

$sPlot$  :  
a statistical tool to unfold data distributions

M. Pivk<sup>a</sup> and F.R. Le Diberder<sup>b</sup>

<sup>a</sup> CERN,  
CH-1211 Geneva 23, Switzerland

<sup>b</sup> Laboratoire de l'Accélérateur Linéaire,  
IN2P3-CNRS et Université de Paris-Sud, F-91898 Orsay, France

**Abstract**

*The paper advocates the use of a statistical tool dedicated to the exploration of data samples populated by several sources of events. This new technique, called  $sPlot$ , is able to unfold the contributions of the different sources to the distribution of a data sample in a given variable. The  $sPlot$  tool applies in the context of a Likelihood fit which is performed on the data sample to determine the yields of the various sources.*

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Basics and definitions</b>	<b>3</b>
2.1	Likelihood method . . . . .	3
2.2	Analysis Validation . . . . .	4
<b>3</b>	<b>First step towards <math>_s\mathcal{Plot}</math>: <math>_{in}\mathcal{Plot}</math></b>	<b>5</b>
<b>4</b>	<b>The <math>_s\mathcal{Plot}</math> technique</b>	<b>6</b>
4.1	The $_s\mathcal{Plot}$ formalism . . . . .	7
4.2	$_s\mathcal{Plot}$ Properties . . . . .	9
4.2.1	Normalization . . . . .	9
4.2.2	Statistical uncertainties . . . . .	10
4.2.3	Merging $_s\mathcal{Plots}$ . . . . .	11
4.3	$_s\mathcal{Plot}$ implementation . . . . .	12
4.4	Illustrations . . . . .	12
4.5	Application: efficiency corrected yields . . . . .	16
<b>5</b>	<b>Conclusion</b>	<b>16</b>
<b>A</b>	<b>Pedagogical examples</b>	<b>18</b>
A.1	Simple cut-and-count analysis . . . . .	18
A.2	Extended cut-and-count analysis . . . . .	21
A.2.1	Generalized cut-and-count analysis: $n_y = N_s$ . . . . .	21
A.2.2	Extended cut-and-count analysis: $n_y > N_s$ . . . . .	22
<b>B</b>	<b>Extended <math>_s\mathcal{Plots}</math>: a species is known (fixed)</b>	<b>24</b>
B.1	Assuming $\mathbf{M}_0$ to be known . . . . .	24
B.2	Assuming $\mathbf{M}_0$ to be unknown . . . . .	26

# 1 Introduction

This paper describes a new technique to explore a data sample when the latter consists of several sources of events merged into a single sample of events. The events are assumed to be characterized by a set of variables which can be split into two components. The first component is a set of variables for which the distributions of all the sources of events are known: below, these variables are collectively referred to as a (unique) *discriminating* variable. The second component is a set of variables for which the distributions of some sources of events are either truly unknown or considered as such: below, these variables are collectively referred to as a (unique) *control* variable.

The new technique, termed *sPlot*, allows to reconstruct the distributions for the control variable, independently for each of the various sources of events, without making use of any *a priori* knowledge on this variable. The aim is thus to use the knowledge available for the discriminating variable to be able to infer the behavior of the individual sources of events with respect to the control variable. An essential assumption for the *sPlot* technique to apply is that the control variable is uncorrelated with the discriminating variable.

The *sPlot* technique is developed in the context of a data sample analyzed using a maximum Likelihood method making use of the discriminating variable. Section 2 is dedicated to the definition of fundamental objects necessary for the following. Section 3 presents an intermediate technique, simpler but inadequate, which is a first step towards the *sPlot* technique. Section 4 is the core of the document where the *sPlot* formalism is developed (Section 4.1) and its properties explained in detail (Section 4.2). Section 4.3 then gives instructions about how to implement and use *sPlot*. Finally, illustrations of *sPlots* are provided with simulated events (Section 4.4) and an application for branching ratios measurements (Section 4.5) is briefly described.

To provide some intuitive understanding of how and why the *sPlot* formalism works, the problem of reconstructing the true distributions is reconsidered in Appendix A, in a simpler analysis framework. An extension of the *sPlot* technique is presented in Appendix B.

## 2 Basics and definitions

A common method used to extract parameters from a data sample is the maximum Likelihood method which is briefly reviewed in Section 2.1 since it constitutes the foundation of the *sPlot* technique. Section 2.2 discusses the need for checks of an analysis based on the Likelihood method and introduces more precisely the goal of the *sPlot* technique.

### 2.1 Likelihood method

One considers an extended Likelihood analysis of a data sample in which are merged several species of events. These species represent various signal components (ie. sources of events in which one is interested) and background components (ie. irrelevant sources of events accompanying the signal components) which all together account for the data

sample. The log-Likelihood is expressed as:

$$\mathcal{L} = \sum_{e=1}^N \ln \left\{ \sum_{i=1}^{N_s} N_i f_i(y_e) \right\} - \sum_{i=1}^{N_s} N_i , \quad (1)$$

where

- $N$  is the total number of events in the data sample,
- $N_s$  is the number of species of events populating the data sample,
- $N_i$  is the number of events expected on the average for the  $i^{\text{th}}$  species,
- $y$  is the set of discriminating variables,
- $f_i$  is the Probability Density Function (PDF) of the discriminating variables for the  $i^{\text{th}}$  species,
- $f_i(y_e)$  denotes the value taken by the PDFs  $f_i$  for event  $e$ , the later being associated with a set of values  $y_e$  for the set of discriminating variables,
- $x$  is the set of control variables which, by definition, do not appear in the above expression of  $\mathcal{L}$ .

The log-Likelihood  $\mathcal{L}$  is a function of the  $N_s$  yields  $N_i$  and, possibly, of implicit free parameters designed to tune the PDFs on the data sample. These parameters as well as the yields  $N_i$  are determined by maximizing the above log-Likelihood.

## 2.2 Analysis Validation

The crucial point for such an analysis of the data sample to be reliable is to use an exhaustive list of sources of events combined with an accurate description of all the PDFs  $f_i$ .

To assess the quality of the fit, one may rely on an evaluation of the goodness of fit based on the actual value obtained for the maximum of  $\mathcal{L}$ , but this is rarely convincing enough. A complementary quality check is to explore further the data sample by examining the distributions of control variables. If the distributions of these control variables are known for at least one of the sources of events, one can compare the expected distribution for this source to the one extracted from the data sample. In order to do so, one must be able to unfold from the distribution of the whole data sample, the contribution arising from the source under scrutiny.

In some instances of control variables, the PDF might even be known for all the sources of events. Such a control variable can be obtained for instance by removing one of the discriminating variables from the set  $y$  before performing again the maximum Likelihood fit, and considering the removed variable as a control variable  $x$ . Another example is provided by a discriminating variable for which the distributions are known for all sources of events, but which does not improve significantly the accuracy fo the fit, and is not incorporated in the set  $y$ , for the sake of simplicity.

In an attempt to have access to the distributions of control variables, a common method consists in applying cuts which are designed to enhance the contributions to the data sample of particular sources of events (typically of signal species). Having enforced this enhancement, the distribution of  $x$  for the reduced data sample can be used to probe

the quality of the fit through a comparison with a Monte Carlo simulated distribution. However, the result is frequently unsatisfactory: firstly because it can be used only if the signal has prominent features to be distinguished from the background, and secondly because of the cuts applied, a sizeable fraction of signal events can be lost, while a large fraction of background events may remain. Therefore, the resulting data distribution concerns a reduced subsample for which statistical fluctuations, or true anomalies, cannot be attributed unambiguously, neither to the signal, nor to the background. For example, one can be tempted to misinterpret an anomaly in the distribution of  $x$  coming from the signal as a harmless background fluctuation.

The aim of the  ${}_s\mathcal{P}lot$  formalism developed in this paper is to provide a convenient method to unfold the overall distribution of a mixed sample of events in a control variable  $x$  into the sub-distributions of the various species which compose the sample. It is a statistical technique which allows to keep all signal events while getting rid of all background events, and keeping track of the statistical uncertainties per bin.

More formally, one is interested in the true distribution (denoted in boldface  $\mathbf{M}_n(x)$ ) of a control variable  $x$  for events of the  $n^{\text{th}}$  species, the later being any one of the  $N_s$  signal and background species. The purpose of this paper is to demonstrate that one can reconstruct  $\mathbf{M}_n(x)$  from the sole knowledge of the PDFs of the discriminating variables  $f_i$ , the first step being to proceed to the maximum Likelihood fit to extract the yields  $N_i$ .

As an introduction, in Section 3, the case is considered where the variable  $x$  actually belongs to the set of  $y$  discriminating variables. That is to say that one makes the assumption opposite to the interesting one:  $x$  is assumed to be totally correlated with  $y$ . Because of this total correlation, there exists a function of the  $y$  parameters which fully determines the 'control' variable,  $x = x(y)$ . In that case, while performing the fit, an *a priori* knowledge of the  $x$ -distributions is implicitly used, thus  $x$  cannot play the role of a control variable. Although the technique presented in the following Section is inadequate, it provides a natural first step towards  ${}_s\mathcal{P}lot$ .

Section 4, dedicated to the  ${}_s\mathcal{P}lot$  formalism, treats the interesting case, where  $x$  is truly a control variable uncorrelated with  $y$ . In that case, while performing the fit, no *a priori* knowledge of the  $x$ -distributions is used.

