# Topographic Mapping of Astronomical Light Curves via a Physically Inspired Probabilistic Model

Nikolaos Gianniotis[1], Peter Tiňo[2], Steve Spreckley[3], and Somak Raychaudhury[3]

[1] Heidelberg Collaboratory for Image Processing,
University of Heidelberg D-69115 Heidelberg, Germany,
[2] School of Computer Science
The University of Birmingham
Edgbaston B15 2TT United Kingdom,
[3] School of Physics and Astronomy
The University of Birmingham
Edgbaston B15 2TT United Kingdom

**Abstract.** We present a probabilistic generative approach for constructing topographic maps of light curves from eclipsing binary stars. The model defines a low-dimensional manifold of local noise models induced by a smooth non-linear mapping from a low-dimensional latent space into the space of probabilistic models of the observed light curves. The local noise models are physical models that describe how such light curves are generated. Due to the principled probabilistic nature of the model, a cost function arises naturally and the model parameters are fitted via MAP estimation using the Expectation-Maximisation algorithm. Once the model has been trained, each light curve may be projected to the latent space as the the mean posterior probability over the local noise models. We demonstrate our approach on a dataset of artificially generated light curves and on a dataset comprised of light curves from real observations.

**Key words:** Topographic mapping, eclipsing binary stars

## 1 Introduction

The Generative Topographic Map algorithm (GTM) [1] has been introduced as a probabilistic analog to SOM [2], seeking to address certain of its limitations such as the absence of a cost function. The GTM formulates a mixture of spherical Gaussians densities constrained on a smooth image of a low-dimensional latent space. Each point in the latent space is mapped via a smooth non-linear mapping to its image in the high-dimensional data space. This image plays the role of the mean of a local spherical Gaussian noise model that is responsible for modelling the density of data points in its vicinity. The GTM can be readily

extended to structured data by adopting alternative formulations of noise models in the place of Gaussian densities. Such extensions have been proposed in [3] for the visualisation of symbolic sequences and in [4] for the visualisation of tree-structured data.

Here we present a further extension of the GTM to a novel data type, namely light curves that originate from eclipsing binary systems. Binary stars are gravitationally bound pairs of stars that orbit a common centre of mass. Astronomical observations suggest that almost half of the stars are binary ones. Thus, studying such systems procures knowledge for a significant proportion of stars. Binary stars are important to astrophysics because they allow calculation of fundamental quantities such as masses and radii, and are important for the verification of theoretical models for stellar formation and evolution. A particular subclass of binary stars are eclipsing binary stars. The luminosity of such stars varies over time and forms a graph called light curve. Light curves are important because they provide information on the characteristics of stars and help in the identification of their type.

## 2   Physical Model For Eclipsing Binaries

The physical model that generates light curves from eclipsing binary systems is described by the following set of parameters: mass $M_1 \in [0.5, 100]$ (in solar masses) of the primary star (star with highest mass of the pair), mass ratio $q \in [0, 1]$ (hence mass of secondary star is $M_2 = qM_1$), eccentricity $e \in [0, 1]$ of the orbit and period $\rho \in [0.5, 100]$ measured in days, all of which specify the shape of the orbit. Furthermore, two angles describing the orientation of the system are necessary [5] which are known as the inclination $\iota \in [0, \frac{\pi}{2}]$ and the argument of periastron $\omega \in [0, 2\pi]$ (see Fig. 1). Inclination describes the angle between the plane of the sky and the orbital plane and periastron is the angle $\omega \in [0, 2\pi]$ that orients the major axis of the elliptic orbit within its plane, that is $\omega$ is measured within the orbital plane. Finally, a third angle known as the longitude of ascending node ($\Omega \in [0, 2\pi]$) is necessary for the complete description of a binary system. However, since it has no effect on the observed light curves, we omit it from the model. We collectively denote these parameters by vector $\boldsymbol{\theta}$.

