# $_s\mathcal{P}lot$: a statistical tool to unfold data distributions

physics/0402083, to be published in *Nucl. Inst. Meth.*

## Muriel Pivk
### CERN

## 30th September 2005
### ROOT Workshop, CERN

---

# 1 Motivation (1)

**Problem to solve when performing an analysis**

Data sample $\equiv$ black box

Few signal events and lots of background

$\Longrightarrow$ How to  - distinguish them ?

- extract the physics of the signal ?

- probe the validity of analysis ?

$\rightarrow$ check the distributions of events !

**The context of BABAR in 2002**

First goal: $\sin 2\beta$ , *Phys. Rev. Lett.*89:201802 (2002)

- "Golden mode" decay analysis: $B^0 \rightarrow J/\psi \, K_s^0$
- Low background

$\Longrightarrow$ No need for a particular tool

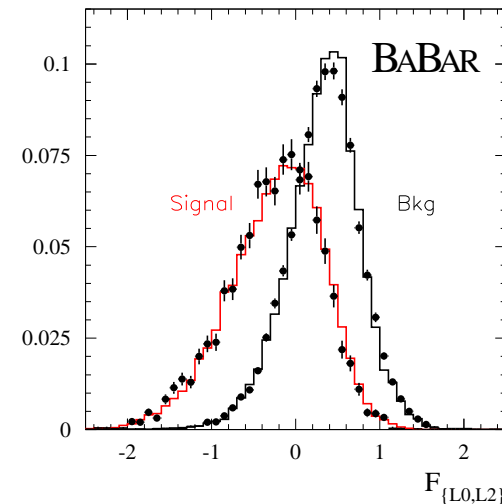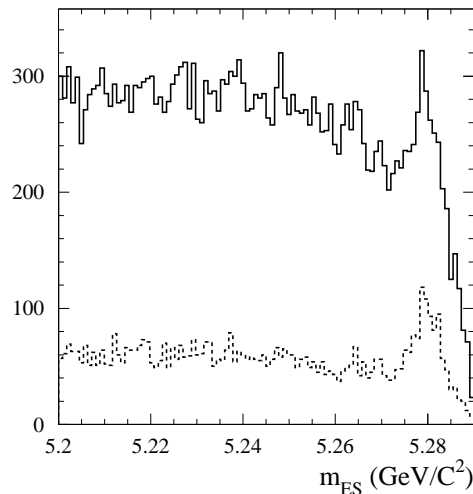**Very rare decay analysis** $\sin 2\alpha$ possible thanks to luminosity

$\Longrightarrow$ Decay channel $B^0 \to h^+ h^-$ $(h = \pi, K)$

Event selection:

- $m_{\mathrm{ES}}$ : reconstructed mass of the $B$ candidate
- $\Delta E$ : difference of energy between $B$ candidate and $\sqrt{s}/2$

Signal/background discrimination:

- Huge $e^+ e^- \to q\bar{q}$ background
- $\mathcal{F}$ : Fisher discriminant, uses topology difference of the events
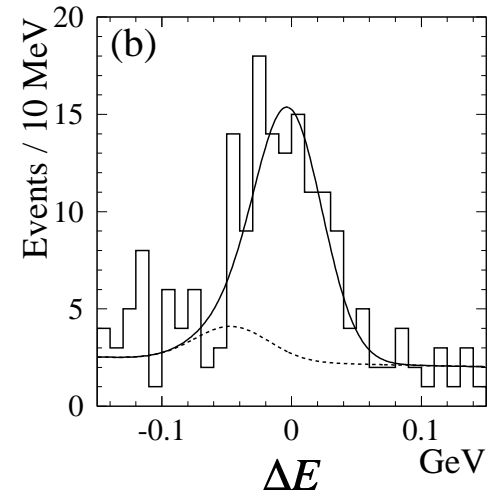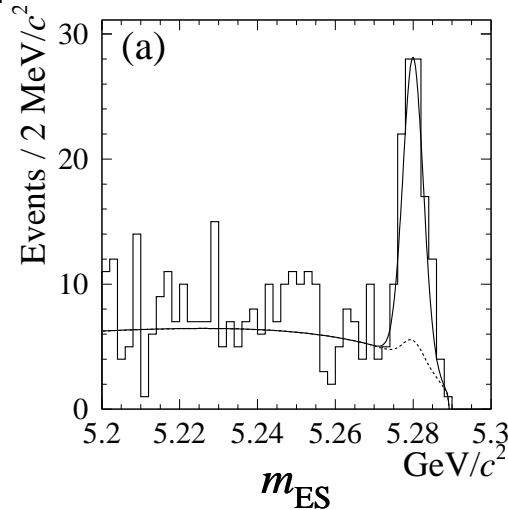


