

# Improved **DATA MODELLING** with Goodness-of-Fit & Likelihood-ratio tests

Terascale Statistics School, DESY, Feb 23, 2018

Olaf Behnke (DESY)

# Content

## Goodness-of-Fit tests for

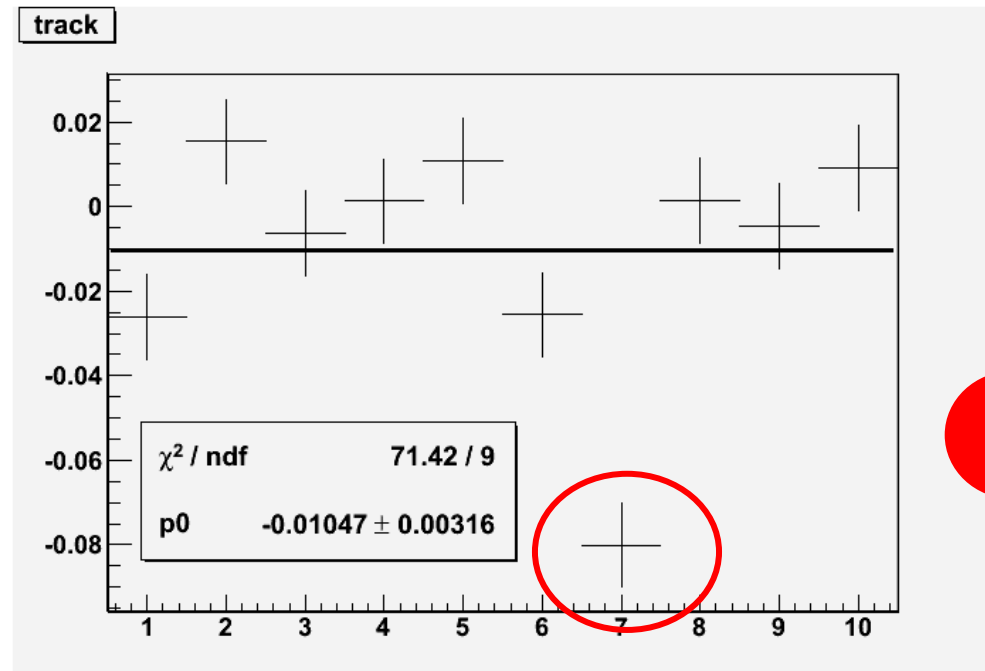
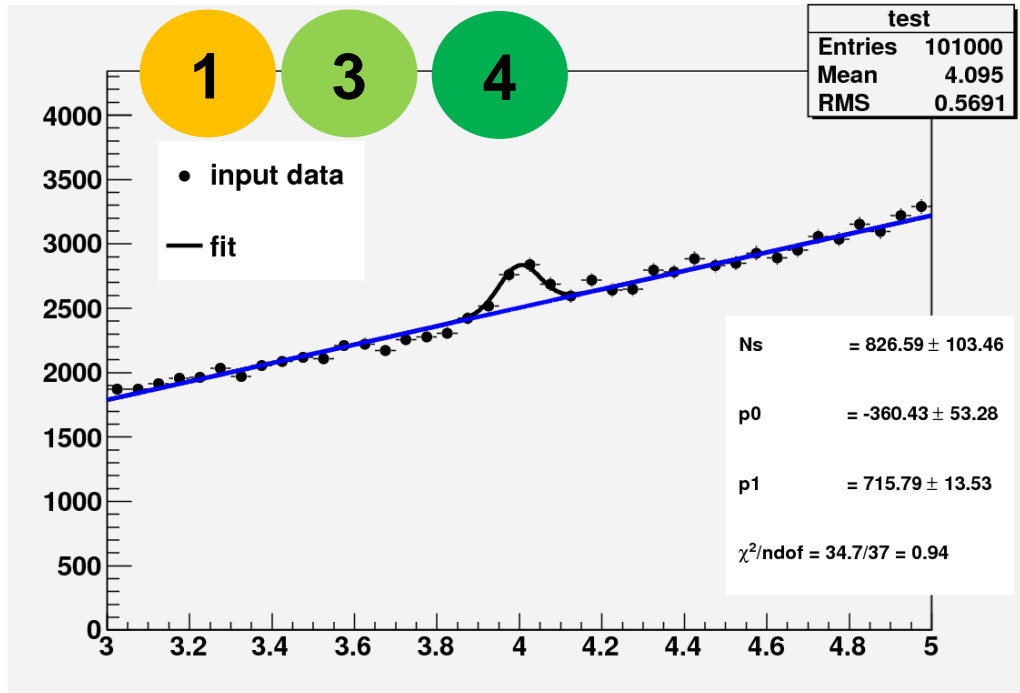
1 Checking data modelling

2 Outlier rejection

## Likelihood ratio tests for background

3 Optimal parametrisation

4 Shape systematics (discrete profiling)



1 GOF tests for  
*checking data modelling*

# Searches with Likelihood ratio

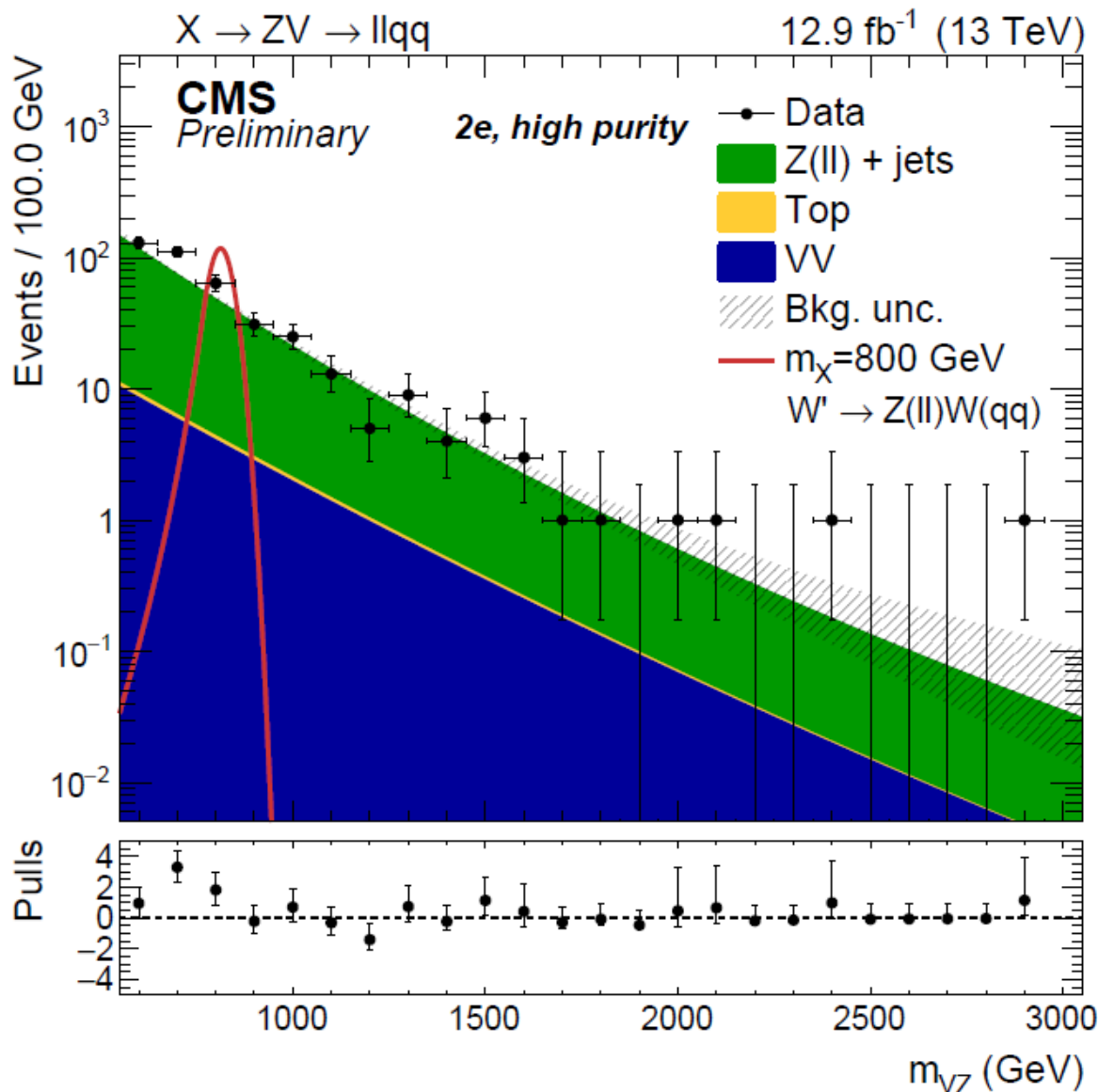
$$L(\mu) = \prod_{j=1}^N \frac{(\mu s_j + b_j)^{n_j}}{n_j!} e^{-(\mu s_j + b_j)}$$

H0:  $\mu=0$ , H1:  $\mu>0$

Neyman-Pearson:  
L(H1)/L(H0)

→ max. power test

→ Use for discovery



# Searches with Likelihood ratio

$$L(\mu) = \prod_{j=1}^N \frac{(\mu s_j + b_j)^{n_j}}{n_j!} e^{-(\mu s_j + b_j)}$$

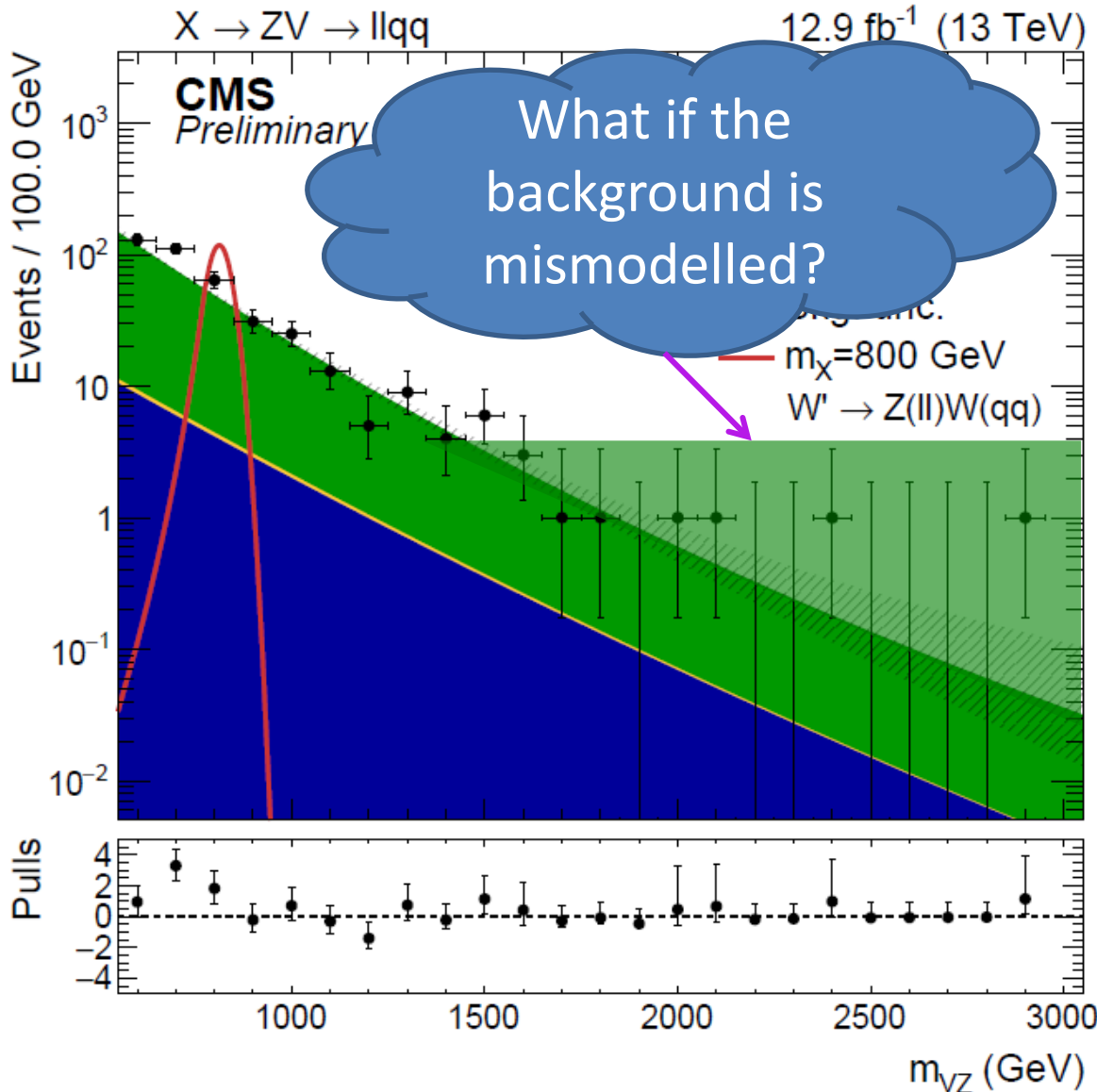
H0:  $\mu=0$ , H1:  $\mu>0$

Neyman-Pearson:  
L(H1)/L(H0)

→ max. power test

→ Use for discovery

⇒ **Don't trust a damn thing !!**



# Searches with Likelihood ratio

$$L(\mu) = \prod_{j=1}^N \frac{(\mu s_j + b_j)^{n_j}}{n_j!} e^{-(\mu s_j + b_j)}$$

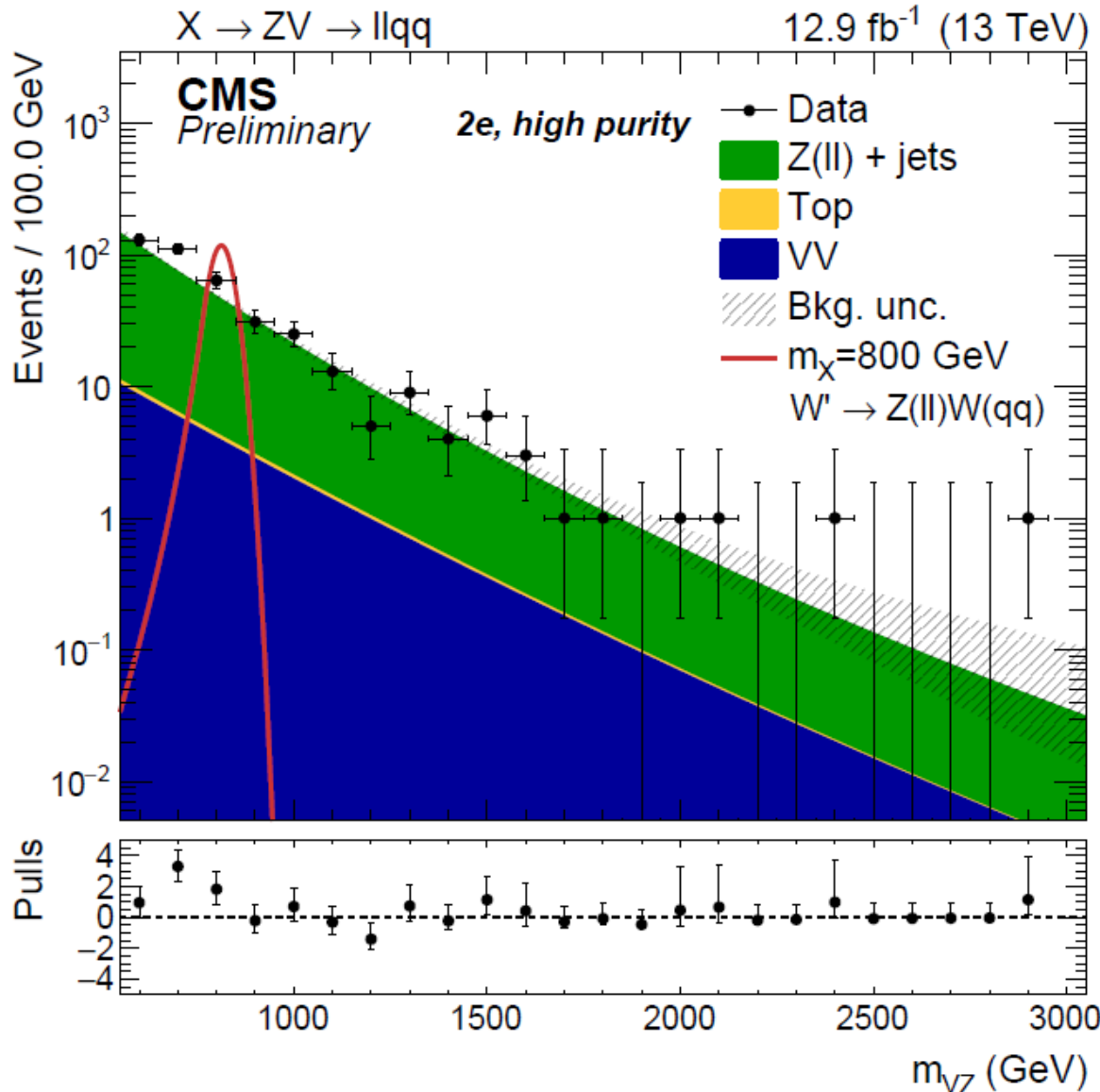
H0:  $\mu=0$ , H1:  $\mu>0$

Neyman-Pearson:  
L(H1)/L(H0)

→ max. power test

→ Use for discovery

→ Proper modelling of background over whole range is essential → do GOF-tests for H0!



# Searches with Likelihood ratio

$$L(\mu) = \prod_{j=1}^N \frac{(\mu s_j + b_j)^{n_j}}{n_j!} e^{-(\mu s_j + b_j)}$$

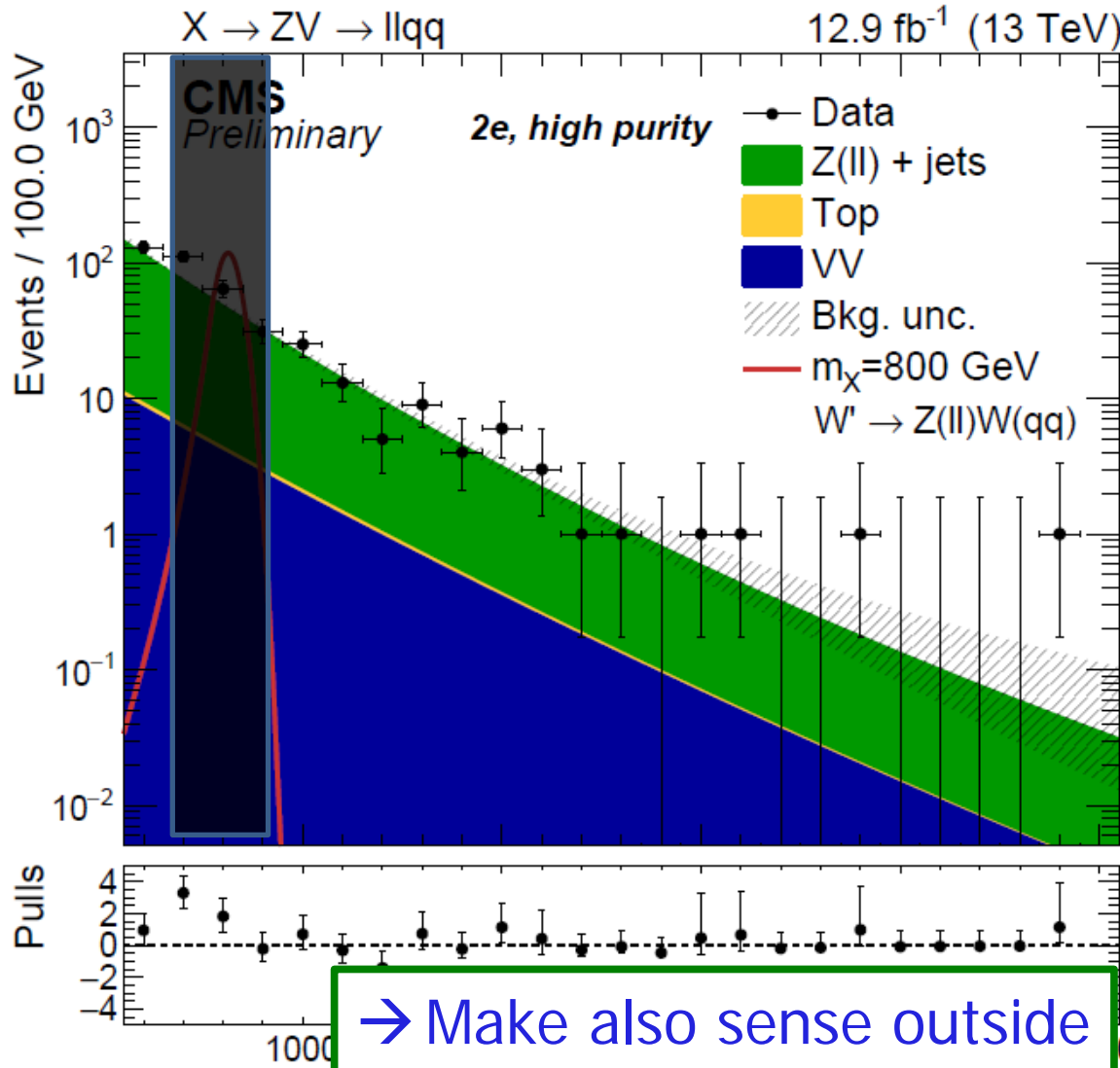
H0:  $\mu=0$ , H1:  $\mu>0$

Neyman-Pearson:  
L(H1)/L(H0)

