

TSARDI: a Machine Learning data rejection algorithm for transiting exoplanet light curves

Mislis D.^{1*}, Pyrzas S.¹, Alsubai K.A.¹

¹*Qatar Environment and Energy Research Institute (QEERI), Hamad Bin Khalifa University (HBKU), Qatar Foundation, P.O. Box 5825, Doha, Qatar*

Accepted ??, Received ??; in original form ??

ABSTRACT

We present TSARDI, an efficient rejection algorithm designed to improve the transit detection efficiency in data collected by large scale surveys. TSARDI is based on the Machine Learning clustering algorithm DBSCAN, and its purpose is to serve as a robust and adaptable filter aiming to identify unwanted noise points left over from data detrending processes. TSARDI is an unsupervised method, which can treat each light curve individually; there is no need of previous knowledge of any other field light curves. We conduct a simulated transit search by injecting planets on real data obtained by the QES project and show that TSARDI leads to an overall transit detection efficiency increase of $\sim 11\%$, compared to results obtained from the same sample, but using a standard sigma-clip algorithm. For the brighter end of our sample (host star magnitude < 12), TSARDI achieves a detection efficiency of $\sim 80\%$ of injected planets. While our algorithm has been developed primarily for the field of exoplanets, it is easily adaptable and extendable for use in any time series.

Key words: Extrasolar planets – transits – survey – algorithm.

1 INTRODUCTION

Large-scale, ground-based surveys for transiting extrasolar planets (e.g. HAT: Bakos et al. 2004; TreS: Alonso et al. 2004; SuperWASP: Pollaco et al. 2006; KELT: Pepper et al. 2007; QES: Alsubai et al. 2013) have been the steady work-horses of the field during the last 15 years. Almost since the beginning of these endeavours, it became readily apparent that the data collected were severely affected by *systematics*, i.e. unwanted flux variations introduced by fixed, ordered trends, such as airmass and seeing variations, colour-dependent extinction, object merging etc.

The answer to the problem came with the development of *detrending* algorithms, with TFA (Kovács et al. 2005) and SysRem (Tamuz et al. 2005) being among the most well-known. While both these two, as well as other similar detrending algorithms (e.g. Mislis et al. 2010; Ofir et al. 2010; Still et al. 2012; Mislis et al. 2017), can effectively remove (or at least minimise) the effects of major trends, they are not necessarily designed to tackle more subtle data irregularities that remain after detrending. Such irregularities can arise from infrequent and/or aperiodic events, e.g. the presence of cirrus, variations in atmospheric transparency and the presence of dust, variations in the sky background etc.

By design, large-scale surveys carry out long campaigns, observing their fields for a given time-period (mainly defined by the field’s visibility in the sky) and returning to them when next visible;

as such, field observations can span years, with considerable time gaps inbetween. Additionally, it is not uncommon for surveys to combine observations from different stations in a multi-longitude mode of observing. The longer a campaign lasts and the more data from different years and/or places are combined, the more susceptible light curves become to the irregular variations described above.

The net effect of these variations is mainly two-fold: (i) randomly distributed nights with higher RMS than the majority and (ii) in the absence of global flux calibration, nights with a mean flux level distinctly different from the overall light curve mean. While individually (i.e. from a single night) the effect on the overall light curve is most likely negligible, it can quickly escalate with additional nights and conceivably reach the 1% level of a typical transit; the end-result, when phase-folding the data to look for periodic transit signals, is a “puffed up” light curve, i.e. a light curve with an RMS higher than it *should* have, which can prove detrimental in identifying transit signals.

In recent years, Machine Learning (ML) algorithms have started becoming popular in a variety of research topics in Astrophysics, with the field of exoplanets prominent among them (e.g. Tornaiainen et al. 2008; Carrasco et al. 2014; Masci et al. 2014; Armstrong et al. 2016; McCauliff et al. 2015; Mislis et al. 2016; Armstrong et al. 2017, 2018). In this paper, we make use of ML, and develop a filtering algorithm, designed to tackle data irregularities that remain after the detrending process.

We present the TSARDI (TimeSeries Analysis for Residual Data Irregularities) algorithm; our approach is generally based on

* E-mail:dmislis@qf.org.qa

2 Mislis D.

the *class identification* methodology, i.e. the goal is to group the data points of a light curve into meaningful subclasses. To achieve this, we use the *clustering* algorithm DBSCAN (Density-Based Spatial Clustering of Applications with Noise), originally developed by Ester et al. (1996), as a more efficient and more effective way of discovering clusters of arbitrary shape, compared to other clustering algorithms (e.g. CLARANS and *k-means/k-medoid* partitioning algorithms). DBSCAN has been used by the K2 mission in order to optimize the photometric aperture size (Barros et al. 2016) and by *ASTERISM* (Tramacere et al. 2016) for galaxy detection and shape classification, but has otherwise attracted little attention in astronomical applications.

