

# SDSS Stripe 82 : quasar variability from forced photometry

Krzysztof Suberlak,<sup>1\*</sup> Željko Ivezić,<sup>1</sup> Yusra AlSayyad,<sup>1</sup>

<sup>1</sup>*Department of Astronomy, University of Washington, Seattle, WA, United States*

Accepted XXX. Received YYY; in original form ZZZ

## ABSTRACT

### 1 INTRODUCTION

### 2 DATA ANALYSIS

We use data from all SDSS runs up to and including run 7202 (Data Release 7), including all 6 SDSS camera columns.

All epochs (individual images) were background-subtracted, and then scaled from the Digital Unit counts to fluxes by comparing standard objects against the Ivezić+2007 catalog (similar to Jiang+2014).

The data were coadded, and all objects detected in the i-band coadds were assigned a deepSourceId (== objectId). For star-galaxy separation, the entire clump was considered as one parent source (with single ParentSourceId). For an object which is a parent (eg. a galaxy), ParentSourceId is null. This amounts to 40 million i-band detections down to  $3\sigma$ . 8 million of those are brighter than  $23^{\text{rd}}$  mag. Thus the total number of photometric measurements is : (40 million i-band detections)  $\times$  (80 epochs)  $\times$  (5 filters) = 16 billion measurements, including (8 million i-band detections  $i < 23$ )  $\times$  (80 epochs)  $\times$  (5 filters) = 3 billion measurements brighter than  $23^{\text{rd}}$  mag.

Forced photometry was performed in u,g,r,i,z on the individual epoch images (NOT difference imaging), in locations specified by i-band coadds. (DIFFERENCE imaging is when photometry is done on coadd - individual-epoch-image). Therefore in some cases the flux reported for a given aperture is negative, because after background subtraction noise oscillates around 0, and when it is scaled up, it can have negative values. [stored in rawDataFPSplit] The background in the optical bands is bright, and if we assume that the measured number of background counts oscillates around the value  $B_0$ , then the measured background count  $B$  is distributed as a Gaussian of width  $\sigma_B$  :  $B - B_0 \sim \mathcal{N}(0, \sigma_B)$ . The noise is Poissonian, i.e. depends on the number of counts, and since for optical measurements the number of counts is large,  $\sigma_B = \sqrt{B}$ . On a 4kx4k CCD, with 16 Mpix,  $5\sigma$  (corresponding to 1 false detection in a million), we would expect about 16 false detections. Now considering the distribution of the probability (likelihood) of flux measurement  $L(F|data)$ , for bright sources it is a very narrow Gaussian centered on the measured  $F_S$ , width  $\sigma_F$  on the level of  $1 - 2\% \approx 0.01 - 0.02$  mag. However, for faint sources the probability, centered around the  $F_S$ , is much wider, so that there is a nonzero probability of negative flux measurement. A Bayesian way to address this issue is to

impose the prior  $p(F)$ , since we understand that physically flux cannot be negative, so that the posterior probability  $p(F|data) \propto L(F|data)p(F)$ . A simple flat prior, being 0 for  $F < 0$  and 1 otherwise, would not affect the measured  $F_S$  for bright sources, but for faint sources it would move the distribution (posterior) to be above zero flux. This would be the upper limit on the flux of that source. Therefore we decided to apply the Bayesian prior in case where  $\langle F_L \rangle < k\sigma$ , with  $k = 2$  (for a Gaussian likelihood this corresponds to 2% probability of  $F_L < 0$ ).

### 2.1 Faint sources

To test our method we generate fiducial lightcurves (DRW / sinusoidal), with a uniform sampling ( $N = 100 \div 1000$ ). Based on the generated flux ( $F_{\text{true}}$ ) we define the  $5\sigma$  level as the robust 25-th percentile (or median) of the ensemble  $F_{\text{true}}$  distribution :  $\sigma_F = (1/5)F_{25\%}$  (in reality,  $\sigma$  increases for fainter observations, but this is a good approximation). Thus we define  $F_{\text{obs}} = F_{\text{true}} + F_{\text{noise}}$ , where the Gaussian noise  $F_{\text{noise}} = \sigma_F \mathcal{N}(0, 1)$  was added to each point. For a weak signal, defined as  $F_{\text{obs}}^i < 2\sigma_F$ , we consider  $p(F)$  - a Gaussian likelihood associated with  $i$ -th measurement:  $p_i(F) = \mathcal{N}(\mu = F_{\text{obs}}^i, \sigma = \sigma_F)$ . Each measurement  $F_{\text{obs}}$  is a mean of this likelihood:  $F_{\text{obs}} = \langle p_i(F) \rangle$ . We call it  $p(F)$  for short :