→ max. power test

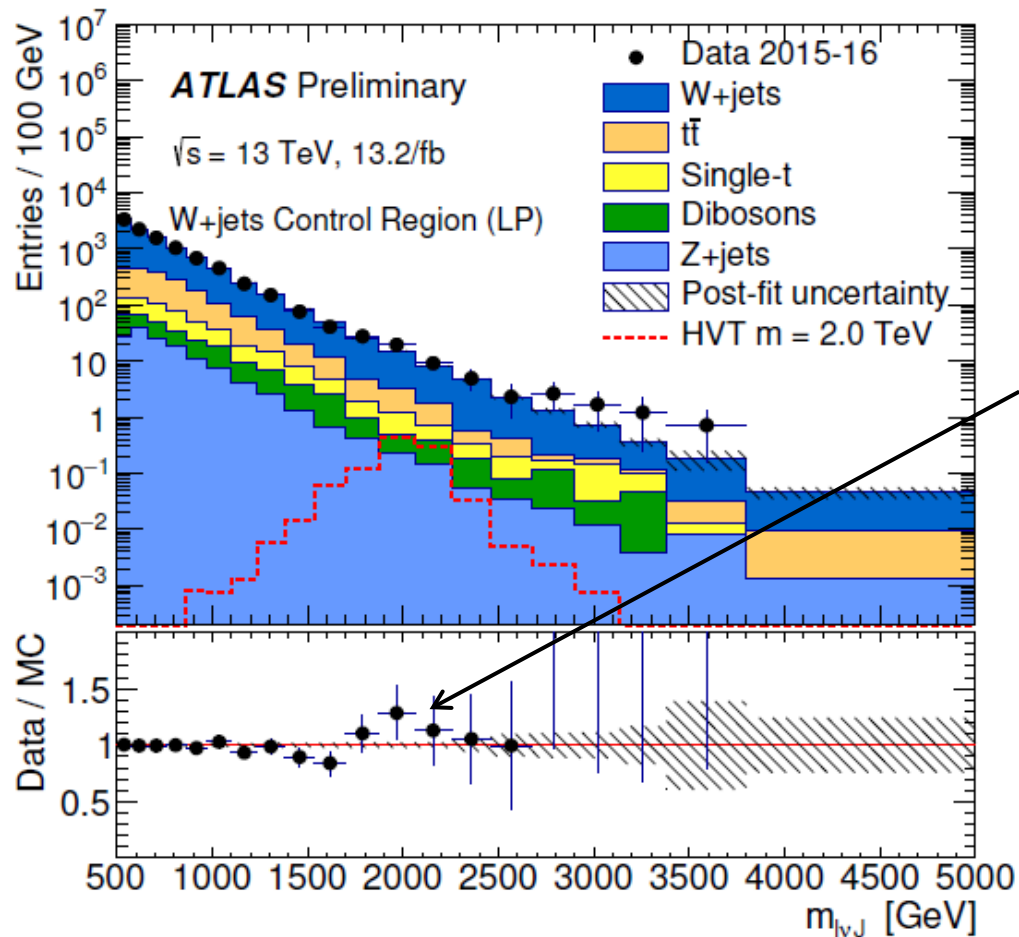
→ Use for discovery

→ Proper modelling of background over whole range is essential → do GOF-tests for H0!



# Goodness-Of-Fit Tests in ATLAS and CMS searches

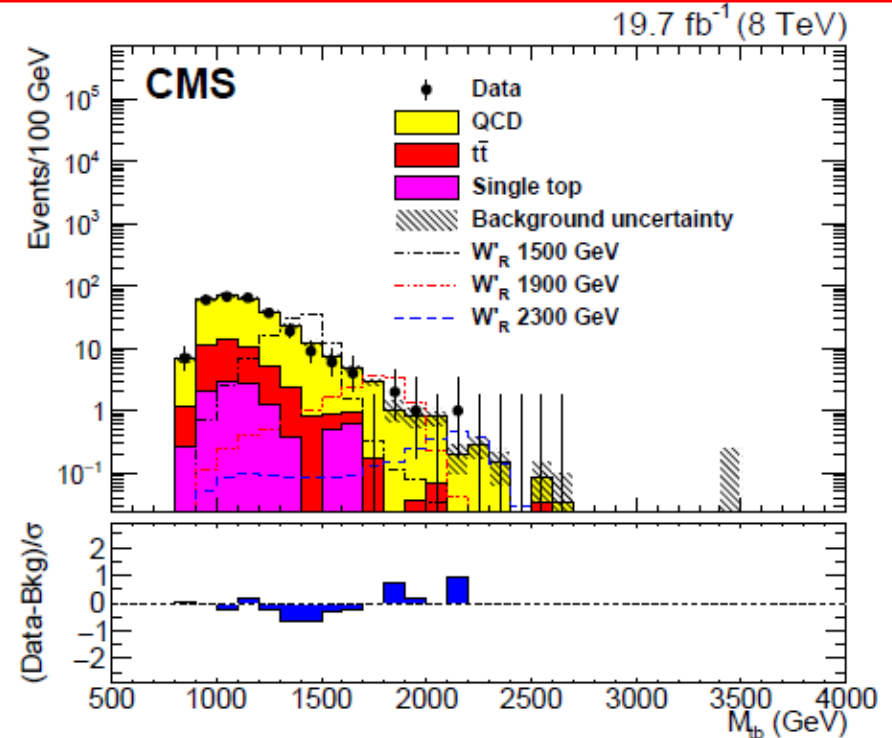
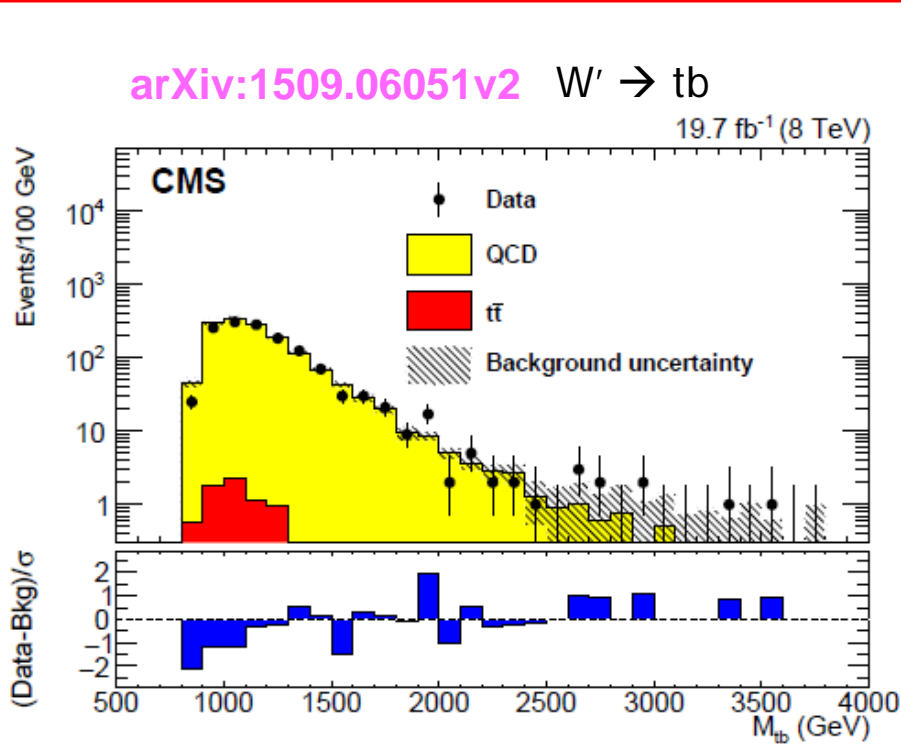
ATLAS-CONF-2016-062 Diboson resonance



“Good agreement  
(optical inspection of ratio)  
between the data and the  
background prediction”



# Goodness-Of-Fit Tests in ATLAS and CMS searches

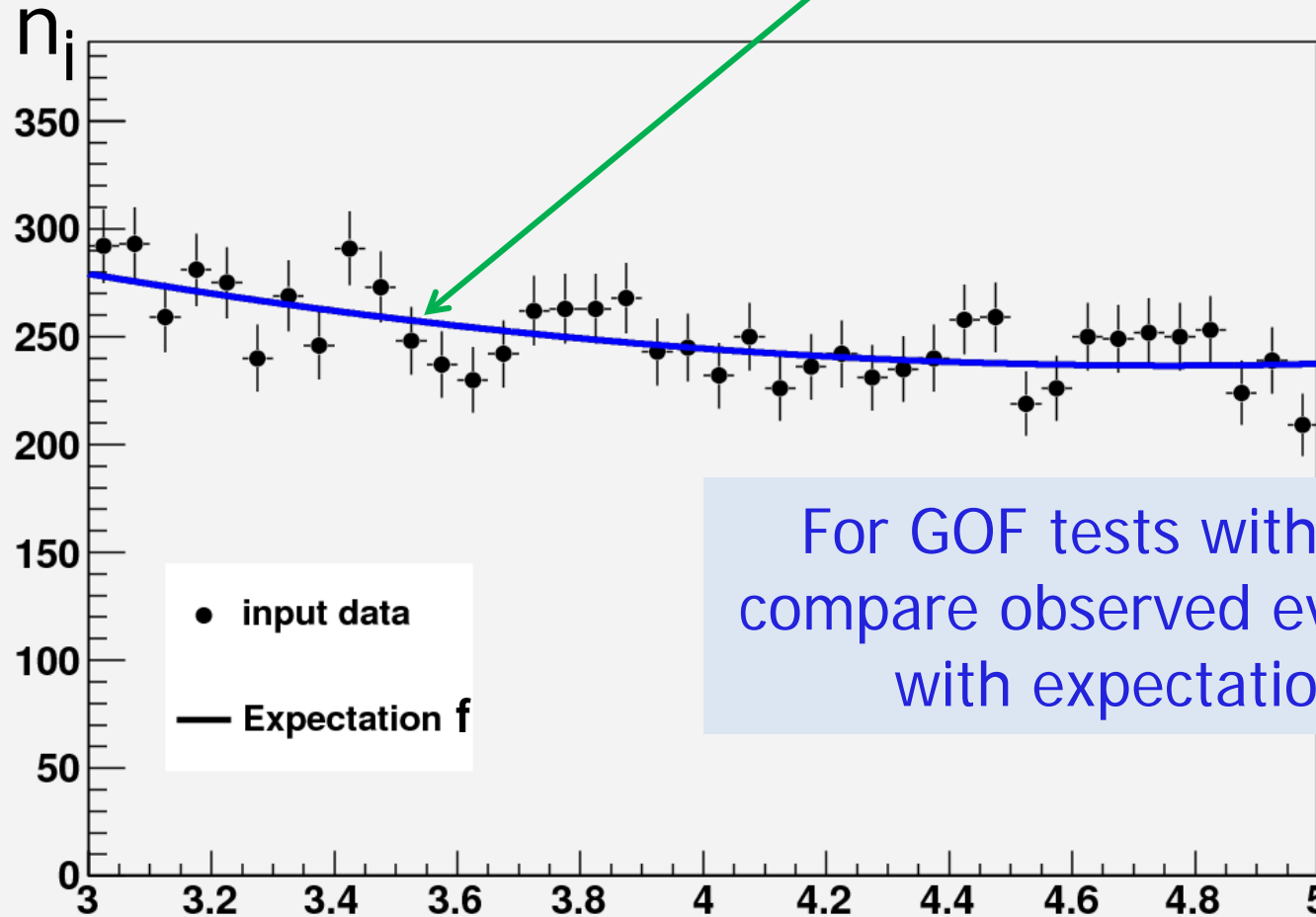


“This test (optical pull inspection) shows good agreement between data and SM ”

→ Optical inspection of bin-wise pulls is crucial but should do also global tests as discussed in the following slides!

# Goodness-Of-Fit Tests – basics

Basic question: how well does  $H_0$  describe the data?



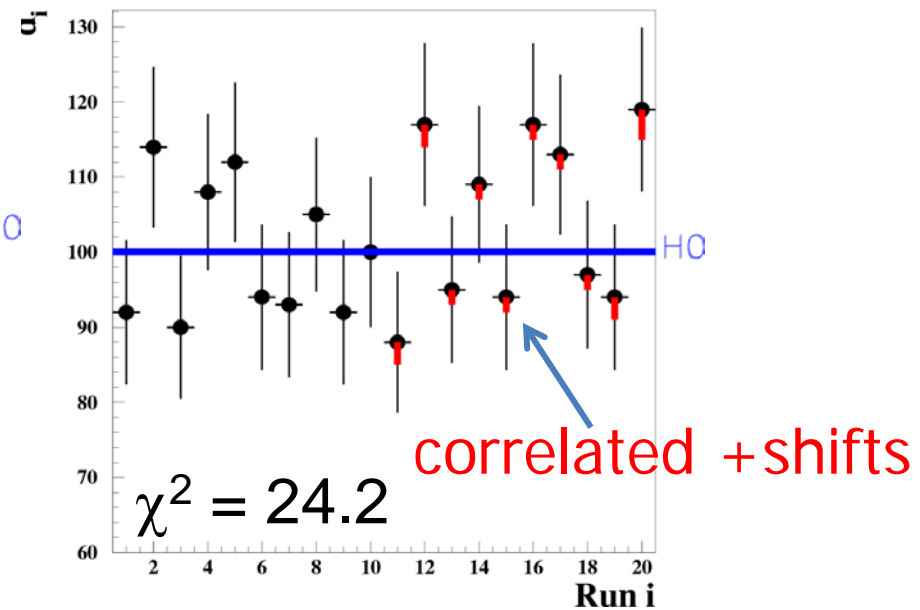
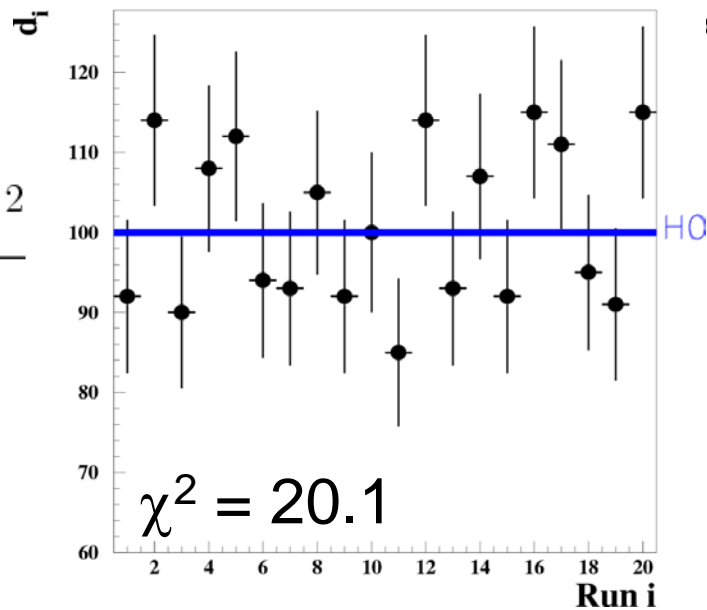
For GOF tests with binned data:  
compare observed event numbers  $n_i$   
with expectation values  $f_i$

$$f_i = \int_{bin\ i} f dm$$

→ Since no  $H_1$  specified → many different GOF tests possible

# Goodness-Of-Fit Test – $\chi^2$ tests

$$\chi^2 = \sum_i \frac{(f_i - n_i)^2}{\sigma_i^2}$$



- $\chi^2$  throws away all sign and order info → not very sensitive to correlated shifts in a certain region.
- apply further GOF tests to check all data/model facets!

Note: p-values for  $\chi^2$ : `TMath::Prob( $\chi^2_{\text{obs}}$ , ndf)`

# GOF-tests: exemplary analysis

Likelihood ratio  
improved  $\chi^2$  test

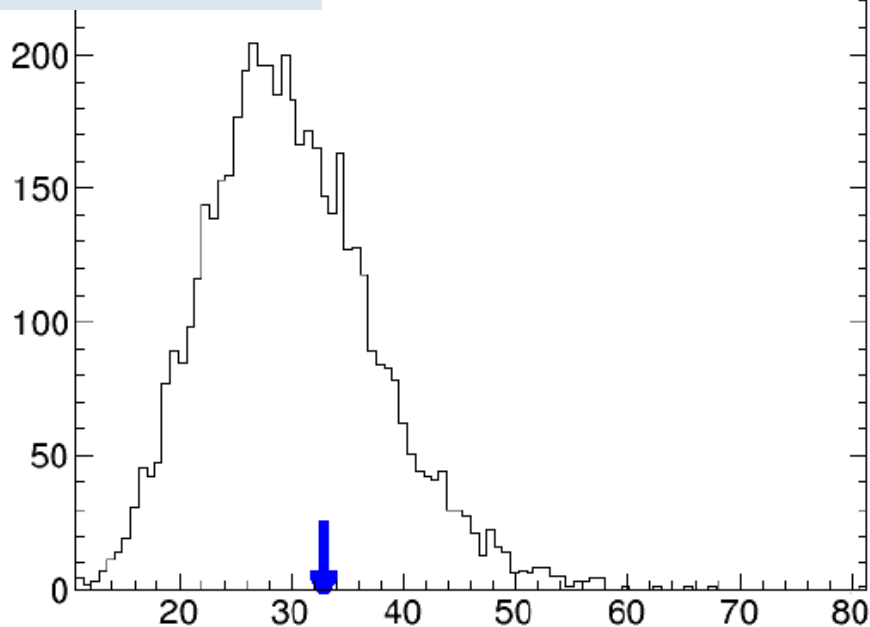
S. Baker & R.D. Cousins,  
NIM 221 (1984) 437

$$\tilde{\chi}^2 = -2 \ln \left( \frac{L(\mathbf{n}|\mathbf{f})}{L(\mathbf{n}|\mathbf{n})} \right) = 2 \sum_i f_i - n_i + n_i \ln(n_i/f_i)$$

