# Forecasting the Outbursts of the Photometry Light Curve of Star V363 Lyr

Alexander Grigorievskiy[1], Maarit Mantere[1] Anton Akusok[1], Emil Eirola[1], and Amaury Lendasse[1,2,3,4]

[1] Aalto University, Department of Information and Computer Science,
PO Box 15400, FI-00076 Aalto , Finland
{alexander.grigorevskiy,anton.akusok,maarit.mantere,emil.eirola,amaury.lendasse}@aalto.fi
[2] IKERBASQUE, Basque Foundation for Science, 48011 Bilbao, Spain
[3] Department of Mechanical and Industrial Engineering, 3131 Seamans Center, The University of Iowa, Iowa City, IA 52242-1527, USA
[4] Arcada University of Applied Sciences, 00550 Helsinki, Finland

**Abstract.** In this paper we investigate the astronomical time series (TS) which is the photometric observations of the variable object V363 Lyr. We perform the spectral analysis of the time series and compare two approaches to forecast the outbursts of this time series. Since the data contain missing values we do the missing values imputation as well. The outbursts occur regularly, but our analysis shows that they are not strictly periodic. Hence, to improve the forecast of outburst position we compare several machine learning techniques for the two main forecasting approaches. Our results show that each approach can be beneficial depending on the starting point of forecast.

**Keywords:** Astronomical Time Series, OP-ELM, TROP-ELM, Recursive, Random Forest, K-NN, Mixture of Gaussians, Missing Data

## 1 Introduction and Problem Motivation

In this work, we study astronomical time series and we are interested in long-term forecasting. Data is obtained from the V363 Lyr star and described precisely in the subsequent paragraphs. Besides general forecasting accuracy, we are interested in forecasting outbursts (peaks) in the time series. In total there are 4201 data points, which are depicted in Fig. 1 (a). The data also contains missing values which prohibits the direct application of many time series forecasting methods. Therefore, as a preprocessing step, we perform the imputation of missing values by the method described in Section 2. Forecasting algorithms are described in Section 4.

The variable object V363 Lyr has been classified as a dwarf novae type variable based on its light curve properties [1], also the spectroscopic properties matching well those of dwarf novae [2]. A dwarf novae is a sub-type of cataclysmic variable stars, that is thought to consist of a close binary star system,

in which one of the components is a white dwarf that accretes matter from its companion. Dwarf novae exhibit outburst states that are thought to be rather regular over time. The mechanism of the excitation of the outburst is believed to result from an instability in the accretion disk, which causes an enhancement in the viscosity of the matter in the accretion disk. As a result of the disk becoming more turbulent and angular momentum being more efficiently transported outwards in the disk, the accretion rate increases, which causes the matter to collapse onto the white dwarf. This release of large amounts of potential energy is then seen as brightening of the object. This disk instability scenario, although quite generally accepted amongst astronomers, lacks its final proof, and competes with the mass-transfer outburst model, which explains the rapid accumulation of matter as being due to more mass suddenly been brought into the accretion disk due to some process occurring in the companion star.

From earlier photometric observations [3], a regular outburst cycle of roughly 22 days has been reported for V363 Lyr, associated with almost symmetrical rise and fade phases of the light curve. The increase in brightness from minimum to maximum has been measured to be of the order of 3 magnitudes or slightly more. Both the length of the outburst cycle and brightness variations roughly fall in the category of typical dwarf novae. The light curve shape and its pronounced stability, however, have been noted to the unusual for dwarf novae. Evidence for short-term periodicity was also looked for, but not found in earlier investigations. One of the problems in the earlier analysis was the rather poor time sampling of the observations (except very short subsamples), and additionally, the time span of the observations did not cover more than two outburst states. Therefore it may well be that all the characteristics of the outburst cycle of this object have yet not been detected due to these limitations.



(a) Original time series      (b) Example of the longest imputation

Fig. 2: Time Series

In this work, we study the dense photometry obtained of V363 Lyr with the KEPLER satellite. In the KEPLER database the object has an identifier KIC

7431243, and roughly 86 days (4201 samples) of recent photometry is available for the object. During this time, the object exhibits six outbursts, immediately casting doubt on the stability of the length of the outburst cycle detected in previous studies [3]. The outbursts do not appear always symmetric, and far less similar to each other than previously thought. Most notably, one super outburst occurred during the KEPLER observations. This dense KEPLER sample also enables us to look for higher frequency oscillations in the time series, and reconsider the possible origin of the outbursts in V363 Lyr.