The mass $M$ of each star relates to the luminosity $L$ radiated by a surface element [6] of the star according to $L = M^{3.5}$ . Moreover, masses relate to the radii $R$ of the stars via:

$$R = \begin{cases} 10^{0.053+0.977\log_{10}(M)}, \text{ if } M < 1.728; \\ 10^{0.153+0.556\log_{10}(M)}, \text{ otherwise.} \end{cases} \tag{1}$$

These relations show that the primary star is the most luminous one and the one with the greatest area of the pair (a star appears as a disc to an observer). Thus, the *observed* area of a star is $A = \pi R^2$ and the *observed* luminosity is $L\pi R^2$. Henceforth, we index quantities related to the primary star by 1 (e.g. primary mass is $M_1$) and 2 for the secondary star.

It is shown from Newton's laws that the orbits of an object in the gravitational field of another object is a conic section of eccentricity $e$. Here we are interested in the case where $0 \leq e < 1$ that corresponds to closed orbits. We formulate two-body systems as systems where one body is fixed and the other is in orbital motion[4].

The position of the orbiting body is calculated by Kepler's equation as the distance $r$ from the fixed companion star on the elliptical orbit [5],

$$r(t) = \frac{a(1 - e^2)}{1 + e \cos \theta(t)},\tag{2}$$

where $t$ is time and $a$ is the semi-major axis of the ellipse calculated by Kepler's third law. Point $\Pi$ in Fig. 1 is the periastron, the point where the distance between the orbiting and fixed body is minimal. Angle $\theta$ is the angle between the radius and the periastron. Knowledge of $\theta$ would allow us to determine the position of the orbiting body. Angle $\theta$ is indirectly inferred via an auxiliary circle centered at the center of the ellipse $O$ and radius equal to semi-major axis. Point $Q$ is the vertical projection of the orbiting body's position $P$ to the auxiliary circle. Angle $E$ is called the eccentric anomaly and is given by Kepler's equation [5]:

$$E(t) = e \sin E(t) + \frac{2\pi}{\rho}(t - \tau),\tag{3}$$

where $\tau$ is the instance of time that the body was at the periastron. Kepler's equation does not admit an analytical solution but can be approximated through the Newton-Raphson method. By geometrical arguments it is shown that the relation between the true and eccentric anomaly reads:

$$\tan \frac{\theta(t)}{2} = [(1 + e)/(1 - e)]^{\frac{1}{2}} \tan(\frac{E(t)}{2})\tag{4}$$

By knowledge of $\theta$ we can determine the position of the second star on the orbit using (2) and (4). These positions correspond to the orbital plane and must be projected to the plane of the observer in the form of Cartesian coordinates [5]:

$$X(t) = r(t)(\cos(\Omega)\cos(\omega + \theta(t)) - \sin(\Omega)\sin(\omega + \theta(t))\cos(\imath)),\tag{5}$$

$$Y(t) = r(t)(\cos(\Omega)\cos(\omega + \theta(t)) + \cos(\Omega)\sin(\omega + \theta(t))\cos(\imath)),\tag{6}$$
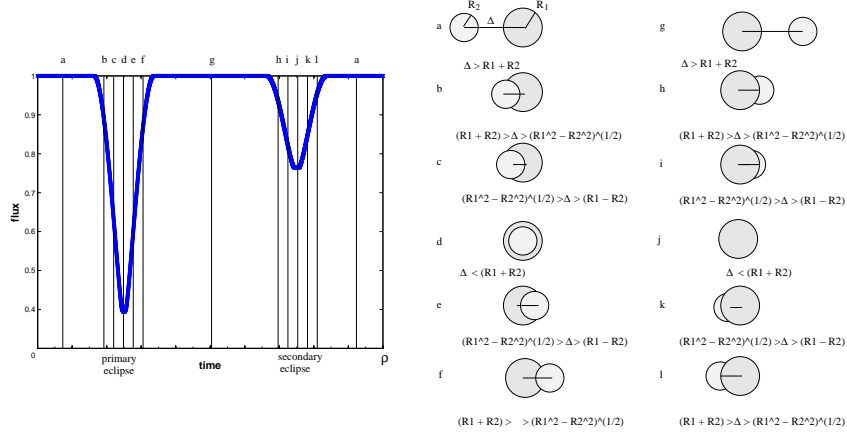
$$Z(t) = r(t)\sin(\omega + \theta(t))\sin(\imath),\tag{7}$$

which concludes the determination [5] of positions of the stars with respect to the observer.