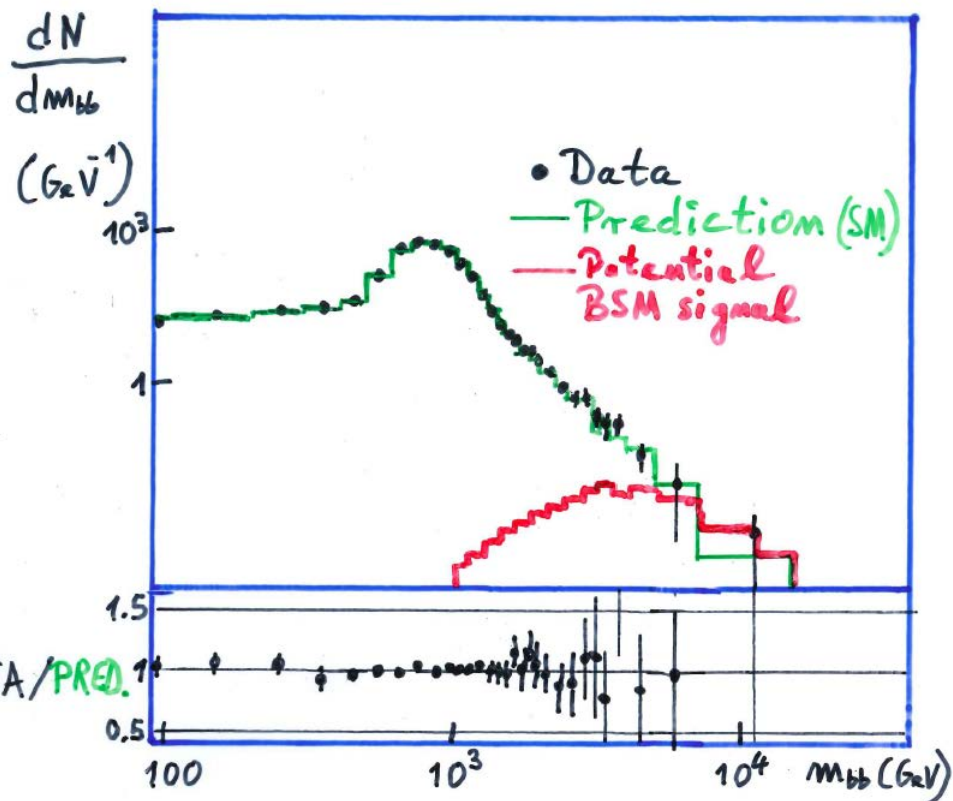
Saturated model

Analysed  
with CMS  
combine tool

saturated, 5000 Toys  
p-value = 0.321



# Hypothetical pp data@100 TeV



→ Test result ok

# GOF-tests: exemplary analysis

## Kolmogorov-Smirnov test

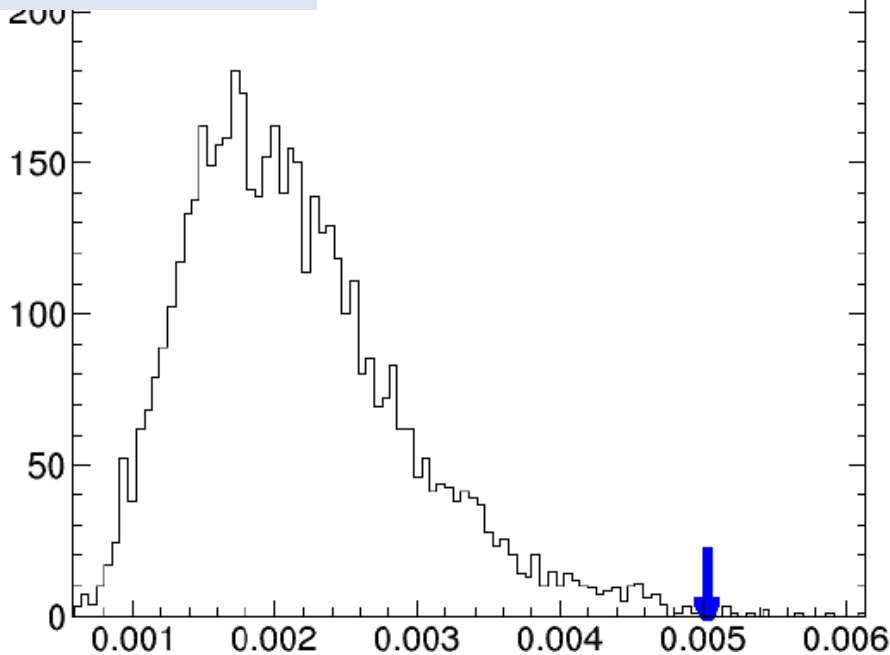
$F_c$ : Cumulative distribution function

$F_e$ : Empirical distribution function

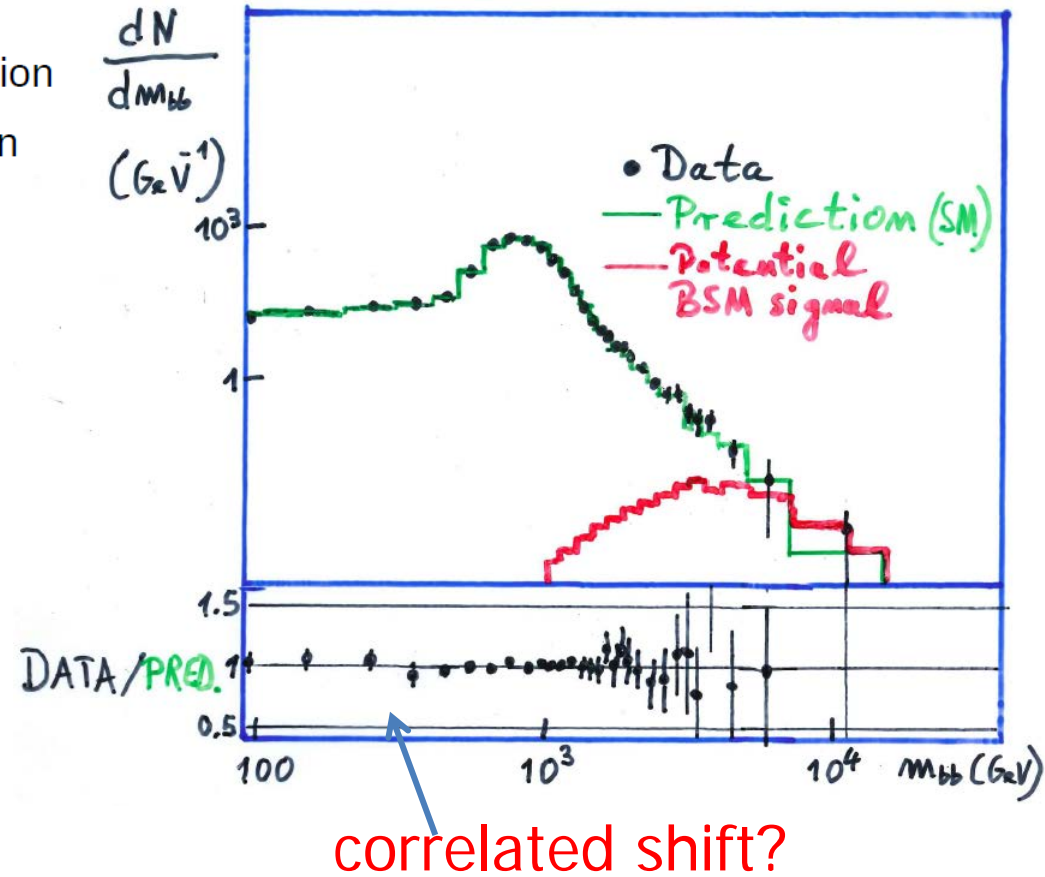
$$Q_{GoF,KS} = \sup |F_c(x) - F_e(x)|$$

Analysed  
with CMS  
combine tool

KS, 5000 Toys  
p-value = 0.002



## Hypothetical pp data@100 TeV



→ Test result **bad!**

# GOF-tests: exemplary analysis

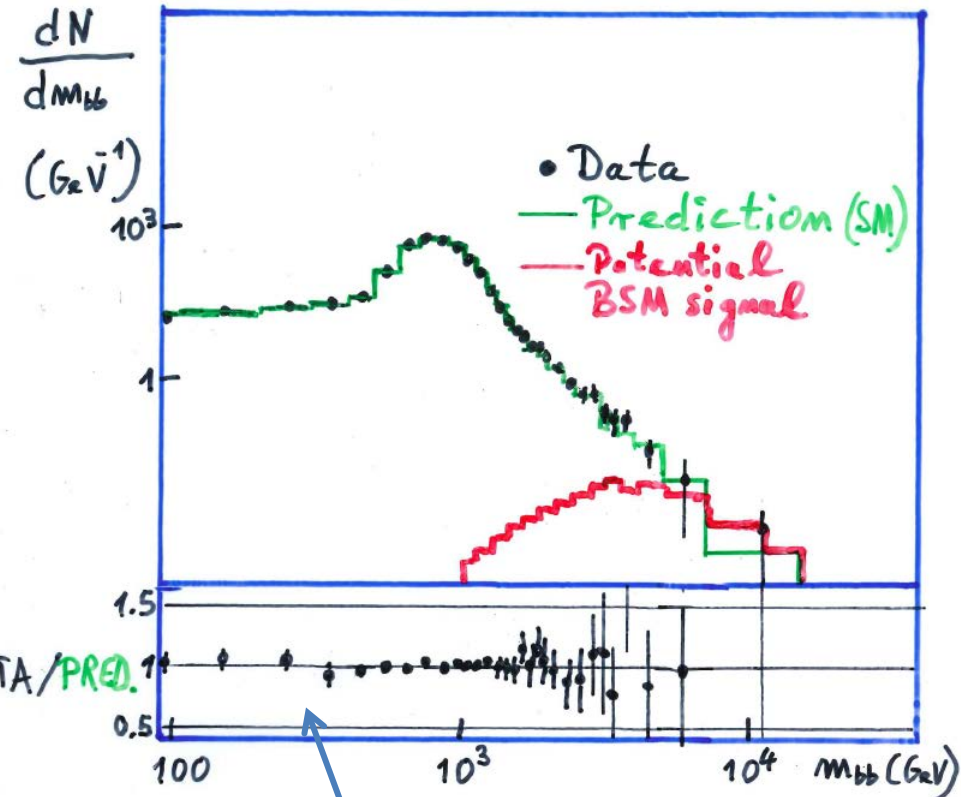
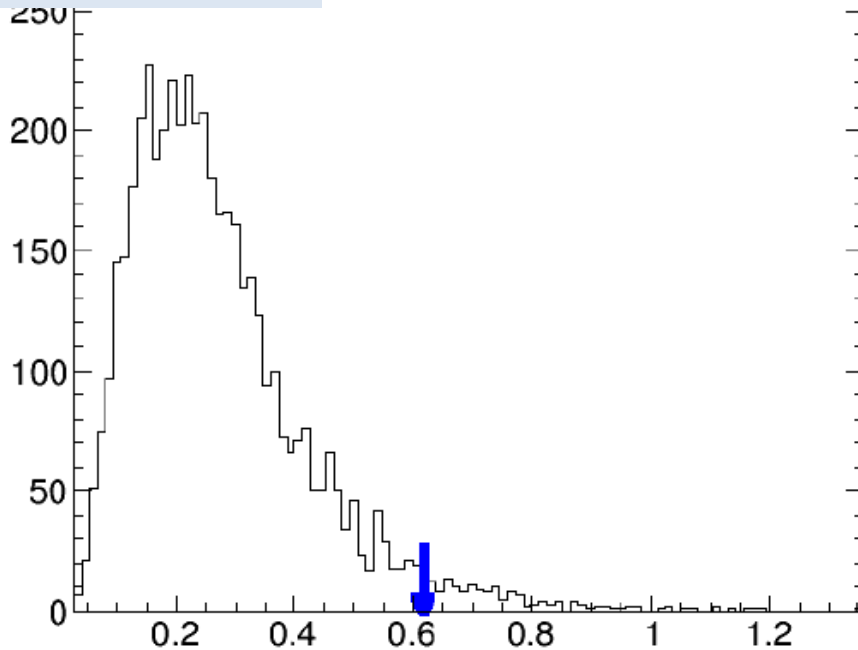
## Hypothetical pp data@100 TeV

### Anderson-Darling test

$$q_{GoF,AD} = n \cdot \int dF_e(x) \frac{(F_c(x) - F_e(x))^2}{F_e(x) \cdot (1 - F_e(x))}$$

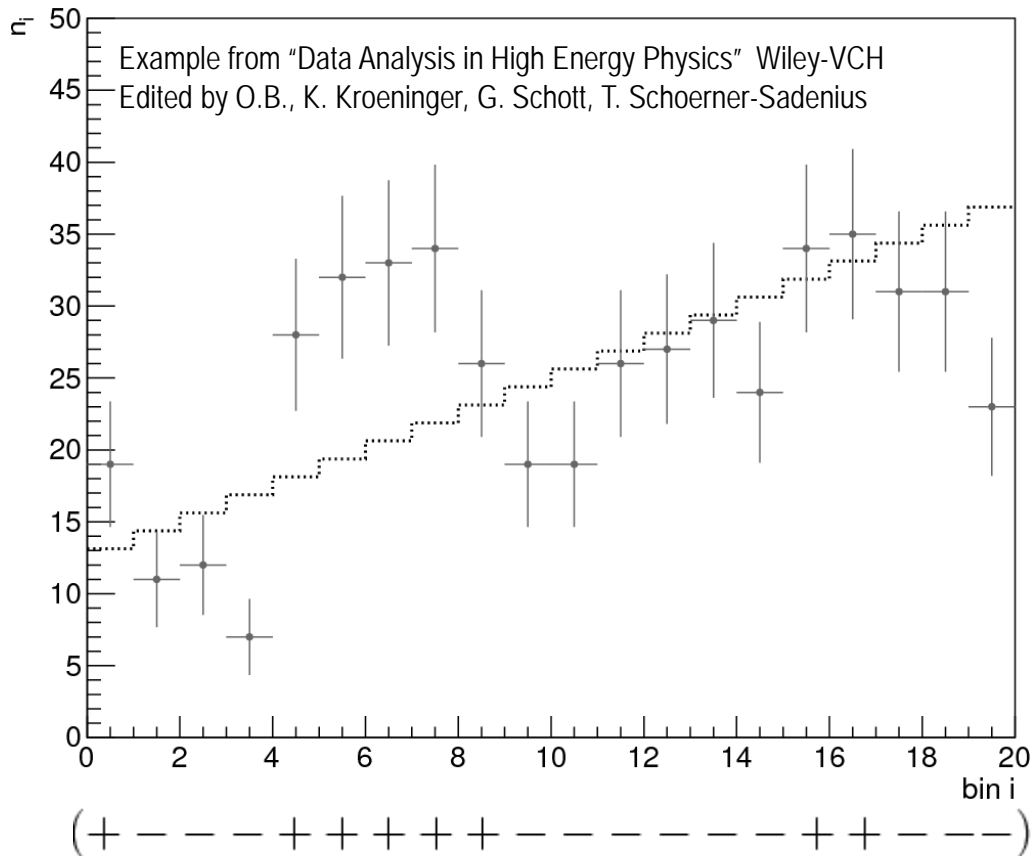
Analysed  
with CMS  
combine tool

AD, 4800 Toys  
p-value = 0.033



→ Test result **bad!**

# Goodness of Fit - Run test



Runs

→ Easy to do test!

Idea: count runs = regions with same sign of deviation

$r$  = #runs

$N_+$  = #bins data > model

$N_-$  = #bins data < model

$r$  should follow Binomial statist.

$$E[r] = 1 + \frac{2N_+N_-}{N_+ + N_-}$$

$$V[r] = \frac{2N_+N_-(2N_+N_- - N_+ - N_-)}{(N_+ + N_-)^2(N_+ + N_- - 1)}$$

$$Z = \frac{r - E[r]}{\sqrt{V[r]}}$$

Approximate Significance

Here:  $r=6$ ,  $E(r)=10.6 \pm 2.1$   
 → p-value = 0.0285

Perform GOF tests through various analysis stages:

✓ **Control plots (!)**

✓ **Signal Extraction (!!)**

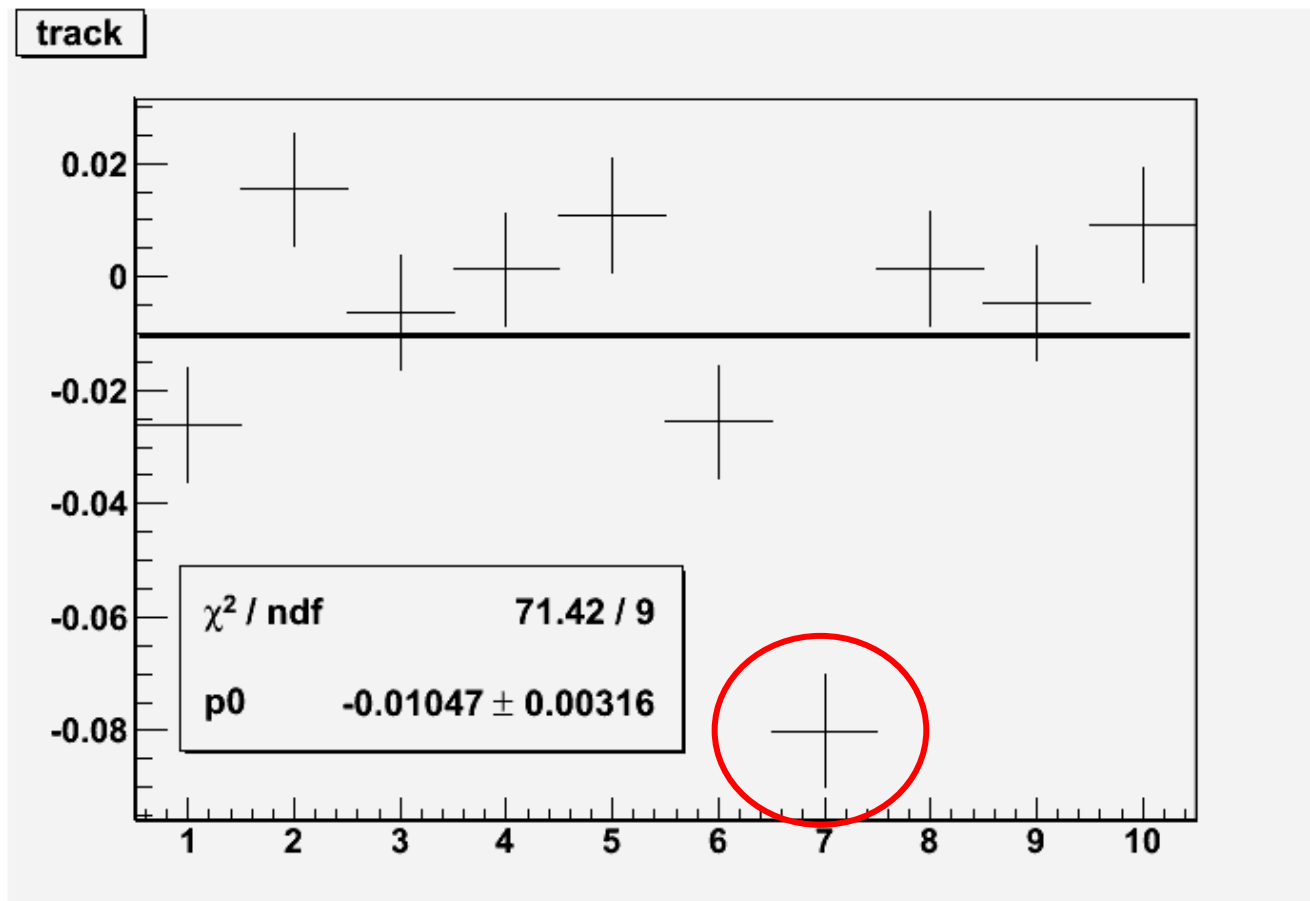
✓ **Comparisons to theory (!!)**

→ essential for understanding/control of analysis results and theory!

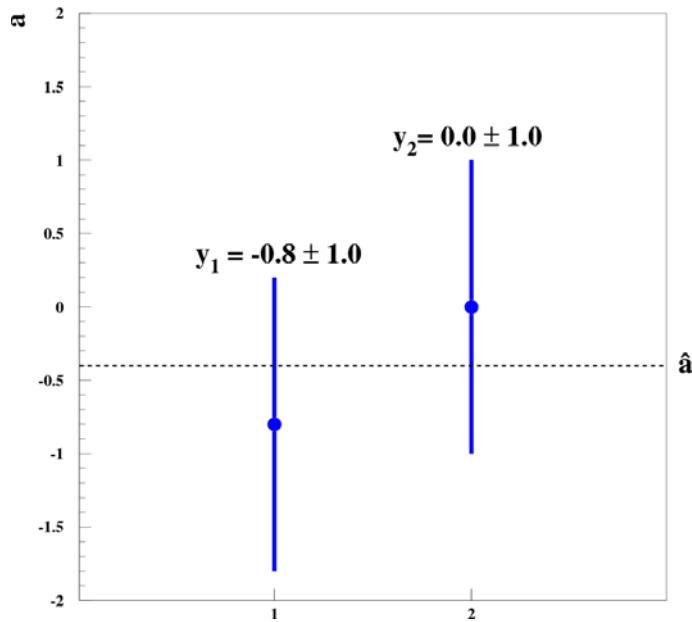
→ Apply  $\geq$  **two** different tests, e.g.  $\tilde{\chi}^2$  and K.S.