In Section 2 we describe the algorithm; Section 3 illustrates the effect of the algorithm on detrended light curves and describes the results of our simulated transit search; and Section 4 contains some concluding remarks.

2 THE ALGORITHM

In what follows, we will first give a brief overview of the DBSCAN algorithm and summarise the necessary definitions; interested readers are referred to Ester et al. (1996) for the complete, in-depth analysis. Subsequently, we will give a detailed description of TSARDI. We note here that TSARDI was built upon the DBSCAN routines as implemented in Python’s `scikit-learn` package¹ (Pedregosa 2011).

2.1 DBSCAN

The main function of DBSCAN is to take a *sample* of points S (in this case, a light curve) and organise all points in S into *clusters*, C . To achieve this, DBSCAN defines (i) a *distance function*², $dist(p, q)$, between points $p, q \in S$; (ii) an upper-limit/maximum value for the distance function, denoted as Eps ; and (iii) a minimum number of points $MinPts$. We can now proceed to the following definitions:

The Eps-neighbourhood of a point p , $N_{Eps}(p)$, is given by $N_{Eps}(p) = \{q \in S \mid dist(p, q) \leq Eps\}$

A point p is directly density-reachable from another point q , if the following two conditions are met: (1) $p \in N_{Eps}(q)$ and (2) $|N_{Eps}(q)| \geq MinPts$. In this case, i.e. when $|N_{Eps}(q)| \geq MinPts$, q is called a *core-point*.

If there is a chain of points $p_1 = q, \dots, p_n = p$ such that p_{i+1} is directly density-reachable from p_i , then p is density-reachable from q .

For three points p, q, w , if both p and q are density-reachable from w , then p and q are density-connected.

With the above definitions, we can define a *cluster* C as a non-empty subset of S , satisfying the conditions: (I) if $p \in C$ and q is density-reachable from p , then $q \in C$; and (II) $\forall p, q \in C$: p is density-connected to q . Finally, we note that points which do not belong to any cluster, are considered *noise points*.

An overview of these definitions is visualised in Figure 1. The

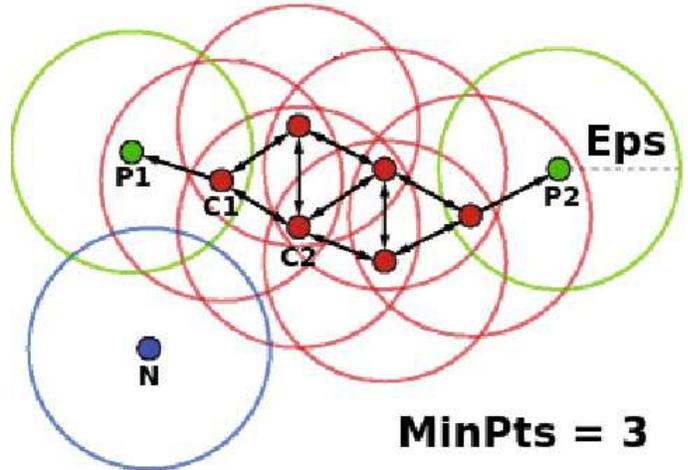


Figure 1. A schematic representation of the DBSCAN definitions. See text for details.

Eps-neighbourhood of each point is shown as a circle with radius Eps (shown as a dashed line), while $MinPts = 3$. All red points are core points, because $|N_{Eps}(q)| \geq MinPts$. Point C1 is directly density-reachable from point C2, because it belongs to the Eps-neighbourhood of C2 and C2 is a core point; the opposite is also true, i.e. point C2 is directly density-reachable from point C1. This pair-wise relation is indicated by the bidirectional arrow. Point P1 is directly density-reachable from point C1; however the opposite is *not* true. This is indicated by the single arrow. Point P1 is density-reachable from point C2 (because $C2 \rightarrow C1 \rightarrow P1$). Points P1 and P2 are density-connected, as they are both density reachable from e.g. C2. All red points together with points P1 and P2 belong to the *same* cluster. Finally, point N does *not* belong to the cluster and is considered a noise point.

It is obvious that the classification of a set into one, or more, clusters and (most importantly) the identification of those points that do not belong to *any* cluster, depends heavily upon the choice of the values for Eps and $MinPts$.

2.2 TSARDI

The TSARDI algorithm consists of two major parts: (1) the core-algorithm part and (2) the external-shell part. A detailed account of both parts follows.

2.2.1 The core-algorithm part

The core-algorithm part is based on four distinct, but chain-linked steps. Each step implements DBSCAN with step-unique distance function and values for Eps and $MinPts$. At each step, the target is to classify the input light curve points into one or more clusters and identify the *noise points*. These noise points are subsequently eliminated, and the resulting “filtered” light curve is used as input for the next step.

In what follows, we assume that our light curve consists of N pairs of time-and-flux values, $P_i = (t_i, f_i)$ with $i = 1, \dots, N$ and spans a total of K nights of observation.

