

LA-UR- 96-4462

CONF-9609317--1

Approved for public release;
distribution is unlimited.

Title:

HEURISTIC ESTIMATES OF WEIGHTED BINOMIAL
STATISTICS FOR USE IN DETECTING RARE
POINT SOURCE TRANSIENTS

Author(s):

James Theiler, NIS-2
Jeff Bloch, NIS-2

RECEIVED

FFR 14 1997

OSTI

Submitted to:

Proceedings of the Meeting on
Astronomical Data Analysis Software and
Systems (ADASS IV), Charlottesville, VA,
September 22-25, 1996.

MASTER

DISTRIBUTION OF THIS DOCUMENT IS UNLIMITED

8

Los Alamos
NATIONAL LABORATORY

Los Alamos National Laboratory, an affirmative action/equal opportunity employer, is operated by the University of California for the U.S. Department of Energy under contract W-7405-ENG-36. By acceptance of this article, the publisher recognizes that the U.S. Government retains a nonexclusive, royalty-free license to publish or reproduce the published form of this contribution, or to allow others to do so, for U.S. Government purposes. Los Alamos National Laboratory requests that the publisher identify this article as work performed under the auspices of the U.S. Department of Energy. The Los Alamos National Laboratory strongly supports academic freedom and a researcher's right to publish; as an institution, however, the Laboratory does not endorse the viewpoint of a publication or guarantee its technical correctness.

DISCLAIMER

**Portions of this document may be illegible
in electronic image products. Images are
produced from the best available original
document.**

Heuristic estimates of weighted binomial statistics for use in detecting rare point source transients

James Theiler and Jeff Bloch

*Astrophysics and Radiation Measurements Group, MS-D436
Los Alamos National Laboratory, Los Alamos, NM 87545
email: {jt,jbloch}@lanl.gov*

Abstract. The ALEXIS¹ (Array of Low Energy X-ray Imaging Sensors) (Priedhorsky *et al.*, 1989) satellite scans nearly half the sky every fifty seconds, and downlinks time-tagged photon data twice a day. The standard science quicklook processing produces over a dozen sky maps at each downlink, and these maps are automatically searched for potential transient point sources. We are interested only in *highly significant* point source detections, and based on earlier Monte-Carlo studies (Roussel-Dupré *et al.*, 1996), only consider $p < 10^{-7}$, which is about 5.2 “sigmas”. Our algorithms are therefore required to operate on the far tail of the distribution, where many traditional approximations break down. Although an exact solution is available for the case of unweighted counts (Lampton, 1994), the problem is more difficult in the case of weighted counts. We have found that a heuristic modification of a formula derived by Li and Ma (1983) provides reasonably accurate estimates of p -values for point source detections even for very low p -value detections.

1. Introduction

We test the null hypothesis of no point source (assuming a spatially uniform background) at a given location by enclosing that location with a source kernel (whose area A_{src} is generally matched to the point-spread-function of the telescope) and then enclosing the source kernel with a relatively large background annulus (area A_{bak}). Given N_{src} photons in the source kernel, and N_{bak} photons in the background annulus, the problem is to determine whether the number of source photons is *significantly* larger than expected under the null.

More sensitive point source detection is obtained by weighting the photons to more precisely match the point-spread function of the telescope. Further enhancements are obtained for ALEXIS data by weighting also according to instantaneous scalar background rate, pulse height, and position on the detector. In this case, we ask whether the weighted sum of photons in the source region is significantly larger than expected under the null.

¹<http://nis-www.lanl.gov/nis-projects/alexis/>

2. Unweighted counts

If counts are unweighted (*i.e.*, all weights are equal), then it is possible to write down an exact, explicit expression for the probability of seeing N_{src} or more photons in the source kernel, assuming $N_{\text{total}} = N_{\text{src}} + N_{\text{bak}}$ is fixed. This is a binomial distribution, and Lampton (1994) showed that the p -value associated with this observation can be expressed in terms of the incomplete beta function: $p = I_f(N_{\text{src}}, N_{\text{bak}} + 1)$, where $f = A_{\text{src}} / (A_{\text{src}} + A_{\text{bak}})$. See also Alexandreas *et al.* (1994), for an alternative derivation of an equivalent expression (the assumption that N_{total} is fixed is replaced by a Bayesian argument).

If the count rate is high (or the exposure long), so that N_{src} and N_{bak} are large, then an appropriate Gaussian approximation can be used. In general, this involves finding a “signal” and dividing it by the square root of its variance.

Case 1u. The most straightforward approach uses the signal $N_{\text{src}} - \alpha N_{\text{bak}}$, where $\alpha = A_{\text{src}} / A_{\text{bak}}$. Under the null hypothesis, this signal has an expected value of zero, and a variance — if N_{src} and N_{bak} are treated as independent Poisson sources — of $N_{\text{src}} + \alpha^2 N_{\text{bak}}$. To get a p -value, use

$$p = \mathcal{S} \left(\frac{N_{\text{src}} - \alpha N_{\text{bak}}}{\sqrt{N_{\text{src}} + \alpha^2 N_{\text{bak}}}} \right), \quad (1)$$

where $\mathcal{S}(s) = \frac{1}{2}(1 - \text{erfc}(s/\sqrt{2}))$ converts “sigmas” of significance into a one-tailed p -value.