$$p(F) = \frac{1}{\sqrt{2\pi}\sigma_F} \exp\left(-\frac{(F - \mu)^2}{2\sigma_F^2}\right) \quad (1)$$

After imposing the Bayesian uniform prior  $p(F)$  becomes a truncated Gaussian (without the negative part), centered on  $F_{\text{obs}}^i$ , with a width  $\sigma_F$ . Thus for the truncated  $p(F)$  the mean ceases to be  $F_{\text{obs}}^i$ , but it can be defined as

$$F_{\text{mean}} = \int_0^\infty F p(F) dF \quad (2)$$

We also define the median as

$$\int_0^{F_{\text{median}}} p(F) dF = \int_{F_{\text{median}}}^\infty p(F) dF \quad (3)$$

Finally, since for a Gaussian distribution the area contained between  $\mu \pm \sigma$  is 95.5% of the total area under the

curve, for the truncated Gaussian we define the  $2\sigma$  level as two areas  $B = 0.05A$ , or :

$$\int_{F_{2\sigma}}^{\infty} p(F) dF = 0.05 * \int_0^{\infty} p(F) dF \quad (4)$$

Fro the mean, we can calculate it as

$$F_{mean} = \int_0^{\infty} \frac{F}{\sqrt{2\pi}\sigma^2} \exp\left(-\frac{(F-\mu)^2}{2\sigma^2}\right) = \frac{\sigma}{\sqrt{2\pi}} \exp\left(-\left(\frac{F_{obs}}{\sigma\sqrt{2}}\right)^2\right) + \frac{F_{obs}}{2} \operatorname{erfc}\left(\frac{F_{obs}}{\sigma\sqrt{2}}\right) \quad (5)$$

using a suitable substitution, and recognizing that

$$\operatorname{erfc} = \int_{t_0}^{\infty} \frac{2}{\sqrt{\pi}} e^{-t^2} \quad (6)$$

is a complementary error function, related to the error function as  $\operatorname{erfc}(x) = 1 - \operatorname{erf}(x)$ .

To find the median and the  $2\sigma$  level we transform from  $F$  space to  $x$  space to  $z$  space, where  $x = F/\sigma_F$  (scale by  $\sigma_F$ ), and  $z = x - x_{obs}$  (shift by  $x_{obs} = F_{obs}/\sigma_F$ ). Thus  $p(x) \sim \mathcal{N}(x_{obs}, 1)$  and  $p(z) \sim \mathcal{N}(0, 1)$ .

In  $z$ -space, the median from

$$\int_0^{x_{med}} p(x) dx = \int_{x_{med}}^{\infty} p(x) dx \quad (7)$$

becomes

$$\int_{x_0}^{z_{med}} p(z) dz = \int_{z_{med}}^{\infty} p(z) dz \quad (8)$$

with  $x_0 = -x_{obs}$

This expression for  $z_{med}$  can be evaluated : rhs is a survival function (sf) = 1 - cumulative density function (cdf) :

$$\int_{z_{med}}^{\infty} p(z) dz = \operatorname{sf}(z_{med}) \quad (9)$$

and the lhs, assuming  $z_{med} > x_0$ , is :

$$\int_{x_0}^{z_{med}} p(z) dz = \int_{-\infty}^{z_{med}} p(z) dz - \int_{-\infty}^{x_0} p(z) dz = \operatorname{cdf}(z_{med}) - \operatorname{cdf}(x_0) \quad (10)$$

Rearranging, and using the percent point function (ppf) - the inverse of the cdf, we find:

$$z_{med} = \operatorname{ppf}\left(\frac{1 + \operatorname{cdf}(x_0)}{2}\right) \quad (11)$$

and transforming back to  $F$  space:

$$F_{med} = F_{obs} + \sigma_F \operatorname{ppf}\left(\frac{1 + \operatorname{cdf}(x_0)}{2}\right) \quad (12)$$

with  $x_0$  and  $x_{obs}$  as above.

