

# Interesting Title...Stellar Flares in Kepler

Nicole Loncke & Lucianne Walkowicz

May 5, 2014

## 1 Introduction

The Kepler Mission is designed to survey our galaxy in the hopes of discovering planets in or near the habitable zone of their stars and determine how many of the billions of stars in our galaxy have such planets. It essentially stares at a relatively small portion of the sky for a long time to gather brightness data about these stars.

In addition to planet detection, this data can be used to gather properties about the stars themselves. In this paper we are concerned with the flaring behaviors of these nearby stars and their effect on any orbiting planets.

## 2 Using the Tools

In this section I describe the tools I’ve written to facilitate the vetting process.

### 2.1 Formatting

The `lightcurves` module makes some assumptions about the format of the input data. The original light curve data should be in a file containing a whitespace-separated table with time in the first column and flux in the second.

<i>808.51470</i>	<i>6338.22</i>
<i>808.53514</i>	<i>6340.73</i>
<i>808.55557</i>	<i>6346.89</i>
<i>808.57601</i>	<i>6341.10</i>
<i>808.59644</i>	<i>6340.22</i>

Table 1: Light curve data sampled from Kepler ID 10068383.

The other fundamental input file is the one containing the potential flare events, or “flags.” This file contains one column listing indices into the time array at the points that mark a suspected event.

Those two files — the light curve data and the flags — are all you need to start using these tools.

## 2.2 Plotting

If your aim is simply to generate arrays of the light curve data — one array for time, another for mean-normalized flux — then use `ltcurve()`. This function takes as its primary argument a string of the name of the file containing the Kepler data and returns the time and brightness arrays. By default it also displays the light curve corresponding to the file on a time vs. brightness plot, but this feature can be switched off by passing the function an optional argument.

If you have *multiple* light curve files and would like to view them one at a time, then `ltcurves()` is more appropriate. Its only required argument is a list or array of filename strings. Note that this function does not return any of the data. In addition, if for each of the Kepler data files you have a set of corresponding event flags<sup>1</sup>, you may use the `flags` kwarg to overplot the potential flares.

## 2.3 Vetting

Instead of cycling through the light curves with overplotted flags, you may find it helpful to inspect and record whether or not the marked events could potentially be stellar flares. In that case, you should use `flareshow()`, which writes user input (either 'y', 'n', 'm') to two files for later retrieval. One file contains a space-separated table of the Kepler IDs and the corresponding user responses to its events. The other file contains information about the length of each event. These two files work in conjunction to gather more information about the potential flares.

```
8848271  n
8908102  n
8953257  n n n n n n n
9002237  n n n y
```

Table 2: Example output.txt file.

```
8848271  3735 03
8908102  1757 03
8953257  1454 6 1610 7 1890 4 2359 3 2516 4 2829 5 2985 6 3265 5
9002237  3337 4 3547 5 3756 3 3967 4
```

Table 3: Corresponding example output.indices.txt file.

Note that before using `flareshow()`, you must have your flags in the proper format, generated by `getflags()`. This helper function outputs a nested list of event indices given a list of the names of the files containing the flags.

---

<sup>1</sup>You can generate these flags using `getflags()` and passing it a list of the names of the files holding the flare flags.

### 3 Data Processing

After evaluating the marked events by eye, quantifying data about the remaining candidates is the next step. Assuming that you used `flareshow()` for vetting, you now have two output files for the set of light curve data. Use the helper function `getEvents()`, which reads from these output files to pare down your list of flags to only those that have been marked with 'y' or 'm' depending on how you set the kwargs.

To calculate the cumulative brightness of each event found within a single light curve, use `intFlare()`. This function returns an array of the integrated brightness over the course of the events, an array of the duration of the events (in hours), and the peak brightnesses of each event. It is important to note that this function assumes you have vetted the flags already and are only providing those about which you would like to find more information (ie, you have used `getEvents()`).

### 4 Machine Learning

After creating the tools for by-eye vetting, the major bottleneck for data analysis was just that—as a human, vetting manually takes a lot of time. The flare detection program that produces the flare flags is trained to recognize some data metrics but does not correctly identify events with high accuracy. Our human brains allow us to do the same thing but with more nuance. If we could write a program to recognize the same patterns that humans so easily detect in the light curves then we could nearly entirely automate the vetting process. To accomplish this goal, we trained a classifier with a handful of metrics from each potential flare event and their respective light curves.