## 2    Missing Values Imputation

There are several long intervals of missing values which are seen in Fig. 1 (a) and many single missing values. The longest missing interval of length 555 is not imputed, it and the data before it is ignored in the subsequent analysis. The remaining time series has a length of 3392 points and about 3% of the values are missing. The longest interval of missing values has a length of 57 points.

After this, the gaps in the time series are filled using a Gaussian mixture model, according to the procedure described in detail in [4]. Gaussian mixture is a very flexible distribution which can fit any other continuous distribution provided enough number of Gaussian components is taken [5]. The exact procedure is the following. A rolling window is used to extract sub-sequences of length $d$, $d = 100$ is used for the final imputation. The next step is a time-delay embedding, where each subsequence is interpreted as a point in $\mathbb{R}^d$. The coordinates are determined by the respective values of the time series. As an output of the method the matrix is obtained. Each row of it is one subsequence of length $d$. Because every missing value appears possibly in several places in this matrix we take their average as the final imputation.

A Gaussian mixture model can be fit to the data in the $d$-dimensional space by the EM algorithm, appropriately marginalizing over any missing values. Additional constraints are applied to ensure that the covariance structure of the mixture model is consistent with the autoregressive time series configuration. After convergence, the resulting model can be used to find the expected value of any missing value, conditional on the nearby known values. This conditional mean imputation is used to fill the gaps. The number of Gaussian components is chosen by trying different number of components and selecting the one with smallest BIC value.

The results of the imputation are not analyzed further in a quantitative way. However, visual investigation shows that the patterns existing in the time series are preserved after imputation. For example, imputation of the longest missing values interval of length 57 is shown on Fig. 1 (b).

## 3    Spectral analysis

Since the time series is taken from the astronomical domain we are also interested in a spectral analysis. Frequency content of the signal can provide astronomers

information about, for example, rotation frequency of the object and the accretion characteristics. Another issue to check is the stationarity of the time series with respect to constituent frequencies. The time series has been divided into three approximately equal parts, each containing two outbursts. We estimate the power spectral density (PSD) of each part and compare them in order to evaluate the stationarity of frequency content of the signal.



(a) Spectrum estimation by Welch method



(b) Frequency estimation by MUSIC method

Fig. 4: Spectral Estimation. Frequencies are relative to the sampling frequency.

Before estimating the spectrum, signal is normalized to have zero mean and no zero frequency in the spectrum. Power spectrum estimation by the Welch method [6, p. 415] is shown in Fig. 2 (a). The Welch method is a nonparametric spectral estimation method whose idea is similar to the basic periodogram method (Fourier transform), but the signal is divided into several overlapping windows, the periodogram of every window is computed and the final PSD estimation is the average of the window periodograms.

As we see from Fig. 2 (a), there are four main spectral components. The lowest one correspond to the periodicity resulted by outbursts. Then there is a double peak around the frequency 0.1 corresponding to the main frequency of oscillations. The highest frequency peak is slightly above 0.2. To precisely estimate frequencies, the MUSIC method [6, p. 463] for frequency estimation has been applied. Before applying the MUSIC method, the trend has been removed from the time series. The trend is calculated by the basic method of rolling averaging of time series. The length of the rolling window is 20. The trend removal is done in order to estimate more precisely the higher frequency components rather than frequency of outbursts. The idea of the MUSIC frequency estimation is to model the signal as a sum of sinusoids and using the eigenvectors of the autocorrelation matrix to compute the pseudospectrum. This pseudospectrum is shown in Fig. 2 (b). Low frequency content is absent in the pseudospectrum because of the trend removal.

