

Correcting $s\mathcal{P}lots$ in the presence of fixed yields

E. Ben Haim, L. Henry

March 16, 2016

Contents

1	Introduction	1
2	$s\mathcal{P}lots$ with fixed yields	1
3	RooStats implementation of the $s\mathcal{P}lots$ method	3
4	Proposed method and test	4
5	Conclusion	9

1 Introduction

Subtracting background from the signal in physics analyses can be performed in several ways. The $s\mathcal{P}lots$ procedure presented in Ref. [1] is an efficient and widespread method to do so; it is implemented in the RooStats package [2].

In this note, we focus on the effect of fixed yields on the results, which is documented in Ref. [1], Annexes B.1 and B.2. In Sec. 2, we shortly remind the $s\mathcal{P}lots$ subtraction procedure and its modifications in the presence of fixed yields. We then review in Sec. 3 the implementation of $s\mathcal{P}lots$ in RooStats and the issues it may introduce. Finally in Sec. 4, we propose an alternative constructor to the RooStats::SPlot class that would solve these issues, along with some tests.

2 $s\mathcal{P}lots$ with fixed yields

We consider a model with N_S event species; the yield of a species k is noted N_k and its normalised PDF f_k . The $s\mathcal{P}lots$ procedure allows to use the information from a fit performed on a discriminating variable X to extract the distributions of the control variable Y for the different species.

16 A main ingredient of the *sPlots* calculation is the covariance matrix V of the fit, which
 17 can be taken from the output of a fit routine (e.g. `TMinuit` [3]). Alternatively, its inverse
 18 can be directly computed using:

$$V_{ij}^{-1} = \sum_{e=1}^N \frac{f_i(e)f_j(e)}{\left(\sum_{k=1}^{N_S} N_k f_k(e)\right)^2}, \quad (1)$$

19 where the sum is running over N events, and $f_i(e)$ designates the value of the PDF f_i for
 20 the event e . We can then use the covariance matrix to compute, for each event species n ,
 21 the per-event sWeight $sP_n(e)$, using:

$$sP_n(e) = \frac{\sum_{j=1}^{N_S} V_{nj} f_j(X_e)}{\sum_{k=1}^{N_S} N_k f_k(X_e)}. \quad (2)$$

22 The distribution of the event species n on the control variable Y is then estimated by the
 23 $sM_n(Y)$ distribution, defined by:

$$N_n \cdot sM_n(Y) \cdot \delta Y = \sum_{e \in [Y-\delta Y, Y+\delta Y]} sP_n(e). \quad (3)$$

24 We now introduce another event species in the model, with a fixed yield N_0 and a normalised
 25 PDF f_0 . The covariance matrix changes, and its inverse becomes:

$$V_{ij}^{-1} = \sum_{e=1}^N \frac{f_i(e)f_j(e)}{\left(\sum_{k=1}^{N_S} N_k f_k(e) + N_0 f_0(e)\right)^2}, \quad (4)$$

26 whereas the per-event sWeight becomes:

$$sP_n(e) = \frac{\sum_{j=1}^{N_S} V_{nj} f_j(X_e)}{\sum_{k=1}^{N_S} N_k f_k(X_e) + N_0 f_0(X_e)} = \frac{\sum_{j=1}^{N_S} V_{nj} f_j(X_e)}{\sum_{k=1}^{N_S} N_k f_k(X_e)} \frac{\sum_{k=1}^{N_S} N_k f_k(X_e)}{\sum_{k=1}^{N_S} N_k f_k(X_e) + N_0 f_0(X_e)}. \quad (5)$$

27 This expression differs from Eq. 2 only by an event-by-event factor that depends on the
 28 yields N_k and the PDFs f_k .

29 In the case where the distribution of the control variable for the species with fixed
 30 yield, $M_0(Y)$, is known, the distribution of the control variable Y for the species n is:

$$N_n \cdot sM_n(Y) \cdot \delta Y = \sum_{e \in [Y-\delta Y, Y+\delta Y]} sP_n(e) + c_n \cdot M_0(Y), \quad (6)$$

31 where

$$c_n = N_n - \sum_{j=1}^{N_S} V_{nj} \quad (7)$$

32 is a coefficient depending uniquely on the considered event species. It quantifies the impact
 33 of the species with fixed yields on $sM_n(Y)$, and vanishes only if $N_0 = 0$.