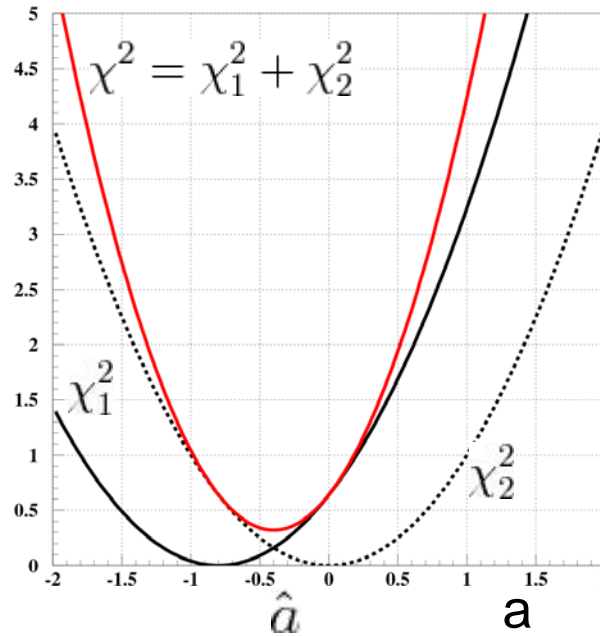
## 2 Use $\chi^2$ tests for outlier rejection



# Role of $\chi^2$ : Combination of two measurements



$$\chi_i^2 = \frac{(a - y_i)^2}{\sigma^2} \quad L \sim e^{-\chi_1^2/2} e^{-\chi_2^2/2} = e^{-\chi^2/2}$$



$$\chi^2(a) = \chi^2(\hat{a}) + H(a - \hat{a})^2$$

$$\hat{a} = (y_1 + y_2)/2; \quad H = 2/\sigma^2 \quad \chi^2(\hat{a}) = \frac{(y_1 - y_2)^2}{2\sigma^2} \quad \leftarrow \text{Squared pull}$$

$$L \sim e^{-\frac{(y_1 - y_2)^2}{2\sigma^2}} \cdot e^{-H(a - \hat{a})^2/2}$$

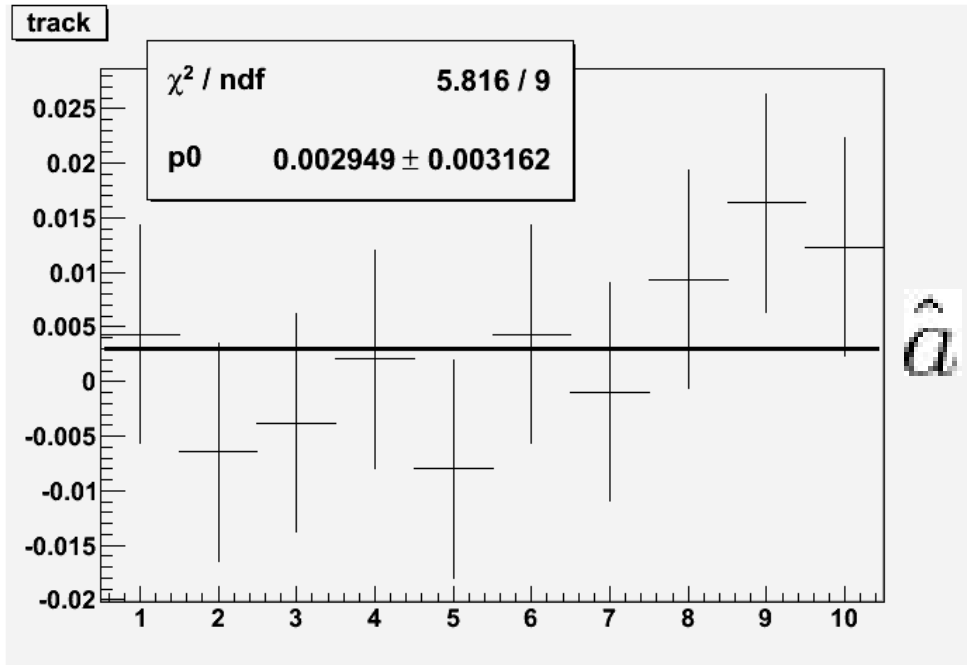
GOF-test

Info on a

→ GOF test and parameter Info decoupled!

# Combination of n measurements

Example: track fit of horizontally flying particle in n detector layers



← Weighted average position

Repeated experiments:  $\chi_{min}^2 = \chi^2(\hat{a})$  follows  $f(\chi^2, n-1)$  distribution

$$f(\chi^2, n) = \frac{1}{\Gamma(n/2)2^{n/2}} \cdot (\chi^2)^{n/2-1} \cdot e^{-\chi^2/2}$$

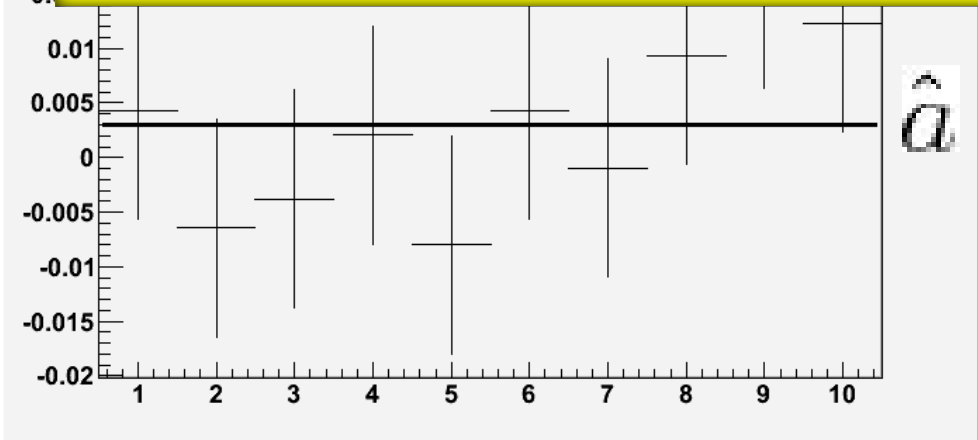
$$\text{with } \Gamma(n/2) = \int_0^\infty dt e^{-t} t^{n/2-1}$$

$$\text{TMath::Prob}(\chi^2, n) = \int_{\chi^2}^\infty f(\chi'^2, n) d\chi'^2$$

uniformly distr. in  $[0, 1]$ , why?

# Track fits to 10 hits – **Interactive work with ROOT, Run 1000 Fits**

- 1) No noise: uncertainty of  $\hat{a}$ , means of  $\chi^2$  and  $\text{prob}(\chi^2, 9)$  distr.
- 2) Repeat with noise hits
  - a) just fitting
  - b) discard fits with bad  $\chi^2$
  - c) outlier rejection + repeat track fit



← Weighted average position

Repeated experiments:  $\chi^2_{min} = \chi^2(\hat{a})$  follows  $f(\chi^2, n-1)$  distribution

$$f(\chi^2, n) = \frac{1}{\Gamma(n/2) 2^{n/2}} \cdot (\chi^2)^{n/2-1} \cdot e^{-\chi^2/2}$$

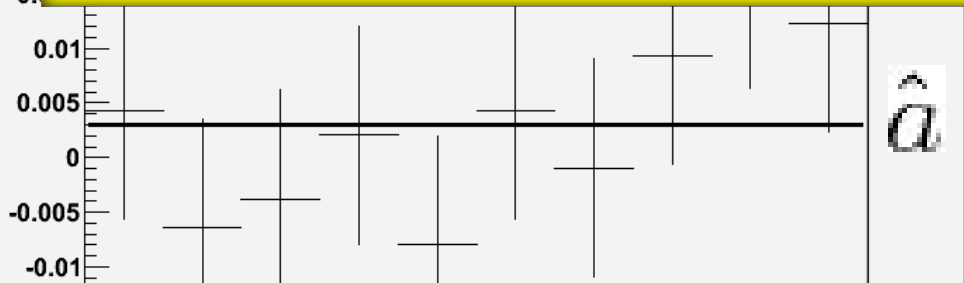
with  $\Gamma(n/2) = \int_0^\infty dt e^{-t} t^{n/2-1}$

$$\text{TMath::Prob}(\chi^2, n) = \int_{\chi^2}^{\infty} f(\chi'^2, n) d\chi'^2$$

uniformly distr. in  $[0, 1]$ , **why?**

## Track fits to 10 hits – **Interactive work with ROOT, Run 1000 Fits**

- 1) No noise: uncertainty of  $\hat{a}$ , means of  $\chi^2$  and  $\text{prob}(\chi^2, 9)$  distr.
- 2) Repeat with noise hits:
  - a) just fitting
  - b) discard fits with bad  $\chi^2$
  - c) outlier rejection + repeat track fit



← Weighted average position

ROOT 6 Macro available at

[www.desy.de/~obehnke/stat/school\\_18feb/p0toyf.C](http://www.desy.de/~obehnke/stat/school_18feb/p0toyf.C)

Instructions available at

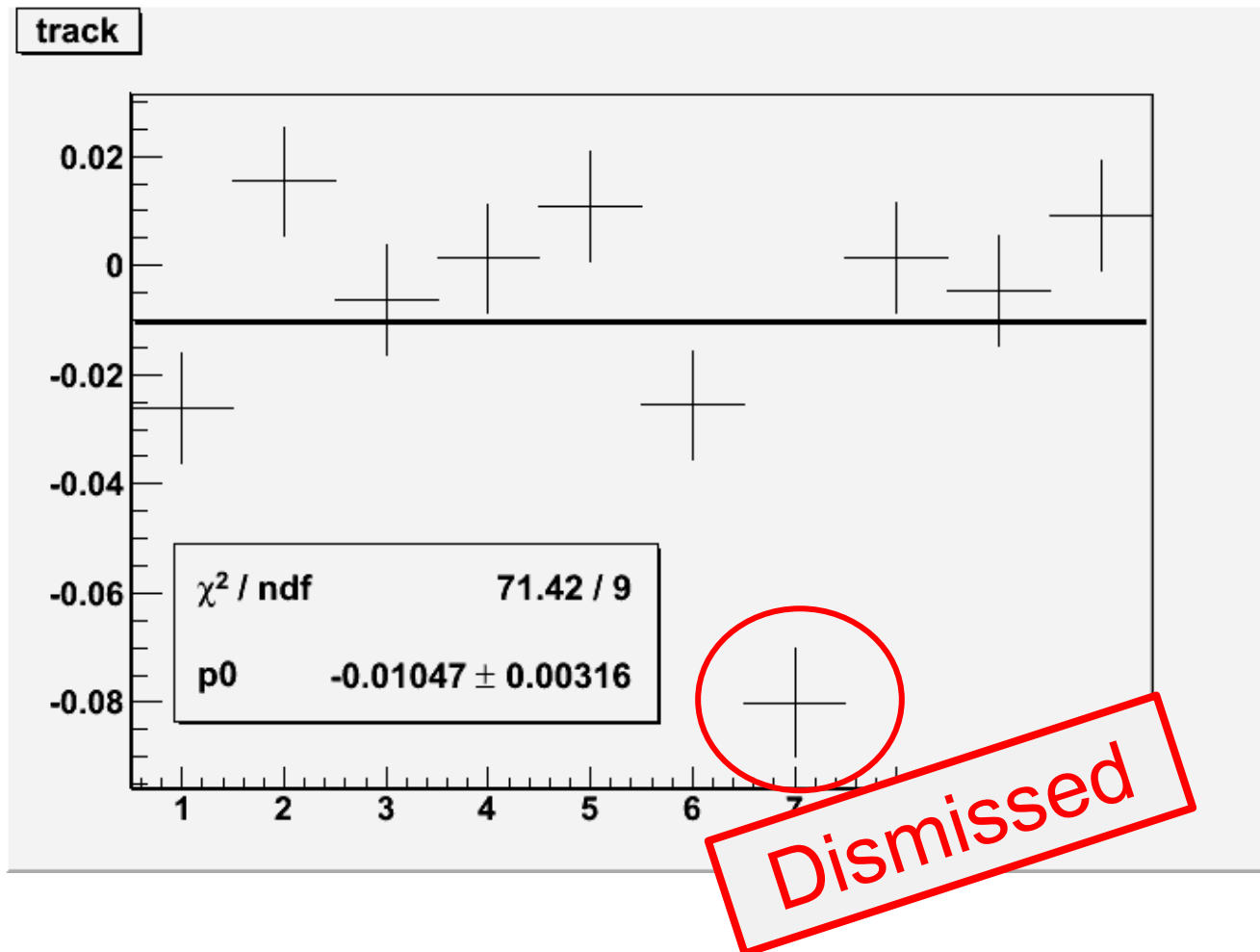
[www.desy.de/~obehnke/stat/school\\_18feb/compueb\\_p0toyf.pdf](http://www.desy.de/~obehnke/stat/school_18feb/compueb_p0toyf.pdf)

$$f(\chi^2, n) = \frac{1}{\Gamma(n/2)2^{n/2}} \cdot (\chi^2)^{n/2-1} \cdot e^{-\chi^2/2}$$

$$\text{with } \Gamma(n/2) = \int_0^\infty dt e^{-t} t^{n/2-1}$$

$$\Gamma\text{Math}::\text{Prob}(\chi^2, n) = \int_{\chi^2} f(\chi'^2, n) d\chi'^2$$

uniformly distr. in  $[0, 1]$ , **why?**



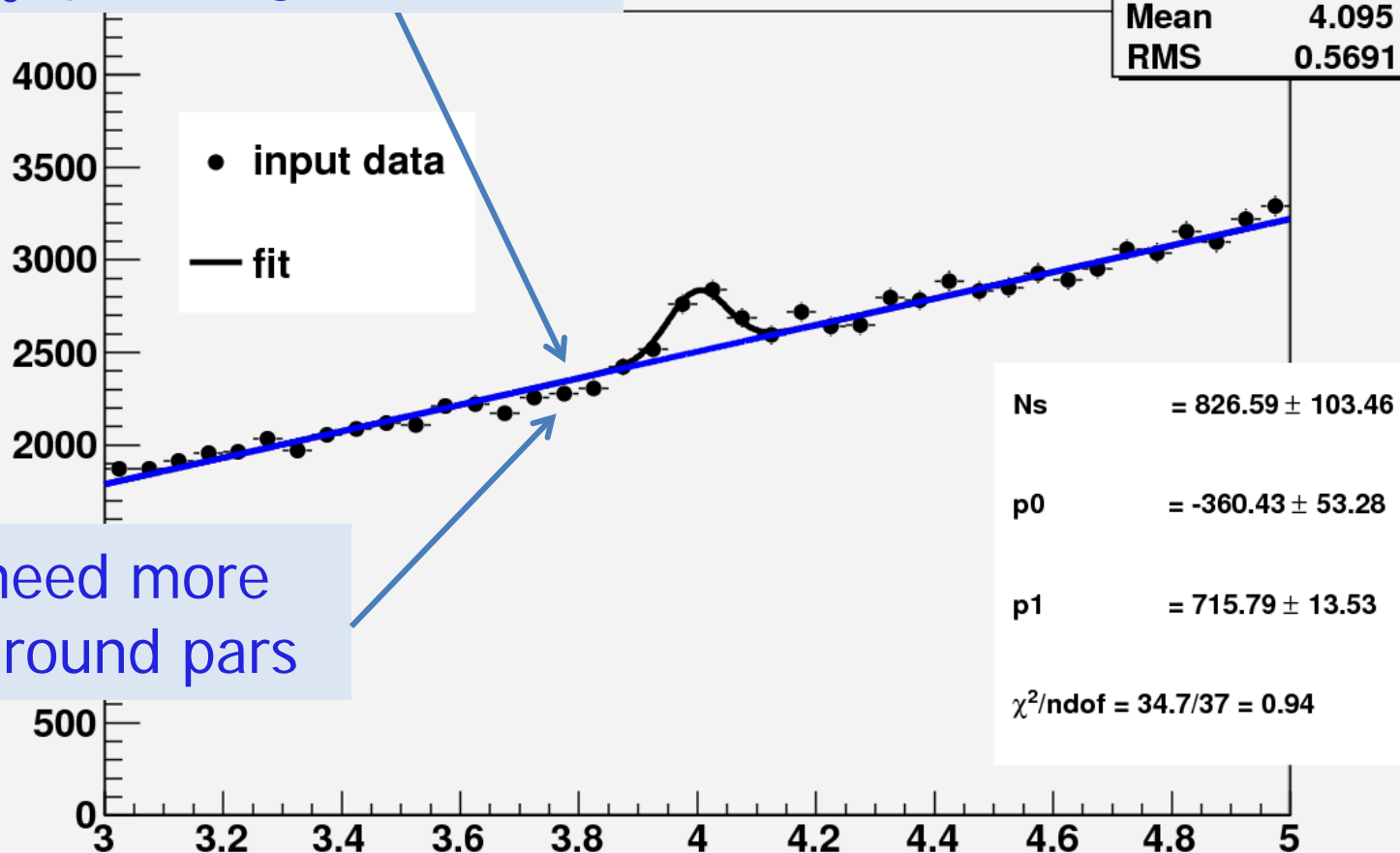
→ Can use  $\chi^2$  tests as powerful tool for Pattern recognition tasks

Note: Rejecting hits with  $\chi^2 > 5$  is hard cut, tune (e.g. try 10)