Table 1: Frequencies estimated by the MUSIC method (relative to the sampling frequency).

|                        | 1st segment | 2nd segment | 3rd segment |
|------------------------|-------------|-------------|-------------|
| **1st peak**           | 0.103625    | 0.103875    | 0.103000    |
| **2nd peak**           | 0.109875    | 0.109625    | 0.109875    |
| **mean of first two peaks** | 0.106750 | 0.106750  | 0.106430    |
| **3rd peak**           | 0.209125    | 0.209375    | 0.209125    |

Frequencies estimation by MUSIC method are summarized in Table 1. For all three segments, frequencies are very close. So, the signal is quite stationary with respect to its frequency content. The corresponding period in hours for the frequency 0.103 is 4.76h and for the frequency 0.109875 is 4.46h. We believe that the double peaks near the frequency $f = 0.1\ Hz$ are a result of amplitude modulation of the signal. According to the elementary trigonometric formula

$$A \cos(2\pi f_m + \phi) \cos(2\pi f_c) = \frac{A}{2} \big[ \cos(2\pi(f_c - f_m) + \phi) + \cos(2\pi(f_c + f_m)) \big]$$

If the signal is amplitude modulated then we obtain two frequencies in the spectrum which are sum and difference of the the main frequency $f_c$ and the modulation frequency $f_m$. Hence, the main frequency equals the mean of the two close peaks and is given in the third row of Table 1. The third peak around 0.209 appears to be the second harmonic of the main frequency.

## 4 Forecasting methodology

### 4.1 Forecasting approaches

The main goal of this paper is to forecast the next outburst of the time series. We define that the outburst starts when the value of the time series exceeds 750. In the forecasting setup, we assume that the last outburst has just ended and the starting point of forecasting is the value in the *valley* of the time series. In other words, the starting point of forecasting can not be on the previous outburst (values higher than 750) it must be one of the time points where time series has *low value*. We have divided our time series into three sets: training, validation, and test, and this division is show in Fig. 3 (a). We want to estimate how well we can predict the outburst depending on the starting time point of forecasting. For instance, consider the validation set. First we assume that no points of validation set are observed and we try the forecast the next outburst. Then we assume that only one point of validation set is observed and again forecast the outburst. Finally, we assume that all points in validation set before the outburst are observed and we check whether our forecasting method can predict the outburst in the next point in time. The marks where the outbursts start and end on the training set are drawn in Fig. 3 (b).

There are two different approaches to forecasting the next outburst which we investigate in this paper:

(a) Data division into traininng, validation and test sets

(b) Outburst beginings and outburst ends for the training set

Fig. 6: Time series forecasting

– **Directly predict the time point where the next outburst happen.** This can done by building a regression model between the last observed time window of length $d$ (this is called regressor size) and the position of the next outburst. This approach we further divided onto two sub-methods:
  - Explicitly include in the model a variable whose value is the distance from the end of the previous outburst. If the time series is strongly periodic then this variable alone may provide a very good estimation of the next outburst.
  - Do not include this *periodicity variable* and build a regression model where regressors are just the previous values of the time series.
– **Conduct time series prediction and monitor where the predictions indicate an outburst.** At first this seems to be a more complex problem because we forecast not only the position of outburst but also the time series values. However, as shown later, this approach may be more beneficial than the first one under certain conditions.

### 4.2    Regression models

All aforementioned approaches intrinsically use some regression methods. We have tried several regression models starting from basic ones and proceeding to the state-of-the-art ones.

The most simple one is the *Linear model* which is linear regression with Tikhonov regularization. The regularization parameter is optimized by leave-one-out (LOO) cross-validation computed by Allen's PRESS statistic [7].

The second model is *Optimally-Pruned Extreme Learning Machine* (OP-ELM) [8]. This is a version of a single layer feed-forward neural network where weights of the hidden layer are not optimized but instead randomly generated. Randomly generated neurons are sorted by Least Angle Regression (LARS) and pruned by the LOO validation. The output weights are calculated by a least

squares method. It has proven to be a fast and accurate **nonlinear** regression model, and its application to time series prediction has been studied in [9]. Note that the output of the model varies from run to run because of the randomly generated weights, however average accuracy is very good. In the experiments we use 100 runs to evaluate forecasting with this model.

The third model is *Tikhonov Regularized Optimally-Pruned Extreme Learning Machine* (TROP-ELM) [10]. This is a variation of OP-ELM where the difference is that the additional Tikhonov regularization is used in the least squares calculation of output weights. The regularization parameter is again optimized by LOO with Press statistic.

The fourth model is Random Forest (RF) [11] which is a well known regression technique which may be considered as the current state-of-the-art. The output of this regression model also vary from run to run because random subsets of variables (and samples) are chosen to train each tree in the forest.