34 One of the issues of extracting sWeights using a tool that treats the varied and fixed
 35 yields on an equal footing is related to the calculation of the covariance matrix. To
 36 illustrate this, we consider a model including three species: a signal S, a combinatorial
 37 background C, and another background B with a fixed yield. Using Eq. 4, we calculate
 38 the inverse of the covariance matrix, and obtain:

$$V^{-1} = \begin{pmatrix} (V^{-1})_{SS} & (V^{-1})_{SC} & (V^{-1})_{SB} \\ (V^{-1})_{CS} & (V^{-1})_{CC} & (V^{-1})_{CB} \\ (V^{-1})_{BS} & (V^{-1})_{BC} & (V^{-1})_{BB} \end{pmatrix}, \quad (8)$$

39 which incorrectly includes terms related to the species with fixed yields. As the terms
 40 $(V^{-1})_{SB}$ and $(V^{-1})_{CB}$ do not vanish *a priori*, the inverse of this matrix has no clear link
 41 with the correct covariance matrix. Also, we notice that in the case where $f_B = f_C$ or
 42 $f_B = f_S$, the matrix is no longer invertible (it has two identical columns and lines), whereas
 43 sWeights should still be calculable. This shows that there is something fundamentally
 44 flawed with this approach.

45 **3 RooStats implementation of the *sPlots* method**

46 The `RooStats::SPlot` method, used to calculate sWeights, computes the inverse of the
 47 covariance matrix using Eq. 1 with a list of yields that the user provides as an argument.
 48 The covariance matrix itself is then obtained from its inverse. As shown in Sec. 2, in
 49 the case where there are some fixed yields in the arguments, this results in an incorrect
 50 covariance matrix.

51 However, building the `RooStats::SPlot` object using only the varied yields is also
 52 incorrect, as it would result in using Eq. 1 and Eq. 2 rather than Eq. 4 and Eq. 5. Correcting
 53 the `sWeights` event-by-event using Eq. 5 is not possible either, as the covariance matrix is
 54 not correctly calculated.

55 4 Proposed method and test

56 It is clear from Sec. 2 that it is necessary to differentiate the fixed yields from the others
 57 in the `RooStats::SPlot` object. To address this requirement, we propose an alternative
 58 constructor to the `RooStats::SPlot` object, shown in Fig. 1 along with the original
 59 constructor¹. We stress that this alternative constructor does not perform the correction
 60 shown in Eq. 6. However, it calculates and stores the c_n coefficients in a new attribute of
 61 the class, as shown in Fig. 2 and Fig. 3.

62 We test this alternative constructor and the c_n extraction tool using a toy model
 63 containing 3 event species: a signal S, a combinatorial background C, and a peaking
 64 background B with a fixed yield. The PDFs of these species on the discriminating variable
 65 X and on the control variable Y are taken as:

- 66 • Signal (S): Gaussian ($\mu = 0, \sigma = 0.1$) for X ; Gaussian ($\mu = 0, \sigma = 0.05$) for Y .
- 67 • Combinatorial background (C): Constant for X ; Gaussian ($\mu = -0.5, \sigma = 0.05$) for
 68 Y .
- 69 • Fixed background. (B): Gaussian ($\mu = 0, \sigma = 0.05$) for X ; Gaussian ($\mu = 0.5, \sigma =$
 70 0.05) for Y .

71 Both variables X and Y are defined in the interval $[-1,1]$. This model is chosen in order to
 72 ensure a sizable nuisance of the species with fixed yield on the discriminating variable,
 73 whereas its impact on the control variable is easy to spot. We generate a sample of
 74 16000 events, including 5000 signal events, 1000 fixed background events, and 10000
 75 combinatorial background events. We then assume a wrong hypothesis on the yield of the
 76 fixed background ($N_B = 1200$), in order to simulate the general case where the value of
 77 the fixed yield is not precisely known. We show the result of the one-dimensional fit on
 78 the variable X using this model, along with the projection of this fit on the variable Y , in
 79 Fig. 4.

80 We consider two approaches, A and B: the former is the calculation of the `sWeights`
 81 with the original `RooStats::SPlot` method, providing only the list of varying yields to
 82 the `RooStats::SPlot` constructor; the latter is the proposed approach, where we use the
 83 alternative `RooStats::SPlot` constructor providing the list of all yields and the list of
 84 fixed yields.