---

[4] It is shown in [6] that in the relative motion system, the eccentricity, period and semi-major of the moving body's orbit are equal to their counterparts in the two-body system, and only the masses transform.

[5] The angle $\Omega$ does influence the position of the orbiting body. However, it does not have an influence on the light curve and thus we treat it as a constant $\Omega = 0$.

**Fig. 1.** Angles orientating the orbital plane with respect to the plane of sky, and angles associated with the orbits. Adapted from [5].

An observer of the binary system receives a variable luminosity from the eclipsing binary system that plotted against time forms a light curve. This variability is due to the eclipses that occur when one body passes in front (in the line of sight of the observer) of the other. This is illustrated in Fig. 2. When no eclipse occurs (positions $a, g$) the luminosity is equal to the sum of the luminosities radiated from the two bodies. The curved parts of the light curve occur due to partial occlusions. Two eclipses take place at each period, one primary eclipse (position $d$), when the most luminous body of the pair is obscured the most, and a secondary eclipse (position $j$), when the most luminous body obscures its companion the most.

Obscured parts of the disks of the stars can be calculated via geometrical arguments. [6] The obscured area of each star at time $t$ is denoted by $\Delta A_1(t)$ and $\Delta A_2(t)$. The luminosity $f_{\boldsymbol{\theta}}(t)$ received by the observer at time $t$ depends on the luminosities $L_i$, areas $A_i$ and obscured areas[7] $\Delta A_i$ via

$$f_{\boldsymbol{\theta}}(t) = L_1(A_1 - \Delta A_1(t)) + L_2(A_2 - \Delta A_2(t)). \tag{8}$$

---

[6] see http://www.physics.sfasu.edu/astro/ebstar/ebstar.html. Last access on 12-0-07.
[7] Recall that $i = 1$ and $i = 2$ index the primary and secondary stars, respectively

**Fig. 2.** Positions of stars (relative to observer's line of sight) and corresponding light curve phases.

## 3 Noise Model for Light Curves

Based on the physical model a probabilistic generative noise model arises naturally. Observed light curves, denoted by $\boldsymbol{O}$, are noisy signals:

$$\boldsymbol{O}(t) = f_{\boldsymbol{\theta}}(t) + \epsilon(t), \tag{9}$$

where $\epsilon$ is i.i.d. Gaussian noise with variance $\sigma^2$. Thus, we regard a light curve $\boldsymbol{O}$ of period $\rho(\mathbf{O})$ sampled at times $t \in \mathcal{T} = \{t_1 = 0, t_2, ..., t_T = \rho(\mathbf{O})\}$ as a realisation drawn from a multivariate spherical normal distribution. We denote the noise model associated with parameters $\boldsymbol{\theta}$ by $p(\boldsymbol{O}|f(.;\boldsymbol{\theta}), \sigma^2)$ or simply by $p(\boldsymbol{O}|\boldsymbol{\theta})$.