STEP 1: The first distance function, df_1 , is the absolute flux difference between two consecutive points, that is $df_1(j) = abs(f_{j+1} - f_j)$ for $j = 1, \dots, N - 1$. Here, the light curve is treated as one “whole”,

¹ <http://scikit-learn.org>

² This can be any appropriate function of a given problem.

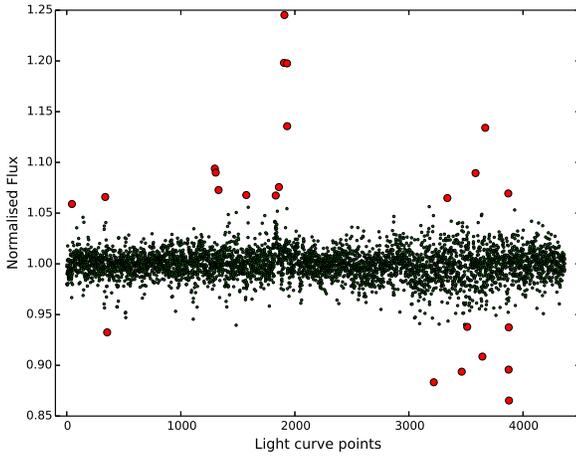


Figure 2. STEP 1 of our algorithm: the single “good points” cluster (smaller, green dots) and the noise points (larger, red dots), identified while treating the light curve as a whole, based on the median absolute flux difference between two consecutive points.

i.e. a continuous timeseries, and “consecutive” is used in an ordinal sense, so that the first point of one night is consecutive to the last point of the previous night. As Eps_1 , we set the median value of all calculated $df_1(j)$ times a *multiplication factor* mf_1 , that is $Eps_1 = mf_1 \times df_1(j)$; the function of mf_1 will become obvious in Sec. 2.2.2 and 2.2.3. $N/100$ is set as $MinPts_1$. Figure 2 visualises the first step.

STEP 2: We now split the light curve into its constituent nights; for each night, we calculate its mean flux value \bar{f}_k , with $k = 1, 2, \dots, K$. The second distance function, df_2 , is the absolute difference of mean fluxes between two consecutive nights³, that is $df_2(m) = abs(\bar{f}_{m+1} - \bar{f}_m)$ for $m = 1, \dots, K - 1$. As before, Eps_2 is set as the median value of all calculated $df_2(m)$ values times *another* multiplication factor mf_N , so $Eps_2 = mf_N \times df_2(m)$. As $MinPts_2$ we set $K/5$. Figure 3 visualises the second step.

STEP 3: This step is almost identical to the previous, only this time we calculate the *standard deviation* of the flux values of a given night, σ_k . The distance function, df_3 , is the absolute difference of standard deviations between two consecutive nights; Eps_3 is taken to be the median value of all $df_3(m)$ times mf_N , i.e. the *same* multiplication factor as in Step 2; and $MinPts_3$ is again $K/5$. Figure 4 visualises the third step.

STEP 4: The fourth, and final, step is similar to Step 1, in that it uses the absolute flux difference between two consecutive points, but in this case, the light curve is once more split into its constituent nights (as in Steps 2 & 3). For a given night k , with n points, we calculate $df_4^k(j) = abs(f_{j+1} - f_j)$ where $j = 1, \dots, n - 1$. Both the Eps and MinPts values are set on a *per night* basis, as the median of the corresponding $df_4^k(j)$ values times mf_1 (the multiplication factor of Step 1), and as $n/5$, respectively. Figure 5 visualises the fourth step.

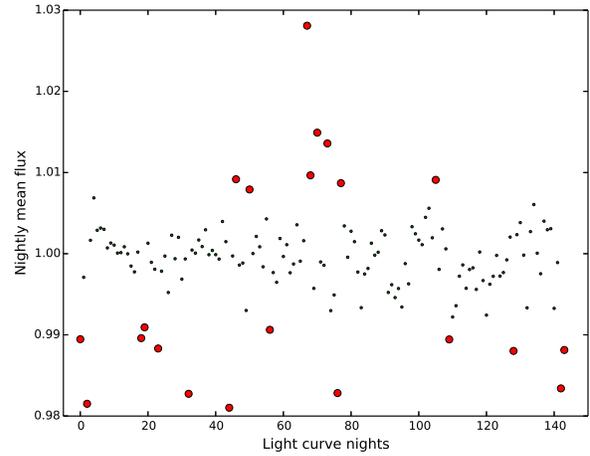


Figure 3. STEP 2 of our algorithm: the single “good points” cluster (smaller, green dots) and the noise points (larger, red dots), identified after binning the light curve per night, based on the mean flux value of a given night.

