

# Error Estimates for the Bin Contents of a Convolution Histogram

Luc Demortier<sup>1</sup> and Luca Lista<sup>2</sup>

<sup>1</sup>*Laboratory of Experimental High-Energy Physics, Rockefeller University*

<sup>2</sup>*INFN, Napoli*

February 12, 2010

## Abstract

We examine a generalization of a Monte Carlo procedure commonly used to model the missing transverse energy spectrum of  $W \rightarrow \tau\nu$  events. This procedure introduces correlations between spectrum bins, thereby invalidating the usual multinomial error structure. We compute the errors on the bin contents from first principles.

## 1 Problem Setup

A common technique for generating the distribution of a random variable  $Z$  that is itself a function  $c(\vec{X}, \vec{Y})$  of two possibly multidimensional, possibly dependent random variables  $\vec{X}$  and  $\vec{Y}$ , is the following:

1. Draw a number  $n$  of variates distributed as  $\vec{X}$ ;
2. For each  $\vec{X}$ :
  - (a) Draw a number  $k$  of variates distributed as  $\vec{Y}$  given  $\vec{X}$ ;
  - (b) Compute  $Z = c(\vec{X}, \vec{Y})$ ;
3. Histogram the  $Z$ 's.

This technique can be applied to the modeling of the missing transverse energy spectrum of events containing  $W \rightarrow \tau\nu$  decays, where one must combine contributions from the  $W$  decay neutrino and from the tau lepton. A typical procedure is to substitute Monte Carlo tau decays for the muons in a data sample of  $W \rightarrow \mu\nu$  events. When the size of the muon data sample is limited, the precision of the result can be increased by substituting several tau decays for each available muon. However, this introduces nontrivial correlations between spectrum bins, because a given  $W \rightarrow \mu\nu$  event can now contribute to more than one bin (only in the special case  $k = 1$  are the correlations purely multinomial).

The question we wish to address is how to calculate the errors on the bin contents of the  $Z$  histogram in general. We derive these errors in the next section.

## 2 Solution

We adopt the standard statistical convention of designating random variables with capital letters and their observed values with the corresponding small-case letters. Let  $f(\vec{x})$  be the probability density of  $\vec{X}$  and  $g(\vec{y}|\vec{x})$  that of  $\vec{Y}$  given  $\vec{X}$ . Define further:

$M_{ab} \equiv$  number of  $Z$  values in bin  $[z_a, z_b]$  coming from a given  $\vec{X}$   
 (observed values of  $M_{ab}$  will be written as  $m_{ab}(\vec{x})$  to indicate  
 their association with observed values of  $\vec{X}$ );

$$N_{ab} \equiv \text{total number of } Z \text{ values in } [z_a, z_b]: \quad n_{ab} = \sum_{i=1}^n m_{ab}(\vec{x}_i); \quad (2.1)$$

$$n'_{ab} \equiv \sum_{i=1}^n [m_{ab}(\vec{x}_i)]^2.$$

The crucial insight is that within a given  $Z$  bin  $[z_a, z_b]$  the  $M_{ab}$  are independent, since each of them is obtained by drawing a new  $\vec{X}$  from  $f(\vec{x})$  and  $k$  new  $\vec{Y}$ 's from  $g(\vec{y}|\vec{x})$ . In addition, the  $M_{ab}$  are identically distributed since the same  $f$  and  $g$  are used throughout. We therefore have, from Eq. (2.1):

$$\text{Var}(N_{ab}) = n \text{Var}(M_{ab}). \quad (2.2)$$

An unbiased estimate of the expectation  $\mathbb{E}(M_{ab})$  is  $n_{ab}/n$ , and the corresponding unbiased estimate of  $\text{Var}(M_{ab})$  is then:

$$\text{Var}(M_{ab}) = \frac{1}{n-1} \sum_{i=1}^n \left( m_{ab}(\vec{x}_i) - \frac{n_{ab}}{n} \right)^2 = \frac{1}{n-1} \left( n'_{ab} - \frac{n_{ab}^2}{n} \right). \quad (2.3)$$

Therefore:

$$\text{Var}(N_{ab}) = \frac{n}{n-1} \left( n'_{ab} - \frac{n_{ab}^2}{n} \right). \quad (2.4)$$

Note that if  $k = 1$  all the  $m_{ab}(\vec{x}_i)$  are either 0 or 1, so that  $n'_{ab} = n_{ab}$  and the above formula reduces to the standard binomial case. More importantly, Eq. (2.4) expresses  $\text{Var}(N_{ab})$  in terms of quantities that are directly accessible when making the histogram; there is no need to run large ensembles of toy experiments to figure out the uncertainties on the bin contents. The histogramming procedure could be something like this:

- 1 Initialize histograms  $H_1$ ,  $H_2$ , and  $H_3$  with  $B$  bins from  $z_{\min}$  to  $z_{\max}$ .
- 2 For  $i = 1, \dots, n$ :
- 3     Generate  $\vec{x}_i \sim f(\vec{x})$ .
- 4     Zero histogram  $H_3$ .
- 5     For  $j = 1, \dots, k$ :
- 6         Generate  $\vec{y}_{ij} \sim g(\vec{y}|\vec{x}_i)$ .
- 7         Set  $z_{ij} = c(\vec{x}_i, \vec{y}_{ij})$ .
- 8         Histogram  $z_{ij}$  in  $H_1$ .
- 9         Histogram  $z_{ij}$  in  $H_3$ .
- 10     Add the square of  $H_3$  to  $H_2$ , bin by bin.
- 11 Compute  $H_4 \equiv \sqrt{\frac{n}{n-1} [H_2 - H_1^2/n]}$  (also bin by bin).