Case 2u. An alternative approach, suggested by Li and Ma (1983), treats the sum $N_{\text{total}} = N_{\text{src}} + N_{\text{bak}}$, as fixed, so that N_{src} and N_{bak} are binomially distributed. In particular, choose the signal $N_{\text{src}} - f N_{\text{total}}$, and note that the variance of N_{src} is given by $f(1-f)N_{\text{total}}$, while the variance of N_{total} is by definition zero. In that case

$$p = \mathcal{S} \left(\frac{N_{\text{src}} - f N_{\text{total}}}{\sqrt{f(1-f)N_{\text{total}}}} \right) = \mathcal{S} \left(\frac{N_{\text{src}} - \alpha N_{\text{bak}}}{\sqrt{\alpha N_{\text{src}} + \alpha N_{\text{bak}}}} \right). \quad (2)$$

Case 3u. By looking at a ratio of Poisson likelihoods, Li and Ma (1983) also derived a more complicated equation

$$p = \mathcal{S} \left(\sqrt{2 \left\{ N_{\text{src}} \ln(N_{\text{src}} / \hat{N}_{\text{src}}) + N_{\text{bak}} \ln(N_{\text{bak}} / \hat{N}_{\text{bak}}) \right\}} \right), \quad (3)$$

where $\hat{N}_{\text{src}} = f N_{\text{total}}$ and $\hat{N}_{\text{bak}} = (1-f)N_{\text{total}}$. This is considerably more accurate than Eqs. (1,2) when N_{src} and N_{bak} are not large, but is still just an approximation to Lampton’s exact formula. Abramowitz and Stegen (1972) provide several approximations to the incomplete beta function, one of which (25.5.19) is an asymptotic series whose first term looks very much like the Li and Ma formula. The left panel of Figure 1 compares these cases, along with the Lampton (1994) formula, using a Monte-Carlo simulation.

3. Weighted counts

Define $W_{\text{src}} = \sum_{i \in \text{src}} w_i$ and $Q_{\text{src}} = \sum_{i \in \text{src}} w_i^2$, where w_i is the weight of the i -th photon. Notice that when all the weights are equal to one, we have $Q_{\text{src}} =$

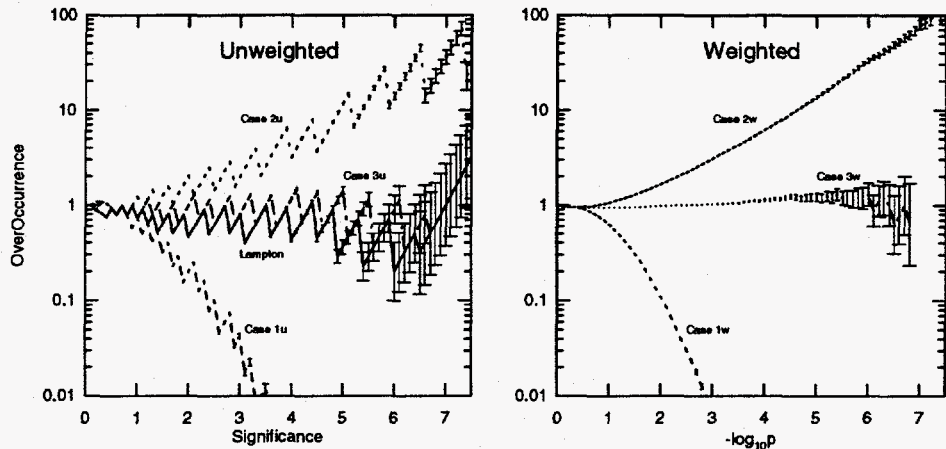


Figure 1. Results of Monte-Carlo experiments with $N = 100$ photons, with $f = 0.1$, and with $T = 10^7$ trials. For the weighted experiment, N weights were uniformly chosen from zero to one, and assigned to the N photons. The photons were randomly assigned to the source kernel or background annulus with probabilities f and $1 - f$ respectively. Values of W_{src} , W_{bak} , Q_{src} , and Q_{bak} were computed, and a p -value was computed using the formulas for the three cases. As the p -values were computed, a cumulative histogram $H(p)$ was built indicating the number of times a p -value less than p was observed. Since we expect $H(p) = pT$, we plotted $H(p)/pT$ as the frequency of “overoccurrence” of that p -value. The plot is this overoccurrence as a function of “significance”, defined by $-\log_{10} p$.

$W_{\text{src}} = N_{\text{src}}$ and $Q_{\text{bak}} = W_{\text{bak}} = N_{\text{bak}}$. Note also that $W_{\text{src}}/N_{\text{src}} = \langle w_i \rangle_{i \in \text{src}}$, and that $Q_{\text{src}}/W_{\text{src}} = \langle w_i^2 \rangle / \langle w_i \rangle$. We do not make any assumptions about weights averaging or summing to unity. (We define W_{bak} and Q_{bak} similarly.)