3

# GOF + Likelihood-ratio tests for optimal background parametrisation

$H_0$ : p1 is a good model

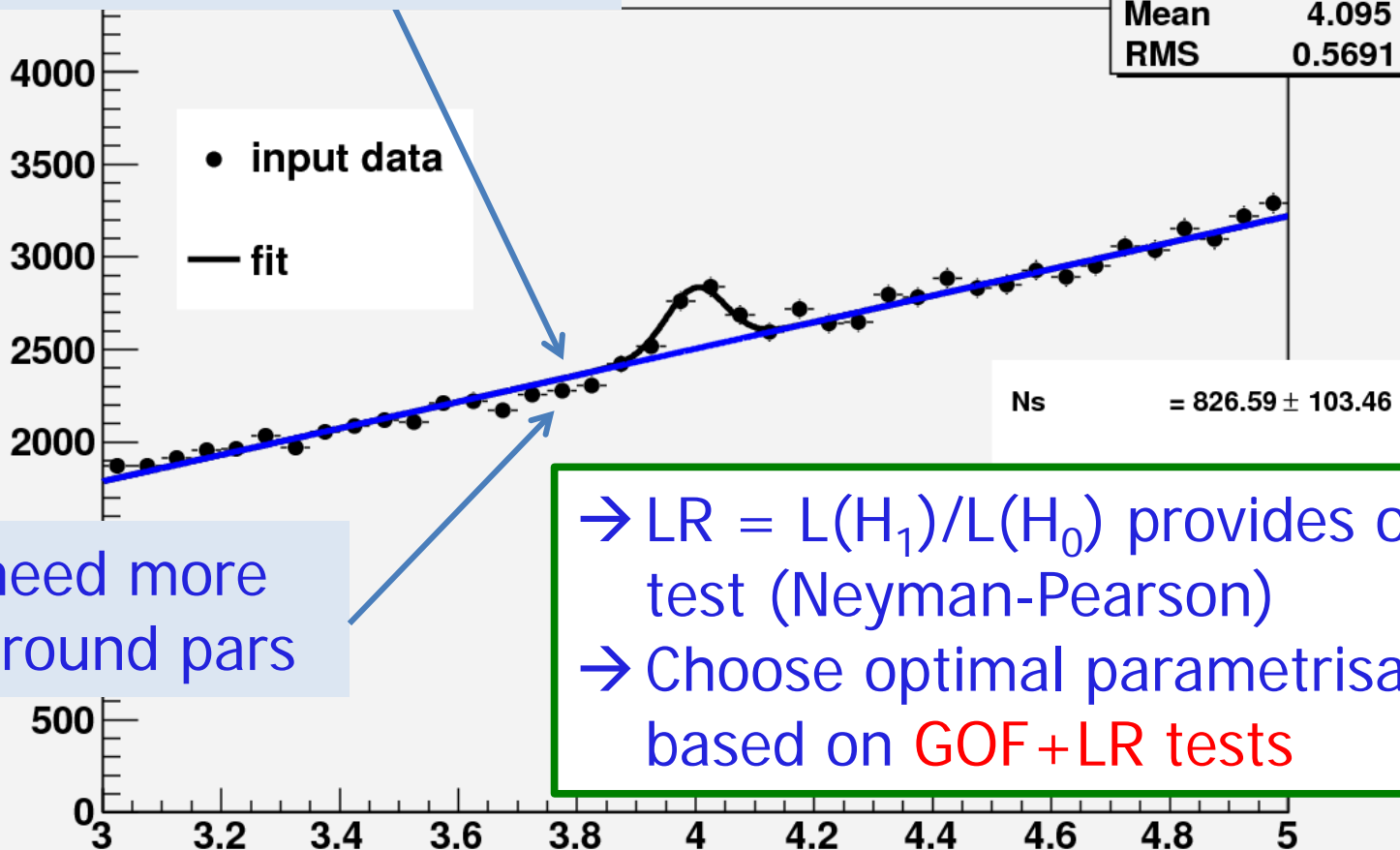


3

# GOF + Likelihood-ratio tests for

## *optimal background parametrisation*

$H_0$ : p1 is a good model

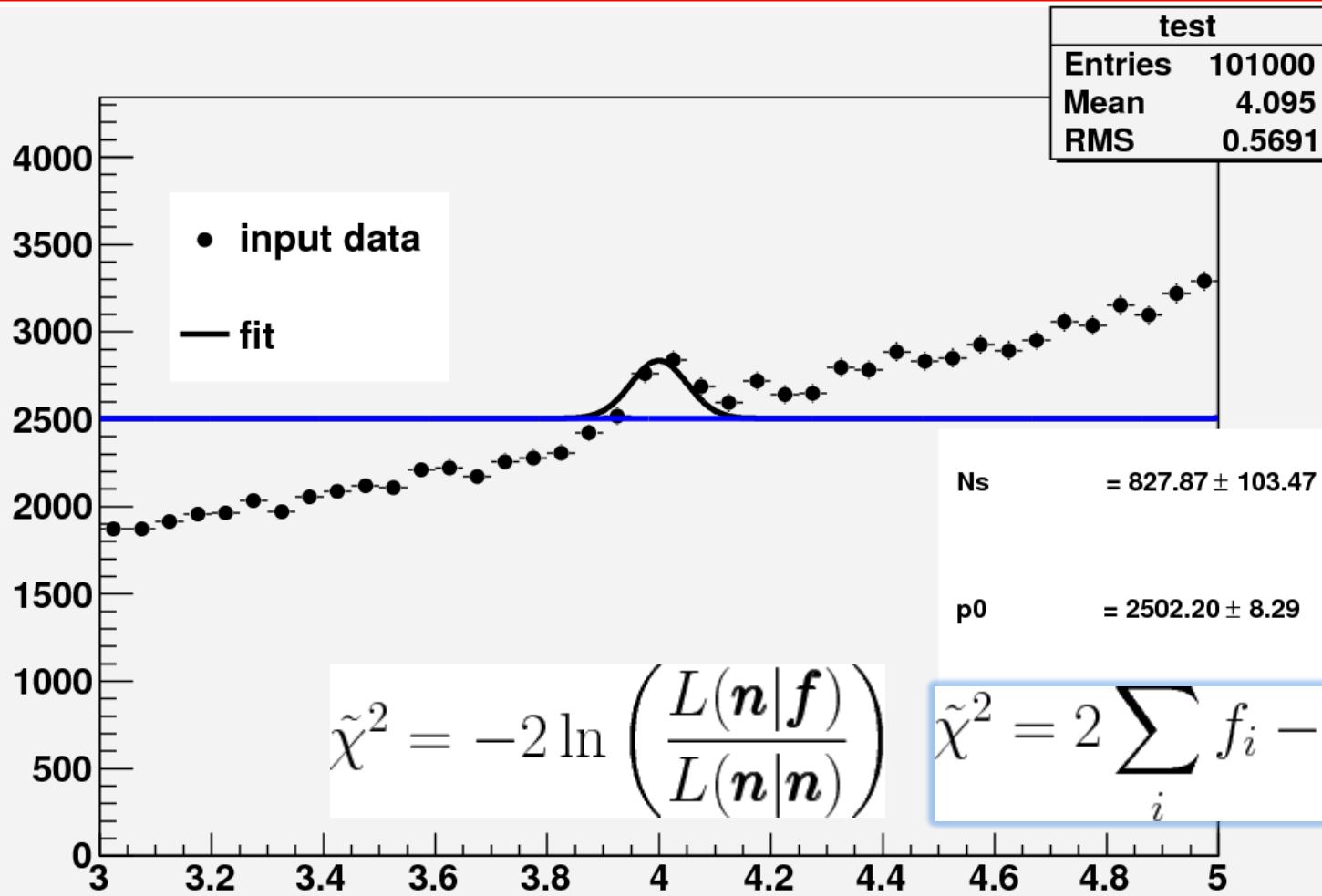


$H_1$ : need more background pars

→  $LR = L(H_1)/L(H_0)$  provides optimal test (Neyman-Pearson)  
→ Choose optimal parametrisation based on **GOF+LR tests**



# #of background fit pars – How many are needed?

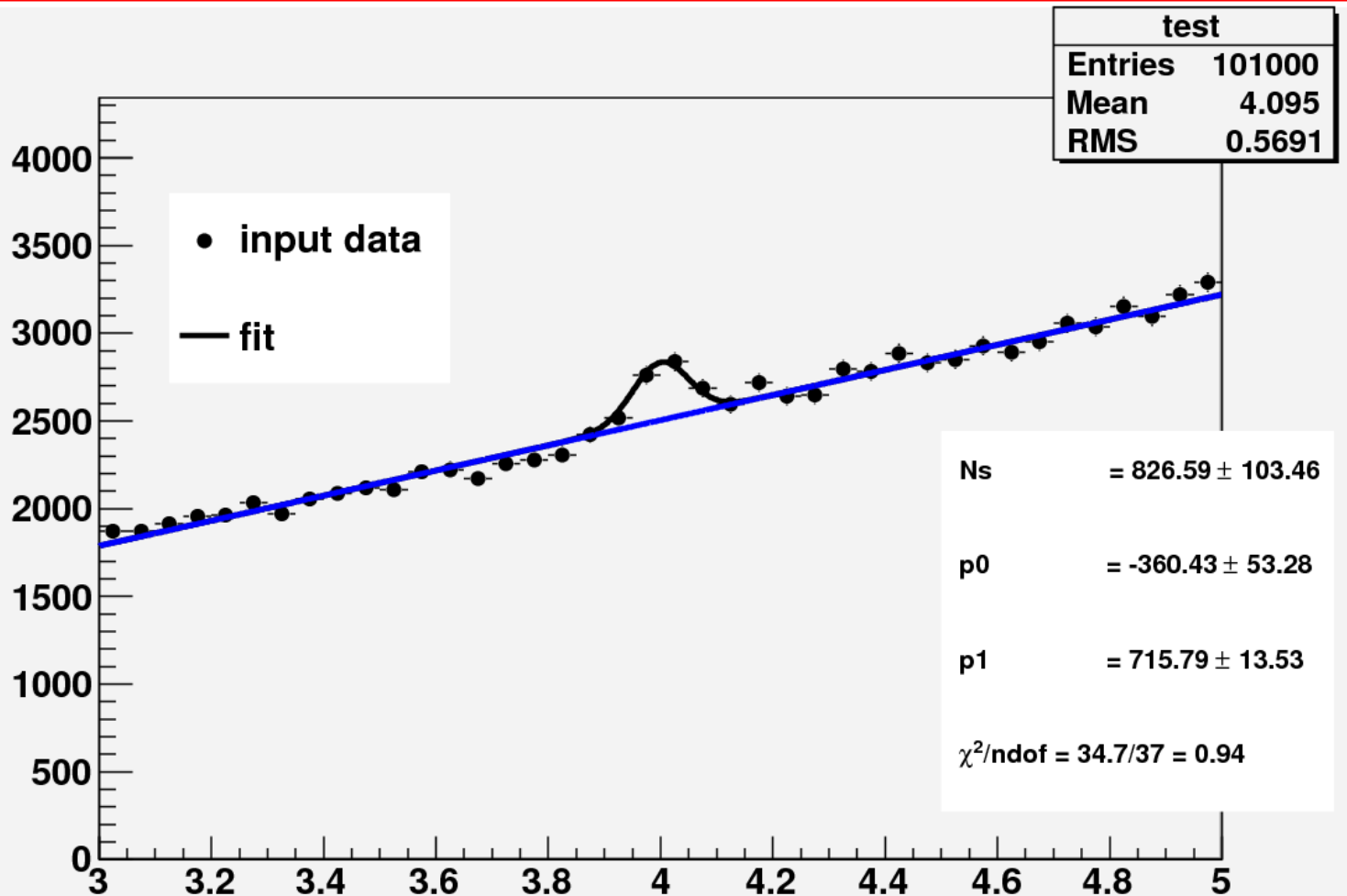


Fit function:  
gauss+p0

$$\tilde{\chi}^2 = 2880$$

→ Very poor fit: TMath::Prob(2880,38) = 0.

# #of background fit pars



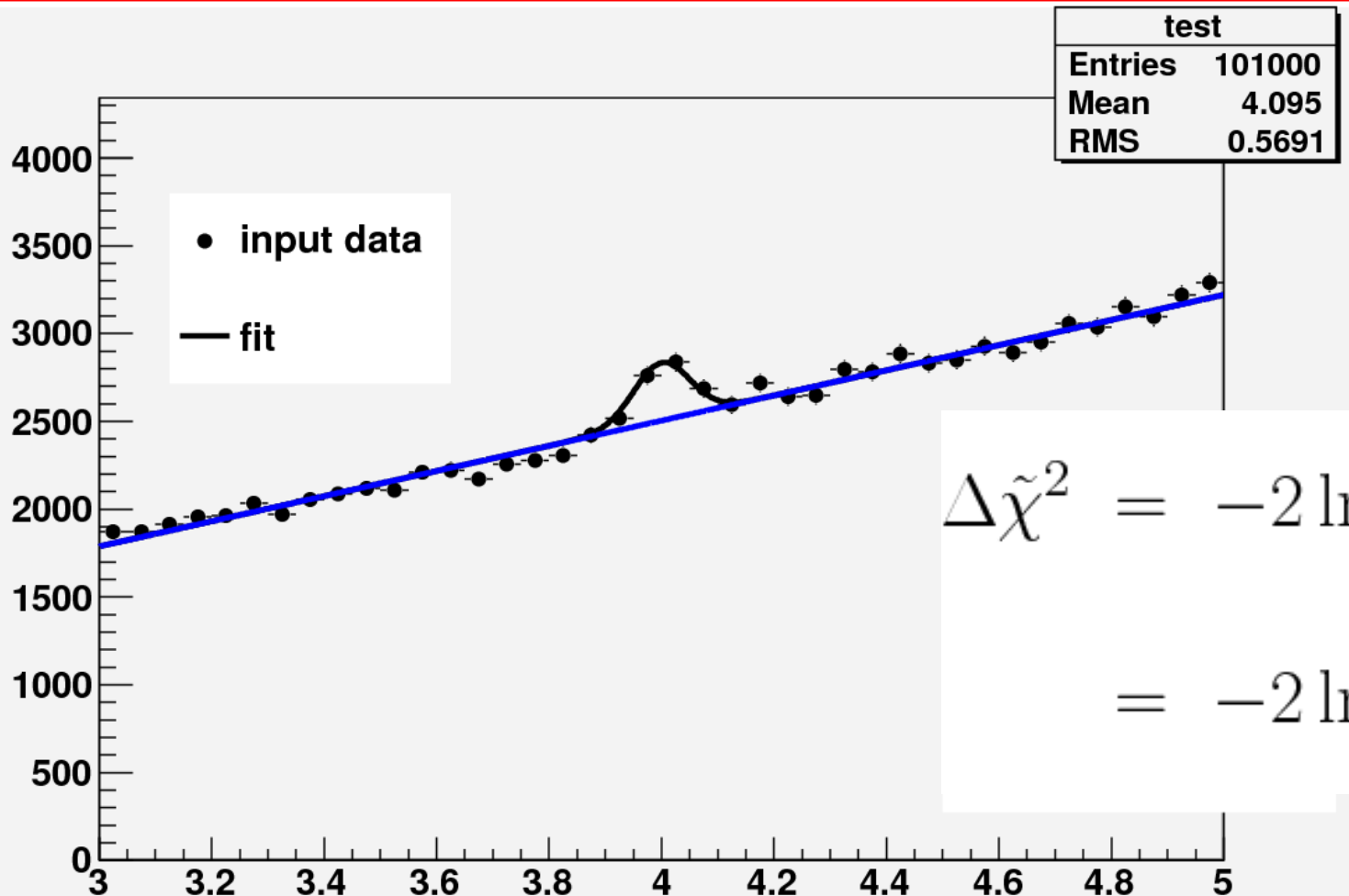
Fit function:  
gauss+p1

$$\tilde{\chi}^2 = 34.7$$

→ Reasonable  $\tilde{\chi}^2$  TMath::Prob(34.7,37) = 0.58

Should we  
stop here?

# #of background fit pars



Fit function:  
gauss+p1

$$\tilde{\chi}^2 = 34.7$$

$$\Delta\tilde{\chi}^2 = -2 \ln \left( \frac{L(H_1)}{L(H_0)} \right)$$

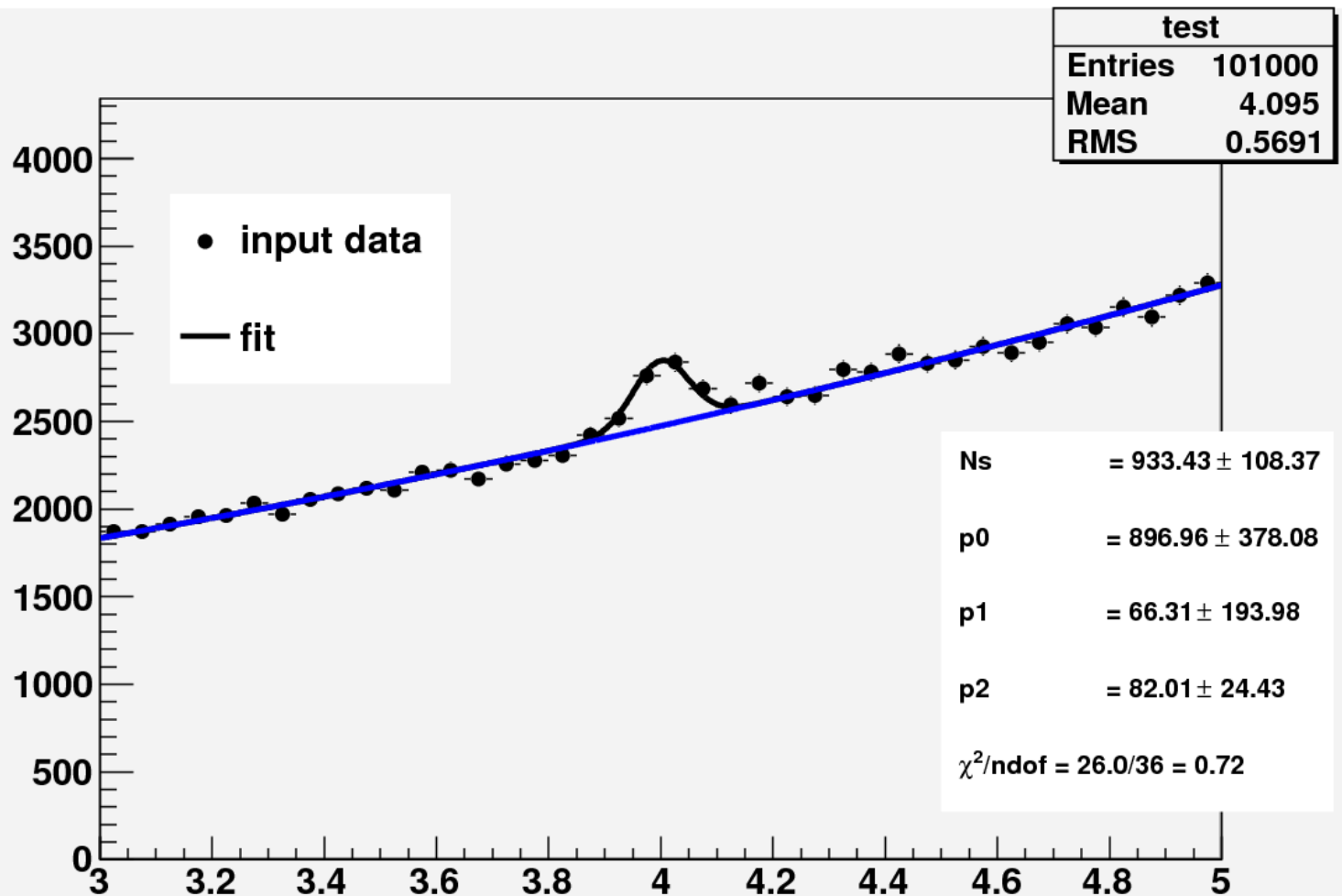
$$= -2 \ln \left( \frac{L(g + p_1)}{L(g + p_0)} \right)$$

$$= -2845.3$$

→ Reasonable  $\tilde{\chi}^2$  TMath::Prob(34.7,37) = 0.58

Should we  
stop here?

# #of background fit pars



Fit function:  
gauss+p2

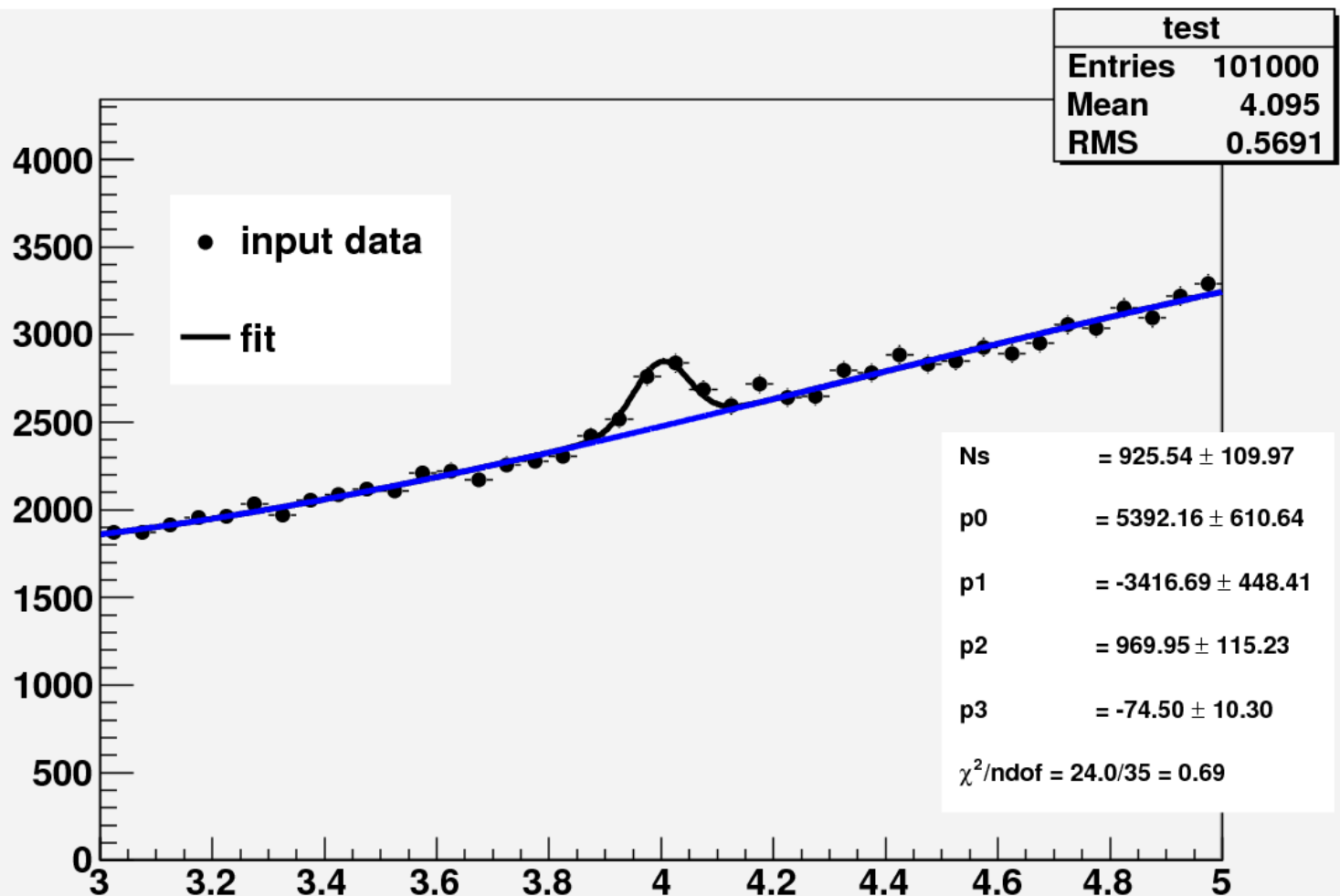
$$\tilde{\chi}^2 = 26.0$$

$$\Delta\tilde{\chi}^2 = -8.7$$

→ Reasonable  $\tilde{\chi}^2$  TMath::Prob(26.0,36) = 0.89

Should we  
stop here?