## 4 Model for Topographic Organisation

The starting point of our model formulation is the form of a mixture model composed of $C$ noise models as described in section 3:

$$p(\boldsymbol{O}|\boldsymbol{\Theta}) = \sum_{c=1}^{C} P(c) \, p(\boldsymbol{O}|\boldsymbol{\theta}_c), \tag{10}$$

where $P(c)$ are the mixing coefficients, $\boldsymbol{\Theta}$ encapsulates all parameter vectors $\{\boldsymbol{\theta}_c\}_{c=1:C}$ and $p(\boldsymbol{O}|\boldsymbol{\theta}_c)$ corresponds to the $c-$th model component with parameter vector $\boldsymbol{\theta}_c$. We simplify notation $p(\boldsymbol{O}|\boldsymbol{\theta}_c)$ to $p(\boldsymbol{O}|c)$. Assuming that dataset $\mathcal{D}$ contains $N$ independently generated fluxes $\boldsymbol{O}^{(n)}$, the posterior of the $\boldsymbol{\Theta}$ is expressed as:

$$p(\boldsymbol{\Theta}|\mathcal{D}) \propto p(\boldsymbol{\Theta}) \prod_{n=1}^{N} p(\boldsymbol{O}^{(n)}|\boldsymbol{\Theta}) = p(\boldsymbol{\Theta}) \prod_{n=1}^{N} \sum_{c=1}^{C} P(c) p(\boldsymbol{O}^{(n)}|c) \tag{11}$$

where the mixing coefficients can be ignored as $P(c) = \frac{1}{C}$.

Topographic organisation is introduced in the spirit of the GTM [1] by requiring that the component parameter vectors $\boldsymbol{\theta}_c$ correspond to a regular grid of points $\boldsymbol{x}_c, c = 1, \ldots, C$, in the two dimensional latent space $\mathcal{V} = [-1, 1]^2$. A smooth nonlinear function $\Gamma$ maps each point $\boldsymbol{x} \in \mathcal{V}$ to a point $\Gamma(\boldsymbol{x})$ that addresses a model $p(\cdot|\boldsymbol{x})$. Points $\Gamma(\boldsymbol{x})$ are constrained on a two-dimensional manifold $\mathcal{M}$ that is embedded in space $\mathcal{H}$, the space of parametrisations of our noise models. Since the neighbourhood of $\Gamma$-images of $\boldsymbol{x}$ is preserved due to continuity of $\Gamma$, a topographic organisation emerges for the models $p(\cdot|\boldsymbol{x})$. Function $\Gamma$ is realised as a RBF network [1]:

$$\Gamma(\boldsymbol{x}) = \boldsymbol{W}\boldsymbol{\phi}(\boldsymbol{x}), \tag{12}$$

where matrix $\boldsymbol{W} \in R^{6 \times K}$ contains the free parameters of the model (6 is the number of parameters in $\{M_1, q, e, \imath, \omega, \rho\}$), and $\boldsymbol{\phi}(.) = (\phi_1(.), ..., \phi_K(.))^T, \phi_k(.) : R^2 \to R$ is an ordered set of K nonlinear smooth basis functions. However, this mapping may produce invalid parameter vectors, since the output of the RBF network is unbounded. We therefore redefine mapping $\Gamma$ as:

$$\Gamma(\boldsymbol{x}) = \boldsymbol{A}g(\boldsymbol{W}\boldsymbol{\phi}(\boldsymbol{x})) + \boldsymbol{v}, \tag{13}$$

where:

- $g$ a vector-valued version of the sigmoid function that "squashes" each element in $[0, 1]$:

$$g(\boldsymbol{y}) = \left[ \frac{1}{1 + \exp(-y_1)}, \ \frac{1}{1 + \exp(-y_2)}, \ \cdots, \ \frac{1}{1 + \exp(-y_Y)} \right]^T, \tag{14}$$

- $\boldsymbol{A}$ is a diagonal matrix that scales parameters to the appropriate range. $\boldsymbol{A}$ has as diagonal elements the length of range $(\theta_i^{max} - \theta_i^{min})$ for each parameter, so that $A = diag((100-0.5), (1-0), (1-0), (2\pi-0), (\frac{\pi}{2}-0), (100-0.5))$.
- vector $\boldsymbol{v}$ shifts the parameters to the appropriate interval. $\boldsymbol{v}$ contains the minimum value $\theta_i^{min}$ for each parameter $\theta_i$: $\boldsymbol{v} = [0.5, \ 0, \ 0, \ 0, \ 0, \ 0.5]^T$.