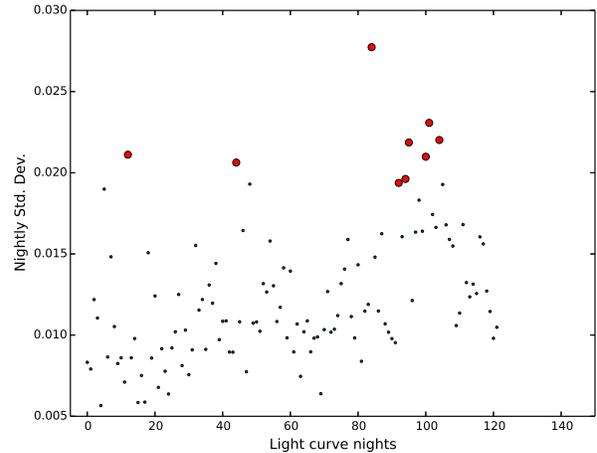


Figure 4. STEP 3 of our algorithm: similar to Step 2, but this time based on the standard deviation of a given night. Notice how the number of nights has decreased, after discarding some nights in the previous step.

In Figure 6 we plot both the original light curve of this example, as well as the resulting TSARDI-filtered light curve; in terms of numbers of points, the original and the final light curve consist of 4,367 and 3,636 respectively, i.e. the filtered light curve retains ~83% of the original number of points.

2.2.2 Large signals and over-rejection

A common caveat of clipping/rejection algorithms is the possibility of rejecting valid points (i.e. true signal) that are found far away from the majority of points in a light curve, as is the case in (deeply) eclipsing binaries and even “large planetary”-sized bodies (e.g. a brown dwarf transiting a late-K or an M-dwarf star).

For a straightforward implementation of the core-algorithm part of TSARDI with rigid Eps values in each Step (in other words, *without* the multiplication factors mf_1 and mf_N) there is a high

³ Here, again, “consecutive” is used in an ordinal sense.

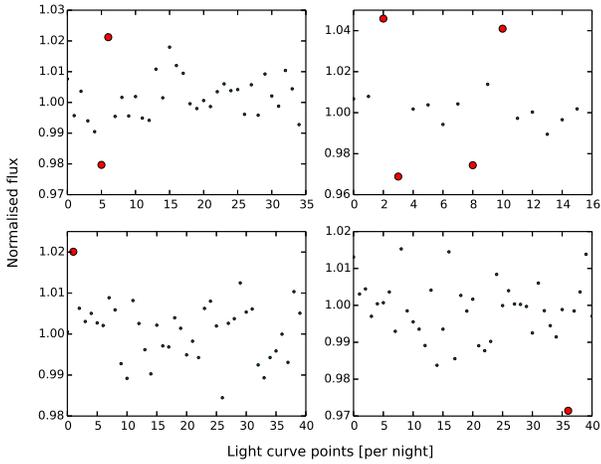


Figure 5. STEP 4 of our algorithm: four representative examples of the output are shown. In this step, each night is treated as a “mini-light curve”, and the identification of good points and noise points is based on the median absolute flux difference between two consecutive points.

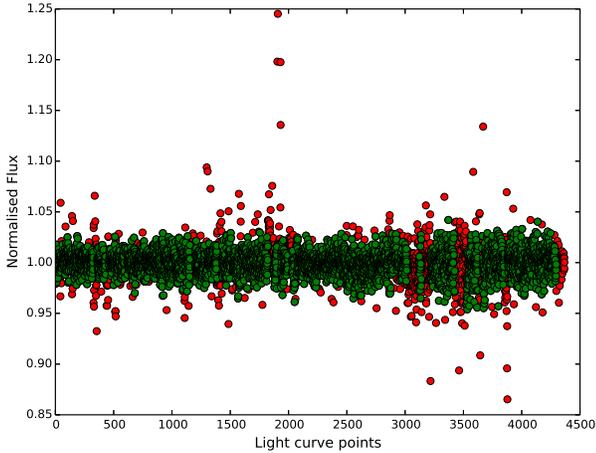


Figure 6. The original light curve (red points) and the TSARDI-filtered, final light curve (green points). For clarity, we plot the light curve as consecutive points, not according to their timestamps.

probability of in-transit points being classified as noise points and rejected from the final light curve, as illustrated in Figure 7.

While this could serve as a fast way to reliably remove large signals from a light curve and re-search the residuals for additional periodic signals, it could also have adverse effects on a survey looking for transiting candidates.

Dealing with the issue of large signals requires an assumption and a caveat. The assumption (in our opinion, quite justified) is that any sufficiently large signal (of the order of 3% and more) will be readily detectable with transit-detection algorithms, such as the BLS algorithm (Kovács et al. 2002), on the *detrended* light curve itself, without the need for any additional clipping or filtering. The caveat is that the presence of a large signal is, of course, not known beforehand, so that a transit-detection algorithm needs to run first.