### 3 First step towards ${}_s\mathcal{P}lot$ : ${}_{in}\mathcal{P}lot$

In this Section, one is considering a variable  $x$  which can be expressed as a function of the discriminating variables  $y$  used in the fit. A fit having been performed to determine the yields  $N_i$  for all species, from the knowledge of the PDFs  $f_i$  and of the values of the  $N_i$ , one can define naively, for all events, the weight <sup>1</sup>

$$\mathcal{P}_n(y_e) = \frac{N_n f_n(y_e)}{\sum_{k=1}^{N_s} N_k f_k(y_e)} , \quad (2)$$

---

<sup>1</sup>It was pointed out to the authors that a weight similar to the naive one of Eq. (2) was introduced long ago in [1].

which can be used to build the  $x$ -distribution  $\tilde{M}_n$  defined by:

$$N_n \tilde{M}_n(\bar{x}) \delta x \equiv \sum_{e \in \delta x} \mathcal{P}_n(y_e) , \quad (3)$$

where the sum  $\sum_{e \in \delta x}$  runs over the  $N_{\delta x}$  events for which  $x_e$  (i.e. the value taken by the variable  $x$  for event  $e$ ) lies in the  $x$ -bin centered on  $\bar{x}$  and of total width  $\delta x$ .

In other words,  $N_n \tilde{M}_n(\bar{x}) \delta x$  is the  $x$ -distribution obtained by histogramming events, using the weight of Eq. (2).

This procedure reproduces, on average, the true distribution  $\mathbf{M}_n(x)$ . In effect, on average, one can replace the sum in Eq. (3) by the integral

$$\left\langle \sum_{e \in \delta x} \right\rangle \longrightarrow \int dy \sum_{j=1}^{N_s} N_j f_j(y) \delta(x(y) - \bar{x}) \delta x . \quad (4)$$

Similarly, identifying the number of events  $N_i$  as determined by the fit to be the expected number of events, one obtains:

$$\begin{aligned} \langle N_n \tilde{M}_n(\bar{x}) \rangle &= \int dy \sum_{j=1}^{N_s} N_j f_j(y) \delta(x(y) - \bar{x}) \mathcal{P}_n(y) \\ &= \int dy \sum_{j=1}^{N_s} N_j f_j(y) \delta(x(y) - \bar{x}) \frac{N_n f_n(y)}{\sum_{k=1}^{N_s} N_k f_k(y)} \\ &= N_n \int dy \delta(x(y) - \bar{x}) f_n(y) \\ &\equiv N_n \mathbf{M}_n(\bar{x}) . \end{aligned} \quad (5)$$

Therefore, the sum over events of the naive weight  $\mathcal{P}_n$  provides a direct estimate of the  $x$ -distribution of events of the  $n^{\text{th}}$  species. Plots obtained that way are referred to as *inPlots*: they provide a correct means to reconstruct  $\mathbf{M}_n(x)$  only insofar as the variable considered is **in** the set of discriminating variables  $y$ . These *inPlots* suffer from a major drawback:  $x$  being correlated to  $y$ , the PDFs of  $x$  enter implicitly in the definition of the naive weight, and as a result, the  $\tilde{M}_n$  distributions cannot be used easily to assess the quality of the fit, because these distributions are biased in a way difficult to grasp, when the PDFs  $f_i(y)$  are not accurate. For example, let us consider a situation where, in the data sample, some events from the  $n^{\text{th}}$  species show up far in the tail of the  $\mathbf{M}_n(x)$  distribution which is implicitly used in the fit. The presence of such events implies that the true distribution  $\mathbf{M}_n(x)$  must exhibit a tail which is not accounted for by  $\tilde{M}_n(x)$ . These events would enter in the reconstructed *inPlot*  $\tilde{M}_n$  with a very small weight, and they would thus escape detection by the above procedure:  $\tilde{M}_n$  would be close to  $\mathbf{M}_n$ , the distribution assumed for  $x$ . Only a mismatch in the core of the  $x$ -distribution can be revealed with *inPlots*. Stated differently, the error bars which can be attached to each individual bin of  $\tilde{M}_n$  cannot account for the systematical bias inherent to the *inPlots*.

## 4 The *sPlot* technique

It was shown in the previous Section that if the 'control' variable  $x$  belongs to the set  $y$  of discriminating variables, one can reconstruct the expected distribution of  $x$  with *inPlots*.

However, the  $_{\text{in}}\mathcal{P}lots$  are not easy to decipher because knowledge of the  $x$  distribution enters in their construction.

In this Section is considered the more interesting case where the variable  $x$  is truly a control variable, i.e. where  $x$  does not belong to  $y$ . More precisely, the two sets of variables  $x$  and  $y$  are assumed to be uncorrelated: hence, the total PDFs  $f_i(x, y)$  all factorize into products  $\mathbf{M}_i(x)f_i(y)$ .

## 4.1 The $_{\text{s}}\mathcal{P}lot$ formalism

One may still consider the above distribution  $\tilde{\mathbf{M}}_n$ , but this time the naive weight is no longer satisfactory: as shown below, Eq. (5) does not hold. This is because, when summing over the events, the  $x$ -PDFs  $\mathbf{M}_j(x)$  appear now on the right hand side of Eq. (4), while they are absent in the Likelihood function. However, a simple redefinition of the weights allows to overcome this difficulty.

Considering the naive weight of Eq. (2):

$$\begin{aligned} \langle N_n \tilde{\mathbf{M}}_n(\bar{x}) \rangle &= \int \int dy dx \sum_{j=1}^{N_s} N_j \mathbf{M}_j(x) f_j(y) \delta(x - \bar{x}) \mathcal{P}_n \\ &= \int dy \sum_{j=1}^{N_s} N_j \mathbf{M}_j(\bar{x}) f_j(y) \frac{N_n f_n(y)}{\sum_{k=1}^{N_s} N_k f_k(y)} \\ &= N_n \sum_{j=1}^{N_s} \mathbf{M}_j(\bar{x}) \left( N_j \int dy \frac{f_n(y) f_j(y)}{\sum_{k=1}^{N_s} N_k f_k(y)} \right) \\ &\neq N_n \mathbf{M}_n(\bar{x}) . \end{aligned} \tag{6}$$

Indeed, as announced, the previous procedure does not apply. In effect, the correction term appearing in Eq. (6)

$$N_j \int dy \frac{f_n(y) f_j(y)}{\sum_{k=1}^{N_s} N_k f_k(y)} \tag{8}$$

is not identical to the kroenecker symbol  $\delta_{jn}$ . The  $_{\text{in}}\mathcal{P}lot$  distribution  $N_n \tilde{\mathbf{M}}_n$  obtained using the naive weight is a linear combination of the true distributions  $\mathbf{M}_j$ . Only if the  $y$  variable was totally discriminating would one recover the correct answer. In effect, for a total discrimination,  $f_{j \neq n}(y)$  vanishes if  $f_n(y)$  is non zero. Thus, the product  $f_n(y) f_j(y)$  is equal to  $f_n^2(y) \delta_{jn}$ , and one gets:

$$N_j \delta_{jn} \int dy \frac{f_n^2(y)}{N_n f_n(y)} = \delta_{jn} . \tag{9}$$

But this is purely academic, because, if  $y$  was totally discriminating, the obtention of  $\mathbf{M}_n(x)$  would be straightforward: one would just apply cuts on  $y$  to obtain a pure sample of events of the  $n^{\text{th}}$  species and plot them to get  $\mathbf{M}_n(x)$ .

However, in the case of interest where  $y$  is not totally discriminating, one observes that the correction term is related to the inverse of the covariance matrix, given by the second derivatives of  $-\mathcal{L}$ , which the analysis minimizes:

$$\mathbf{V}_{nj}^{-1} = \frac{\partial^2(-\mathcal{L})}{\partial N_n \partial N_j} = \sum_{e=1}^N \frac{f_n(y_e) f_j(y_e)}{(\sum_{k=1}^{N_s} N_k f_k(y_e))^2} . \tag{10}$$

On average, replacing the sum over events by an integral (Eq. (4)) the variance matrix reads:

$$\begin{aligned}
\langle \mathbf{V}_{nj}^{-1} \rangle &= \int \int dy dx \sum_{l=1}^{N_s} N_l \mathbf{M}_l(x) f_l(y) \frac{f_n(y) f_j(y)}{(\sum_{k=1}^{N_s} N_k f_k(y))^2} \\
&= \int dy \sum_{l=1}^{N_s} N_l f_l(y) \frac{f_n(y) f_j(y)}{(\sum_{k=1}^{N_s} N_k f_k(y))^2} \int dx \mathbf{M}_l(x) \\
&= \int dy \frac{f_n(y) f_j(y)}{\sum_{k=1}^{N_s} N_k f_k(y)} .
\end{aligned} \tag{11}$$

Therefore, Eq. (6) can be rewritten:

$$\langle \tilde{\mathbf{M}}_n(\bar{x}) \rangle = \sum_{j=1}^{N_s} \mathbf{M}_j(\bar{x}) N_j \langle \mathbf{V}_{nj}^{-1} \rangle . \tag{12}$$

Inverting this matrix equation, one recovers the distribution of interest:

$$N_n \mathbf{M}_n(\bar{x}) = \sum_{j=1}^{N_s} \langle \mathbf{V}_{nj} \rangle \langle \tilde{\mathbf{M}}_j(\bar{x}) \rangle . \tag{13}$$

Hence, if the control variable  $x$  is uncorrelated with the discriminating variable, the true distribution of  $x$  can still be reconstructed using the naive weight of Eq. (2), through a linear combination of the  ${}_n\mathcal{P}lots$ . This result is better restated as follows. When  $x$  does not belong to the set  $y$ , the appropriate weight is not given by Eq. (2), but is the covariance-weighted quantity (thereafter called sWeight) defined by:

$$\boxed{{}_s\mathcal{P}_n(y_e) = \frac{\sum_{j=1}^{N_s} \mathbf{V}_{nj} f_j(y_e)}{\sum_{k=1}^{N_s} N_k f_k(y_e)}} . \tag{14}$$

With this sWeight, the distribution of the control variable  $x$  can be obtained from the  ${}_s\mathcal{P}lot$  histogram:

$$N_n {}_s\tilde{\mathbf{M}}_n(\bar{x}) \delta x \equiv \sum_{e \in \delta x} {}_s\mathcal{P}_n(y_e) , \tag{15}$$

which reproduces, on average, the true distribution:

$$\langle N_n {}_s\tilde{\mathbf{M}}_n(x) \rangle = N_n \mathbf{M}_n(x) . \tag{16}$$

If the control variable  $x$  exhibits significant correlation with the discriminating variable  $y$ , the  ${}_s\mathcal{P}lots$  obtained with Eq. (15) cannot be compared directly with the pure distributions of the various species. In that case, one must proceed to a Monte-Carlo simulation of the procedure to obtain the expected distributions to which the  ${}_s\mathcal{P}lots$  should be compared with.