The redefined mapping $\Gamma$ now takes a point $\boldsymbol{x}$ in space $\mathcal{V}$ to a valid parameter vector $\Gamma(\boldsymbol{x})$ that addresses a noise model in $\mathcal{M}$. Thus, $\boldsymbol{\Theta}$ has become a function of the weight matrix $\boldsymbol{W}$ of the RBF network, $\boldsymbol{\Theta}(\boldsymbol{W})$. Hence. the logarithm of the posterior from (11) now reads:

$$\log p(\boldsymbol{\Theta}(\boldsymbol{W})|\mathcal{D}) \propto \log p(\boldsymbol{\Theta}(\boldsymbol{W})) + \sum_{n=1}^{N} \log \sum_{c=1}^{C} p(\boldsymbol{O}^{(n)}|\boldsymbol{x}_c). \tag{15}$$
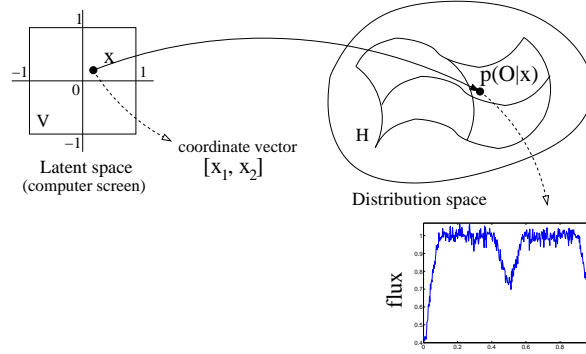
Figure 3 summarises the model formulation. Each point $\boldsymbol{x}$ of the visualisation space $\mathcal{V}$ is non-linearly and smoothly mapped via $\Gamma$ to model parameters that identify the corresponding noise model $p(\cdot|\boldsymbol{x})$. These parameters are constrained on a two-dimensional manifold $\mathcal{M}$ embedded in $\mathcal{H}$, the space of all possible

parametrisations of our noise model. In the spirit of [1], the model can be used to visualise observed fluxes $\boldsymbol{O}$ by calculating the posterior probability of each grid point $\boldsymbol{x}_c \in \mathcal{V}$, given $\boldsymbol{O}$:

$$p(\boldsymbol{x}_c|\boldsymbol{O}) = \frac{P(\boldsymbol{x}_c)p(\boldsymbol{O}|\boldsymbol{x}_c)}{p(\boldsymbol{O})} = \frac{P(\boldsymbol{x}_c)p(\boldsymbol{O}|\boldsymbol{x}_c)}{\sum_{c'=1}^{C} P(\boldsymbol{x}_{c'})p(\boldsymbol{O}|\boldsymbol{x}_{c'})} = \frac{p(\boldsymbol{O}|\boldsymbol{x}_c)}{\sum_{c'=1}^{C} p(\boldsymbol{O}|\boldsymbol{x}_{c'})}. \quad (16)$$

Each observed flux $\boldsymbol{O}$ is then represented in the visualisation space $\mathcal{V}$ by a point $proj(\boldsymbol{O}) \in \mathcal{V}$ given by the expectation of the posterior distribution over the grid points:

$$proj(\boldsymbol{O}) = \sum_{c=1}^{C} p(\boldsymbol{x}_c|\boldsymbol{O})\boldsymbol{x}_c. \quad (17)$$



**Fig. 3.** Formulation of the topographic mapping model.

We train our model in the MAP estimation framework with a physically motivated prior $p(\boldsymbol{\Theta})$ obtained from relevant literature [7,8,9,10]. To that purpose we employ the EM algorithm. Note that, due to the nature of the physical model formulation in sections 2 and 3, the M-step cannot be carried out analytically, nor can the derivatives of expected complete-data log-posterior with respect to the RBF network parameters $\boldsymbol{W}$ be analytically obtained. However, the EM algorithm does not necessarily require that an optimum is achieved in the M-step; it is sufficient that the likelihood is merely improved [11]. For our purposes we resort to numerical optimisation by employing a $(1+1)$ evolutionary strategy described in [12]. The fitness function for the evolutionary strategy is the expected complete-data log-posterior.