# #of background fit pars



Fit function:  
gauss+p3

$$\tilde{\chi}^2 = 24.0$$

$$\Delta\tilde{\chi}^2 = -2.0$$

→ Reasonable  $\tilde{\chi}^2$  TMath::Prob(24.0,35) = 0.92

Lets stop  
here 😊

# #of background fit pars

g+p0

$$\tilde{\chi}^2 = 2880$$

See also: [www.pd.infn.it/~dorigo/rolkelrvsftest.pdf](http://www.pd.infn.it/~dorigo/rolkelrvsftest.pdf)

g+p1

$$\tilde{\chi}^2 = 34.7 \quad \Delta\tilde{\chi}^2 = -2845.3$$

g+p2

$$\tilde{\chi}^2 = 26.0 \quad \Delta\tilde{\chi}^2 = -8.7$$

g+p3

$$\tilde{\chi}^2 = 24.0 \quad \Delta\tilde{\chi}^2 = -2.0$$

## When to stop adding further parameters?

H<sub>0</sub> Hypo: Additional parameter not needed (= zero)

If H<sub>0</sub> correct then according to Wilks' theorem:  $-\Delta\tilde{\chi}^2$  should follow  $\chi^2$  function with ndf=1 (in asymptotic regime of large n)

TMath::Prob(8.7,1) = 0.003 → g+p2 favoured over g+p1

Tmath::Prob(2.0,1) = 0.15 → g+p3 not favoured over g+p2

Gaussian z-scores:  $\sqrt{8.7} \sim 3$  and  $\sqrt{2} = 1.4$

# Wilks' theorem

Samuel S. Wilks  
(1906-1964)



- $H_0$ : Additional parameters (as predicted by  $H_1$ ) not needed (= zero)
- If  $H_0$  correct then according to Wilks' theorem:  
 $-\Delta\tilde{\chi}^2 = -2\ln[L(H_1)/L(H_0)]$  should follow for  $n \rightarrow \infty$   
 $\chi^2$  function with  $\text{ndf} = \# \text{added parameters}$   
(e.g.  $\text{ndf} = 3$  for  $p_2 \rightarrow p_5$ )

Wilks' theorem only applies for nested hypotheses:

- ✓  $H_0$ : 1<sup>st</sup> order polynomial  $\rightarrow$   $H_1$ : 2<sup>nd</sup> order polynomial
- ✗  $H_0$ : 1<sup>st</sup> order polynomial  $\rightarrow$   $H_1$ :  $a \cdot \exp(bx + cx^2)$



Stop adding parameters  $k \rightarrow k+1$  when

- $\text{TMath}::\text{Prob}(\tilde{\chi}_k^2, ndf) > 5\%$
- $\text{TMath}::\text{Prob}(\tilde{\chi}_{k+1}^2 - \tilde{\chi}_k^2, 1) > 5\%$

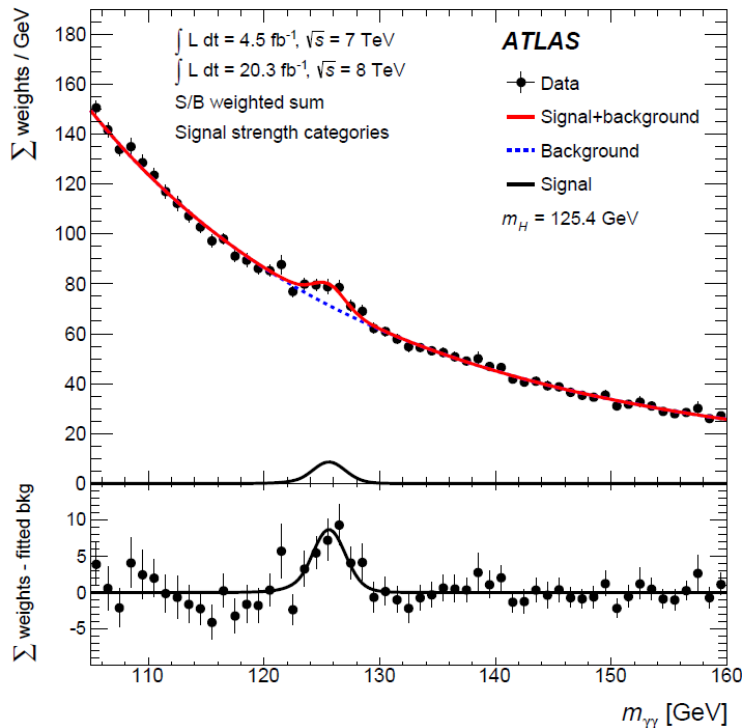
Equivalent  $\tilde{\chi}_{k+1}^2$  vs  $\tilde{\chi}_k^2$  test:

- Fisher F-test

**What about background shape systematics?**

→ Discuss next





† arXiv:1408.7084  
Phys. Rev. D90, 112015 (2014)

Conventional shape systematics:

- repeat fits with different functions (e.g. polynomials, exponential)
- changes on signal strength  $\mu \rightarrow \Delta \mu_{sys}^{bgr}$

Spurious signal idea†: absorb systematics in fit function  $f = \mu \cdot signal + \mu' \cdot signal + bgr$  with  $\mu' =$  extra fit par. for spurious signal

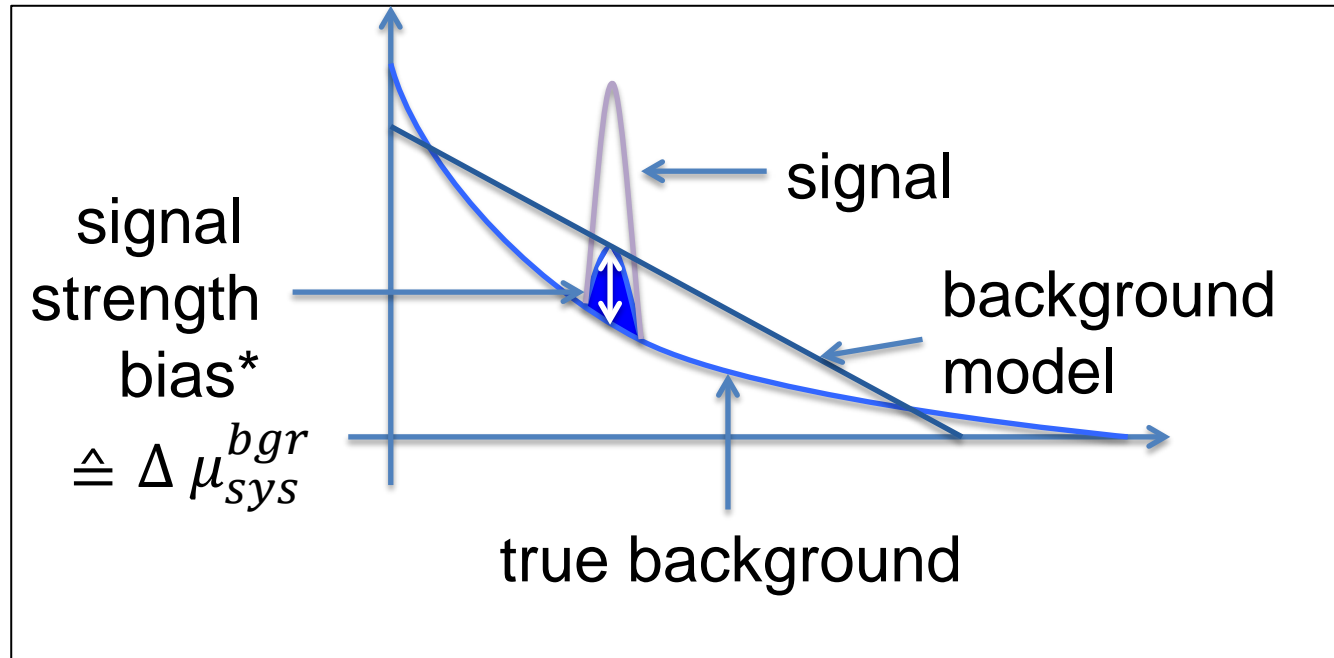
constraint on  $\mu'$

$$\tilde{\chi}^2 = 2 \left[ \sum_i f_i - n_i - n_i \ln \frac{f_i}{n_i} \right] + \left[ \frac{\mu'}{\Delta \mu_{sys}^{bgr}} \right]^2$$

- effective way of treating systematics as statistical uncertainty
- Perhaps looks a bit 'ugly'?

Determine  $\Delta \mu_{sys}^{bgr}$  from MC background toys

- generate with one function  $\rightarrow$  fit with another function + signal



\*too little signal measured in this case

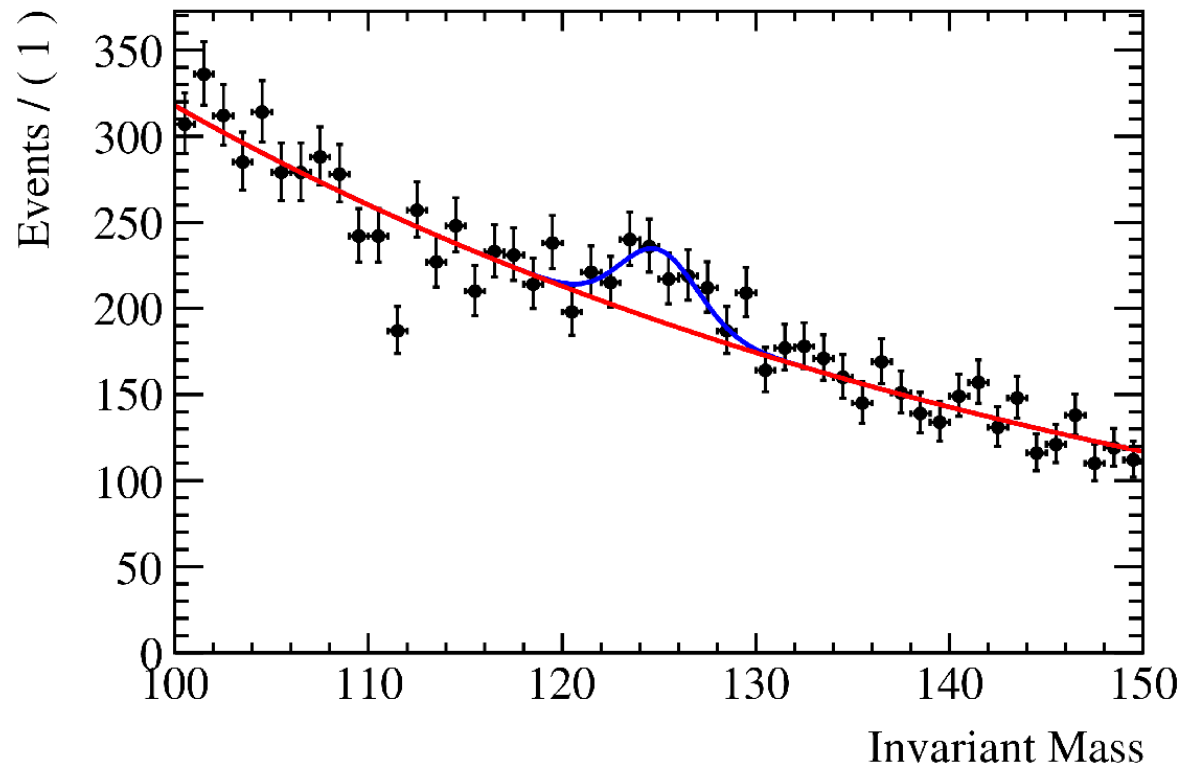
$\rightarrow$  Lets look now at another method used in CMS

# Discrete profiling method

P. Dauncey, M. Kenzie, N. Wardle and G. Davies

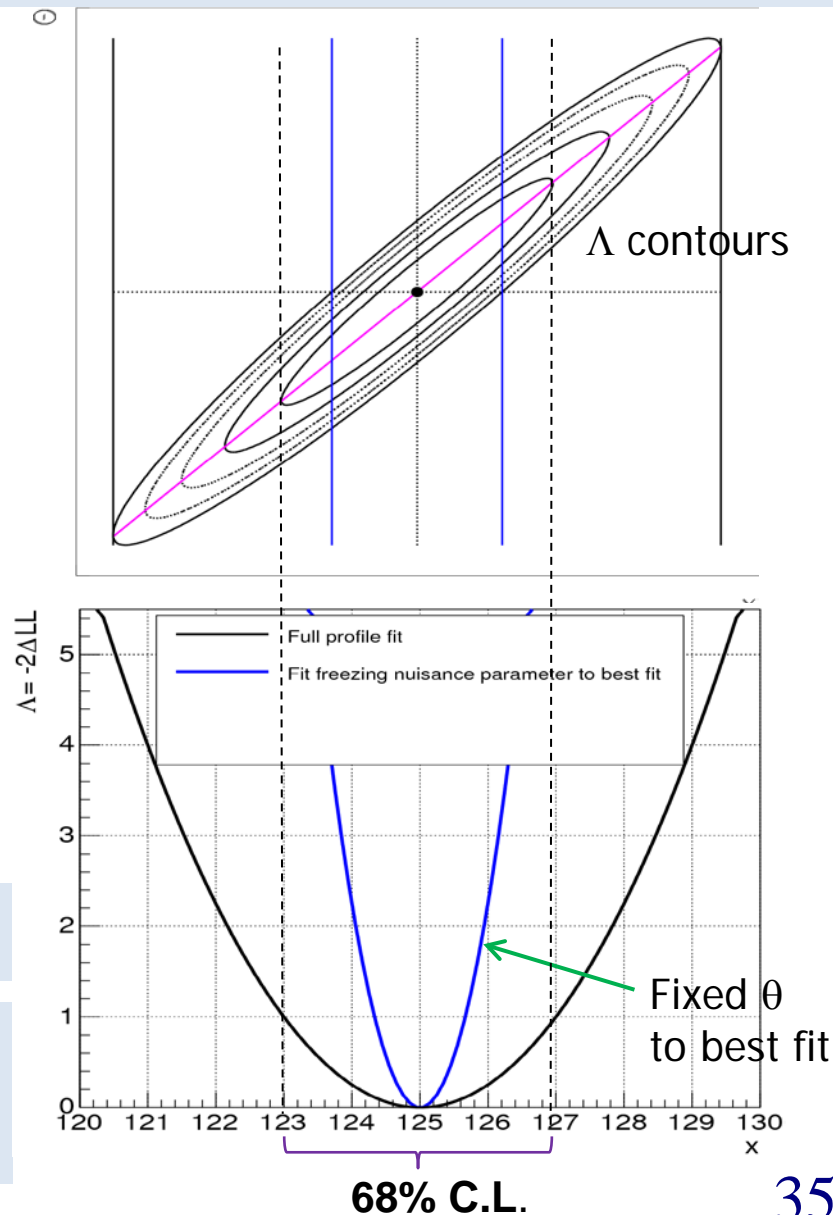
JINST 10 P04015 [arXiv:1408.6865]

Standard Profile likelihood:  
Scan  $\Lambda = -2\Delta\ln(L)$  vs  $x$ ; profiling  $\theta$



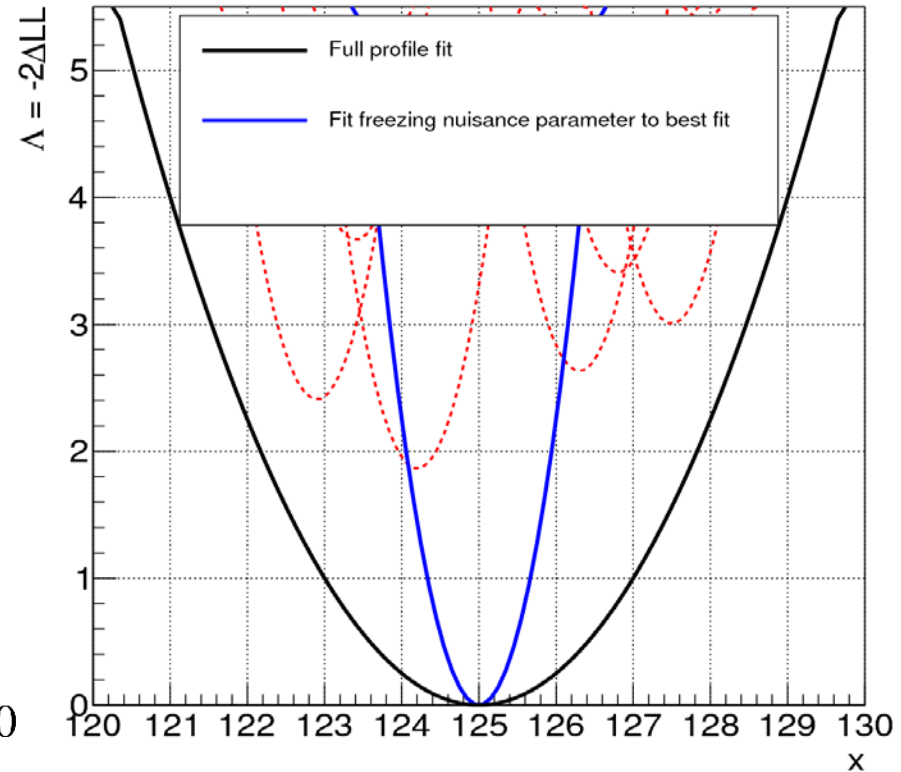
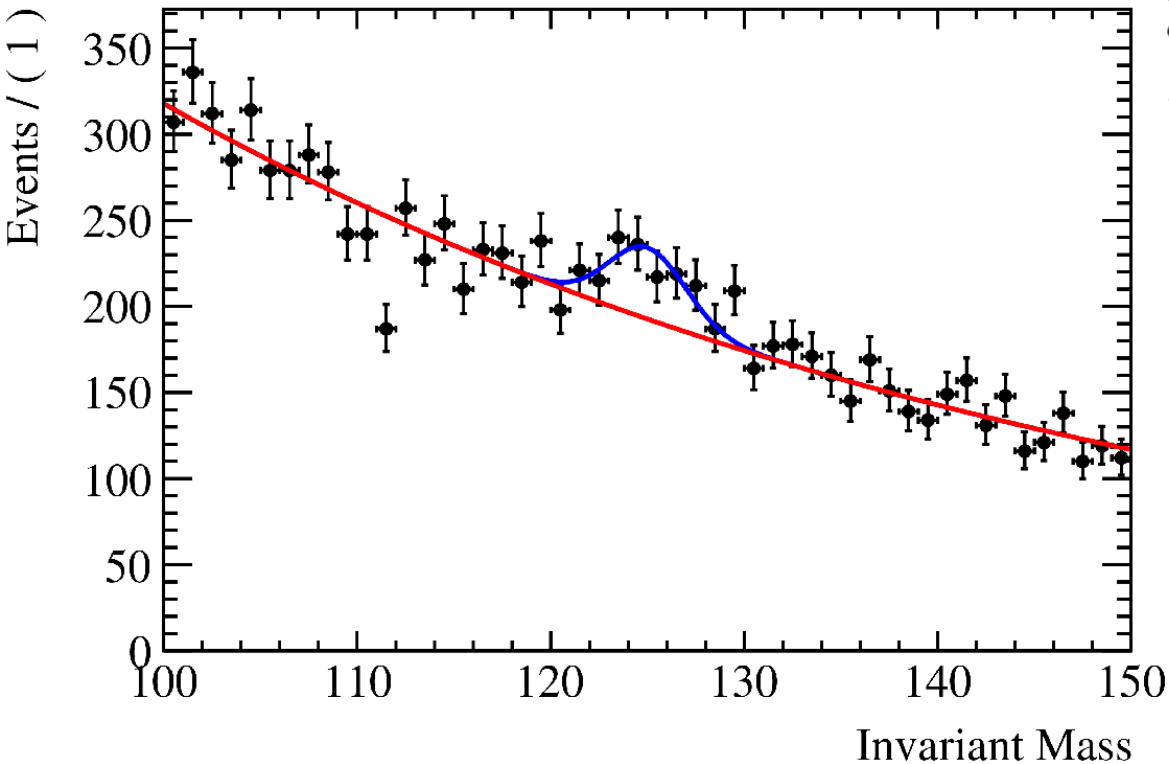
Fit gaussian signal + exponential background

Parameters:  $x$  = signal mass;  $\theta$  = background exponential slope (nuisance)