Among 88 million of $B\overline{B}$ pairs

$\Longrightarrow$ 156 $\pi^+\pi^-$ and 588 $K^+\pi^-$ among 26k events

The question is: how to check the distributions of events ?

**Solution** ? "Projection plots"

Cut applied on the $\mathcal{L}$ ratio to reduce background



1. subset of sample only

2. signal and background events mixed

3. hard (impossible) if distributions not really different (Fisher ?)

**Solution** ! $_s\mathcal{P}lot$

New tool: firstly meant as projection plots optimization

1. keep all data

2. separate signal and background

3. applicable for ANY variable

## $\boxed{\textbf{2.1}}$ **Likelihood analyses**

### Extended log-likelihood

$$\mathcal{L} = \sum_{e=1}^{N} \ln \left\{ \sum_{i=1}^{N_{\mathrm{s}}} N_i \mathrm{f}_i(y_e) \right\} - \sum_{i=1}^{N_{\mathrm{s}}} N_i \tag{1}$$

- $N$ : number of events in the data sample
- $e$ : event number
- $N_{\mathrm{s}}$ : number of species in the data sample
- $i$ : species number (signals, backgrounds)
- $y$ : discriminating variables
- $\mathrm{f}_i(y_e)$ : distribution of variables $y$ of species $i$ for event $e$, normalized to unity

### Analysis $B^0 \to h^+h^-$

- $N_{\mathrm{s}}$ : three species
- $i$ : signal $\pi^+\pi^-$ $(N_{\pi\pi})$, signal $K^+\pi^-$ $(N_{K\pi})$, background $q\bar{q}$ $(N_{q\bar{q}})$
- $y$ : $m_{\mathrm{ES}}$, $\Delta E$, $\mathcal{F}$ $(\dots)$

**Distribution of $x$ for species n, $x \in y$, using the (naive) weight**

$$\mathcal{P}_{\mathrm{n}}(y_e) = \frac{N_{\mathrm{n}}\mathrm{f}_{\mathrm{n}}(y_e)}{\sum_{k=1}^{\mathrm{N_s}} N_k\mathrm{f}_k(y_e)} \tag{2}$$

The reconstructed distribution $\tilde{\mathrm{M}}_{\mathrm{n}}$ of variable $x$ is defined by:

$$N_{\mathrm{n}}\tilde{\mathrm{M}}_{\mathrm{n}}(x)\delta x \;\equiv\; \sum_{e \subset \delta x}^{N} \mathcal{P}_{\mathrm{n}}(y_e) \tag{3}$$

Replacing $\sum_{e \subset \delta x}^{N}$ by $\int dy$ (total pdf) $\delta(x(y) - x)\delta x$:

$$N_{\mathrm{n}}\tilde{\mathrm{M}}_{\mathrm{n}}(x) \;=\; \int dy \sum_{i=1}^{\mathrm{N_s}} N_i\mathrm{f}_i(y)\delta(x(y) - x)\frac{N_{\mathrm{n}}\mathrm{f}_{\mathrm{n}}(y)}{\sum_{k=1}^{\mathrm{N_s}} N_k\mathrm{f}_k(y)} \tag{4}$$

$$\;=\; N_{\mathrm{n}} \int dy\,\delta(x(y) - x)\mathrm{f}_{\mathrm{n}}(y) \tag{5}$$

$$\;\equiv\; N_{\mathrm{n}}\mathbf{M}_{\mathrm{n}}(x) \tag{6}$$

where $\mathbf{M}_{\mathrm{n}}(x)$ is the TRUE distribution of variable $x$ for species n
$\Longrightarrow$ Not a clean test:
the Pdf of x is implicitly used to reconstruct itself ... can we avoid it ?