In  $z$  space, the  $2\sigma$  areas  $A$  and  $B$  are :

$A = \operatorname{sf}(x_0)$  and  $B = \operatorname{sf}(z_B)$ , so to find  $z_B$  we use the inverse survival function  $\operatorname{isf} : z_B = \operatorname{isf}(0.05A)$ . Thus transforming back to  $F$ -space we have:

$$F_{2\sigma} = F_{obs} + \sigma_F (\operatorname{isf}(0.05 \operatorname{sf}(x_0))) \quad (13)$$

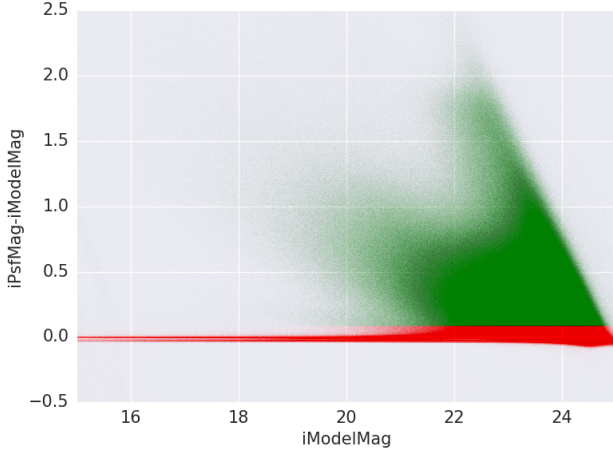
## 2.2 Photometric colors

The average brightness of an object in a given filter can be found in two ways. The median of the forced photometry values (over all epochs, including the negative fluxes) will better reflect the actual brightness of a variable source. This may mean that the median is negative, i.e. the median flux is negative. Since the magnitude is undefined for negative fluxes, we then revert to the lightcurve, and for each exposure with negative flux we find the zero point magnitude ( $m_1$ ) - the magnitude for a source with a flux of 1 count per sec, different for each exposure. The zero point magnitude for each exposure with negative flux is calculated from the Flux of 0 magnitude source,  $F_0$ , as  $m_1 = 2.5 \log_{10} F_0$ . For that object the new median magnitude in that filter will be the upper limit.

Colors can be calculated in two ways: using the median of forced photometry over all epochs (object detected in coadded i-band has photometry in all epochs), or the median over single-epoch detections (only when an object was above the detection threshold for a single epoch). Thus the median over all detections will be biased (especially for faint sources) towards higher brightness. On the other hand, the median over all epochs will be more representative of the true brightness of an object in a given filter. If a median brightness is negative, we use zero point magnitudes and in such cases median over all epochs will be an upper limit on brightness, but still less biased than median over all detections. Therefore we choose to use median over all epochs to calculate colors (see Fig. 3 for an example).

Since the reported fluxes are not extinction-corrected, we use a table of E(B-V) in a direction of a given source to correct for the galactic extinction. We use the formula  $x_{corr} = x_{obs} + A_x * E(B-V)$ , where  $x$  is u,g,r,i,z, and  $A_x$  is 5.155, 3.793, 2.751, 2.086, 1.479 for each filter respectively [Schlegel 98, Av are for RV = 3.1, also suggested by Eddie Schlafly]