The fact that the matrix  $\mathbf{V}_{ij}$  enters in the definition of the sWeights is enlightening, and, as discussed in the next Section, this confers nice properties to the  ${}_s\mathcal{P}lots$ . But this is not the key point. The key point is that Eq. (6) is a matrix equation which can be inverted



using a numerical evaluation of the matrix based only on data, thanks to Eq. (10). Rather than computing the matrix by this direct sum over the events, one can use the covariance matrix resulting from the fit, but this option is numerically less accurate than the direct computation<sup>2</sup>.

## 4.2 ${}_s\mathcal{P}lot$ Properties

Beside satisfying, on the average, the essential asymptotic property Eq. (16),  ${}_s\mathcal{P}lots$  bear properties which hold even under non-asymptotic conditions.

### 4.2.1 Normalization

The distribution  ${}_s\tilde{M}_n$  defined by Eq. (15) is guaranteed to be normalized to unity and the sum over the species of the  ${}_s\mathcal{P}lots$  reproduces the data sample distribution of the control variable. These two properties are not obvious because, from expression Eq. (14), neither is it obvious that the sum over the  $x$ -bins of  $N_n {}_s\tilde{M}_n \delta x$  is equal to  $N_n$ , nor is it obvious that, in each bin, the sum over all species of the expected numbers of events equates to the number of events actually observed. The demonstration uses the three sum rules below.

#### 1. Maximum Likelihood Sum Rule

The Likelihood Eq. (1) being extremal for  $N_j$ , one gets the first sum rule:

$$\sum_{e=1}^N \frac{f_j(y_e)}{\sum_{k=1}^{N_s} N_k f_k(y_e)} = 1, \quad \forall j. \quad (17)$$

#### 2. Variance Matrix Sum Rule

From Eq. (10) and Eq. (17) one derives:

$$\sum_{i=1}^{N_s} N_i \mathbf{V}_{ij}^{-1} = \sum_{i=1}^{N_s} N_i \sum_{e=1}^N \frac{f_i(y_e) f_j(y_e)}{(\sum_{k=1}^{N_s} N_k f_k(y_e))^2} = \sum_{e=1}^N \frac{f_j(y_e)}{\sum_{k=1}^{N_s} N_k f_k(y_e)} = 1. \quad (18)$$

#### 3. Covariance Matrix Sum Rule

Multiplying both sides of Eq. (18) by  $\mathbf{V}_{jl}$  and summing over  $j$  one gets the sum rule:

$$\sum_{j=1}^{N_s} \mathbf{V}_{jl} = \sum_{j=1}^{N_s} \mathbf{V}_{jl} \sum_{i=1}^{N_s} N_i \mathbf{V}_{ij}^{-1} = \sum_{i=1}^{N_s} \left( \sum_{j=1}^{N_s} \mathbf{V}_{ij}^{-1} \mathbf{V}_{jl} \right) N_i = \sum_{i=1}^{N_s} \delta_{il} N_i = N_l. \quad (19)$$

It follows that:

- Each  $x$ -distribution is properly normalized (cf. Eq. (17) and Eq. (19)):

$$\sum_{[\delta x]} N_n {}_s\tilde{M}_n(x) \delta x = \sum_{e=1}^N {}_s\mathcal{P}_n(y_e) = \sum_{e=1}^N \frac{\sum_{j=1}^{N_s} \mathbf{V}_{nj} f_j(y_e)}{\sum_{k=1}^{N_s} N_k f_k(y_e)} = \sum_{j=1}^{N_s} \mathbf{V}_{nj} = N_n. \quad (20)$$

---

<sup>2</sup>Furthermore, when parameters are fitted together with the yields  $N_j$ , in order to get the correct matrix, one should take care to perform a second fit, where these parameters are frozen.

- The contributions  ${}_s\mathcal{P}_j(y_e)$  add up to the number of events actually observed in each  $x$ -bin. In effect, for any event (cf. Eq. (19)) :

$$\sum_{l=1}^{N_s} {}_s\mathcal{P}_l(y_e) = \sum_{l=1}^{N_s} \frac{\sum_{j=1}^{N_s} \mathbf{V}_{lj} f_j(y_e)}{\sum_{k=1}^{N_s} N_k f_k(y_e)} = \frac{\sum_{j=1}^{N_s} N_j f_j(y_e)}{\sum_{k=1}^{N_s} N_k f_k(y_e)} = 1. \quad (21)$$

Therefore, an  ${}_s\text{Plot}$  provides a consistent representation of how all events from the various species are distributed in the control variable  $x$ . The contributions to the data sample distribution in  $x$  from the various species are disentangled according to a fit based on the discriminating variable  $y$ , provided  $x$  and  $y$  are uncorrelated. Summing up the  $N_s$   ${}_s\text{Plots}$ , one recovers the data sample distribution in  $x$ , and summing up the number of events entering in a  ${}_s\text{Plot}$  for a given species, one recovers the yield of the species, as it is provided by the fit.

For instance, if one observes an excess of events for a particular  $n^{\text{th}}$  species, in a given  $x$ -bin, this excess is effectively accounted for in the number of event  $N_n$  resulting from the fit. To remove these events (for whatever reason and by whatever means) implies a corresponding decrease in  $N_n$ . It remains to gauge how significant is an anomaly in the  $x$ -distribution of the  $n^{\text{th}}$  species. This is the subject of the next Section.

#### 4.2.2 Statistical uncertainties

The statistical uncertainty on  $N_n {}_s\tilde{M}_n(x)\delta x$  can be defined in each bin by

$$\sigma[N_n {}_s\tilde{M}_n(x)\delta x] = \sqrt{\sum_{e \in \delta x} ({}_s\mathcal{P}_n)^2}. \quad (22)$$

The proof that Eq. (22) holds asymptotically goes as follows:

$$\begin{aligned} \left\langle \left( \sum_{e \in \delta x} {}_s\mathcal{P}_n \right)^2 \right\rangle - \left\langle \sum_{e \in \delta x} {}_s\mathcal{P}_n \right\rangle^2 &= \langle N_{\delta x} \rangle \langle {}_s\mathcal{P}_n^2 \rangle + \langle N_{\delta x} (N_{\delta x} - 1) \rangle \langle {}_s\mathcal{P}_n \rangle^2 - \langle N_{\delta x} \rangle^2 \langle {}_s\mathcal{P}_n \rangle^2 \\ &= \langle N_{\delta x} \rangle \langle {}_s\mathcal{P}_n^2 \rangle + \left( \langle N_{\delta x}^2 \rangle - \langle N_{\delta x} \rangle \right) \langle {}_s\mathcal{P}_n \rangle^2 - \langle N_{\delta x} \rangle^2 \langle {}_s\mathcal{P}_n \rangle^2 \\ &= \langle N_{\delta x} \rangle \langle {}_s\mathcal{P}_n^2 \rangle \\ &\quad + \left( \langle N_{\delta x} \rangle + \langle N_{\delta x} \rangle^2 - \langle N_{\delta x} \rangle \right) \langle {}_s\mathcal{P}_n \rangle^2 - \langle N_{\delta x} \rangle^2 \langle {}_s\mathcal{P}_n \rangle^2 \\ &= \langle N_{\delta x} \rangle \langle {}_s\mathcal{P}_n^2 \rangle = \left\langle \sum_{e \in \delta x} ({}_s\mathcal{P}_n)^2 \right\rangle = \left\langle \sigma^2[N_n {}_s\tilde{M}_n\delta x] \right\rangle. \quad (23) \end{aligned}$$

The above asymptotic property is completed by the fact that the sum in quadrature of the uncertainties Eq. (22) reproduces the statistical uncertainty on the yield  $N_n$ , as it is provided by the fit:  $\sigma[N_n] \equiv \sqrt{\mathbf{V}_{nn}}$ . The sum over the  $x$ -bins reads:

$$\begin{aligned} \sum_{[\delta x]} \sigma^2[N_n {}_s\tilde{M}_n\delta x] &= \sum_{[\delta x]} \sum_{e \in \delta x} ({}_s\mathcal{P}_n)^2 = \sum_{e=1}^N \left( \frac{\sum_{j=1}^{N_s} \mathbf{V}_{nj} f_j(y_e)}{\sum_{k=1}^{N_s} N_k f_k(y_e)} \right)^2 \\ &= \sum_{j=1}^{N_s} \sum_{l=1}^{N_s} \mathbf{V}_{nl} \mathbf{V}_{nj} \sum_{e=1}^N \frac{f_l(y_e) f_j(y_e)}{(\sum_{k=1}^{N_s} N_k f_k(y_e))^2} \end{aligned}$$

$$\begin{aligned}
&= \sum_{j=1}^{N_s} \sum_{l=1}^{N_s} \mathbf{V}_{nl} \mathbf{V}_{nj} \mathbf{V}_{lj}^{-1} = \sum_{l=1}^{N_s} \mathbf{V}_{nl} \delta_{nl} \\
&= \mathbf{V}_{nn} ,
\end{aligned} \tag{24}$$

and more generally, the whole covariance matrix is reproduced:

$$\sum_{e=1}^N ({}_s\mathcal{P}_i)({}_s\mathcal{P}_j) = \mathbf{V}_{ij} . \tag{25}$$