### 4.3   Forecasting the outburst position directly

Before presenting the result of direct outburst forecasting let us consider how to measure the quality of forecast. For both validation and test sets we know the position of true outburst. So, for every starting point of forecasting, a regression model predicts the position of the next outburst and we compute the absolute error between the prediction and the true outburst. Since the starting point roll from the beginning of the validation (or test) set up until one point before the true outburst, we average all the predictions. So, the error is the Mean Absolute Error (MAE). Because all regression models except linear involve randomness we repeat the experiment 100 times (20 for Random Forest) and again average the result.

Table 2: MAE of outburst position forecast

| Periodicity variable | Data | Linear Model | OP-ELM | TROP-ELM | Random Forest |
|---|---|---|---|---|---|
| Excluded | **Valid.** | 80.37 | $44.69 \pm (1.05)$ | $43.59 \pm (0.87)$ | $\mathbf{42.42 \pm (0.002)}$ |
|  | **Test** | 66.86 | $35.28 \pm (1.38)$ | $32.27 \pm (0.95)$ | $\mathbf{19.76 \pm (0.13)}$ |
| Included | **Valid.** | **3.27** | $3.41 \pm (0.05)$ | $3.62 \pm (0.05)$ | $3.53 \pm (0.01)$ |
|  | **Test** | **29.0** | $29.66 \pm (0.04)$ | $29.33 \pm (0.05)$ | $29.25 \pm (0.02)$ |

Results of direct outburst forecasting are presented in Table 2. Hyper parameters of all models are optimized on the validation set and the best values are used to forecast on the test set. For the OP-ELM and TROP-ELM models, the hyper parameter is the regressor size. The optimal value is 20 which equals the two periods of the main frequency of the time series. For Random Forests, hyper parameters include the regressor size (again 20), minimum number of samples in leafs and the number of randomly selected variables used to build trees. For

the linear model, the only hyper parameter is again the regressor size. It turns out that for the linear model which includes the periodicity variable regressor size = 0, so only the periodicity variable is used. For linear model without the periodicity variable regressor size = 5.

Looking at Table 2, one can notice that when the periodicity variable is included, the best model is the Linear model. This indicates that outbursts are quite periodic. However, test error is much larger than validation error because the last outburst (which belongs to the test set) is closer to previous outburst than the outburst distances in the rest time series. Hence, even though the outburst periodicity is quite apparent, it is not constant. When the periodicity variable is not included, the best model on validation and test sets is Random Forest. This agrees with the fact that it is a state-of-the-art regression model. It produces the lowest MAE on the test set – significantly lower than the linear model with periodicity variable. Therefore, Random Forest model without periodicity variable is compared in the following section with the approach when outburst predicted by forecasting time series values.

An interesting observation about the forecasts is that the accuracy is on average higher the closer the next outburst is. For instance, the MAE of the last 30 predictions before the outburst in the test set is only 12.65 for Random Forest. This is demonstrated in more details below.

### 4.4   Forecasting the outburst by time series prediction

In this approach, time series prediction is performed and when the predicted values become larger than 750 (outburst detection threshold) we assume that the outburst is starting. We can quantify the error between the true outburst position and forecasted position similarly as in the previous case - by absolute error of difference of two positions. However, it may happen that predicted values are never larger than the threshold as shown in Fig. 4 (b). Therefore, we must treat this special case separately. So, we assume that the outburst forecasting error is infinity and denoted by *inf*.

A typical outburst detection error vector is presented in Table 3. Different values there correspond to different starting points of time series forecasting. More precisely, the top left value correspond to the case when the prediction starts from the first point of the validation set (or test set). Then row-wise are the outburst errors for starting points which go further into the validation set. The last value in the table is the outburst error when we start prediction just one point before the true outburst. In the validation set the outburst happens at sample 270, so the are 269 values in Table 3. We can summarize this table by four numbers:

- **Total non-*inf* number**. This number count how many non-infinity numbers there are in the table.
- **Last non-*inf* sequence length**. This measures how many non-*inf* values in a row are at the end of the error vector. This is a measure of the accuracy and stability of the forecasting model because the closer we are to the true outburst the more accurate (in theory) we are able to forecast it.

(a) Outburst beginning is estimated (non-infinite error)

(b) Forecast does not indicate an outburst (infinite error)

Fig. 8: Successful and unsuccessful time series forecasts

– **Mean absolute error (MAE)**. This is the average absolute error (MAE) of the non-infinity values in the table.
– **First non-*inf* index**. This quantity measure how early we can forecast the outburst relatively correctly.