### Distribution of $x$, $x \notin y$

$$N_\mathrm{n}\tilde{\mathrm{M}}_\mathrm{n}(x) = \int dy \sum_{i=1}^{\mathrm{N_s}} N_i \mathbf{M}_i(x) \mathrm{f}_i(y) \frac{N_\mathrm{n}\mathrm{f}_\mathrm{n}(y)}{\sum_{k=1}^{\mathrm{N_s}} N_k \mathrm{f}_k(y)} \tag{7}$$

$$= N_\mathrm{n} \sum_{i=1}^{\mathrm{N_s}} \mathbf{M}_i(x) \left( N_i \int dy \frac{\mathrm{f}_\mathrm{n}(y)\mathrm{f}_i(y)}{\sum_{k=1}^{\mathrm{N_s}} N_k \mathrm{f}_k(y)} \right) \tag{8}$$

$$\neq N_\mathrm{n}\mathbf{M}_\mathrm{n}(x) \tag{9}$$

### But but but ... !

Variance matrix:

$$\mathbf{V}_{\mathrm{n}i}^{-1} = \frac{\partial^2(-\mathcal{L})}{\partial N_\mathrm{n} \partial N_i} = \sum_{e=1}^{N} \frac{\mathrm{f}_\mathrm{n}(y_e)\mathrm{f}_i(y_e)}{(\sum_{k=1}^{\mathrm{N_s}} N_k \mathrm{f}_k(y_e))^2} \tag{10}$$

$$= \int dy \frac{\mathrm{f}_\mathrm{n}(y)\mathrm{f}_i(y)}{\sum_{k=1}^{\mathrm{N_s}} N_k \mathrm{f}_k(y)} \tag{11}$$

Eq. (8) becomes $\tilde{\mathrm{M}}_\mathrm{n}(x) = \sum_{i=1}^{\mathrm{N_s}} \mathbf{M}_i(x) N_i \mathbf{V}_{\mathrm{n}i}^{-1}$

$\implies$ By inversion:

$$N_\mathrm{n}\mathbf{M}_\mathrm{n}(x) = \sum_{i=1}^{\mathrm{N_s}} \mathbf{V}_{\mathrm{n}i}\tilde{\mathrm{M}}_i(x) \tag{12}$$

**New tool $_s\mathcal{P}lot$: weight computed for each event and each species**
$N_s$ species in the sample, discriminating variables $y$, $f_i(y)$ their pdfs.

For species $n$ :

$$\boxed{\;{_s}\mathcal{P}_n(y_e) = \frac{\sum_{i=1}^{N_S} \mathbf{V}_{ni} f_i(y_e)}{\sum_{k=1}^{N_S} N_k f_k(y_e)}\;} \qquad (13)$$

with $\mathbf{V}_{ni}$ the covariance matrix of the fit (number of events)
The TRUE distribution of $x$ $(x \notin y)$ is:

$$N_n \mathbf{M}_n(x) \equiv \sum_{e \subset \delta x} {_s}\mathcal{P}_n(y_e) \qquad (14)$$

**NB**

- The most discriminating the variables are, the most powerful $_s\mathcal{P}lot$ is.
- The variables must be uncorrelated (already necessary with the $\mathcal{L}$).

## Normalization

1. Each $x$-distribution is properly normalized:

$$\sum_{e=1}^{N} {}_s\mathcal{P}_{\rm n}(y_e) = N_{\rm n} \qquad (15)$$

2. The contributions ${}_s\mathcal{P}_{\rm n}(y_e)$ add up to the number of events actually observed in each $x$-bin. For any event:

$$\sum_{{\rm n}=1}^{{\rm N_s}} {}_s\mathcal{P}_{\rm n}(y_e) = 1 \qquad (16)$$

## Uncertainties

3. For each species:

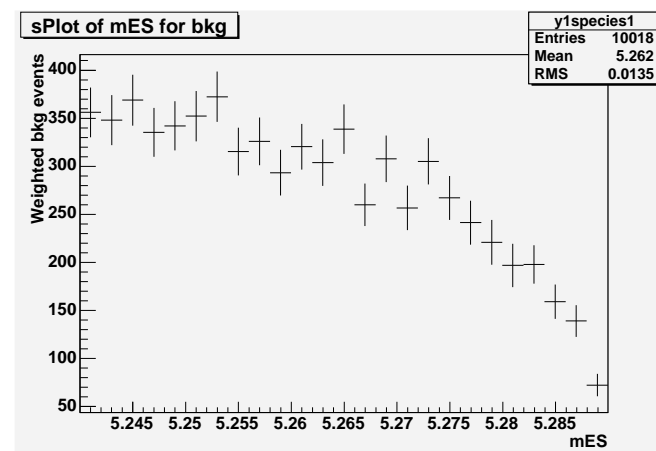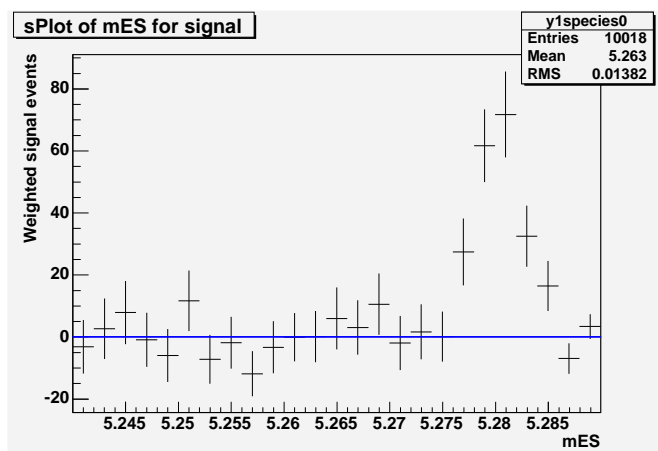$$\sum_{e=1}^{N} ({}_s\mathcal{P}_{\rm n}(y_e))^2 = \sigma^2(N_{\rm n}) \qquad (17)$$

as given by the fit

# 3 Easy implementation

## The way to follow

1. Perform the fit to obtain the $N_{\mathrm{n}}$ of each n species present in the data sample without the variable one wants to get the distribution of

2. Compute the sWeights $_s\mathcal{P}$ following Eq. 13, using the covariance matrix given by Minuit or computed directly

3. Fill histograms with the value of the variable $x$ weighted with the sWeights $_s\mathcal{P}$ for each species present in the data sample
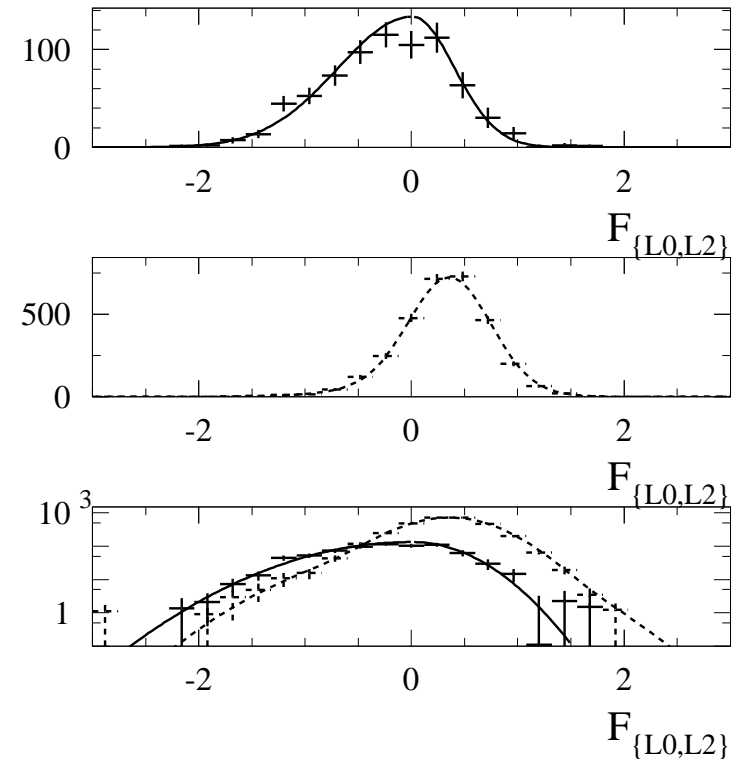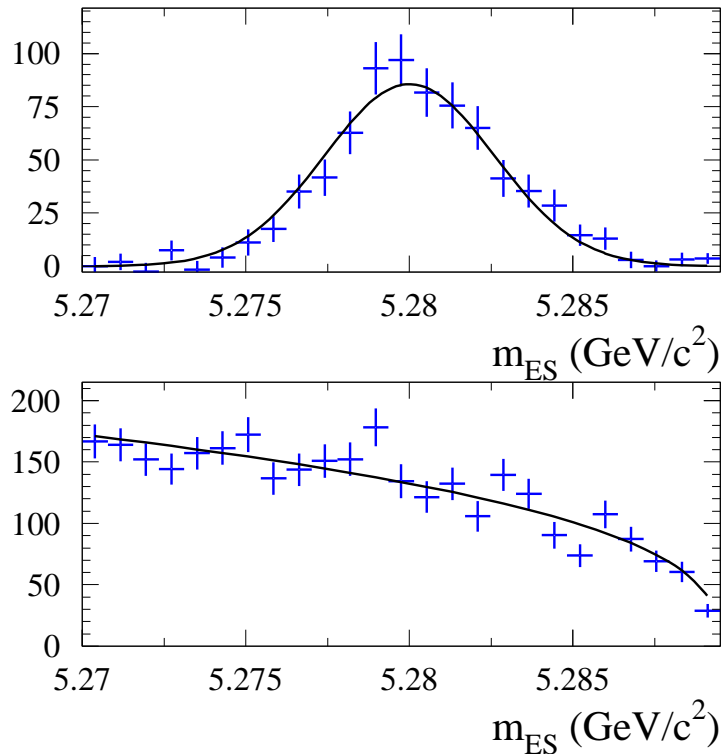
## Tool $_s\mathcal{P}lot$ in ROOT