In TSARDI, we implement a safeguard against over-rejection by setting a strict lower limit for the percentage of points in the fi-

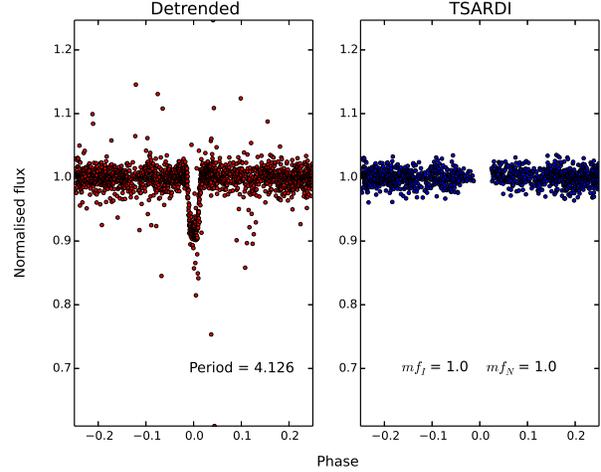


Figure 7. An example transit of a $2R_J$ object with an orbital period of 4.126 d, yielding a depth of 0.132. On the left-hand side panel, the detrended light curve. Running the core-algorithm part of TSARDI in a straightforward fashion, i.e. with $m_{f_I} = m_{f_N} = 1$ (effectively without multiplication factors) results in the in-transit points being classified as noise points and, therefore, rejected.

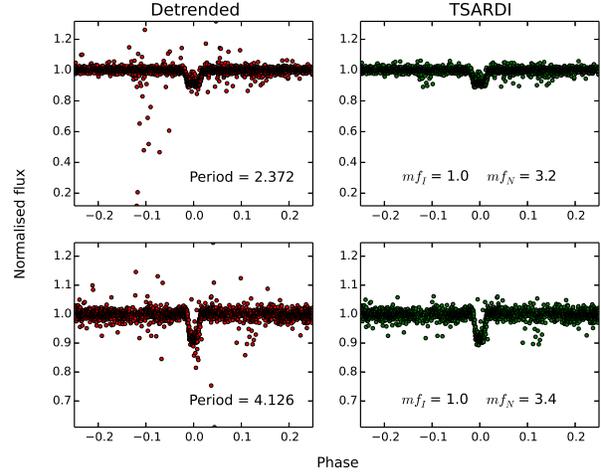


Figure 8. Two examples of transits of a $2R_J$ object (with orbital periods indicated), yielding depths of 0.088 and 0.132 (top and bottom, respectively). We show the detrended light curve on the left-hand side panels; and the final TSARDI filtered light curves on the right-hand side panels. We also indicate the corresponding (m_{f_I}, m_{f_N}) values.

nal light curve compared to the original number of points, that is, $NPTS_{LIM} = N_{fin}/N_{org}$. Depending on the tentative depth returned by the transit-detection algorithm, $NPTS_{LIM}$ is set to 82% for signals up to 2%; 95% for signals up to 5%; and 98% for signals larger than 5%.

The $NPTS_{LIM}$ is also the reason for the presence of the multiplication factors m_{f_I} and m_{f_N} first mentioned in Sec. 2.2.1, as they allow for easy adjustment of the Eps value at each Step of the core algorithm, depending on the current rejection rate compared to $NPTS_{LIM}$. Varying the (m_{f_I}, m_{f_N}) values allows TSARDI to retain deep signals, as indicated in Figure 8.

The interaction between all components is governed by the external-shell part of TSARDI, detailed below.

2.2.3 The external-shell part

The external-shell part serves as a wrapper for the core-algorithm part. At its centre, it hosts a double loop designed to run consecutive iterations of the core-algorithm part, while safeguarding against over-rejection. The external-shell (and by extension, TSARDI itself) runs in the following fashion:

- Two sets of values are defined for the multiplication factors m_{f_I} and m_{f_N} ; both sets range from 1 to 5, but with steps of 1 and 0.1 for m_{f_I} and m_{f_N} respectively.
- BLS is run on the detrended light curve, and $NPTS LIM$ is set according to the result.
 - A *first* loop begins for each m_{f_I} value.
 - A *second* loop begins for each m_{f_N} value.
 - For the given pair of (m_{f_I}, m_{f_N}) values, the core algorithm runs on the detrended light curve.
 - The percentage of remaining points (PRP) in the output light curve is recorded and compared against $NPTS LIM$.
 - If $PRP \geq NPTS LIM$ both loops break, otherwise the algorithm continues with the next pair of (m_{f_I}, m_{f_N}) values.

If the loops reach their end (i.e. $PRP < NPTS LIM$ for every pair of (m_{f_I}, m_{f_N}) values), then the maximum recorded value of PRP is compared to $NPTS LIM$. If the difference is less than 0.5%, TSARDI is re-run with that pair of (m_{f_I}, m_{f_N}) values giving $max(PR P)$; else, the detrended light curve remains untouched.

Finally, we should note that we have settled on these specific choices for the values of different variables, such as the range and step of m_{f_I} and m_{f_N} , the values for $NPTS LIM$ and the Step-unique values of $MinPTs$, for consistently yielding the best results based on extensive tests carried out on data by the Qatar Exoplanet Survey (QES, Alsubai et al. 2013). Some adjustment might be required to these parameters and/or the Steps themselves when applying TSARDI to different sets of data; for example, splitting the light curve into night segments isn't really applicable on continuous space-based data sets (but, perhaps, splitting into *some* form of segments is).

