minimizers.

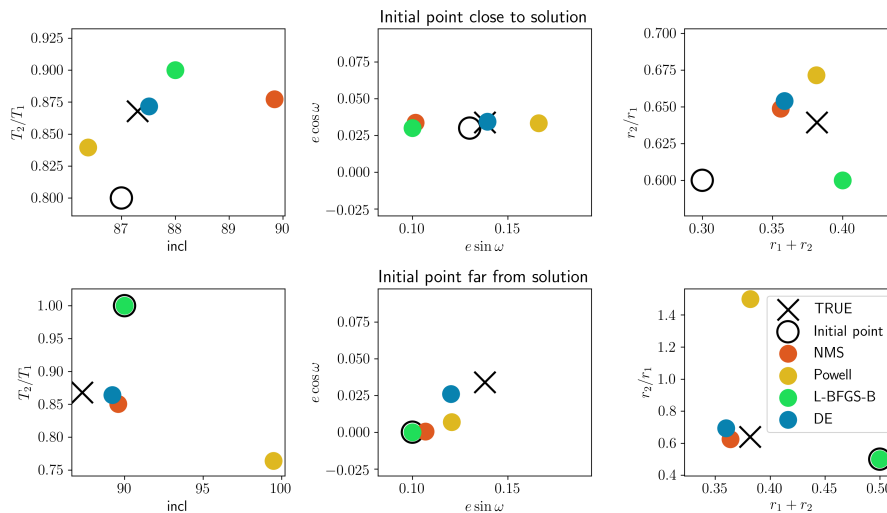


Figure 1. Cross-sections of the minimizer solutions in two-parameter space, with the chosen initial point relatively close to the true solution (top) and relatively far from the true solution (bottom).

3. Samplers

As we demonstrated in Section 2, global minimizers, at the cost of computational time, can yield a solution close to the global minimum. However, samplers are more robust in terms of exploring the topology of the parameter space around the global minimum and, thus, yield more reliable parameter uncertainties. Using MCMC for this purpose has become very common in our field, but it comes

with certain caveats. MCMC is not a search algorithm or optimizing algorithm (Hogg & Foreman-Mackey, 2018) and as such, only performs well once initialized close to the global minimum. Otherwise, we risk abuse of the algorithm at the cost of prohibitively long convergence times and incorrect interpretation of the results. Fig. 2 shows the results of sampling the parameter space of our synthetic light curve with the package `emcee` (Foreman-Mackey et al., 2013), with both bad use of MCMC (using it as a search algorithm for the global minimum) and good use of MCMC (using it to sample the posterior around the global minimum).

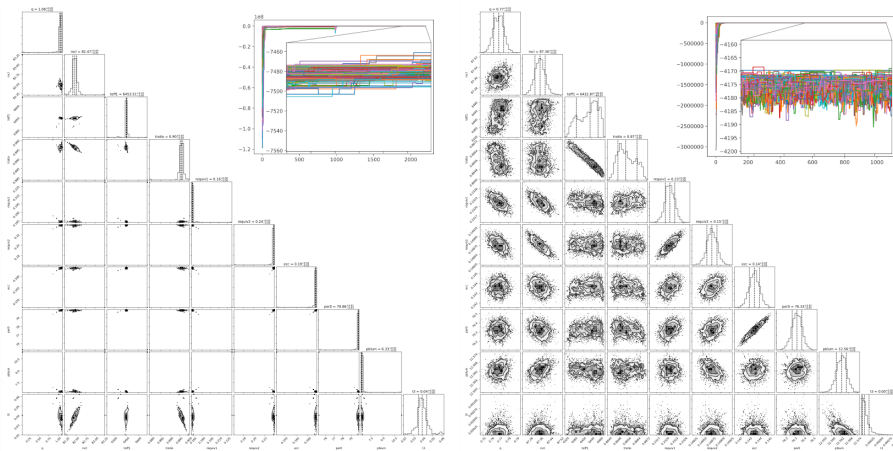


Figure 2. Examples of sampling the posterior with MCMC using `emcee`. Left: wide initial sampling range and use as a search algorithm to find the global minimum. The $\log p$ plot in the top right corner shows the solution is still converging. Right: initial sampling range in a tight ball around the global minimum. The $\log p$ plot in the top right corner shows oscillations around a constant value, which is a good indicator of a converged sampling.

Fortunately, other sampling algorithms can yield more robust results if we are not completely certain of the position of the global minimum. Based on our synthetic light curve test, nested sampling (Skilling, 2004) begins to reveal structure around the true global minimum in the likelihood of some parameters after several hundred iterations. Fig. 3 shows the trace plots of the position of the live points in all parameters for a nested sampling run with the package `dynesty` (Speagle, 2019). The sampling has not completely converged, but the values of the parameters that the light curve is sensitive to quickly move towards the global optimum.