Histogram  $H_4$  then contains the estimated errors on the bins of histogram  $H_1$ .

### 3 Example

We illustrate our result with an example that is somewhat inspired by the  $W \rightarrow \tau\nu$  problem. Suppose  $\vec{X}$  and  $\vec{Y}$  are two-dimensional vectors generated as follows:

$$\vec{X} \equiv (X_1, X_2) \equiv (U \cos S, U \sin S) : \quad U \sim \frac{1}{\tau} e^{-u/\tau}, \quad S \sim \mathcal{U}[0, 2\pi[, \quad (3.1)$$

$$\vec{Y} \equiv (Y_1, Y_2) \equiv (V \cos T, V \sin T) : \quad V | U \sim \mathcal{N}(u, \sigma^2), \quad T \sim \mathcal{U}[0, 2\pi[, \quad (3.2)$$

where  $\tau$  and  $\sigma$  are known constants. In words,  $\vec{X}$  has an exponentially distributed magnitude,  $\vec{Y}$  a Gaussian distributed magnitude conditional on the magnitude of  $\vec{X}$ , and both have uniform azimuth. One could think of  $\vec{X}$  as the missing transverse energy of the neutrino from  $W$  decay and  $\vec{Y}$  as that of the tau lepton, although the above distributions are motivated by convenience rather than by any resemblance with real physics processes. The combination rule  $c(\vec{X}, \vec{Y})$  is the magnitude of  $\vec{X} + \vec{Y}$ :

$$Z = c(\vec{X}, \vec{Y}) = \sqrt{(X_1 + Y_1)^2 + (X_2 + Y_2)^2} = \sqrt{U^2 + V^2 + 2UV \cos(S - T)}. \quad (3.3)$$

We generated 50,000 toy experiments with  $\tau = \sigma = 1$ ,  $n = 10$ , and  $k = 100$ . The results are displayed in Fig. 1. The left panel shows agreement between the toy experiment averages of the estimated variances of the bin contents and the toy experiment variances of the bin contents, indicating that the variance estimated with eq. (2.4) is unbiased. The right panel shows the same with standard deviation instead of variance. Here we do see some bias at low bin contents, where the estimated standard deviation underestimates the “true” standard deviation. This behavior is expected since bias is not invariant under nonlinear transformations such as the square root. Furthermore, the square root is a concave function, so that by Jensen’s inequality the estimated standard deviation should indeed *underestimate*.

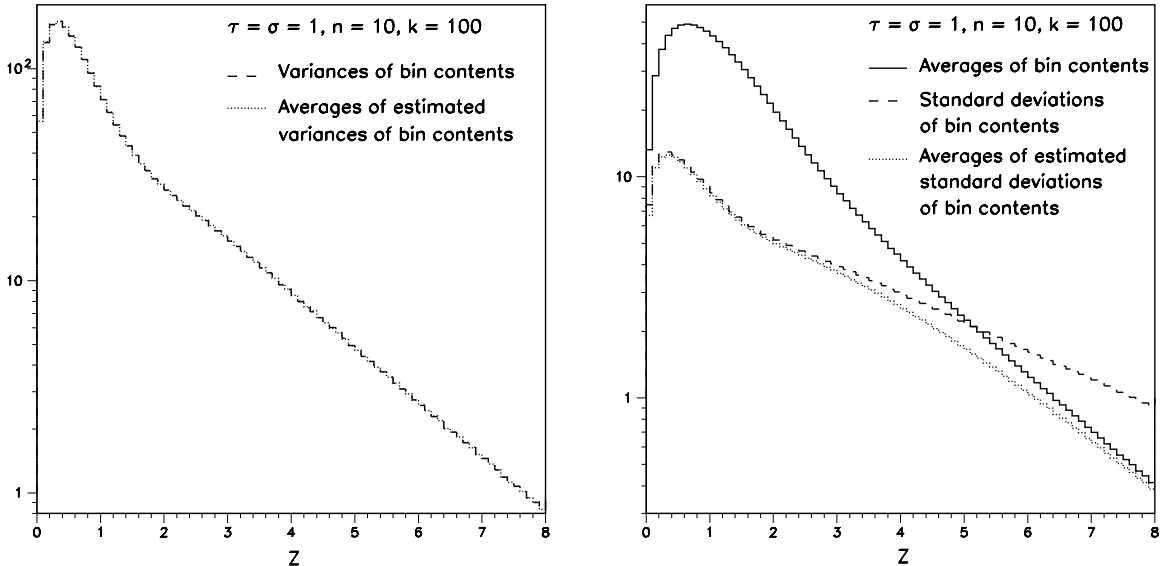


Figure 1: Result of 50,000 toy experiments with  $\tau = \sigma = 1$ ,  $n = 10$ , and  $k = 100$ .