#### 4.1 Training

Our first task was to gather quantitative data about the stellar flares to feed into the classifier. In total we use 10 metrics.

- (a) *amplitude*: the range of the entire light curve. Stars with great stellar variability tend to be more magnetically active than those without. We expect high light curve amplitude to correlate with real flares.
- (b) *number of events*: Light curves that have many flagged events tend to have real flares, so we expect a high number of events to correlate with real flares.
- (c) *standard deviation*: The standard deviation of the entire light curve with stellar variability subtracted. **Though I'm not exactly sure of the direct bearing of this metric on the likelihood of the event, it seemed like it could be a useful metric.**
- (d) *consecutive points*: Sometimes there are gaps in the Kepler data. Kepler must rotate and point its antenna towards Earth to send its light curve data roughly every month. When the satellite begins recording again, there may be a sudden increase in brightness that resembles a flare but isn't. In order to avoid marking these as true flares we check whether the time intervals are evenly spaced across the event.

- (e) *kurtosis*: The kurtosis measures the “peakedness” of a flare event. A sharp increase and decrease in brightness is likely to indicate a true flare, though the decay ought to be more gradual than the incline.
- (f) *midpoint check*: A stellar flare typically requires a monotonic increase then monotonic decrease in brightness. Ensuring that the middle point is higher than the beginning and end points of the event is one way to rule out falsely marked events.
- (g) *second derivative*: Smoothing over the flagged event, is the light curve locally concave up or down? The second derivative of the window around the potential flare can capture the shape of a light curve in the neighborhood of an event.
- (h) *skew*: Skewness is a measure of the asymmetry of the event brightness—is the flare left-leaning or right-leaning? Because flares are characterized by very quick increases in brightness followed by a slow decay, left-leaning events (and therefore those with negative skew) are more likely to be true flares.
- (i) *slope*: Is the brightness of the star generally increasing or decreasing at the time of the event? This metric measures the slope of the line formed by connecting the point at the beginning of the flare window to the point at the end of the flare window. time of the event?
- (j) *slope ratio*: We also compute the ratio of the light curve’s slope just before the event begins and the slope just after it ends. We hope to capture more information about the local shape of the light curve with this metric.

These data were gathered for potential flaring events that we labelled by-eye to form a training set.

## 4.2 Initial Classification Performance

We tried a few different classification techniques. We used Python’s `sklearn` package for our machine learning framework.

### 4.2.1 Support Vector Classification

For our first attempt we used support vector classification as packaged in `sklearn.svm.SVC`. We initially used a radial basis function (RBF) kernel.

We also tried using a linear kernel, but this performed worse.

### 4.2.2 Random Forest Classifier

Next we attempted the same task using random forest classification, as packaged in `sklearn.ensemble`. While it performed superbly on the training set, it performed about average on the training set.

	precision	recall	f1-score	support
n	0.73	0.81	0.77	57
y	0.74	0.89	0.81	61
m	0.71	0.31	0.43	32
avg / total	0.73	0.73	0.71	150

Table 4: Reconstructing the training set (150 flares total) with RBF kernel.

	precision	recall	f1-score	support
n	0.61	0.88	0.72	60
y	0.62	0.60	0.61	65
m	0.33	0.12	0.18	40
avg / total	0.55	0.59	0.55	165

Table 5: Classification report for predicting the testing set with RBF kernel.

	precision	recall	f1-score	support
n	1.00	0.92	0.99	57
y	0.95	1.00	0.98	61
m	1.00	0.94	0.97	32
avg / total	0.98	0.98	0.98	150

Table 6: Reconstructing the training set (150 flares total) with Random Forest classification method.

	precision	recall	f1-score	support
n	0.63	0.77	0.69	60
y	0.60	0.72	0.66	65
m	0.29	0.10	0.15	40
avg / total	0.54	0.59	0.55	165

Table 7: Classification report for predicting the testing set with the Random Forest classification method.

## 5 Conclusion

Ultimately it seems the human element cannot be removed from the vetting process, but with the aid of machine learning,...