## 5    Experiments

### 5.1    Datasets

We performed experiments on two datasets. Dataset *1* is a synthetic dataset that consists of 200 light curves (fluxes). A common set of model parameters, $\{M_1 = 5, q = 0.8, e = 0.3, \imath = \frac{\pi}{2}\}$ was defined. However, two distinct values $\rho_1 = 2, \rho_2 = 5$ of period and $\omega_1 = 0, \omega_2 = \frac{5}{6}\pi$ of argument of periastron were used, to create 4 classes of light curves (50 in each class) by the combinations of these values, $\{\rho_1, \rho_2\} \times \{\omega_1, \omega_2\}$. The discerning characteristic of each class is the position of each secondary eclipse and the widths of the eclipses. Each light curve was then generated from these four "prototypical" parameter settings corrupted by a Gaussian noise. Gaussian noise was also subsequently added to the generated light curves to simulate observational errors.

Dataset *2* consists of light curves from real observations obtained from two resources available[8] on the WWW: the *Catalogue and Archive of Eclipsing Binaries* at http://ebola.eastern.edu/ and the *All Sky Automated Survey*. Dataset *2* was preprocessed before training using local linear interpolations. Preprocessing is necessary as one needs to account for gaps in the monitoring process and for overlapping observations. Light curves must also be phase-shifted so that their first point is the primary eclipse and resampled to equal length as described in section 3. Finally, the light curves were resampled at $T = 100$ regular intervals which was judged an adequate sample rate.
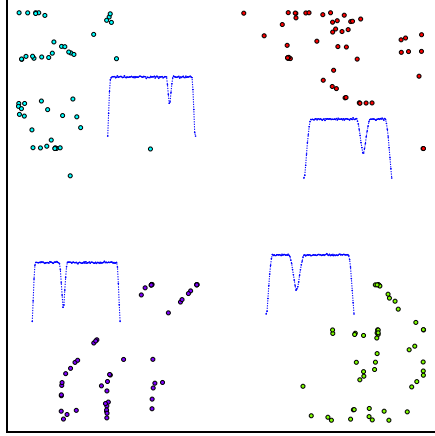
### 5.2    Training

The lattice was a $10 \times 10$ regular grid (i.e. $C = 100$) and the RBF network consisted of $M = 17$ basis functions; 16 of them were Gaussian radial basis functions of variance $\sigma^2 = 1$ centred on a $4 \times 4$ regular grid in $\mathcal{V} = [0,1]^2$, and one was a bias term. The variance of the observation noise in the local models $p(\boldsymbol{O}|\boldsymbol{x})$ was set to $\sigma^2 = 0.075$.

### 5.3    Results

Fig. 4 presents the topographic map constructed for the synthetic dataset. Each point stands for a light curve projected to latent visualisation space $\mathcal{V}$ and is coloured according to class membership. The class memberships of synthetic fluxes were not used during the training process. Also, next to each cluster, a typical light curve has been plotted. The classes have been identified and organised appropriately, each occupying one of the four corners of the plot.

Fig. 5 presents the topographic map constructed for the dataset of real observed light curves. The red curves are the data projected against the underlying local noise models displayed in black. Several interesting observations can be made about the topographic formation of the light curves on the resulting map.

---

[8] Last accessed on the 12th September 2007.