¹Providing an empty `RooArgSet` as the `fixedYields` argument of the alternative constructor gives the same results as calling the original constructor with the same arguments

Original constructor:

```
SPlot::SPlot(const char* name, const char* title ,
            RooDataSet& data, RooAbsPdf* pdf,
            const RooArgList &yieldsList , const RooArgSet &projDeps ,
            bool includeWeights, bool cloneData, const char* newName):
    TNamed(name, title)
```

Implementation of the alternative constructor:

```
SPlot::SPlot(const char* name, const char* title ,
            RooDataSet& data, RooAbsPdf* pdf,
            const RooArgList &allYieldsList ,const RooArgList &fixedYields ,
            const RooArgSet &projDeps ,
            bool includeWeights, bool cloneData, const char* newName):
    TNamed(name, title)
{
    /* Skipped: Original body of the SPlot constructor. */
    // Add check that yieldsList contains only varying yields
    // Check that fixed yields are in the allYields arguments
    iter = fixedYields.createIterator();
    while((arg=(RooAbsArg*)iter->Next()))
        if (!(allYieldsList.contains(*arg)))
        {
            // Throw exception and error message.
        }
    // Call new method to build sWeights, with fixed yields
    this->AddSWeight(pdf, allYieldsList , fixedYields , projDeps , includeWeights);
}
```

Figure 1: Snippets of code showing the original (top) and the alternative (bottom) constructors for the `RooStats::SPlot` class.

```

void SPlot::AddSWeight( RooAbsPdf* pdf, const RooArgList &allYieldsList,
                      const RooArgList &fixedYields,
                      const RooArgSet &projDeps, bool includeWeights)
{
    /* Skipped: store the constant parameters (other than yields) values */
    /* Store indexes of the allYieldsList corresponding to variables yields */
    TIterator *it = allYieldsList.createIterator();
    RooAbsArg* arg;
    unsigned int iArg(0);
    std::vector<unsigned int> varIndexes;
    while ((arg = (RooAbsArg*) it->Next()) != NULL){
        if (!(fixedYields.find(arg->GetName())))
            varIndexes.push_back(iArg);
        iArg++;}
    /* We now have two indexes over which we iterate */
    Int_t nAllSpec = allYieldsList.getSize();
    Int_t nVarSpec = allYieldsList.getSize() - fixedYields.getSize();

    /* Skipped:
    -store the initial yield parameters
    -calculate the value of the component pdf for each event and
    species. */
    /* Inverse of the covariance matrix */
    TMatrixD covInv(nVarSpec, nVarSpec);
    /* Skipped: Initialisation to 0 of each covInv(i,j) */
    for (Int_t ievt = 0; ievt < numevents; ++ievt){
        fSData->get(ievt) ;
        /* Sum for the denominator */
        Double_t dsum(0);
        for(Int_t k = 0; k < nAllSpec; ++k)
            dsum += pdfvalues[ievt][k] * yieldvalues[k] ;

        for(Int_t n=0; n<nVarSpec; ++n)
            for(Int_t j=0; j<nVarSpec; ++j)
                if(includeWeights == kTRUE)
                    covInv(n,j) += fSData->weight()
                        *pdfvalues[ievt][varIndexes[n]]
                        *pdfvalues[ievt][varIndexes[j]]/(dsum*dsum) ;
                else
                    covInv(n,j) +=
                        pdfvalues[ievt][varIndexes[n]]
                        *pdfvalues[ievt][varIndexes[j]]/(dsum*dsum) ;}
    /* Invert to get the covariance matrix */
    TMatrixD covMatrix(TMatrixD::kInverted, covInv);
}

```

Figure 2: Proposed addition to the `RooStats::SPlot` class to compute `sWeights` in the presence of fixed yields (first part).

```

// Create and label the variables used to store the SWeights
for(Int_t k=0; k<nVarSpec; ++k){
  /* Skipped: usual variables definitions, including sWeights
  variables */

  // c_n coefficients
  wname = std::string(yieldvars[varIndexes[k]]->GetName()+"_c");
  var = new RooRealVar(wname.c_str(),wname.c_str(),0);
  double cVal = yieldvalues[varIndexes[k]];
  for (Int_t n = 0 ; n<nVarSpec; ++n)
    cVal -= covMatrix[k][n];
  var->setVal(cVal);
  fSWeightCoefs.add(*var); // new attribute of the class.
}
// Create and fill a RooDataSet with the SWeights
RooDataSet* sWeightData = new RooDataSet("dataset", "", sweightset);
for(Int_t ievt = 0; ievt < numevents; ++ievt){
  fSData->get(ievt) ;
  // sum for denominator
  Double_t dsum(0);
  for(Int_t k = 0; k < nAllSpec; ++k)
    dsum += pdfvalues[ievt][k] * yieldvalues[k] ;
  // covariance weighted pdf for each species
  for(Int_t n=0; n<nVarSpec; ++n){
    Double_t nsum(0) ;
    for(Int_t j=0; j<nVarSpec; ++j)
      nsum += covMatrix(n,j) * pdfvalues[ievt][varIndexes[j]] ;
    if(includeWeights == kTRUE)
      sweightvec[n]->setVal(fSData->weight() * nsum/dsum) ;
    else
      sweightvec[n]->setVal( nsum/dsum) ;
  }
  /* Skipped: fill the dataset with sweightvec */
}
// Add the SWeights to the original dataset
fSData->merge(sWeightData);
/* Skipped: reinitialise all parameters and yields */
return;
}