3 RESULTS

To assess the performance of the algorithm, we selected a field from the QES, observed with one of the 400 mm lenses ($f/2.8$, $FOV 5.24^\circ \times 5.24^\circ$). This particular data set was collected over a period of two years, from Jan. 2013 to Jan. 2015, and consists of ~ 4500 points, with an exposure time of 60 sec. The data were reduced with the QES pipeline, described in detail in Alsubai et al. (2013).

We limited the sample by imposing a cut on stellar magnitude of $V < 14$, resulting in 2022 stars. Following a similar procedure to the one described in Collier et al. (2007), for each star, we used the available V and K magnitudes, together with theoretical (and/or empirical) colour-temperature, temperature-radius and mass-radius relations to obtain a first estimate of the stellar masses and radii.

Subsequently, we injected a simulated transit signal of an $R_p = 1.0 R_J$ planet, generated using the Pál (2008) model, in all the *raw* light curves of the sample. We did this for two different orbital periods $P_1 = 2.37217$ d and $P_2 = 4.12669$ d. The transit ephemeris was chosen so that an adequate number of transits

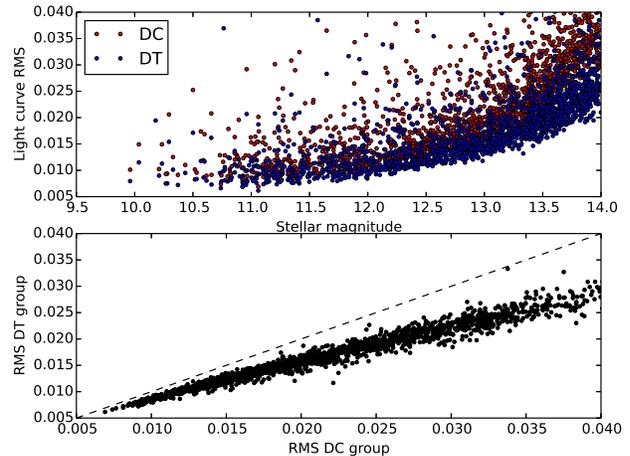


Figure 9. Top panel: RMS diagram for the DC (red) and the DT (blue) groups; see text for definition of the groups. Bottom panel: DC RMS versus DT RMS. In both panels, the overall RMS improvement is obvious, particularly for magnitudes below 12.0.

would be sampled in the given data set, to ensure a large number of detections for statistical purposes.

Subsequently, we detrended the light curves using the DOHA algorithm (Mislis et al. 2017) and processed them further with TSARDI; we will refer to these light curves as the *detrend & TSARDI*, DT group. We also created a “control” group, the DC group, by processing the detrended light curves with a more general sigma-clip algorithm, rejecting (i) points that were more than 8σ from the overall light curve mean; (ii) points that were more than 8σ from individual nightly means; and (iii) nights whose standard deviation was more than 5σ from the average standard deviation.

3.1 Overall RMS improvement

The first test illustrating the efficiency of our algorithm is a straightforward comparison of the light curve RMS between the DC and DT groups. In the top panel of Figure 9 we plot a typical RMS diagram for both groups. The overall RMS improvement is obvious, indicating that our algorithm not only clips obvious outliers (the scattered points in the upper left diagonal), but also that the additional steps of filtering out nights with comparatively high RMS can indeed improve the overall RMS of the main locus. This RMS improvement becomes more evident in the lower panel of Fig. 9, where we plot the RMS of the DC group versus that of the DT group.

3.2 Transit detection efficiency

The major test for our algorithm was to investigate whether it can indeed (positively) affect the transit detection efficiency. For that, each light curve was subjected to the BLS algorithm (Kovács et al. 2002); this was done for both the DC and the DT groups. As “successful recovery”, we consider the identification of the input planet period as the dominant peak in the BLS periodogram. The results for $P_1 = 2.37217$ d are shown in Figure 10.

In the top panel of Figure 10 we plot a histogram of the entire stellar sample in bins of 0.5 mag, together with the successful BLS

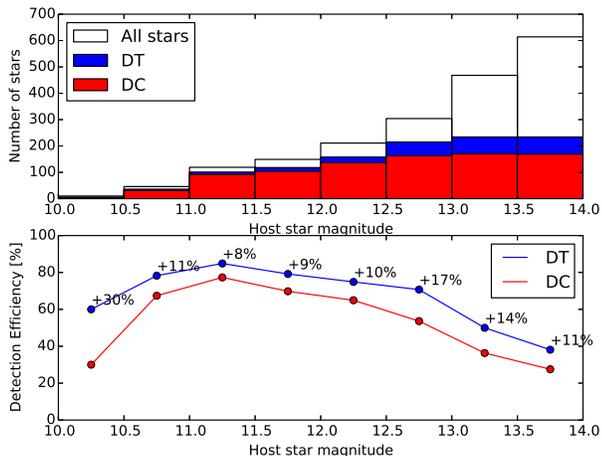


Figure 10. Top panel: Histogram of successful detections for the DC (red) and the DT (blue) groups per magnitude bin of the entire sample (white); **Bottom panel:** Transit detection efficiency per bin. The numbers above the DT points indicate the difference in efficiency from the corresponding DC points.