Class TSPlot: implemented by Anna Kreshuk, to be released soon

BABAR **data:** $_s\mathcal{P}lots$ **of** $m_{\mathrm{ES}}$ **and** $\mathcal{F}$

Distributions used in the fit are superimposed



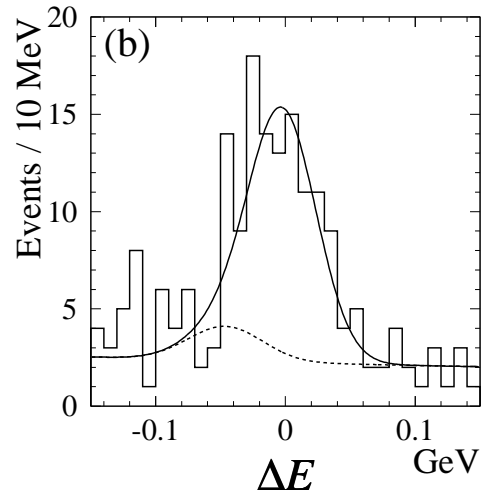- $\Delta\mathrm{E}$ and $\mathcal{F}$ only
- $m_{\mathrm{ES}}$ not in the fit

- $m_{\mathrm{ES}}$ and $\Delta\mathrm{E}$ only
- $\mathcal{F}$ not in the fit

$\Longrightarrow$ Very good agreement

$\Longrightarrow$ Optimal tool to validate an analysis ! Still for Fisher !
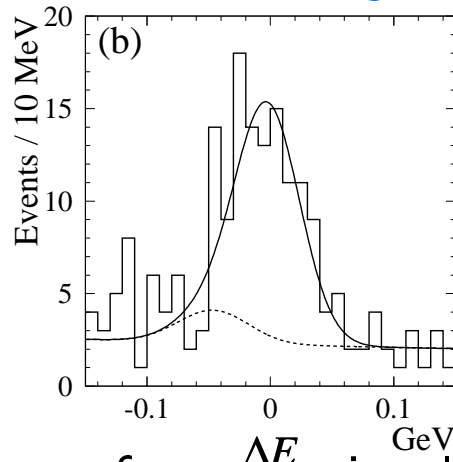
## Comparison with "projection plots"



Projection plot :
- Cut on the $\mathcal{L}$ ratio: signal loss and remaining background
- Uncertainties related to signal + background

$\Longrightarrow$ Excess of events: signal ? background ?