```

Figure 3: Proposed addition to the `RooStats::SPlot` class to compute sWeights in the presence of fixed yields (second part).

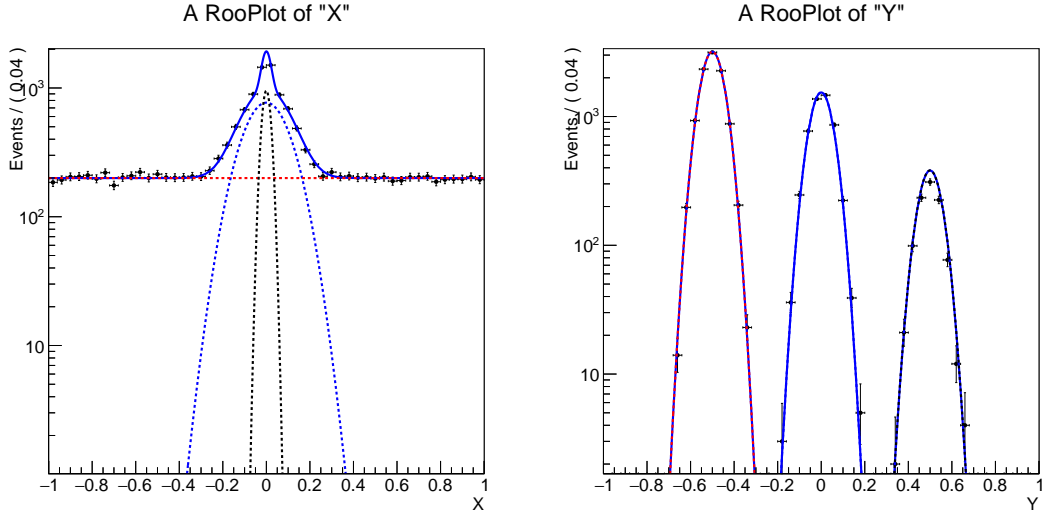


Figure 4: Result of the fit on X , projected on the X dimension (left) and the Y dimension (right). The signal is displayed in blue, the combinatorial background in red, and the peaking background in black. As expected, the projection of the result on the Y dimension shows that we overestimated (on purpose) the number of background events.

85 For each of these two methods, we show in Fig. 5 the samples with signal and background
 86 weights applied, both before and after the correction of Eq. 6. The results are satisfactory
 87 for approach B after the c_n correction, for both the signal and the combinatorial background
 88 distributions. On the contrary, approach A provides an acceptable description of the signal
 89 shape, but shows large discrepancies for the combinatorial background.

90 For each event species, we also show in Table. 1 the sum of sWeights and of its associated
 91 c_n , compared to the fitted yield. According to Ref. [1], the sum has to be compatible
 92 with the yield, which is clearly not the case for approach A.

Species	$\sum_e sP_A(e)+c_A$	$\sum_e sP_B(e)+c_B$	Fitted yield
Signal	-7336.15	4975.27	4975.35
Comb.	16361.5	10007.5	10007.5
Bkg	-	-	1200.

Table 1: Summary of yields and sum of sWeights for the two approaches. As shown in Ref. [1], in both approaches the sum of all yields extracted from the fit is not equal to the total number of events.