Therefore, for the expected number of events per  $x$ -bin indicated by the  ${}_s\mathcal{P}lots$ , the statistical uncertainties are straightforward to compute using Eq. (22). The later expression is asymptotically correct, and it provides a consistent representation of how the overall uncertainty on  $N_n$  is distributed in  $x$  among the events of the  $n^{\text{th}}$  species. Because of Eq. (25), and since the determination of the yields is optimal when obtained using a Likelihood fit, one can conclude that the  ${}_s\mathcal{P}lot$  technique is itself an optimal method to reconstruct distributions of control variables.<sup>3</sup>

#### 4.2.3 Merging ${}_s\mathcal{P}lots$

As a result of the above, two species  $i$  and  $j$  can be merged into a single species  $(i+j)$  without having to repeat the fit and recompute the sWeights. The  ${}_s\mathcal{P}lot$  of the merged species is just the sum of the two  ${}_s\mathcal{P}lots$  obtained by adding the sWeights on an event-by-event basis:

$$N_{(i+j)} \tilde{\mathbf{M}}_{(i+j)} \delta x = \sum_{e \in \delta x} ({}_s\mathcal{P}_i + {}_s\mathcal{P}_j) . \tag{28}$$

The resulting  ${}_s\mathcal{P}lot$  has the proper normalization and the proper error bars (Eqs. (20) and (25)):

$$N_{(i+j)} = \sum_{e=1}^N ({}_s\mathcal{P}_i + {}_s\mathcal{P}_j) = N_i + N_j \tag{29}$$

$$\begin{aligned}
\sigma^2[N_{(i+j)}] &= \sum_{e=1}^N ({}_s\mathcal{P}_i + {}_s\mathcal{P}_j)^2 \\
&= \mathbf{V}_{ii} + \mathbf{V}_{jj} + 2\mathbf{V}_{ij} = \mathbf{V}_{(i+j)(i+j)} .
\end{aligned} \tag{30}$$

---

<sup>3</sup>This is not the case for  ${}_{\text{in}}\mathcal{P}lots$  for which one gets:

$$\sum_{e=1}^N (\mathcal{P}_i)(\mathcal{P}_j) = N_i N_j \mathbf{V}_{ij}^{-1} . \tag{26}$$

Hence, using the fact, that contrary to sWeights, the  ${}_{\text{in}}\mathcal{P}lot$  weights of Eq. (2) are positive definite, one gets:

$$\sum_{e=1}^N (\mathcal{P}_i)^2 \leq \sum_{e=1}^N (\mathcal{P}_i) \left( \sum_{j=1}^{N_s} \mathcal{P}_j \right) = N_i \sum_{j=1}^{N_s} N_j \mathbf{V}_{ij}^{-1} = N_i \leq \mathbf{V}_{ii} . \tag{27}$$

That is to say that the statistical uncertainties attached to the  ${}_{\text{in}}\mathcal{P}lots$  are always not only smaller than the ones resulting from the fit, but even smaller than the statistical uncertainties obtained in a background free situation.

### 4.3 ${}_s\mathcal{P}lot$ implementation

This Section is meant to show that using  ${}_s\mathcal{P}lot$  is indeed easy. The different steps to implement the technique are the following:

1. One is dealing with a data sample in which several species of events are present.
2. A maximum Likelihood fit is performed to obtain the yields  $N_i$  of the various species. The fit relies on a discriminating variable  $y$  uncorrelated with a control variable  $x$ : the later is therefore totally absent from the fit.
3. The sWeights  ${}_s\mathcal{P}$  are calculated using Eq. (14) where the covariance matrix is obtained by inverting the matrix given by Eq. (10).
4. Histograms of  $x$  are filled by weighting the events with the sWeights  ${}_s\mathcal{P}$ . The sum of the entries are equal to the yields  $N_i$  provided by the fit.
5. Error bars per bin are given by Eq. (22). The sum of the error bars squared are equal to the uncertainties squared  $\mathbf{V}_{ii}$  provided by the fit.
6. The  ${}_s\mathcal{P}lots$  reproduce the true distributions of the species in the control variable  $x$ , within the above defined statistical uncertainties.

The  ${}_s\mathcal{P}lot$  method has been implemented in the ROOT framework under the class TSPlot [2].

### 4.4 Illustrations

To illustrate the technique, one considers in this Section an example derived from the analysis where  ${}_s\mathcal{P}lots$  have been first used [3] and [4] (but see also [5]). One is dealing with a data sample in which two species are present: the first is termed signal and the second background. A maximum Likelihood fit is performed to obtain the two yields  $N_1$  and  $N_2$ . The fit relies on two discriminating variables collectively denoted  $y$  which are chosen within three possible variables denoted (following the notations of [3])  $m_{ES}$ ,  $\Delta E$  and  $\mathcal{F}$ . The variable which is not incorporated in  $y$  is used as a control variable  $x$ . The six distributions of the three variables are assumed to be the ones depicted in Fig. 1.

A data sample being built through a Monte Carlo simulation based on the distributions shown in Fig. 1, one obtains the three distributions of Fig. 2. Whereas the distribution of  $\Delta E$  clearly indicates the presence of the signal, the distribution of  $m_{ES}$  and  $\mathcal{F}$  are less obviously populated by signal.

Choosing  $\Delta E$  and  $\mathcal{F}$  as discriminating variables to determine  $N_1$  and  $N_2$  through a maximum Likelihood fit, one builds, for the control variable  $m_{ES}$  which is unknown to the fit, the two  ${}_s\mathcal{P}lots$  for signal and background shown in Fig. 3. For comparison, the PDFs of  $m_{ES}$  taken from Fig. 1 are superimposed on the  ${}_s\mathcal{P}lots$ . One observes that the  ${}_s\mathcal{P}lot$  for signal reproduces correctly the PDF even where the latter vanishes, although the error bars remain sizeable. This results from the almost complete cancellation between positive and negative sWeights: the sum of sWeights is close to zero in the tails while the sum of sWeights squared is not. The occurrence of negative sWeights is provided through the appearance of the covariance matrix, and its negative components, in the definition of Eq. (14).

A word of caution is in order with respect to the error bars. Whereas their sum in quadrature is identical to the statistical uncertainties of the yields determined by the

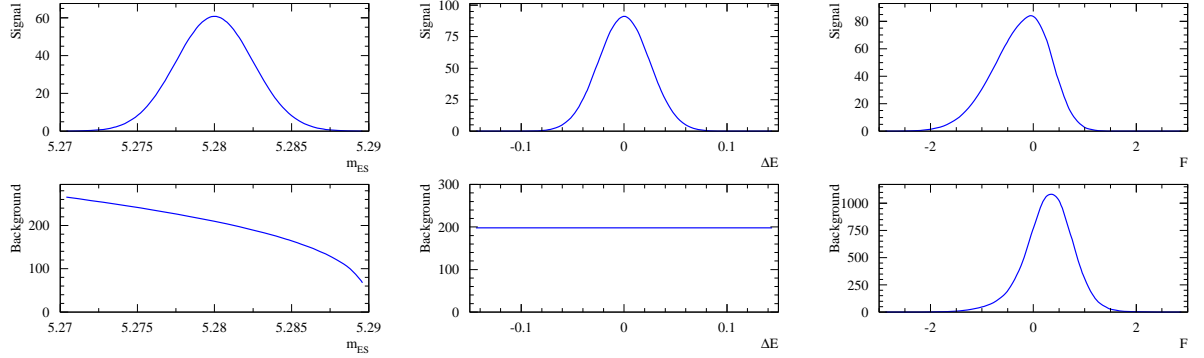


Figure 1: Distributions of the three different discriminating variables available to perform the Likelihood fit:  $m_{\text{ES}}$ ,  $\Delta E$ ,  $\mathcal{F}$ . Among the three variables, two are used to perform the fit while one is kept out of the fit to serve the purpose of a control variable. The three distributions on the top (resp. bottom) of the figure correspond to the signal (resp. background). The unit of the vertical axis is chosen such that it indicates the number of entries per bin, if one slices the histograms in 25 bins.

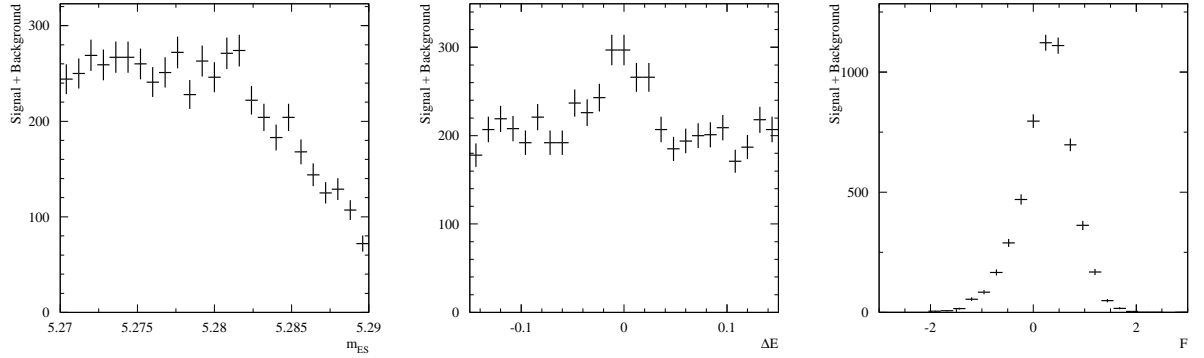


Figure 2: Distributions of the three discriminating variables for signal plus background. The three distributions are the one obtained from a data sample obtained through a Monte Carlo simulation based on the distributions shown in Fig. 1. The data sample consists of 500 signal events and 5000 background events.