The SDSS Stripe 82 data was processed in two data centers : NCSA (National Center for Supercomputing Applications, University of Illinois at Urbana-Champaign, IL) and IN2P3 (Institut national de physique nucléaire et de physique des particules in Paris, France). NCSA processed data with  $-40 (+320) < RA < +10$  and IN2P3 with  $+5 < RA < +55$  degrees. In the NCSA data there are 20978391 sources (`iCoaddAll.csv`). Removing those that overlap with IN2P3, we are left with 16520093 sources in the coadd photometry (`iCoaddPhotometryAll.csv`), and 16514187 sources (5906 less) in `DeepSourceNCSA_i_1t300.csv`. There are 12373162 sources with median photometry, matched with E(B-V) data (`medianPhotometry.csv`), and of these, 5892054 brighter than 23 mag, with calculated median flux and colors (`ugrizMetrics.csv`). Before individual bands are aggregated into one, we have individual bands treated separately, with metrics calculated for each band, eg. `i_metrics.csv`. In this file we have both annual aggregate metrics and full lightcurve aggregate metrics, including the Butler & Bloom classifier, which can for high S/N objects, where it has a good discriminating power. It's advantage over the full DRW analysis for each lightcurve is that by assuming a range of  $\tau$ , amplitude, expected for a DRW for a QSO, we calculate the likelihood of a given lightcurve belonging to a QSO.



**Figure 1.** A plot showing NCSA sources detected in coadds, removing the outliers beyond the edges of the plot. The coloring corresponds to the `extendedness` parameter calculated in the pipeline based on the `iPsFMag-iModelMag` : red being 0 (compact), and green being 1 (extended). As `iModelMag` increases, the separation becomes less certain, as more distant galaxies are more compact.