The tuple of these four numbers for Table 3 is (45, 33, 8.87, 157). The main drawback of the outburst detection by time series forecasting is that the set of non-*inf* MAEs is not sequential. This means that the outburst forecast started from earlier starting point may be more accurate (non-infinite) than the forecast obtained from later starting point. This situation is clearly seen in the table. However, the more the starting point is closer to the true outburst, the less infinite errors we have. Therefore, the most important parameter out of four is the first one, *Total non-inf number*, as it characterizes the forecasting as more stable and reliable.

Table 3: Absolute error of outburst location for TROP-ELM model on the validation set. Every value correspond to different starting point of forecasting.

```
inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf
inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf
inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf
inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf
inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf
inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf
inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf
inf, inf, inf, 12, inf, 20., inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf,
inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, 28, inf, 26, inf, inf, 27,  9,   5,  28,
inf, inf, 18, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, inf, 38, inf, inf,
inf, inf, inf, inf, inf, inf, inf, inf, 36, 39, inf, inf, inf, inf, inf, inf, 11,  1,   9,  10, 11,  9,
 9,  8,  7,  7,  3,  0,  3,  2,  3,  0,  4,  1,  1,  0,  1,  2,  1,  2,  1,  1,  1,  1,
 1,  1,  1,  1,  0
```

We have tried several time series prediction methods. Most of them are based on regressing the future values of a time series on the previous values. In particular, we construct a time delay embedding matrix of width 20 (20 was selected using the validation set) and for each time window we define the next time series value as a value to estimate. Using training data we make a regression model and then we use *Recursive prediction strategy* [9] (sometimes called rolling forecast) for long-term forecasting. There are other long-term forecasting strategies such as *Direct* and *DirRec* but again by experimenting on the validation set we observed that recursive strategy works the best for this time series. Regression models which are used for TS prediction are the same as in Section 4.3, namely *Linear*, OP-ELM, TROP-ELM and Random Forest.

In addition, K-Nearest Neighbor (K-NN) method for time series forecasting which represent another class of methods is taken for the comparison. It is a simple forecasting method which has shown decent accuracy in forecasting competition [12]. The idea of the method is that for the last available 20 values of the time series we search the closest segment (or several closest neighbors) with respect to some simple transformation (scaling and adjusting the mean). Then the forecast is the application of this transformation to the values following the closest segment. It is worth to mention that in contrast to other TS forecasting methods where we predict future values one by one, in K-NN method we predict values by batches of 20 values. Actually, 20 is taken by analogy with other methods, also 10 has been evaluated on validation set, but 20 performed better, so we leave only it in the following. Different number of nearest neighbors have been evaluated on the validation set but forecasting with one nearest neighbor performs the best, so only it is left for the comparison with other models.

The validation performance of outburst forecasting of various time series prediction methods is shown in the Table 4. For the methods which outcomes are probabilistic average values over 100 iterations are presented. Each model is characterized by four numbers which are described earlier. There is no model which is the best by all four parameters, but we suppose that TROP-ELM model is optimal because it has large average *Total non-inf number* and relatively good average MAE.

Table 4: Summary of outburst detection of different models on the validation set

| Models | Total non-*inf* number | Last non-*inf* sequence length | Total average error (MAE) | First non-*inf* index |
|:---:|---|---|---|---|
| **1-NN** | 167 | 12 | 67.80 | 6 |
| **Linear Model** | 269 | 269 | 65.29 | 0 |
| **OP-ELM** | 66.16 | 53.39 | 21.35 | 171.72 |
| **TROP-ELM** | **100.18** | **63.71** | **32.13** | **95.24** |

## 5    Comparison of direct outburst forecasting and forecasting by TS prediction

We have two different approaches for the forecasting the outburst. We have selected the best models from each approach based on validation results. Now we want to compare this two approaches on the test set. It is worth to mention again that the time series prediction approach can not provide outburst forecast at all for some starting points in which case can we assume that the forecasting error is infinite. So, the error we provide for this approach is calculated only for those values for which the forecast is available. The MAE of outburst forecasting is presented in Table 5.