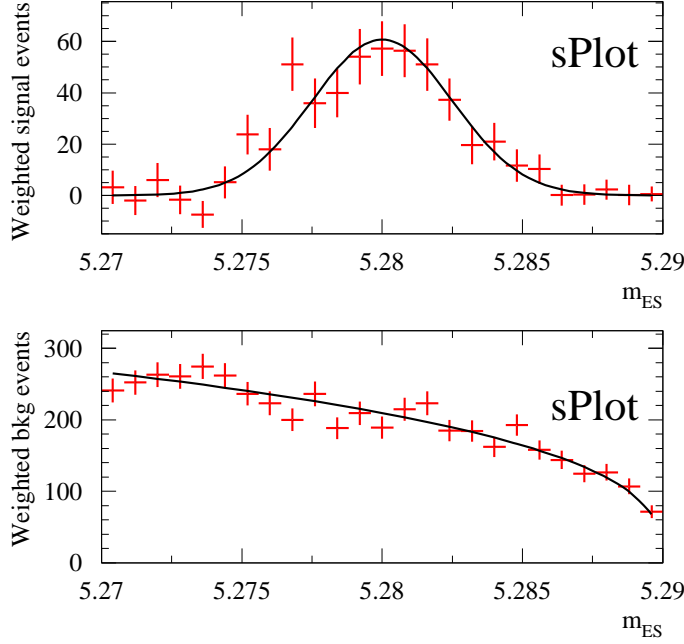


Figure 3: The  $sPlots$  (signal on top, background on bottom) obtained for  $m_{ES}$  are represented as dots with error bars. They are obtained from a fit using only information from  $\Delta E$  and  $\mathcal{F}$ . The black curves are the PDFs of  $m_{ES}$  of Fig. 1: these PDFs are unknown to the fit.

fit, and if, in addition, they are asymptotically correct (cf. Section 4.2.2) the error bars should be handled with care for low statistics and/or for too fine binning. This is because the error bars do not incorporate two known properties of the PDFs: PDFs are positive definite and can be non-zero in a given  $x$ -bin, even if in the particular data sample at hand, no event is observed in this bin. The latter limitation is not specific to  $sPlots$ , rather it is always present when one is willing to infer the PDF at the origin of an histogram, when, for some bins, the number of entries does not guaranty the applicability of the Gaussian regime. In such situations, a satisfactory practice is to attach allowed ranges to the histogram to indicate the upper and lower limits of the PDF value which are consistent with the actual observation, at a given confidence level. Although this is straightforward to implement, even when dealing with  $sWeighted$  events, for the sake of simplicity, this subject is not discussed further in the paper.

Choosing  $m_{ES}$  and  $\Delta E$  as discriminating variables to determine  $N_1$  and  $N_2$  through a maximum Likelihood fit, one builds, for the control variable  $\mathcal{F}$  which is unknown to the fit, the two  $sPlots$  for signal and background shown in Fig. 4. For comparison, the PDFs of  $\mathcal{F}$  taken from Fig. 1 are superimposed on the  $sPlots$ . In the  $sPlot$  for signal one observes that error bars are the largest in the  $x$  regions where the background is the largest.

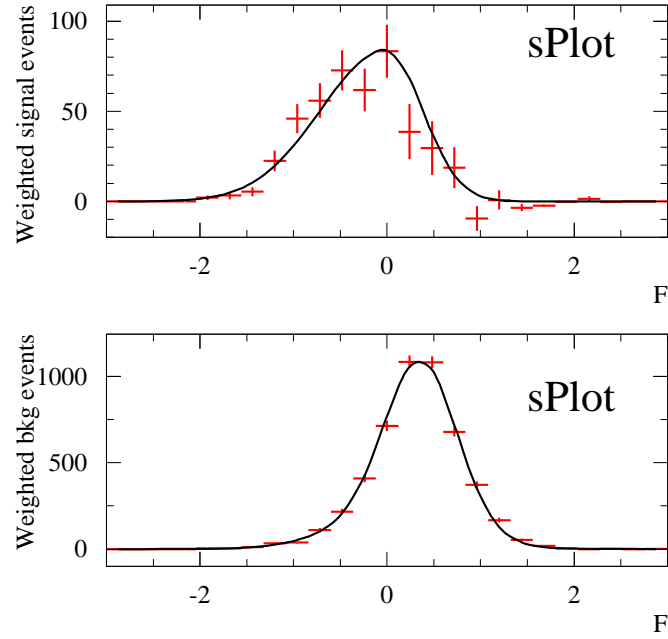


Figure 4: The  $sPlots$  (signal on top, background on bottom) obtained for  $\mathcal{F}$  are represented as dots with error bars. They are obtained from a fit using only information from  $m_{ES}$  and  $\Delta E$ . The black curves are the PDFs of  $\mathcal{F}$  of Fig. 1: these PDFs are unknown to the fit.

## 4.5 Application: efficiency corrected yields

Beside providing a convenient and optimal tool to cross-check the analysis by allowing distributions of control variables to be reconstructed and then compared with expectations, the *sPlot* formalism can be applied also to extract physics results, which would otherwise be difficult to obtain. For example, one may be willing to explore some unknown physics involved in the distribution of a variable  $x$ . Or, one may be interested to correct a particular yield provided by the Likelihood fit from a selection efficiency which is known to depend on a variable  $x$ , for which the PDF is unknown.

To be specific, one can take the example of a three body decay analysis of a species, the signal, polluted by background, while the signal PDF inside the two-dimensional Dalitz plot is not known, because of unknown contributions of resonances, continuum and an interference pattern. Since the  $x$ -dependence of the selection efficiency  $\epsilon(x)$  can be computed without *a priori* knowledge of the  $x$ -distributions, one can build the efficiency corrected two-dimensional *sPlots* (cf. Eq. (15)):

$$\frac{1}{\epsilon(\bar{x})} N_n \tilde{M}_n(\bar{x}) \delta x = \sum_{e \in \delta x} \frac{1}{\epsilon(x_e)} {}_s\mathcal{P}_n(y_e) , \quad (31)$$

and compute the efficiency corrected yields:

$$N_n^\epsilon = \sum_{e=1}^N \frac{{}_s\mathcal{P}_n(y_e)}{\epsilon(x_e)} . \quad (32)$$

Analyses can then use the *sPlot* formalism for validation purposes, but also, using Eq. (31) and Eq. (32), to probe for resonance structures and to measure branching ratios.

## 5 Conclusion

The technique presented in this paper applies when one examines a data sample originating from different sources of events: using a set  $y$  of discriminating variables, a Likelihood fit is performed on the data sample to determine the yields of the sources. By building *sPlots*, one can reconstruct the distributions of variables, separately for each source present in the data sample, provided the variables are uncorrelated with the set  $y$  used in the fit. Although no cut is applied (hence, the *sPlot* of a given species represents the whole statistics of this species) the distributions obtained are pure (free from the potential background arising from the other species) in a statistical sense. The more discriminating the discriminating variables  $y$ , the clearer the *sPlot* is. The technique is straightforward to implement and features several nice properties: both the normalizations and the statistical uncertainties of the *sPlots* reflect the fit outputs.

## References

- [1] P. E. Condon and P. L. Cowell, Phys. Rev. D 9 (1974) 2558-2562
- [2] <http://root.cern.ch/>



- [3] M. Pivk, Thèse de l'Université Paris VII, BABAR-THESIS-03/012 (2003),  
available (in French) at <http://tel.ccsd.cnrs.fr> (ID 00002991)
- [4] The BABAR Collaboration, Phys.Rev.Lett.93 (2004) 131801
- [5] The BABAR Collaboration, Phys.Rev.Lett.93 (2004) 181805;  
The BABAR Collaboration, Phys.Rev.D70 (2004) 091103;  
The BABAR Collaboration, Phys.Rev.Lett.94 (2005) 181802

## A Pedagogical examples

The purpose of this Appendix is to detail in simplified situations how and why *sPlot* works. One begins with the simplest situation and proceed to more complex ones.

### A.1 Simple cut-and-count analysis

In this Section, a very simple situation is considered where the proper way to reconstruct signal and background distributions for a control variable  $x$  is obvious from the start. The purpose is to observe the *sPlot* technique at work, when one knows beforehand what the outcome should be.