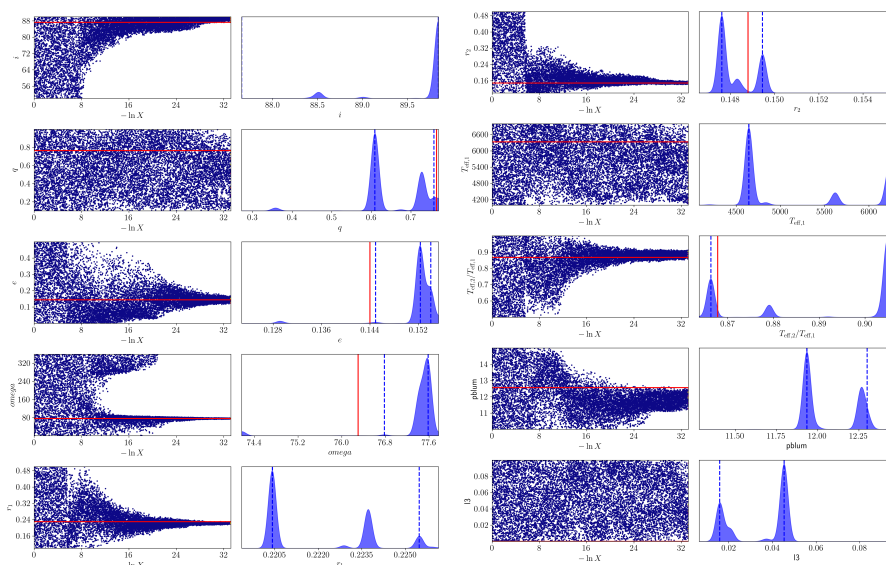


Figure 3. Trace plots of the positions of the live points in each parameter for a non-converged run of nested sampling with `dynesty`. The true parameter values are represented by red lines.

4. Machine Learning

Finding an initial solution for the light curve parameters usually involves a lot of manual work (adjusting the parameters and comparing the model with the data, initializing minimizers from different starting points, etc.). This approach becomes ineffective when dealing with large data sets. To address this, we have explored simple approaches using pre-computed synthetic databases and an algorithm based on nearest-neighbors search. Estimates of the model parameters found in this way are based on light curve similarity between the fitted light curves and a pre-computed database. The parameter estimates are computed as a distance-weighted mean (dw-mean) from the parameter values of the light curve’s nearest neighbors, while the range of possible values (min/max) is taken as the minimum and maximum of the parameter values across the nearest neighbors. This can be useful for providing the boundaries of the prior distributions used in MCMC or nested sampling.

Table 1 shows the parameter estimates from a nearest-neighbors distance-weighted computation for our test light curve. Because our data constrain the model well, our results are very close to the true parameter values. This is not always the case due to the parameter degeneracies in our model, but is a useful first step towards finding a better fit that eliminates the need for manual fitting.

Table 1. Distance-weighted estimates of the parameter values from a nearest neighbors search algorithm and their respective minimum and maximum values, compared to the true parameter values used to generate the test light curve.

	<i>min</i>	<i>dw-mean</i>	<i>max</i>	<i>TRUE</i>
q	0.5117	0.764	0.993	0.765
incl	80.915	85.925	89.996	87.3
r1+r2	0.323	0.397	0.489	0.3817
r2/r1	0.558	0.678	0.798	0.6394
Teff1	5013	6053	6992	6332
T2/T1	0.771	0.879	0.962	0.8678
esinw	0.014	0.145	0.246	0.1379
ecosw	0.009	0.031	0.052	0.03339
pblum	10.17	11.75	13.29	12.5664
l3	0.001	0.045	0.1	0

In some cases, the solution is not as well constrained because the nearest neighbors algorithm is not as sensitive to changes in certain parameters. To visualize this, we can use a dimensionality reduction technique, like t-SNE (van der Maaten & Hinton, 2008), to demonstrate the parameter value distributions across the parameter space in terms of light curve similarity. Fig. 4 showcases the value distributions of inclination, temperature ratio, mass ratio and fillout factor for a data set of synthetic contact binary light curves. It is clear that the light curve similarity is driven by the inclination and temperature ratio (top row), while the fillout factor and mass ratio values (bottom row) are distributed relatively uniformly across the map. Fitting for these two parameters with a nearest-neighbor search would thus yield results that are not as reliable as those for inclination and temperature ratio.