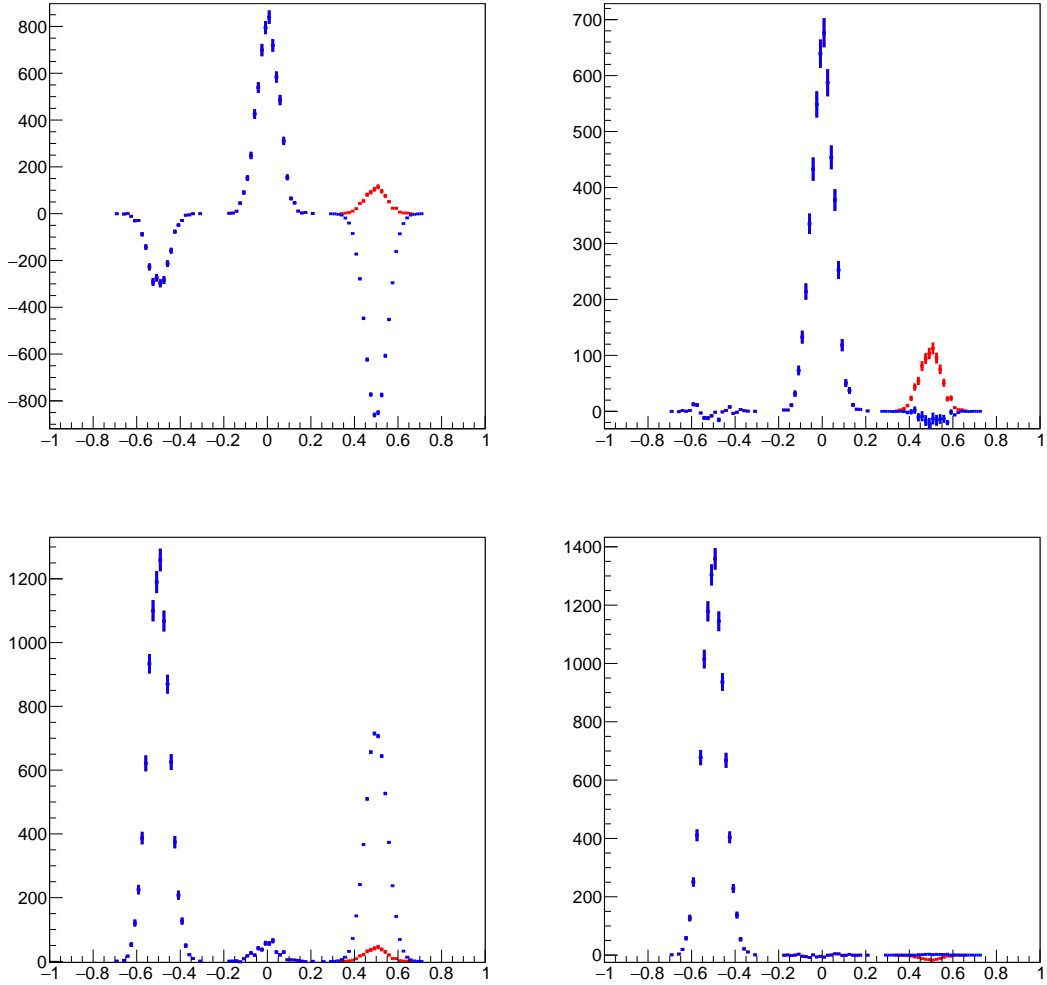


Figure 5: Distribution of Y in data samples $sWeighted$ according to signal (combinatorial) weights are shown on top (bottom). Results of the approach A(B) are shown on the left (right). The red points correspond to the distribution of the sum of $sWeights$, whereas the blue points represent the final distribution (after correction using Eq. 6).

93 5 Conclusion

94 Models that include an event species with fixed yields require a specific treatment when
 95 using the $sPlots$ method, that the `RooStats::SPlot` implementation does not provide.
 96 In this note, we propose a straightforward modification of this class that allows to extract
 97 the correct $sPlots$ in the case of fixed yields. These modifications do not remove the need
 98 to correct the distributions using Eq. 6, but allow to calculate the c_n coefficients inside
 99 the `RooStats::SPlot` object in a way that is coherent with the $sWeights$ extraction.

100 We tested this additional code, both in terms of compatibility with the former imple-
101 mentation (not shown here), and in terms of expected results. The results are satisfactory,
102 and show a clear improvement compared to the original approach, especially in terms of
103 normalisation properties.

104 **References**

- 105 [1] M. Pivk and F. R. Le Diberder, *sPlot: a statistical tool to unfold data distributions*, .
- 106 [2] *RooStats package webpage*, [https://root.cern.ch/root/html/ROOFIT_ROOSTATS_](https://root.cern.ch/root/html/ROOFIT_ROOSTATS_Index.html)
107 [Index.html](https://root.cern.ch/root/html/ROOFIT_ROOSTATS_Index.html).
- 108 [3] *TMinuit webpage*, <https://root.cern.ch/doc/v606/classTMinuit.html>.