One considers a data sample consisting of  $N_s = 2$  species: species 1 is referred to as the signal and species 2 as the background. A unique discriminating variable  $y \in [0, 1]$  is used in the fit. One further assumes that:

- the signal distribution is the step-function:

$$f_1(y < y_0) = 0 \quad (33)$$

$$f_1(y \geq y_0) = (1 - y_0)^{-1} , \quad (34)$$

- the background distribution is uniform in the full range:

$$f_2(y) = 1 . \quad (35)$$

Therefore, one is dealing with a cut-and-count analysis: there is a pure background sideband for  $y < y_0$ , and the shapes of the signal and background distributions offer no discriminating power in the region where the signal is present, for  $y \geq y_0$ . Denoting  $N$  the total number of events present in the data sample,  $N_<$  the number of events located below  $y_0$ , and  $N_>$  the number of events located above  $y_0$ :

1. the expected number of background and signal events can be deduced without any fit, from the sideband:

$$N_2 = \frac{1}{y_0} N_< \quad (36)$$

$$N_1 = -\frac{1 - y_0}{y_0} N_< + N_> , \quad (37)$$

2.  $N_<$  and  $N_>$  being two independent numbers of events, the covariance matrix can be deduced directly from Eqs. (36)-(37):

$$\mathbf{V} = \begin{pmatrix} N_> + \left(\frac{1 - y_0}{y_0}\right)^2 N_< & -\frac{1 - y_0}{y_0^2} N_< \\ -\frac{1 - y_0}{y_0^2} N_< & \frac{1}{y_0^2} N_< \end{pmatrix} , \quad (38)$$

3. denoting  $\delta N_{<}^x$  the number of events in a given  $x$ -bin, with  $y \leq y_0$ , the background distribution  $M_2(x)$  can also be deduced by a mere rescaling of  $\delta N_{<}^x$ , as in Eq. (36):

$$N_2 M_2(x) \delta x = \frac{\delta N_{<}^x}{y_0} . \quad (39)$$

Similarly to Eq. (37), the signal distribution is given by:

$$N_1 M_1(x) \delta x = -(1 - y_0) N_2 M_2(x) + \delta N_{>}^x , \quad (40)$$

that is to say, one can obtain the signal distribution from the (mixed) events populating the domain  $y \geq y_0$ , if one subtracts the contribution of background events, which is known from Eq. (39). Stated differently, one is lead to assign the negative weight  $-(1 - y_0)/y_0$  to those events in the  $x$ -bin which satisfy  $y \leq y_0$ .

Whereas in such a simple situation the use of the  ${}_s\mathcal{P}lot$  formalism would be awkward, the latter should reproduce the above obvious results, and indeed it does. The proof goes as follows:

1. denoting  $f_i(0)$  (resp.  $f_i(1)$ ) the value taken by the PDF of species  $i$  for  $y \leq y_0$  (resp.  $y > y_0$ ), Eq. (17) reads:

$$\begin{aligned} 1 &= \sum_{e=1}^N \frac{f_1(y_e)}{\sum_{k=1}^{N_s} N_k f_k(y_e)} = N_{<} \frac{f_1(0)}{N_1 f_1(0) + N_2 f_2(0)} + N_{>} \frac{f_1(1)}{N_1 f_1(1) + N_2 f_2(1)} \\ &= \frac{N_{>} (1 - y_0)^{-1}}{N_1 (1 - y_0)^{-1} + N_2} \end{aligned} \quad (41)$$

$$\begin{aligned} 1 &= \sum_{e=1}^N \frac{f_2(y_e)}{\sum_{k=1}^{N_s} N_k f_k(y_e)} = N_{<} \frac{f_2(0)}{N_1 f_1(0) + N_2 f_2(0)} + N_{>} \frac{f_2(1)}{N_1 f_1(1) + N_2 f_2(1)} \\ &= \frac{N_{<}}{N_2} + \frac{N_{>}}{N_1 (1 - y_0)^{-1} + N_2} . \end{aligned} \quad (42)$$

The first equation yields:

$$N_1 (1 - y_0)^{-1} + N_2 = N_{>} (1 - y_0)^{-1} \quad (43)$$

and thus, for the second equation:

$$1 = \frac{N_{<}}{N_2} + 1 - y_0 , \quad (44)$$

which leads to Eqs. (36)-(37).

2. similarly, Eq. (10) yields

$$\mathbf{V}^{-1} = \begin{pmatrix} \frac{1}{N_{>}} & \frac{1 - y_0}{N_{>}} \\ \frac{1 - y_0}{N_{>}} & \frac{(1 - y_0)^2}{N_{>}} + \frac{y_0^2}{N_{<}} \end{pmatrix} . \quad (45)$$

For example, using Eq. (43), the  $\mathbf{V}_{11}$  component is computed as follows:

$$\mathbf{V}_{11}^{-1} = \sum_{e=1}^N \frac{f_1(y_e)f_1(y_e)}{(\sum_{k=1}^{N_s} N_k f_k(y_e))^2} = N_{>} \frac{(1-y_0)^{-2}}{(N_1(1-y_0)^{-1} + N_2)^2} = \frac{1}{N_{>}}. \quad (46)$$

And similarly for the other components. Inverting  $\mathbf{V}^{-1}$  one gets Eq. (38).

3. Eq. (15) then reproduces Eqs. (39)-(40). Namely:

$$\begin{aligned} N_1 M_1(x)\delta x &= \sum_{e \in \delta x} \frac{\mathbf{V}_{11}f_1(y_e) + \mathbf{V}_{12}f_2(y_e)}{\sum_{k=1}^{N_s} N_k f_k(y_e)} \\ &= \delta N_{<}^x \frac{\mathbf{V}_{11}f_1(0) + \mathbf{V}_{12}f_2(0)}{N_1f_1(0) + N_2f_2(0)} + \delta N_{>}^x \frac{\mathbf{V}_{11}f_1(1) + \mathbf{V}_{12}f_2(1)}{N_1f_1(1) + N_2f_2(1)} \\ &= \delta N_{<}^x \frac{\mathbf{V}_{12}}{N_2} + \delta N_{>}^x \frac{\mathbf{V}_{11}(1-y_0)^{-1} + \mathbf{V}_{12}}{N_1(1-y_0)^{-1} + N_2} \\ &= \delta N_{<}^x \frac{-\frac{1-y_0}{y_0^2} N_{<}}{N_{<} y_0^{-1}} \\ &\quad + \delta N_{>}^x \frac{(N_{>} + (\frac{1-y_0}{y_0})^2 N_{<})(1-y_0)^{-1} - \frac{1-y_0}{y_0^2} N_{<}}{N_{>}(1-y_0)^{-1}} \\ &= -\frac{1-y_0}{y_0} \delta N_{<}^x + \delta N_{>}^x \end{aligned} \quad (47)$$

and:

$$\begin{aligned} N_2 M_2(x)\delta x &= \sum_{e \in \delta x} \frac{\mathbf{V}_{21}f_1(y_e) + \mathbf{V}_{22}f_2(y_e)}{\sum_{k=1}^{N_s} N_k f_k(y_e)} \\ &= \delta N_{<}^x \frac{\mathbf{V}_{21}f_1(0) + \mathbf{V}_{22}f_2(0)}{N_1f_1(0) + N_2f_2(0)} + \delta N_{>}^x \frac{\mathbf{V}_{21}f_1(1) + \mathbf{V}_{22}f_2(1)}{N_1f_1(1) + N_2f_2(1)} \\ &= \delta N_{<}^x \frac{\mathbf{V}_{22}}{N_2} + \delta N_{>}^x \frac{\mathbf{V}_{21}(1-y_0)^{-1} + \mathbf{V}_{22}}{N_1(1-y_0)^{-1} + N_2} \\ &= \delta N_{<}^x \frac{\frac{1}{y_0^2} N_{<}}{N_{<} y_0^{-1}} \\ &\quad + \delta N_{>}^x \frac{-\frac{1-y_0}{y_0^2} N_{<}(1-y_0)^{-1} + \frac{1}{y_0^2} N_{<}}{N_{>}(1-y_0)^{-1}} \\ &= \frac{\delta N_{<}^x}{y_0}. \end{aligned} \quad (48)$$

4. it can be shown as well that Eqs. (18)-(19)-(20)-(21)-(25) hold.

Therefore, in this very simple situation where the problem of reconstructing the distributions of signal and background events is glaringly obvious, the  ${}_s\mathcal{P}lot$  formalism reproduces the expected results.

## A.2 Extended cut-and-count analysis

The above example of the previous Section A.1 is a very particular case of a more general situation where the  $y$ -range is split into  $n_y$  slices inside which one disregards the shape of the distributions of the species, whether these distributions are the same or not. Using greek letters to index the  $y$ -slices, this amounts to replacing the  $f_i(y)$  PDFs by step functions with constant values. For each  $y$ -bin  $F_i^\alpha$ , these constant values are defined by the integral over the  $y$ -bin  $\alpha$ :

$$f_i(y) \rightarrow F_i^\alpha = \int_\alpha f_i(y) dy \quad (49)$$

$$\sum_{\alpha=1}^{n_y} F_i^\alpha = 1. \quad (50)$$

With this notation, the number of events  $\bar{N}_\alpha$  expected in the slice  $\alpha$  is given by:

$$\bar{N}_\alpha = \sum_{i=1}^{N_s} N_i F_i^\alpha. \quad (51)$$

To make particularly obvious what must be the outcome of the  $sPlot$  technique, in the previous Section it was assumed that  $n_y = N_s = 2$ , and that the signal was utterly absent in one of the two  $y$ -slices:  $F_1^1 = 0$ ,  $F_1^2 = 1$ ,  $F_2^1 = y_0$  and  $F_2^2 = 1 - y_0$ .

Below one proceeds in two steps, first considering the more general case where only  $n_y = N_s$  is assumed (Section A.2.1), then considering the extended cut-and-count analysis where  $n_y > N_s$  (Section A.2.2). Since the general case discussed in the presentation of the  $sPlot$  formalism corresponds to the limit  $n_y \rightarrow \infty$ , what follows amounts to a step-by-step new derivation of the technique.