## Comparison with "projection plots"



Projection plot :
- Cut on the $\mathcal{L}$ ratio: signal loss and remaining background
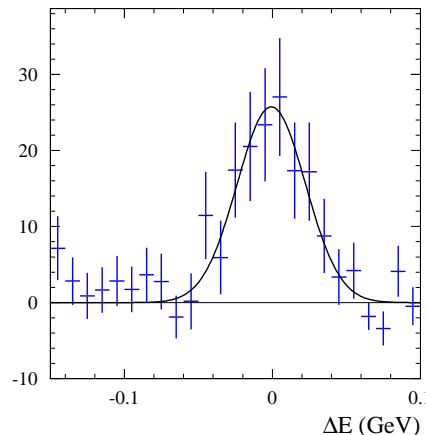- Uncertainties related to signal + background

$\Longrightarrow$ Excess of events: signal ? background ?



$_s\mathcal{P}lot$ : Can reveal subtle effects
- No cut applied: keep all the signal events and get rid of all the background ones (statistically)
- Uncertainties related to the signal only

$\Longrightarrow$ Signal **!** radiative events $(B^0 \rightarrow \pi^+\pi^-\gamma)$ ignored in the analysis

$\Longrightarrow$ $\mathcal{B}(B^0 \rightarrow h^+h^-)$ under-estimated by about $10\%$ (**!!**)
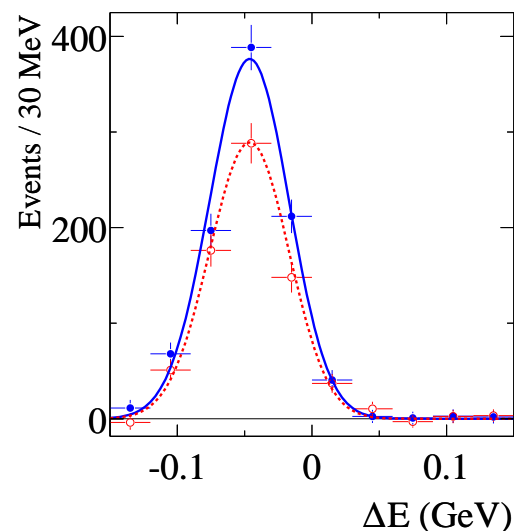
Confirmed later for different charmless *BABAR* analyses (hep-ex/0508046)

## Only BABAR so far ...

1. *Branching fractions and CP asymmetries in $B^0 \to K^+ K^- K^0_S$ and $B^+ \to K^+ K^0_S K^0_S$, Phys. Rev. Lett.93:181805, 2004*
2. *Measurement of neutral $B$ decay branching fractions to $K^0_S \pi^+ \pi^-$ final states, Phys. Rev. D70:091103, 2004*
3. *BF and CP asymmetries in $B^0 \to \pi^0 \pi^0$, $B^+ \to \pi^+ \pi^0$ and $B^+ \to K^+ \pi^0$ decays and isospin analysis of the $B \to \pi\pi$ system, Phys. Rev. Lett.94:181802, 2005*
4. *Measurement of $CP$ asymmetries in $B^0 \to \phi K^0_S$ and $B^0 \to K^+ K^- K^0_S$ decays, Phys. Rev. D71:091102, 2005*
5. . . .

## Observation of direct CP violation in $B^0 \to K^+ \pi^-$

*Phys. Rev. Lett..93:131801 (2004)*

- $N_{K^+\pi^-} + N_{K^-\pi^+} = 1606 \pm 51$

  $N_{K^+\pi^-} = 910$

  $N_{K^-\pi^+} = 696$

- $A_{K\pi} = -0.133 \pm 0.030 \pm 0.009$

# 5 Summary and conclusion

**New tool** $_s\mathcal{P}lot$ **:** optimal for information !

1. Only data involved
2. No bias ($_s\mathcal{P}lot$ted variable not in the fit)
3. Shows signal and background separately
4. Statistical uncertainties
5. Easy to use ! Moreover class TSPlot in ROOT very soon

$\Longrightarrow$ Excellent tool to validate an analysis

Reveal subtle effects : $B^0 \to h^+ h^- (\gamma)$

$\Longrightarrow$ Excellent tool to perform an analysis in Dalitz

**More in the reference**

- Detailed explanations
- Case where species fixed in the fit

**Shall be useful beyond $B$ physics**
Higgs searches, SUperSYmetry, . . .