detections in both the DC and the DT groups. It is clear that the DT detections are more than the DC ones, in every bin. To quantify the improvement, we plot the detection efficiency in each magnitude bin in the lower panel of Fig. 10, where the numbers correspond to the actual percentage difference between the two groups in each bin.

In terms of absolute numbers, out of the possible 2022 planets, the DC group had a 45.3% overall success rate (916 planets) versus a 56.5% rate for the DT group (1142 planets). There were 22 unique detections for the DC group, and 248 unique detections for the DT group, resulting in a very favourable 11:1 ratio for the latter. Selecting the subsample of the moderately bright end ($V < 12.0$), out of the maximum 329 planets, the DC group had a success rate of 73.3% (241 planets), while the DT group had a success rate of 81.5% (268 planets). For the fainter end ($V > 12.0$), where the RMS improvement with TSARDI becomes readily obvious (see again Fig. 9), out of the maximum 1693 planets, the DC group had a success rate of 39.9% (675 planets), while the DT group successfully identified more than half the planets, with a 51.6% success rate (874 planets).

For a more detailed look into the workings of the algorithm, we plot in Figure 11 the expected transit depth D (based on our initial estimate of R_*) versus RMS for both the DC and the DT groups, differentiating between successful and unsuccessful detections. We also plot the $D = RMS$ and $D = 2 * RMS$ lines. It is evident that the majority of unsuccessful detections have small depth ($D < 1\%$) and large RMS ($RMS > 2 * D$); most of these stars have $V > 13.0$ and the photometric accuracy of the survey itself becomes the dominant factor. Notice again the improvement in light curve RMS and how much “tighter” the RMS of the DT group becomes.

To further illustrate the difference between the more “generic” sigma-clip algorithm and TSARDI, we once again plot in Figure 12 the expected transit depth D versus the light curve RMS, but this time, only for the 248 unique TSARDI detections. The ability of TSARDI to improve the overall RMS and pick up small signals (the majority of unique detections have transit depths less than 1%) is evident.

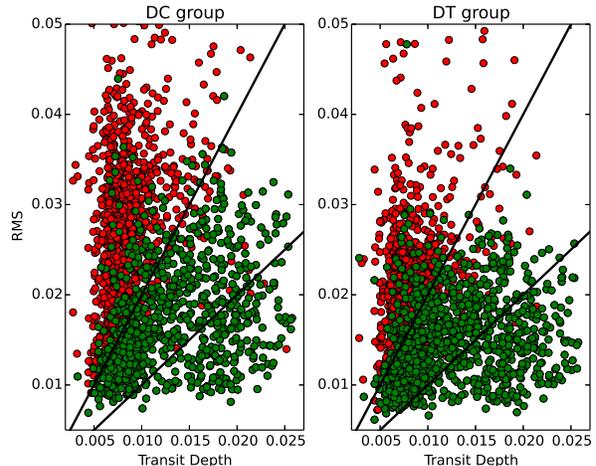


Figure 11. Transit depth versus RMS for the DC (left panel) and the DT (right panel) groups; successful and unsuccessful detections in both groups are plotted as green and red points, respectively. To aid the eye, we also plot the $D = RMS$ and $D = 2 * RMS$ lines. The vast majority of planets not detected have $D < 1\%$ and $RMS > 2 * D$, and correspond to the fainter end ($V > 13.0$) of the stars in our sample.

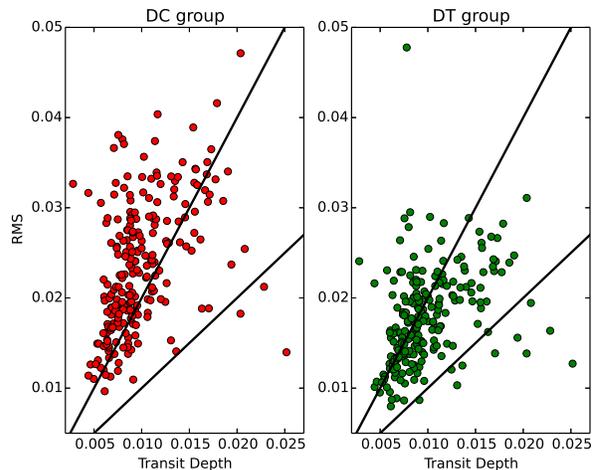


Figure 12. Same as Fig. 11, but collecting only the unique TSARDI detections. Note again the RMS improvement, and that the majority of these systems have transit depths smaller than 1%.