### A.2.1 Generalized cut-and-count analysis: $n_y = N_s$

When the number of  $y$ -slices equals the number of species, the solution remains obvious, if the  $N_s \times N_s$  matrix  $F_i^\alpha$  is invertible (if not, the  $N_i$  cannot be determined). In that case, one can identify the expected numbers of events  $\bar{N}_\alpha$  with the observed number of events  $N_\alpha$ , and thus:

1. one recovers the expected number of events  $N_i$  from the numbers of events  $N_\alpha$  observed in the  $n_y$  slice, by inverting Eq. (51):

$$N_i = \sum_{\alpha=1}^{N_s} N_\alpha (F^{-1})_i^\alpha, \quad (52)$$

2. the number  $N_\alpha$  being statistically independent, one obtains directly the covariance matrix:

$$\mathbf{V}_{ij} = \sum_{\alpha=1}^{N_s} N_\alpha (F^{-1})_i^\alpha (F^{-1})_j^\alpha, \quad (53)$$

3. similarly to Eq. (51), the number of events  $\delta N_\alpha^x$  observed in the  $y$ -slice  $\alpha$  and in the bin  $x$  of width  $\delta x$  is given by:

$$\delta N_\alpha^x = \sum_{i=1}^{N_s} N_i \mathbf{M}_i(x) \delta x \mathbf{F}_i^\alpha \quad (54)$$

and thus, the  $x$ -distribution of species  $i$  is:

$$\delta N_i^x \equiv N_i \mathbf{M}_i(x) \delta x = \sum_{\alpha=1}^{N_s} \delta N_\alpha^x (\mathbf{F}^{-1})_i^\alpha . \quad (55)$$

It remains to be shown that Eq. (55) is reproduced using the  $sPlot$  formalism. First, using Eq. (49) and Eq. (51), one observes that:

$$\sum_{e=1}^N \frac{f_i(y_e)}{\sum_{k=1}^{N_s} N_k f_k(y_e)} \rightarrow \sum_{\alpha=1}^{N_s} N_\alpha \frac{\mathbf{F}_i^\alpha}{\sum_{k=1}^{N_s} N_k \mathbf{F}_k^\alpha} = \sum_{\alpha=1}^{N_s} N_\alpha \frac{\mathbf{F}_i^\alpha}{N_\alpha} = \sum_{\alpha=1}^{N_s} \mathbf{F}_i^\alpha = 1 , \quad (56)$$

which shows that the obvious solution Eq. (52) is the one which maximizes the extended log-Likelihood. Similarly:

$$\mathbf{V}_{ij}^{-1} = \sum_{e=1}^N \frac{f_i(y_e) f_j(y_e)}{(\sum_{k=1}^{N_s} N_k f_k(y_e))^2} \rightarrow \sum_{\alpha=1}^{N_s} N_\alpha \frac{\mathbf{F}_i^\alpha \mathbf{F}_j^\alpha}{N_\alpha^2} = \sum_{\alpha=1}^{N_s} \frac{1}{N_\alpha} \mathbf{F}_i^\alpha \mathbf{F}_j^\alpha , \quad (57)$$

which inverse is given by Eq. (53), and thus:

$$N_i {}_s\tilde{\mathbf{M}}_i(x) \delta x \rightarrow \sum_{\alpha=1}^{N_s} \delta N_\alpha^x \frac{\sum_j \mathbf{V}_{ij} \mathbf{F}_j^\alpha}{N_\alpha} = \sum_{\alpha=1}^{N_s} \delta N_\alpha^x (\mathbf{F}^{-1})_i^\alpha . \quad (58)$$

The  $sPlot$  formalism reproduces Eq. (55).

### A.2.2 Extended cut-and-count analysis: $n_y > N_s$

In the more general situation where the number of  $y$ -slices is larger than the number of species, there is no blatant solution neither for determining the  $N_i$ , nor for reconstructing the  $x$ -distribution of each species (in particular, Eq. (52) is lost). Because of this lack of an obvious solution, what follows is a rephrasing of the derivation of the  $sPlots$ , but taking a different point of view, and in the case where the  $y$ -distributions are binned.

The best determination of the  $N_i$  (here as well as in the previous simpler situations) is provided by the Likelihood method which yields (cf. Eq. (17)):

$$\sum_{\alpha=1}^{n_y} \frac{N_\alpha \mathbf{F}_i^\alpha}{\sum_{k=1}^{N_s} N_k \mathbf{F}_k^\alpha} = 1 \quad , \forall i \quad (59)$$

with a variance matrix (cf. Eq. (10)):

$$\mathbf{V}_{ij}^{-1} = \sum_{\alpha=1}^{n_y} N_\alpha \frac{\mathbf{F}_i^\alpha \mathbf{F}_j^\alpha}{(\sum_{k=1}^{N_s} N_k \mathbf{F}_k^\alpha)^2} , \quad (60)$$

from which one computes the covariance matrix  $\mathbf{V}_{ij}$ . Instead of Eq. (52) the number of events  $N_i$  provided by Eq. (59) is shown below to satisfy the equality (cf. Eq. (20)):

$$N_i = \sum_{\alpha=1}^{n_y} N_{\alpha} ({}_s\mathcal{P})_i^{\alpha} , \quad (61)$$

where the matrix element  $({}_s\mathcal{P})_i^{\alpha}$  is the sWeight (Eq. (14)) for species  $i$  of events with  $y_e$  lying in the  $y$ -slice  $\alpha$ , namely:

$$({}_s\mathcal{P})_i^{\alpha} = \frac{\sum_{j=1}^{N_s} \mathbf{V}_{ij} \mathbf{F}_j^{\alpha}}{\sum_{k=1}^{N_s} N_k \mathbf{F}_k^{\alpha}} . \quad (62)$$

The identity of Eq. (61) is not asymptotic, it holds even for finite statistics, since the contractions with  $\mathbf{V}_{li}^{-1}$  of both the left- and right-hand sides yield the same result. Indeed (Eq. (18)):

$$\sum_{i=1}^{N_s} N_i \mathbf{V}_{li}^{-1} = \sum_{i=1}^{N_s} \sum_{\alpha=1}^{n_y} \frac{N_{\alpha} \mathbf{F}_l^{\alpha} N_i \mathbf{F}_i^{\alpha}}{(\sum_{k=1}^{N_s} N_k \mathbf{F}_k^{\alpha})^2} = \sum_{\alpha=1}^{n_y} \frac{N_{\alpha} \mathbf{F}_l^{\alpha} (\sum_{i=1}^{N_s} N_i \mathbf{F}_i^{\alpha})}{(\sum_{k=1}^{N_s} N_k \mathbf{F}_k^{\alpha})^2} = \sum_{\alpha=1}^{n_y} \frac{N_{\alpha} \mathbf{F}_l^{\alpha}}{\sum_{k=1}^{N_s} N_k \mathbf{F}_k^{\alpha}} = 1 , \quad (63)$$

which is identical to:

$$\sum_{i=1}^{N_s} \mathbf{V}_{li}^{-1} \sum_{\alpha=1}^{n_y} N_{\alpha} \frac{\sum_{j=1}^{N_s} \mathbf{V}_{ij} \mathbf{F}_j^{\alpha}}{\sum_{k=1}^{N_s} N_k \mathbf{F}_k^{\alpha}} = \sum_{\alpha=1}^{n_y} N_{\alpha} \frac{\sum_{j=1}^{N_s} (\sum_{i=1}^{N_s} \mathbf{V}_{li}^{-1} \mathbf{V}_{ij}) \mathbf{F}_j^{\alpha}}{\sum_{k=1}^{N_s} N_k \mathbf{F}_k^{\alpha}} = \sum_{\alpha=1}^{n_y} \frac{N_{\alpha} \mathbf{F}_l^{\alpha}}{\sum_{k=1}^{N_s} N_k \mathbf{F}_k^{\alpha}} = 1 . \quad (64)$$

Since Eq. (61) holds for the complete sample of events, it must hold as well for any sub-sample, provided the splitting into sub-samples is not correlated with the variable  $y$ . Namely, for all  $x$ -bin, one is guaranteed to observe, on average, the same relationship between the numbers of events  $\delta N_i^x$  and  $\delta N_{\alpha}^x$ . The  ${}_sPlot$  obtained from the weighted sum

$$\delta N_i^x = \sum_{\alpha=1}^{n_y} \delta N_{\alpha}^x ({}_s\mathcal{P})_i^{\alpha} , \quad (65)$$

is an unbiased estimator of the true distribution of  $x$  for species  $i$ . One can provide a direct proof that the above  ${}_sPlot$  of Eq. (65) reproduces the true distribution by following the same line which leads to Eq. (12). On average, using successively:

$$\langle \delta N_{\alpha}^x \rangle = \sum_{l=1}^{N_s} N_l \mathbf{M}_l(x) \mathbf{F}_l^{\alpha} \delta x \quad (66)$$

and hence:

$$\langle N_{\alpha} \rangle = \sum_x \langle \delta N_{\alpha}^x \rangle = \sum_x \mathbf{M}_l(x) \sum_{k=1}^{N_s} N_k \mathbf{F}_k^{\alpha} = \sum_{k=1}^{N_s} N_k \mathbf{F}_k^{\alpha} , \quad (67)$$

one gets:

$$\left\langle \sum_{\alpha=1}^{n_y} \delta N_{\alpha}^x ({}_s\mathcal{P})_i^{\alpha} \right\rangle = \sum_{\alpha=1}^{n_y} \left( \sum_{l=1}^{N_s} N_l \mathbf{M}_l(x) \mathbf{F}_l^{\alpha} \delta x \right) \frac{\sum_{j=1}^{N_s} \mathbf{V}_{ij} \mathbf{F}_j^{\alpha}}{\sum_{k=1}^{N_s} N_k \mathbf{F}_k^{\alpha}}$$

$$\begin{aligned}
&= \delta x \sum_{l=1}^{N_s} N_l \mathbf{M}_l(x) \left( \sum_{j=1}^{N_s} \mathbf{V}_{ij} \sum_{\alpha=1}^{n_y} \frac{F_j^\alpha F_l^\alpha}{\sum_{k=1}^{N_s} N_k F_k^\alpha} \right) \\
&= \delta x \sum_{l=1}^{N_s} N_l \mathbf{M}_l(x) \left( \sum_{j=1}^{N_s} \mathbf{V}_{ij} \sum_{\alpha=1}^{n_y} N_\alpha \frac{F_j^\alpha F_l^\alpha}{(\sum_{k=1}^{N_s} N_k F_k^\alpha)^2} \right) \\
&= \delta x \sum_{l=1}^{N_s} N_l \mathbf{M}_l(x) \left( \sum_{j=1}^{N_s} \mathbf{V}_{ij} \mathbf{V}_{jl}^{-1} \right) \\
&= N_i \mathbf{M}_i(x) \delta x \equiv \langle \delta N_i^x \rangle , \tag{68}
\end{aligned}$$

which concludes the discussion of the situation where the  $y$ -distributions are step functions.