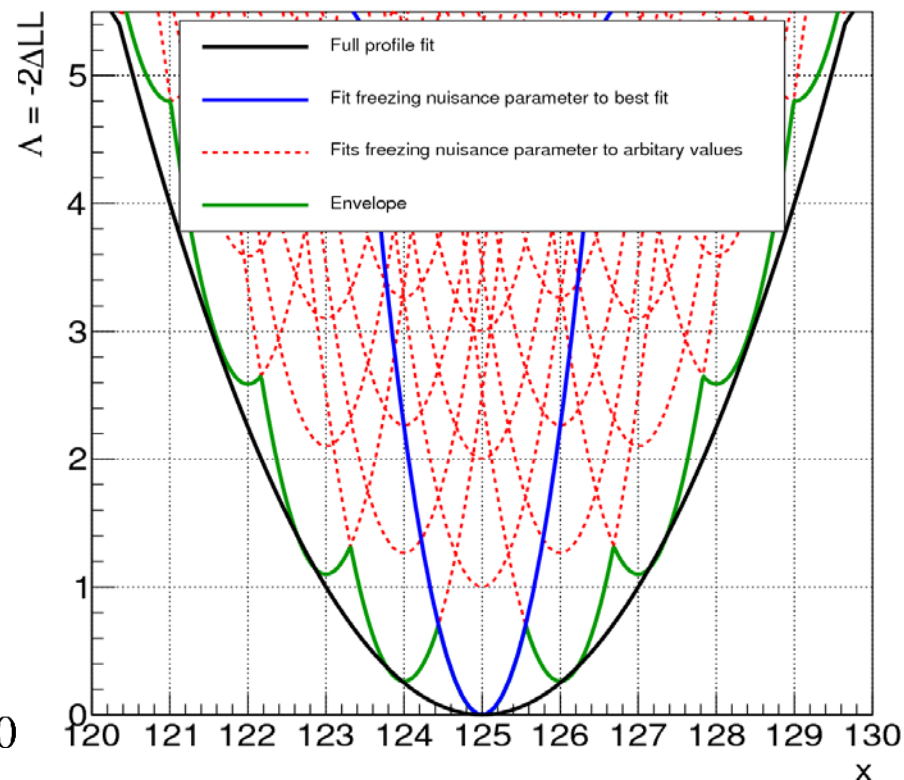
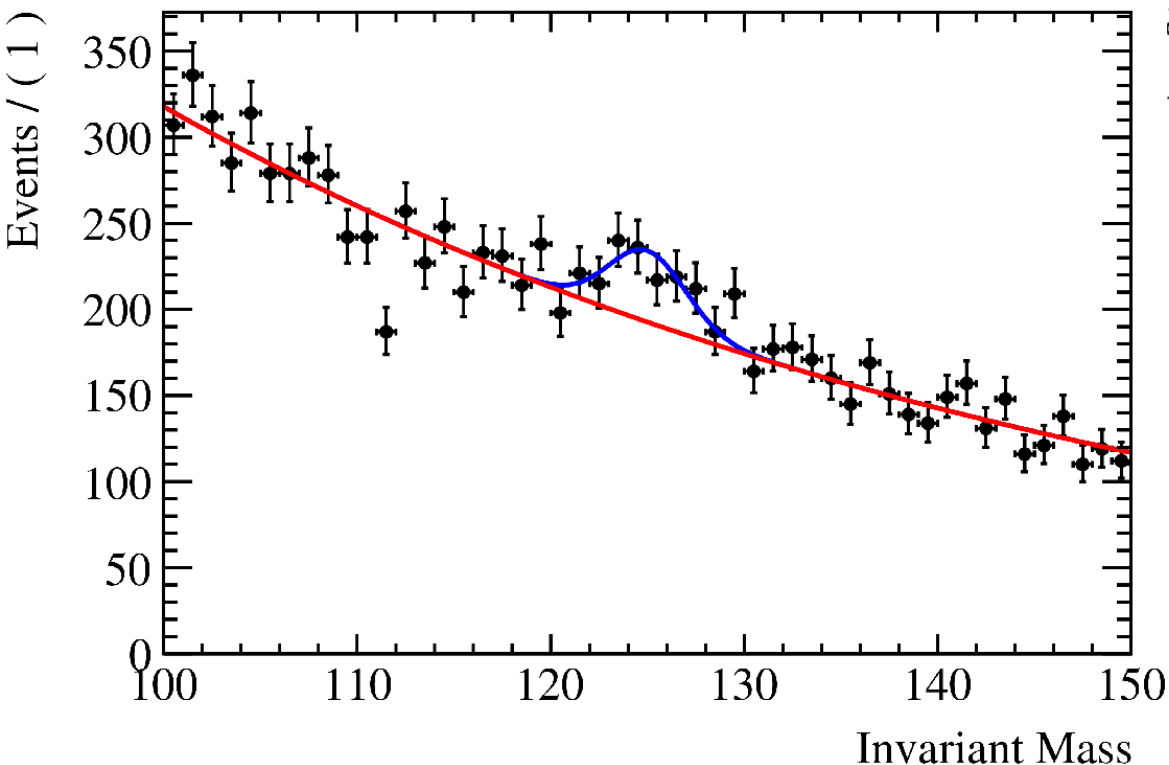
# Play around with nuisance parameter

Fix  $\theta$  to a few random values  $\rightarrow$  red dashed lines



# Play around with nuisance parameter

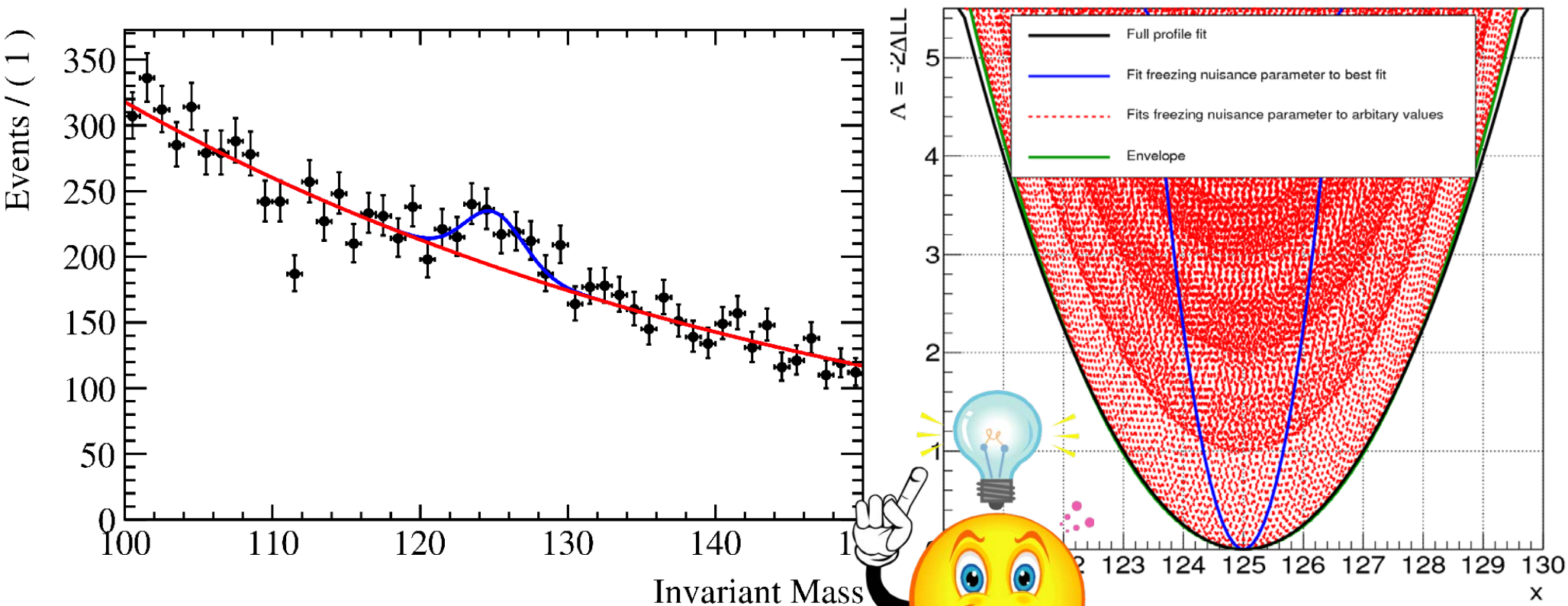
Fix  $\theta$  to **many** random values  
→ more red dashed lines



Draw minimum envelope (green line)  
→ begin to recover original curve

# Play around with nuisance parameter

Fix  $\theta$  to **huge number** of random values  $\rightarrow$  **more red dashed lines**



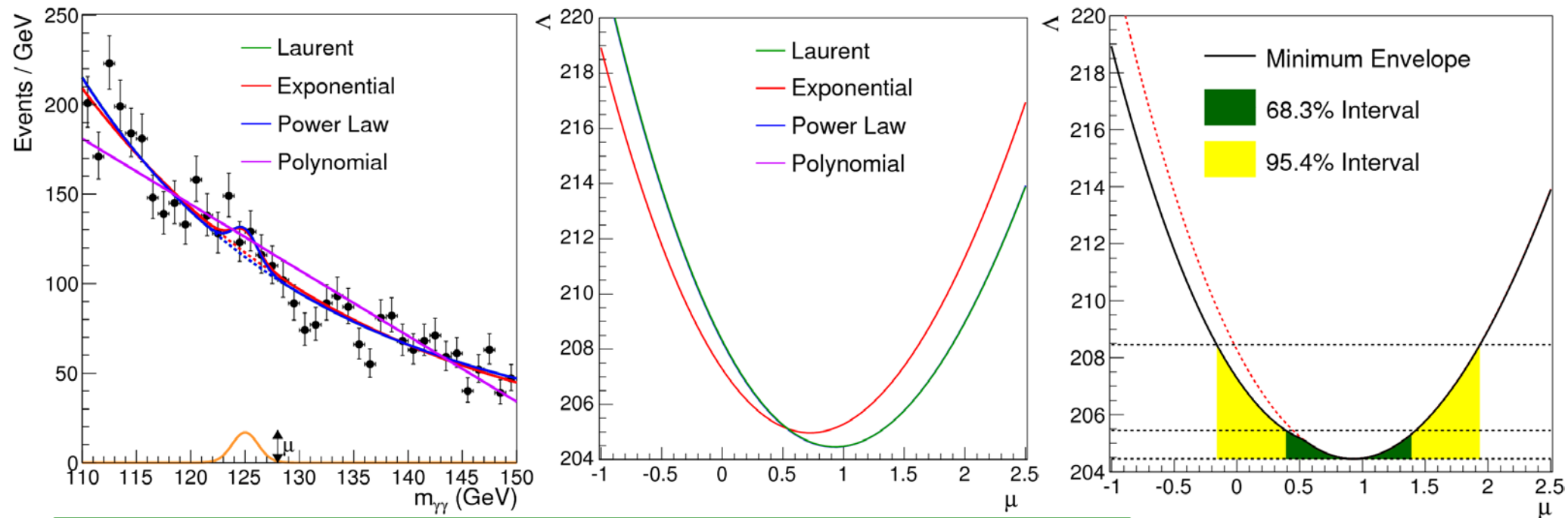
Minimum envelope = original curve  
 $\rightarrow$  One can **mix** discrete nuisance parameters with continuous ones

# A more realistic example

Fit  $\mu$  · signal-model + background

$$\Lambda = 2 \left[ \sum_i f_i - n_i - n_i \ln \frac{f_i}{n_i} \right] \quad (\text{Baker-Cousins } \tilde{\chi}^2)$$

Test background functions with same #fit parameters



→ Minimum envelope provides:

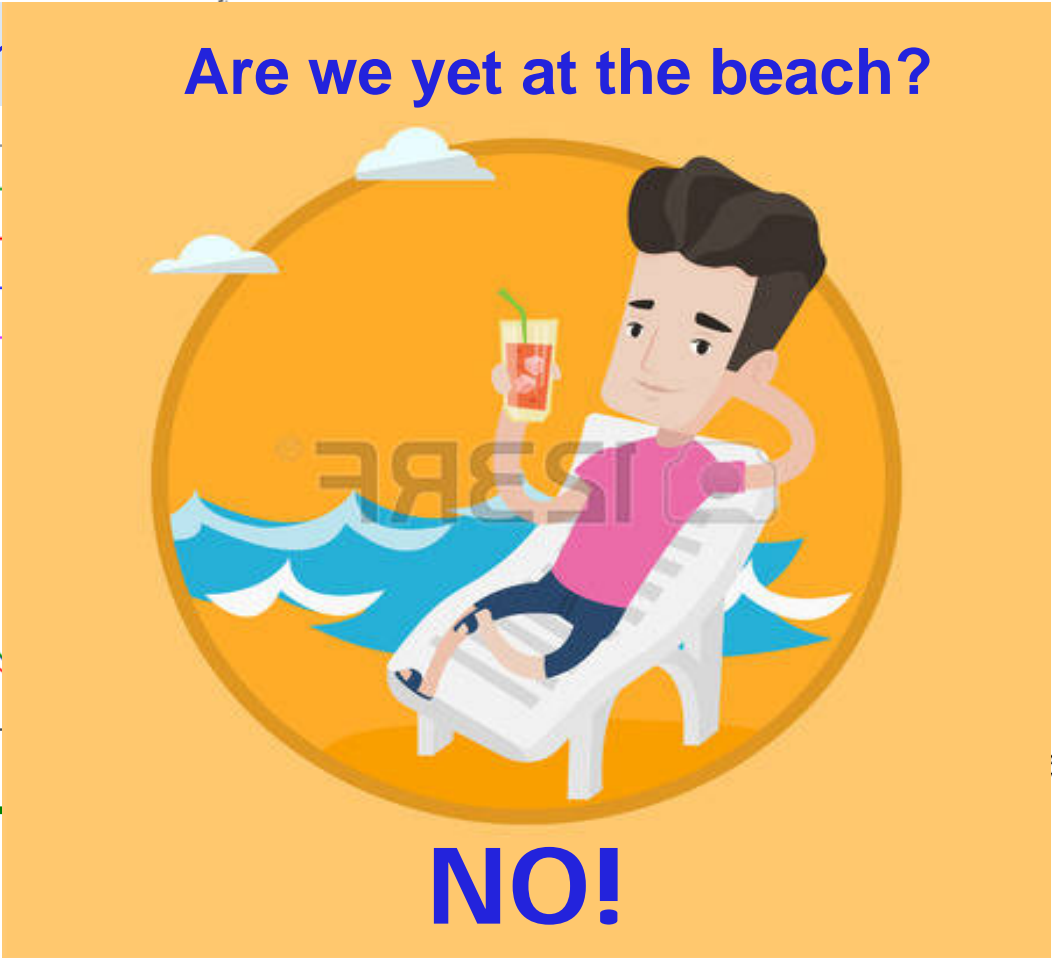
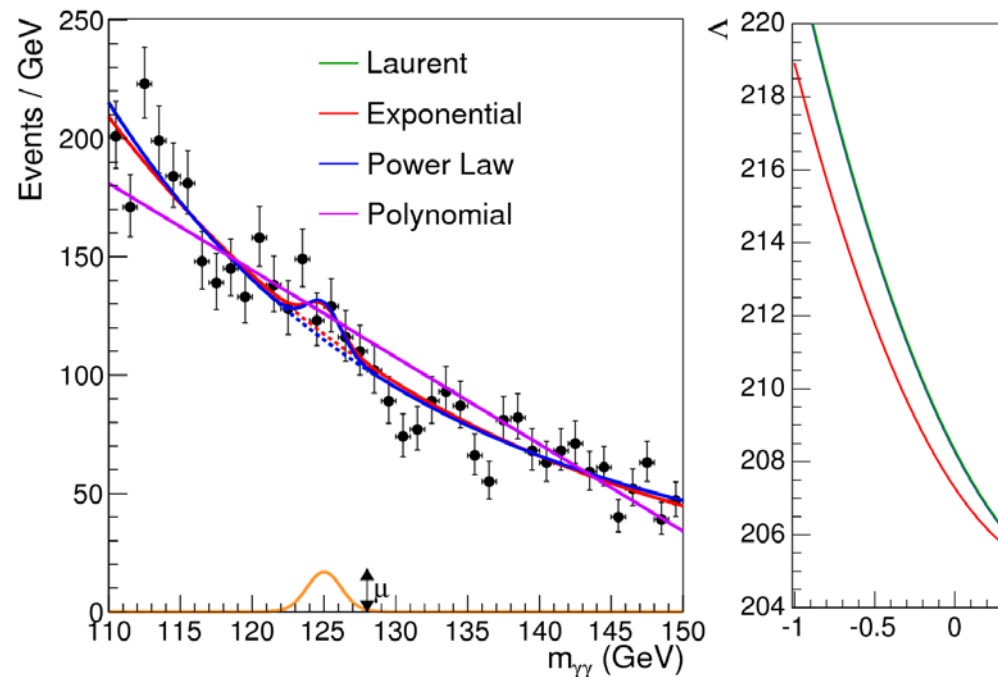
- best fit value  $\hat{\mu}$
- Confidence interval ( $\Delta\Lambda \leq 1$ )
- Systematic from background model choice

# A more realistic example

Fit  $\mu$  · signal-model + background

$$\Lambda = 2 \left[ \sum_i f_i - n_i - n_i \ln \frac{f_i}{n_i} \right] \quad (\text{Baker-Cousins } \tilde{\chi}^2)$$

Test background functions with same



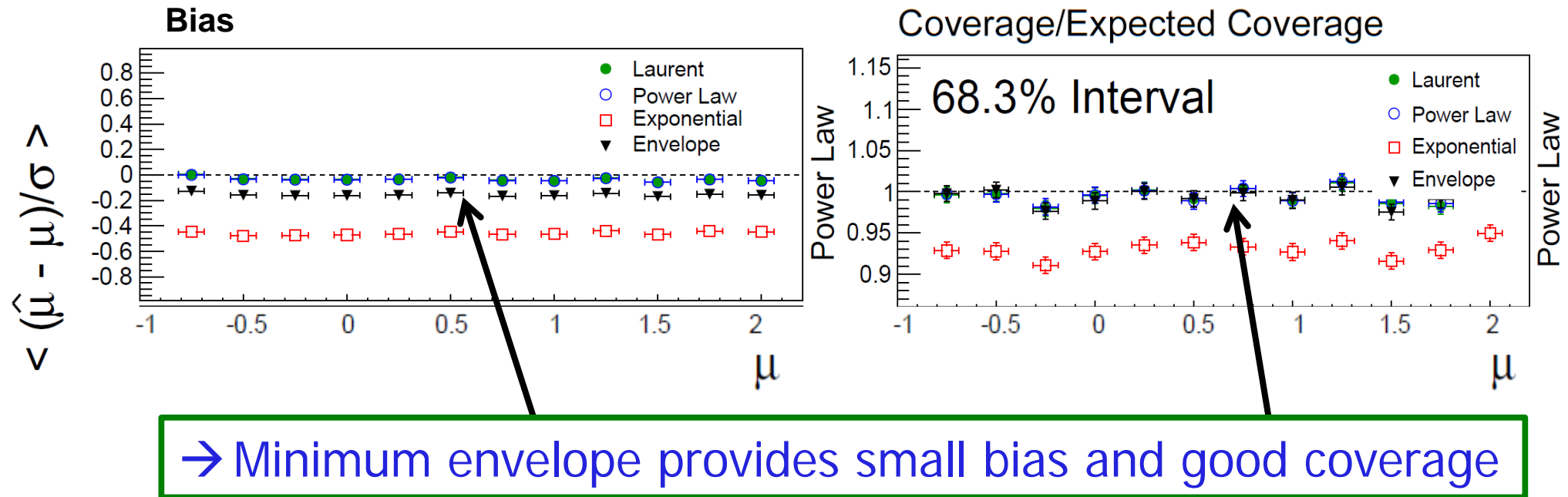
→ Minimum envelope provides:

- best fit value  $\hat{\mu}$
- Confidence interval ( $\Delta\Lambda \leq 1$ )
- Systematic from background model choice



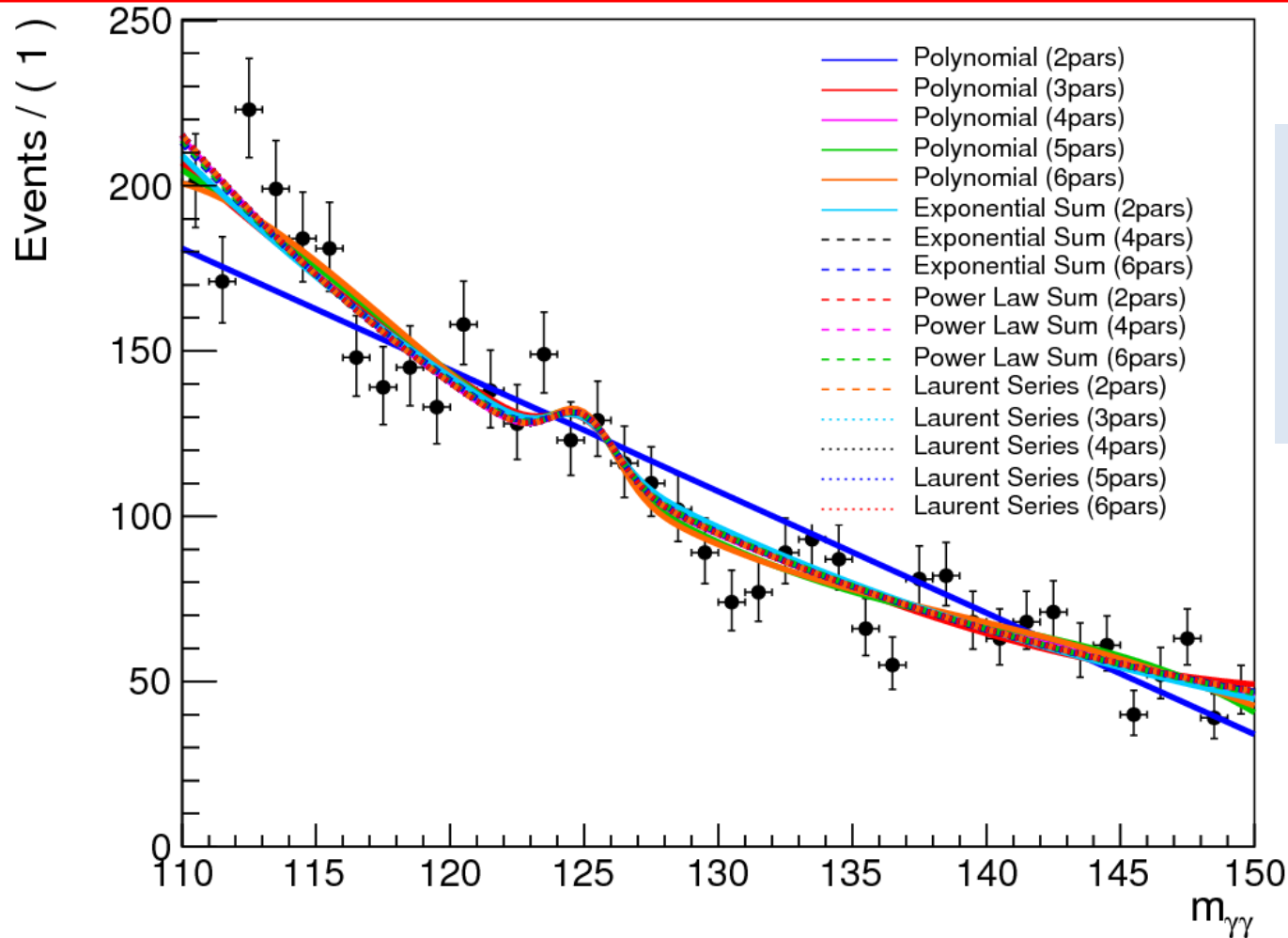
# Bias and Coverage

Generate toy MC from various background hypotheses and study bias and coverage<sup>†</sup> of fitted  $\hat{\mu}$  as function of generated true  $\mu$



<sup>†</sup> Coverage: correct coverage means that in 68.3% of repeated experiments the true parameter value is contained within the estimated  $\pm 1$  sigma region.