The results for the second orbital period, $P_2 = 4.12669$ d, are very similar, both qualitatively and quantitatively. The DC group had an overall success rate of 40.9% versus 51.9% for the DT group; for the moderately bright subsample, the success rates were 71.7% versus 77.5% for the DC and DT groups respectively; and finally the unique detection ratio was 9:1 in favour of the DT group (252 versus 29 planets).

4 CONCLUSIONS

We have developed TSARDI, a time-series analysis algorithm aiming to identify and remove residual data irregularities that remain in light curves, even after a detrending process. TSARDI is built

on the clustering algorithm DBSCAN, and uses the latter’s density-based notion, via an appropriately selected set of distance functions, to find outlying noise points; in our implementation, these noise points can be both “traditional” individual-point outliers, as well as individual *nights* that are distinct and differ significantly from the majority of nights in a long-term light curve.

Based on the results of a search for (simulated) transits on a real data set, we demonstrate that TSARDI can lead to a substantial improvement of the transit detection efficiency; compared to light curves filtered with a straightforward sigma-clip algorithm, TSARDI-processed light curves showed an overall increase of ~10% in the number of detections. Taking into account the accuracy of the data used, and limiting the sample in terms of host star magnitude ($m < 12$), leads to a detection rate of 80% after using TSARDI.

TSARDI was conceived and tailor-built to deal with light curves from ground-based, large-scale surveys of transiting exoplanets. However, due to the flexibility of DBSCAN’s density-based clustering, and with appropriate choices for the key algorithm parameters, it can be easily adapted and extended to essentially any time series.

ACKNOWLEDGMENTS

We thank the anonymous referee for insightful comments and suggestions, which improved not only the original manuscript, but the algorithm itself as well. This publication was made possible by NPRP grant # X-019-1-006 from the Qatar National Research Fund (a member of Qatar Foundation). The statements made herein are solely the responsibility of the author.

REFERENCES

- Alonso R. et al, 2004, ApJ, 613L, 153
 Armstrong D. J., Osborn H. P., Brown D. J. A., et al. 2014, MNRAS, 444, 1873
 Armstrong D. J., Kirk J., Lam K. W. F. et al. 2016, MNRAS, 456, 2260
 Armstrong, D. J.; Pollacco, D.; Santerne, A., 2017, MNRAS, 465, 2634A
 Armstrong, D. J. et al, 2018, MNRAS, in-press
 Alsubai K. et al. 2013, Acta Astron., 63, 465
 Bakos G., Noyes R. W., Kovács G., Stanek K. Z., Sasselov D. D. and Domsa I., 2004, PASP, 116, 266
 Barros S. C. C., Demangeon O. & Deleuil M., 2016, A&A, 594A, 100B
 Carrasco K., & Brunner R. J., 2014, MNRAS, 442, 3380
 Collier Cameron A. et al. 2007, MNRAS, 380, 1230C
 Ester, M., Kriegel, H. P., Sander, J., Xu, X., 1996, AAAI Press, KDD’96 Proceedings of the Second International Conference on Knowledge Discovery and Data Mining, pp. 226-231
 Kovács G., Zucker S. & Mazeh T., 2002, A&A, 391, 369
 Kovács G., Bakos G. and Noyes R., 2005, MNRAS, 356, 557
 Masi F. J., Hoffman D. I., Grillmair C. J., et al. 2014, AJ, 148, 21M
 McCauliff, S. D. et al, 2015, ApJ, 806, 6M
 Mislis D., Schmitt J. H. M. M., Carone L. et al. 2010 A&A, 522, 86
 Mislis D., Bachelet E., Alsubai K. A., et al. 2016, MNRAS, 455, 626
 Mislis, D., Pyrzas, S., Alsubai, K. A. et al. 2017, MNRAS, 465, 3759M
 Ofir A., Alonso R., Bonomo A., et al. 2010, MNRAS, 404L, 990
 Pál A., 2008, MNRAS, 390, 281
 Scikit-learn: Machine Learning in Python, Pedregosa et al., JMLR 12, pp. 2825-2830, 2011
 Pepper, J., Siverd, R. J., Beatty, T. G., et al. 2013, ApJ, 773, 64
 Pollaco, D. L. et al, 2006, PASP, 118, 1407
 Shallue, C. J., Vanderburg, A., 2018, AJ, 155, 94S
 Still M & Barclay T., 2012, Astrophysics Source Code Library, ascl.soft08004S
 Tamuz, O., Mazeh, T., Zucker, S., 2005, MNRAS, 356, 1466
 Tornaiainen I., Tornikoski M., Turunen M. et al. 2008, A&A, 482, 483
 Tramacere A., Paraficz D., Dubath P., et al. 2016, MNRAS, 463, 2939T

This paper has been typeset from a $\text{\TeX}/\text{\LaTeX}$ file prepared by the author.