Generalizing Case 1u, we define the signal as $W_{\text{src}} - \alpha W_{\text{bak}}$ and then treating source and background as independent, we can write the variance as $Q_{\text{src}} + \alpha^2 Q_{\text{bak}}$. We can similarly generalize Case 2u and obtain:

$$\text{Case 1w: } p = S \left(\frac{W_{\text{src}} - \alpha W_{\text{bak}}}{\sqrt{Q_{\text{src}} + \alpha^2 Q_{\text{bak}}}} \right). \quad (4)$$

$$\text{Case 2w: } p = S \left(\frac{W_{\text{src}} - \alpha W_{\text{bak}}}{\sqrt{\alpha Q_{\text{src}} + \alpha Q_{\text{bak}}}} \right). \quad (5)$$

Case 3w: It is not as straightforward to generalize Eq. (3), but we have tried the following heuristic:

$$p = S \left(\sqrt{\left(\frac{2W_{\text{total}}}{Q_{\text{total}}} \right) \left\{ W_{\text{src}} \ln(W_{\text{src}}/\hat{W}_{\text{src}}) + W_{\text{bak}} \ln(W_{\text{bak}}/\hat{W}_{\text{bak}}) \right\}} \right), \quad (6)$$

where $\hat{W}_{\text{src}} = fW_{\text{total}}$ and $\hat{W}_{\text{bak}} = (1-f)W_{\text{total}}$. The Monte-Carlo results shown in Figure 1 indicate that this heuristic provides reasonably accurate p -values even for very small values of p .

4. Limit of precisely known background

An interesting limit occurs as the background annulus becomes large. Here, $A_{\text{bak}} \rightarrow \infty$, and the expected backgrounds \hat{N}_{src} , \hat{W}_{src} , etc. are all precisely known. For the unweighted counts, the exact p -value can be expressed in terms of the incomplete gamma function: $p = 1 - \Gamma(N_{\text{src}}, \hat{N}_{\text{src}}) / \Gamma(N_{\text{src}})$. The Gaussian estimate of significance is straightforward² both for the unweighted case, $p = \mathcal{S}\left(\frac{N_{\text{src}} - \hat{N}_{\text{src}}}{\sqrt{N_{\text{src}}}}\right)$, and for the weighted case: $p = \mathcal{S}\left(\frac{W_{\text{src}} - \hat{W}_{\text{src}}}{\sqrt{\hat{Q}_{\text{src}}}}\right)$. In this limit, Eq. (6) becomes

$$p = \mathcal{S}\left(\sqrt{2\left(\hat{W}_{\text{src}}/\hat{Q}_{\text{src}}\right)}\left(W_{\text{src}} \ln(W_{\text{src}}/\hat{W}_{\text{src}}) - (W_{\text{src}} - \hat{W}_{\text{src}})\right)\right). \quad (7)$$

Marshall (1994) has suggested an empirical formula $p = \mathcal{S}\left(\frac{W_{\text{src}} - \hat{W}_{\text{src}} + \Delta}{\sqrt{\hat{Q}_{\text{src}} + \Delta}}\right)$, where $\Delta = 0.7\hat{Q}_{\text{src}}/\hat{W}_{\text{src}}$, which produced reasonable results in his simulations, but does not appear to be well suited for p -values at the far tail of the distribution.

Acknowledgments. This work was supported by the United States Department of Energy.

References

- Abramowitz, M., & Stegun, I. A. 1972, Handbook of Mathematical Functions (Dover, New York), 945
- Alexandreas, D. E., *et al.* 1993 Nucl. Instr. Meth. Phys. Res. A328, 570
- Babu, G. J., & Feigelson, E. D. 1996, Astrostatistics (Chapman & Hall, London), 113
- Lampton, M. 1994 ApJ 436, 784
- Li, T.-P., & Ma Y.-Q. 1983, ApJ 272, 317
- Marshall, H. L. 1994, in ASP Conf. Ser., Vol. 61, Astronomical Data Analysis Software and Systems III, ed. Dennis R. Crabtree, R. J. Hanisch & Jeannette Barnes (San Francisco: ASP), 403
- Priedhorsky, W. C., Bloch, J. J., Cordova, F., Smith, B. W., Ulibarri, M., Chavez, J., Evans, E., Seigmund, O., H. W., Marshall, H., & Vallergera, J. 1989, in Berkeley Colloquium on Extreme Ultraviolet Astronomy, Berkeley, CA, vol 2873, 464
- Roussel-Dupré, D., Bloch, J. J., Theiler, J., Pfafman, T., & Beauchesne, B. 1996, in ASP Conf. Ser., Vol. 101, Astronomical Data Analysis Software and Systems V, ed. George H. Jacoby & Jeannette Barnes (San Francisco: ASP), 112

²Babu and Feigelson (1996) incorrectly suggest $p = \mathcal{S}\left((N_{\text{src}} - \hat{N}_{\text{src}})/\sqrt{N_{\text{src}} + \hat{N}_{\text{src}}}\right)$.