### 3 RESULTS

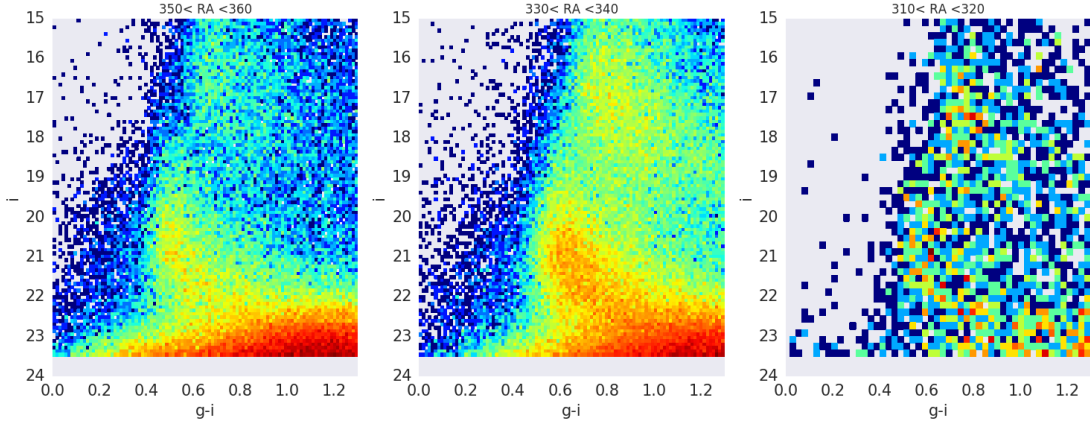
### 4 CONCLUSIONS

### ACKNOWLEDGEMENTS

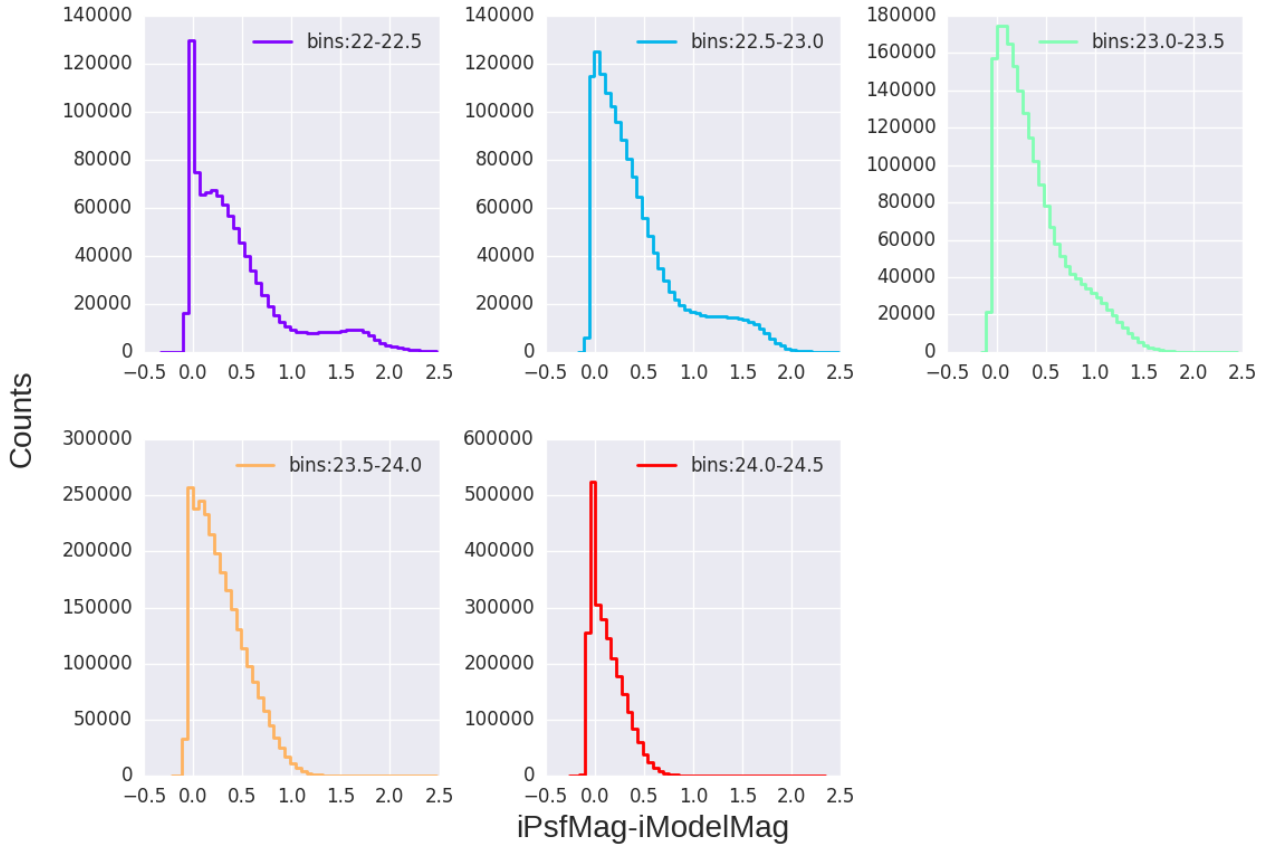
Funding for the SDSS and SDSS-II has been provided by the Alfred P. Sloan Foundation, the Participating Institutions, the National Science Foundation, the U.S. Department of Energy, the National Aeronautics and Space Administration, the Japanese Monbukagakusho, the Max Planck Society, and the Higher Education Funding Council for England. The SDSS Web Site is <http://www.sdss.org/>.

The SDSS is managed by the Astrophysical Research Consortium for the Participating Institutions. The Participating Institutions are the American Museum of Natural History, Astrophysical Institute Potsdam, University of Basel, University of Cambridge, Case Western Reserve University, University of Chicago, Drexel University, Fermilab, the Institute for Advanced Study, the Japan Participation Group, Johns Hopkins University, the Joint Institute for Nuclear Astrophysics, the Kavli Institute for Particle Astrophysics and Cosmology, the Korean Scientist Group, the Chinese Academy of Sciences (LAMOST), Los Alamos National Laboratory, the Max-Planck-Institute for Astronomy (MPIA), the Max-Planck-Institute for Astrophysics (MPA), New Mexico State University, Ohio State University, University of Pittsburgh, University of Portsmouth, Princeton University, the United States Naval Observatory, and the University of Washington.

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



**Figure 2.** A color-magnitude plot, reproducing the results of Sesar+2010, Fig.23. We show here only NCSA-processed sources, which is why certain RA ranges are omitted or have less sources. We only select sources with `extendedness=0` parameter (stars). The scale is showing the  $\log_{10}$  of count. All sources have their colors corrected for extinction. On first two panels the features of Sagittarius Stream are clearly visible.



**Figure 3.** The histograms show the count of sources in 5 magnitude bins, corresponding to the vertical cut through Fig. 2.2. It helps to verify how well can the extended and compact sources be separated based solely on the `iPsMag-iModelMag`