Table 5: Average outburst MAE

|  | average MAE (over all starting pints) | average MAE (over last 30 starting points) |
|---|---|---|
| **Direct outburst forecasting (Random Forest)** | **19.76 ± (0.13)** | 12.65 ± 0.17 |
| **Forecasting outburst by TS prediction (TROP-ELM)** | 26.86 ± 5.84 | **4.14 ± 1.37** |

The difference between two error columns in the table is that in the first one average MAE over all the starting points is shown. In the last one over only the last 30 starting points which are adjacent to true outburst. From Table 5 we see that that direct forecasting of outburst position by Random Forest is more accurate than forecasting by time series prediction. However, the third column of Table 5 shows us that in the vicinity of the true outburst the time series forecasting model becomes more accurate with noticeable gap.

Hence, it is established that for the starting points which are relatively far away from the true outburst forecasting directly the outburst position works better. On the other hand, if the starting point is close to the true outburst then forecasting by time series prediction becomes more accurate. In real situation we don't know when the true outburst is, so we could use mixed approach. First forecast by the direct approach, if the predicted outburst is going to happen in the next 30 values then perform time series forecasting and use its result.

## 6    Conclusions

This paper studies the densely sampled time series of photometry of the V363 Lyr star, and consider it from the perspective of predicting the outbursts. The outbursts occur regularly, but our analysis shows that they are not strictly periodic. Only considering the time since that last outbursts leads to a decent prediction, but the various prediction strategies in Section 4 are able to improve

on this by including information of the current state of the time series. In particular, the autoregressive forecasting models can identify the time of the upcoming outburst with high accuracy when it is imminent. This fact, combined with the smaller fluctuations and slight increasing trend apparent in each valley between the outbursts, seems to support the disk instability mechanism believed to be responsible for the phenomenon.

Separately from the main outbursts, the spectral analysis in Section 3 identifies faster and lower amplitude component with a consistent period of about 4.5–5 hours. This fluctuation seems to originate from the angular rotation of the object.

In the Section 4 two approaches for outburst forecasting are investigated, namely direct outburst forecasting and forecasting by predicting values of the time series. From the first look the second approach seems less feasible because the problem is solved not directly but through predicting time series values. Indeed our analysis shows that the first approach is better when the true outburst is far away. However, we show that when the true forecast is within 30 time steps the second approach is preferable. To use the advantages of both methods we can combine them as written in Section 4.4.

# References

1. Hoffmeister, C. Astronomische Nachrichten **289** (1967) 205
2. Liu, W., Hu, J.Y., Zhu, X.H., Li, Z.Y.: Spectroscopic confirmation of 55 northern and equatorial cataclysmic variables. i. 27 confirmed cataclysmic variables. The Astrophysical Journal Supplement Series **122**(1) (1999) 243
3. Kato, T., Nogami, D., Baba, H., Masuda, S.: Outburst cycle of v363 lyr. Information Bulletin on Variable Stars **5118** (2001)  1
4. Eirola, E., Lendasse, A.: Gaussian mixture models for time series modelling, forecasting, and interpolation. Volume 8207 of Lecture Notes in Computer Science. (2013) 162–173
5. McLachlan, G.J., Peel, D.: Finite Mixture Models. Wiley Series in Probability and Statistics. John Wiley & Sons, New York (2000)
6. Hayes, M.H.: Statistical Digital Signal Processing and Modeling. 1st edn. John Wiley & Sons, Inc., New York, NY, USA (1996)
7. Allen, D.M.: The Relationship between Variable Selection and Data Agumentation and a Method for Prediction. Technometrics **16**(1) (February 1974) 125–127
8. Miche, Y., Sorjamaa, A., Bas, P., Simula, O., Jutten, C., Lendasse, A.: OP-ELM: Optimally-pruned extreme learning machine. IEEE Transactions on Neural Networks **21**(1) (January 2010) 158–162
9. Grigorievskiy, A., Miche, Y., Ventela, A.M., Sverin, E., Lendasse, A.: Long-term time series prediction using op-elm. Neural Networks **51**(0) (2014) 50 – 56
10. Miche, Y., van Heeswijk, M., Bas, P., Simula, O., Lendasse, A.: TROP-ELM: a double-regularized ELM using LARS and tikhonov regularization. Neurocomputing **74**(16) (2011) 2413–2421
11. Breiman, L.: Random forests. Machine Learning **45**(1) (2001)
12. Crone, S.F., Hibon, M., Nikolopoulos, K.: Advances in forecasting with neural networks? empirical evidence from the nn3 competition on time series prediction. International Journal of Forecasting **27**(3) (2011) 635–660