**Fig. 4.** Visualisation of synthetic dataset. A representative light curve is plotted next to each cluster.

In the lower right-hand corner binary systems of large periods are found. The median period of the systems in our sample is 2.7 days, and binaries like *V459 Cas*, with a period of 8.45 days lie in this corner. Systems with short period have the appearance of a wide V-shaped eclipse in the shape of their light curve, and inhabit the top and left edges of the map, e.g. WY Hya (Period: 0.7 days) and RT And (Period: 0.6 days). At the lower left of the map, we find systems with high eccentricity, e.g. V1647 Sgr. High eccentricity causes the light curve to appear assymetric, so that the period of the eclipse occurs further and further away from the center. On the other hand, very symmetric curves indicate orbits of low eccentricity (more circular) and low mass-ratio (stars of similar mass), and indeed we find systems like DM Vir ($e = 0.03$, mass ratio=1) and CD Tau ($e = 0.0$, mass ratio=1.05) in the cluster in the lower-right hand corner of the map. Finally, low-inclination systems, occupy the top left-hand corner of the map, and these orbits will have very shallow eclipses as the companion star barely eclipses the primary star.

## 6   Conclusions

We have presented a model-based probabilistic approach for the visualisation of eclipsing binary systems. The model is formulated as a constrained-mixture of physically motivated noise models. As a consequence, a clear cost function naturally arises which drives the optimisation of the model. In our experiments we have demonstrated that the resulting maps can be interpreted in a transparent way by inspecting the underlying local noise models. Furthermore, modification and refinement of the local noise models is possible, to account for greater physical fidelity by incorporating physical aspects for non-spherical stars and even more sophisticated phenomena such as gravity darkening.

**Fig. 5.** Visualisation of dataset *2* of real data. Light curves in red are the projected real data and light curves in black are the light curves of the underlying local noise models.

# References

1. Bishop, C.M., Svensén, M., Williams, C.K.I.: GTM: The generative topographic mapping. Neural Computation **10**(1) (1998) 215–234
2. Kohonen, T.: The self-organizing map. Proceedings of the IEEE **78**(9) (September 1990) 1464–1480
3. Tiňo, P., Kaban, A., Sun, Y.: A generative probabilistic approach to visualizing sets of symbolic sequences. In: KDD '04: Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining, ACM Press (2004) 701–706
4. Gianniotis, N., Tiňo, P.: Visualisation of tree-structured data through generative probabilistic modelling. In Verleysen, M., ed.: European Symposium on Artificial Neural Networks, D-Facto (2007) 97–102
5. Hilditch, R.W.: An introduction to close binary stars. Cambridge University Press (2001)
6. Karttunen, H., Krger, P., Oja, H., Poutanen, M., Donner, K.J., eds.: Fundamental astronomy. Springer-Verlag (1996)
7. Devor, J.: Solutions for 10,000 eclipsing binaries in the bulge fields of ogle ii using debil. The Astrophysical Journal **628**(1) (2005) 411–425
8. Halbwachs, J.L., Mayor, M., Udry, S., Arenou, F.: Multiplicity among solar-type stars. iii. statistical properties of the f7-k binaries with periods up to 10 years. Astronomy and Astrophysics **397** (2003) 159–175
9. Miller, G.E., Scalo, J.M.: The initial mass function and stellar birthrate in the solar neighborhood. Astrophysical Journal Supplement Series **41** (1979) 513–547

10. Paczyński, B., Szczygieł, D.M., Pilecki, B., Pojmański, G.: Eclipsing binaries in the All Sky Automated Survey catalogue. Monthly Notices of the Royal Astronomical Society **368** (2006) 1311–1318
11. Ng, S., Krishnan, T., McLachlan, G.: The em algorithm. In Gentle, J., Hardle, W., Mori, Y., eds.: Handbook of Computational Statistics. Volume 1. Springer-Verlag (2004) 137–168
12. Rowe, J.E., Hidović, D.: An evolution strategy using a continuous version of the gray-code neighbourhood distribution. In Deb, K., Poli, R., Banzhaf, W., Beyer, H.G., Burke, E., etal, eds.: Genetic and Evolutionary Computation – GECCO-2004, Part I. Volume 3102 of Lecture Notes in Computer Science., Springer-Verlag (2004) 725–736