# Fits with background functions of different orders



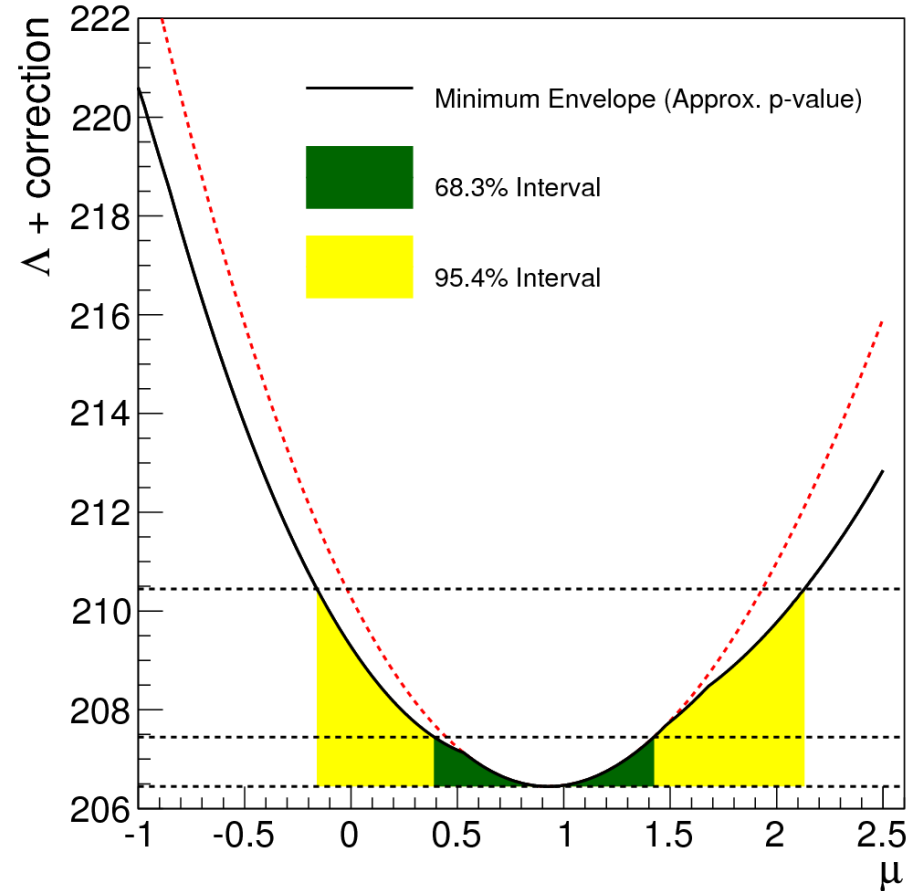
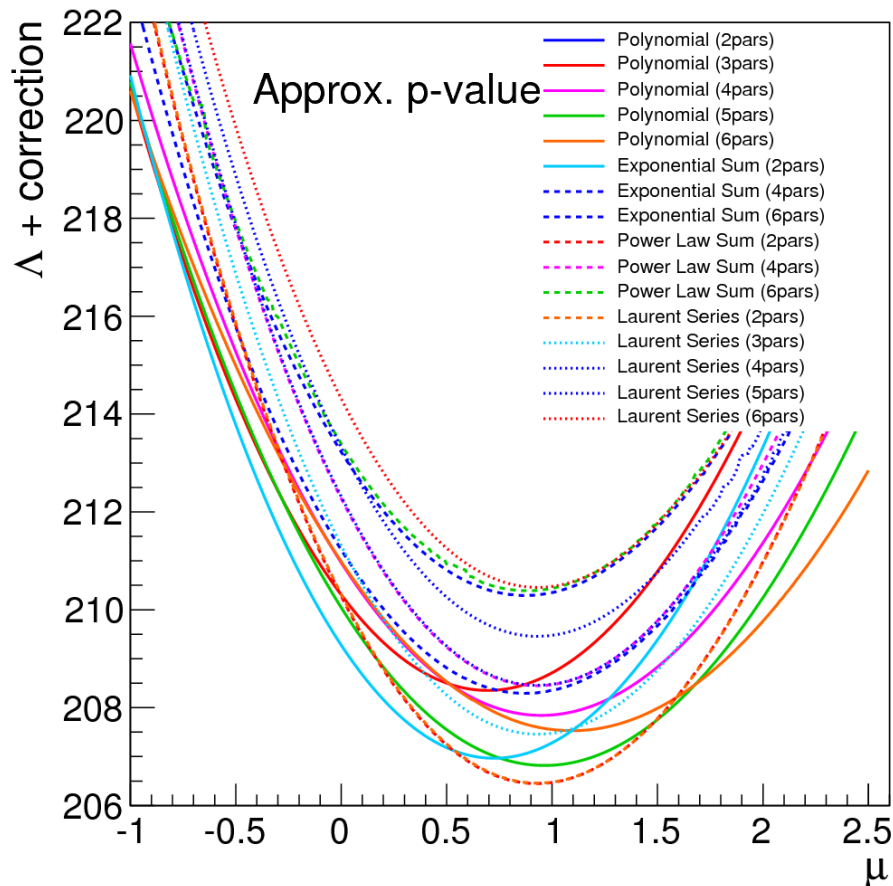
Minimum  $\Delta$  envelope:  
functions with large  
#fit parameters  
(npars) yield lower  $\Delta$

→ need to correct  $\Delta$  for different npars

$$\Delta = -2\ln(L) + c \text{ npars}$$

$c=1 \triangleq$  "approximate p-value correction"

# $\Lambda$ scans and minimum envelope $\Lambda = -2\ln(L) + c \text{ npars}; \quad c=1$



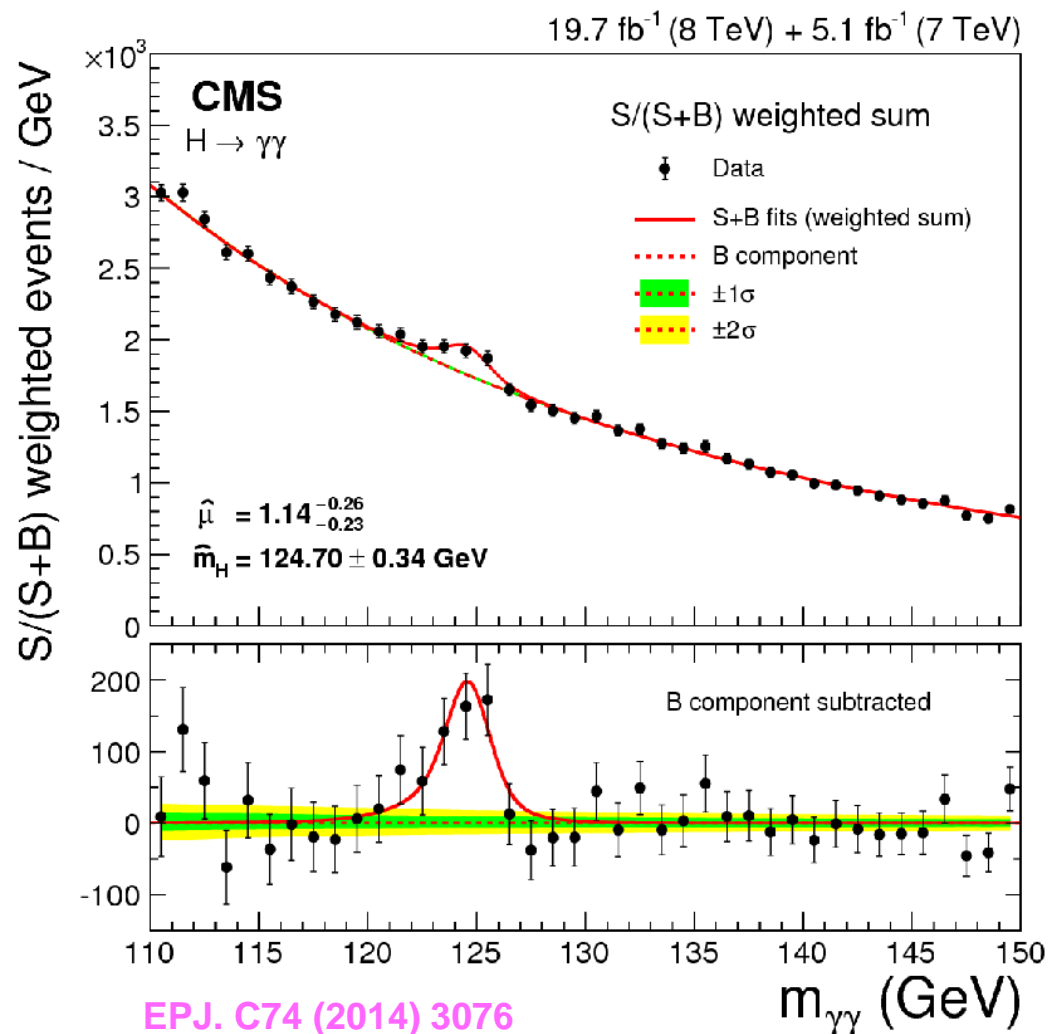
Best fit: 2 parameter power law

## Choice of $c$ :

- Large, e.g. 5  $\rightarrow$  prefer lower order functions  $\rightarrow$  potential biases
- Small, e.g. 0.1  $\rightarrow$  prefer higher order functions  $\rightarrow$  blow up  $\sigma_{\text{stat}}$

# Summary of Discrete Profile Likelihood method

- Method shows good coverage in toy experiments → perform toy experiments for your specific analyses
- Choices (open questions):
  - function models to include?
  - c term
- Method has been used e.g. in CMS  $H \rightarrow \gamma\gamma$  analysis



EPJ. C74 (2014) 3076

# Summary

## Goodness-of-Fit tests for

1 Checking data modelling

→ Perform  $\geq 2$  tests, e.g.  $\tilde{\chi}^2$  and K.S.

2 Outlier rejection

→  $\chi^2$  tests powerful tool

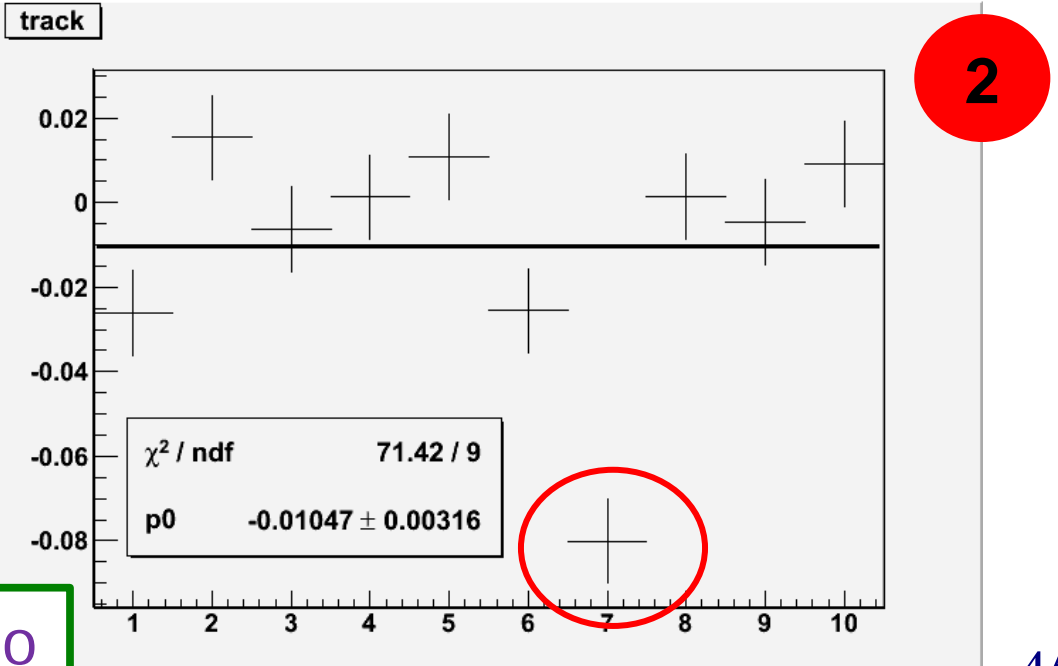
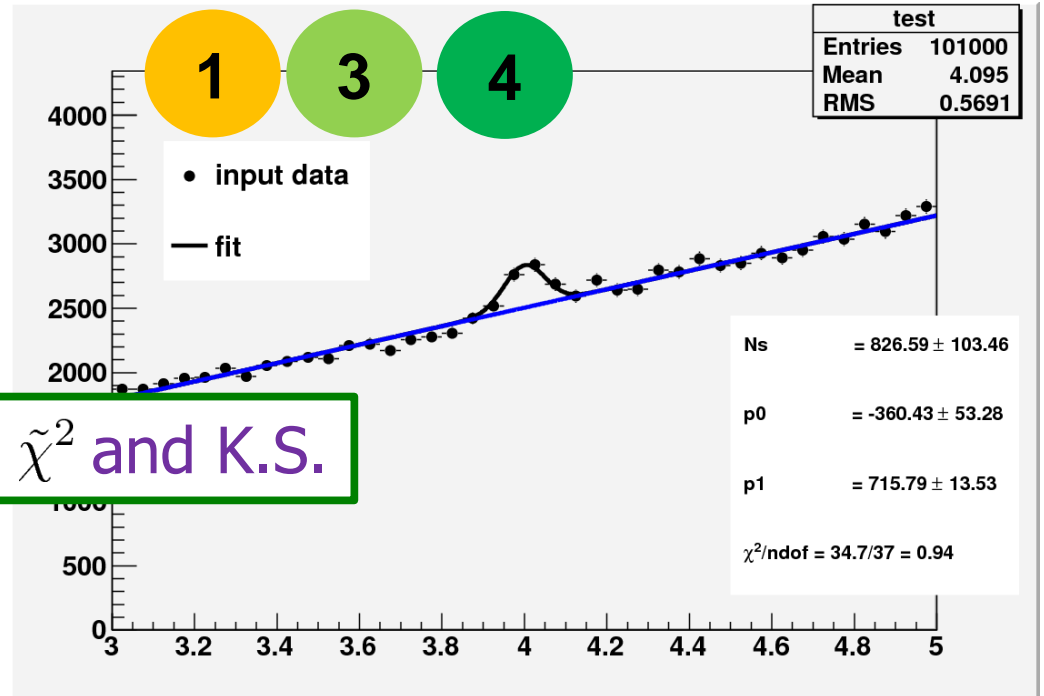
## Likelihood ratios for Background

3 Optimal parametrisation

→ Stop  $k \rightarrow k+1$  when:  
 $\text{TMath::Prob}(\tilde{\chi}_{k+1}^2 - \tilde{\chi}_k^2, 1) > 5\%$  &  
 $\text{TMath::Prob}(\tilde{\chi}_k^2, ndf) > 5\%$

4 Shape systematics  
(discrete profiling)

→ stat+syst error in one go



Final riddle – part I

Meggie has two children and  
the older one is a girl

→ What is the probability that the  
other child is also a girl?

Martin Gardner, “the two-child  
problem,” *Scientific American*, 1959

## Final riddle – part II

From all families with two children, at least one of whom is a girl, a family is chosen at random → What is the probability that both children are girls?

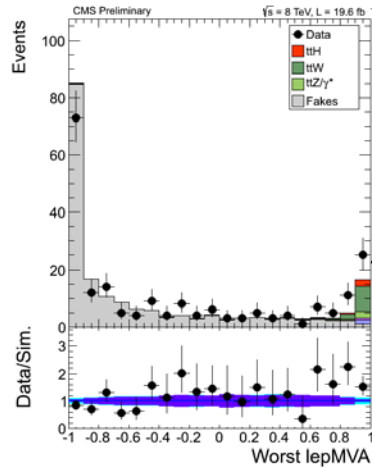
# Backup slides

---



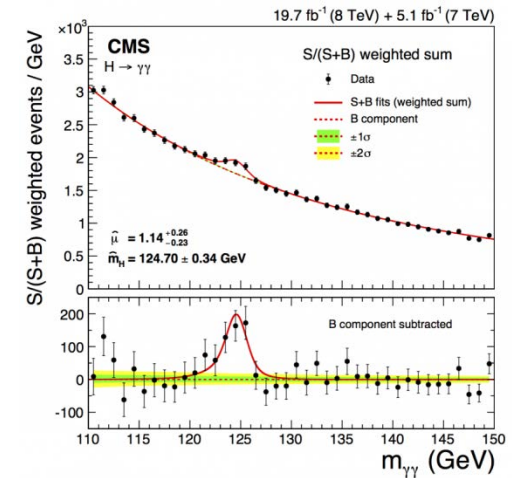
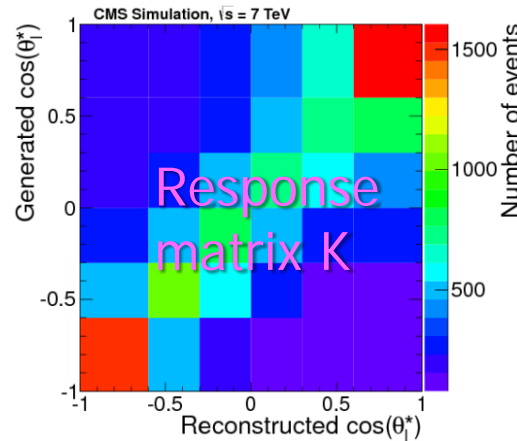
# Statistical Data Analysis – typical tasks

1. Optimal S/B separation



2. Signal searches, fits of all kind of interesting physics parameters to data and limits

3. Unfolding differential  $\sigma$



4. Systematic uncertainties

5. Data combination