## B Extended $s\mathcal{P}lots$ : a species is known (fixed)

It may happen that the yields of some species are not derived from the data sample at hand, but are taken to be known from other sources of information. Here, one denotes collectively as species '0' the overall component of such species. The number of expected events for species '0',  $N_0$ , being assumed to be known, is held fixed in the fit. In this Section, the indices  $i, j \dots$  run over the  $N_s$  species for which the yields  $N_i$  are fitted, the fixed species '0' being excepted ( $i, j \dots \neq 0$ ).

One can meet various instances of such a situation. Two extreme cases are:

1. the species '0' is very well known, such that the information on it contained by the data sample at hand is irrelevant. Not only is  $N_0$  already pinned down by other means, but  $\mathbf{M}_0(x)$ , the marginal distribution of the fixed species, is available,
2. the species '0' is poorly known, and the data sample at hand is unable to resolve its contribution. This is the case if the  $y$  variables cannot discriminate between species '0' against any one of the other  $N_s$  species. Stated differently, if  $N_0$  is left free to vary in the fit, the covariance matrix blows up for certain species and the measurement is lost. To avoid that, one is lead to accept an *a priori* value for  $N_0$ , and to compute systematics associated to the choice made for it. In that case, the worst case scenario is met if  $\mathbf{M}_0(x)$  is unknown as well.

It is shown below that the  $s\mathcal{P}lot$  formalism can be extended to deal with this situation, whether or not  $\mathbf{M}_0(x)$  is known, although in the latter case the statistical price to pay can be prohibitive.

### B.1 Assuming $\mathbf{M}_0$ to be known

Here, it is assumed that  $\mathbf{M}_0(x)$ , is taken for granted. Then, it is not difficult to show that the Extended  $s\mathcal{P}lot$ , which reproduces the marginal distribution of species  $n$ , is now given by:

$$N_n \tilde{M}_n(\bar{x}) \delta x = c_n \mathbf{M}_0(x) \delta x + \sum_{e \subset \delta x} s\mathcal{P}_n , \tag{69}$$

where:



- ${}_s\mathcal{P}_n$  is the previously defined sWeight of Eq. (14):

$${}_s\mathcal{P}_n = \frac{\sum_j \mathbf{V}_{nj} f_j}{\sum_k N_k f_k + N_0 f_0} , \quad (70)$$

where the covariance matrix  $\mathbf{V}_{ij}$  is the one resulting from the fit of the  $N_{i \neq 0}$  expected number of events, that is to say the inverse of the matrix:

$$\mathbf{V}_{ij}^{-1} = \sum_{e=1}^N \frac{f_i f_j}{(\sum_k N_k f_k + N_0 f_0)^2} , \quad (71)$$

- $c_n$  is the species dependent coefficient:

$$c_n = N_n - \sum_j \mathbf{V}_{nj} . \quad (72)$$

Some remarks deserve to be made:

- The Likelihood is now written:

$$\mathcal{L} = \sum_{e=1}^N \ln \left\{ \sum_{i=1}^{N_s} N_i f_i(y_e) + N_0 f_0(y_e) \right\} - \left\{ \sum_{i=1}^{N_s} N_i + N_0 \right\} . \quad (73)$$

Because  $N_0$  is held fixed, in general, its assumed value combined with the fitted values for the  $N_i$ , does not maximize it:

$$\frac{\partial \mathcal{L}}{\partial N_0} = \sum_{e=1}^N \frac{f_0}{\sum_k N_k f_k + N_0 f_0} - 1 \neq 0 . \quad (74)$$

- It follows that the sum over the number of events per species does not equal the total number of events in the sample:

$$\sum_i N_i = N - N_0 \left( \sum_{e=1}^N \frac{f_0}{\sum_k N_k f_k + N_0 f_0} \right) \neq N - N_0 . \quad (75)$$

- Similarly, the Variance Matrix Sum Rule Eq. (18) holds only for  $N_0 = 0$ :

$$\sum_i N_i \mathbf{V}_{ij}^{-1} = 1 - N_0 v_j , \quad (76)$$

where the vector  $v_j$  is defined by:

$$v_j \equiv \sum_{e=1}^N \frac{f_0 f_j}{(\sum_k N_k f_k + N_0 f_0)^2} . \quad (77)$$

- Accordingly, Eq. (19) becomes:

$$\sum_j \mathbf{V}_{jl} = N_l + N_0 \sum_j \mathbf{V}_{lj} v_j . \quad (78)$$

- Thus, as they should, the  $c_n$  coefficients vanish only for  $N_0 = 0$ :

$$c_n = -N_0 \sum_j \mathbf{V}_{nj} v_j . \quad (79)$$

- The above defined Extended  ${}_s\mathcal{P}lots$  share the same properties as the  ${}_s\mathcal{P}lots$ :
  1. They reproduce the true marginal distributions, as in Eq. (16).
  2. In particular, they are properly normalized, as in Eq. (20).
  3. The sum of  ${}_s\mathcal{P}_n^2$  reproduces  $\sigma^2[N_n]$ , as in Eq. (24).

## B.2 Assuming $\mathbf{M}_0$ to be unknown

In the above treatment, because one assumes that a special species '0' enters in the sample composition, the sWeights per event do not add up to unity, as in Eq. (21). Instead one may define the sWeights for species '0' as:

$${}_s\mathcal{P}_0 \equiv 1 - \sum_i {}_s\mathcal{P}_i \quad (80)$$

and introduce the reconstructed  ${}_s\tilde{\mathbf{M}}_0$  distribution (normalized to unity):

$${}_s\tilde{\mathbf{M}}_0(x)\delta x = \left( N - \sum_{i,j} \mathbf{V}_{ij} \right)^{-1} \sum_{e \in \delta x} {}_s\mathcal{P}_0 , \quad (81)$$

which reproduces the true distribution  $\mathbf{M}_0(x)$  if (by chance) the value assumed for  $N_0$  is the one which maximizes the Likelihood.

Taking advantage of  ${}_s\tilde{\mathbf{M}}_0(x)$ , one may redefine the Extended  ${}_s\mathcal{P}lots$  by:

$$N_n {}_s\tilde{\mathbf{M}}_n(\bar{x})\delta x = c_n {}_s\tilde{\mathbf{M}}_0(x)\delta x + \sum_{e \in \delta x} {}_s\mathcal{P}_n = \sum_{e \in \delta x} {}_{es}\mathcal{P}_n , \quad (82)$$

where the redefined sWeight which appears on the right hand side is given by:

$${}_{es}\mathcal{P}_n \equiv {}_s\mathcal{P}_n + \frac{N_i - \sum_j \mathbf{V}_{ij}}{N - \sum_{i,j} \mathbf{V}_{ij}} {}_s\mathcal{P}_0 . \quad (83)$$

It does not rely on *a priori* knowledge on the true distribution  $\mathbf{M}_0(x)$ . With this redefinition, the following properties hold:

- The set of reconstructed  $x$ -distributions  $N_i \tilde{\mathbf{M}}_i$  of Eq. (82) completed by  $(N - \sum_i N_i) \tilde{\mathbf{M}}_0$  of Eq. (81) are such that they add up in each  $x$ -bin to the number of events observed.
- The normalization constant of the  $\tilde{\mathbf{M}}_0$  distribution vanishes quadratically with  $N_0$ . It can be rewritten in the form:

$$N - \sum_{i,j} \mathbf{V}_{ij} = N_0^2 \left( v_0 - \sum_{i,j} \mathbf{V}_{ij} v_i v_j \right) , \quad (84)$$

where  $v_0$  is defined as  $v_j$  (cf. Eq. (77)) and where the last term is regular when  $N_0 \rightarrow 0$ .

- Whereas the normalization of the redefined extended  ${}_s\mathcal{Plots}$  remains correct, the sum of the redefined sWeights Eq. (83) squared is no longer equal to  $\mathbf{V}_{nn}$ . Instead:

$$\sum ({}_es\mathcal{P}_n)^2 = \mathbf{V}_{nn} + \frac{(N_n - \sum_j \mathbf{V}_{nj})^2}{N - \sum_{i,j} \mathbf{V}_{ij}} = \mathbf{V}_{nn} + \frac{\sum_{ij} \mathbf{V}_{ni} \mathbf{V}_{nj} v_i v_j}{v_0 - \sum_{ij} \mathbf{V}_{ij} v_i v_j} . \quad (85)$$

Since the expression on the right hand side is regular when  $N_0 \rightarrow 0$ , it follows that there is a price to pay to drop the knowledge of  $\mathbf{M}_0(x)$ , even though one expects a vanishing  $N_0$ . Technically, this feature stems from

$$\sum ({}_s\mathcal{P}_0)^2 = \sum {}_s\mathcal{P}_0 = N - \sum_{i,j} \mathbf{V}_{ij} . \quad (86)$$

Hence, the sum in quadrature of the  ${}_s\tilde{\mathbf{M}}_0$  uncertainties per bin diverges with  $N_0 \rightarrow 0$ . This just expresses the obvious fact that no information can be extracted on species '0' from a sample which contains no such events.