5. Conclusion

The best way of modeling eclipsing binaries remains a careful, hands-on approach. However, the era of big data is urging us to explore automated methods and advanced treatment of our data and model uncertainties. With the development of open source packages for modeling in wide-spread programming languages like `python`, a plethora of methods become available to us, but not all are suitable for our particular problem. The findings presented here are an initial step towards learning more about the fitting methods available to us and how to best apply them to our problem. Using a test synthetic light curve and PHOEBE2, we have demonstrated that: global optimizers outperform local optimizers at the cost of large computation times, simple machine learning approaches can give us initial estimates of the parameter values and potential prior ranges, while samplers come in many flavors, some of which are suitable

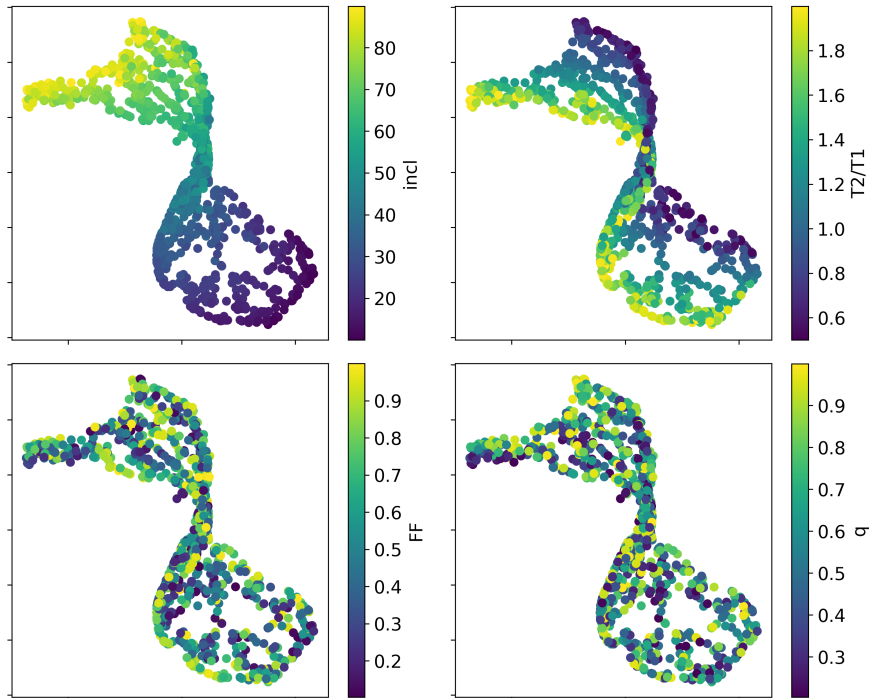


Figure 4. Parameter distributions over a 3D projection of the 4D parameter space in a contact binary light curve database generated with PHOEBE 2.

for posterior estimation near the global minimum, like MCMC, and others, like nested sampling, can reveal the underlying structure of the likelihood when we do not have a good estimate of the position of the global minimum.

References

- Foreman-Mackey, D., Hogg, D. W., Lang, D., & Goodman, J., emcee: The MCMC Hammer. 2013, *Publ. Astron. Soc. Pac.*, **125**, 306, DOI: 10.1086/670067
- Hogg, D. W. & Foreman-Mackey, D., Data Analysis Recipes: Using Markov Chain Monte Carlo. 2018, *Astrophys. J., Suppl.*, **236**, 11, DOI: 10.3847/1538-4365/aab76e
- Prša, A., Conroy, K. E., Horvat, M., et al., Physics Of Eclipsing Binaries. II. Toward the Increased Model Fidelity. 2016, *Astrophys. J., Suppl.*, **227**, 29, DOI: 10.3847/1538-4365/227/2/29

- Prša, A. & Zwitter, T., A Computational Guide to Physics of Eclipsing Binaries. I. Demonstrations and Perspectives. 2005, *Astrophys. J.*, **628**, 426, DOI: 10.1086/430591
- Skilling, J., Nested Sampling. in , *American Institute of Physics Conference Series*, ed. R. Fischer, R. Preuss, & U. V. Toussaint, Vol. **735**, 395–405
- Speagle, J. S., dynesty: A Dynamic Nested Sampling Package for Estimating Bayesian Posteriors and Evidences. 2019, *arXiv e-prints*, arXiv:1904.02180
- van der Maaten, L. & Hinton, G., Visualizing High-Dimensional Data Using t-SNE. 2008, *Journal of Machine Learning Research*, **9**, 2579