# Statistical methods in HEP "Overview"
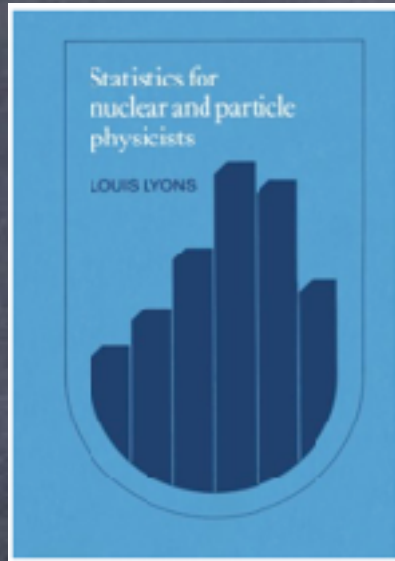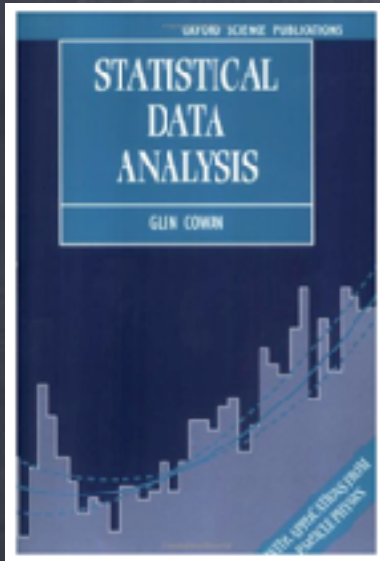
January 2018

Spåtind 2018

A. Read (U. Oslo)
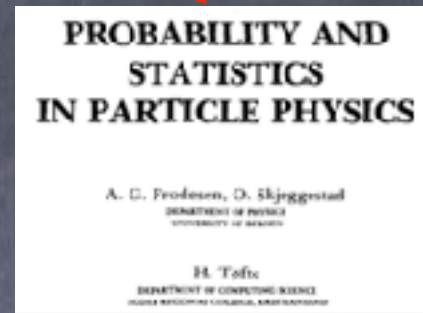
UiO **:** **University of Oslo**

The Research Council of Norway

Note: Many clickable links to documentation!

# Real (!) overviews

# 2 main approaches

- Bayesian - probability(theory|data) $p(\theta|x)$
  - well-defined accounting for beliefs
  - prior-probability for the theory must be given
  - prior-dependence should be studied

- Frequentist/classical - probability(data|theory) $p(x|\theta)$
  - says nothing about probability of theory
  - typically used in HEP to report experimental results "objectively" (as possible)
  - can lead to subset of individual results which are obviously wrong but consistent with methodology

# Bayesian credible intervals

$$P(A|B) = P(B|A)\frac{P(A)}{P(B)}$$

$$p(\boldsymbol{\theta}|\boldsymbol{x}) = \frac{L(\boldsymbol{x}|\boldsymbol{\theta})\pi(\boldsymbol{\theta})}{\int L(\boldsymbol{x}|\boldsymbol{\theta}')\pi(\boldsymbol{\theta}')\,d\boldsymbol{\theta}'}$$

Posterior density
for parameter

$$p(\boldsymbol{\theta}|\boldsymbol{x}) = \int p(\boldsymbol{\theta}, \boldsymbol{\nu}|\boldsymbol{x})\,d\boldsymbol{\nu}$$

Marginalizing nuisance
parameters (e.g. data-driven
backgrounds, systematics)

$$1 - \alpha = \int_{\theta_{\mathrm{lo}}}^{\theta_{\mathrm{up}}} p(\theta|\boldsymbol{x})\,d\theta$$

Interval:

Minimum interval
Highest density
Physical boundry (e.g.
m≥0)

# Confidence intervals (Neyman construction)



parameter $\theta$

$D(\alpha)$

$x_2(\theta), \theta_2(x)$

$\theta_0$

$x_1(\theta), \theta_1(x)$

$x_1(\theta_0)$   $x_2(\theta_0)$
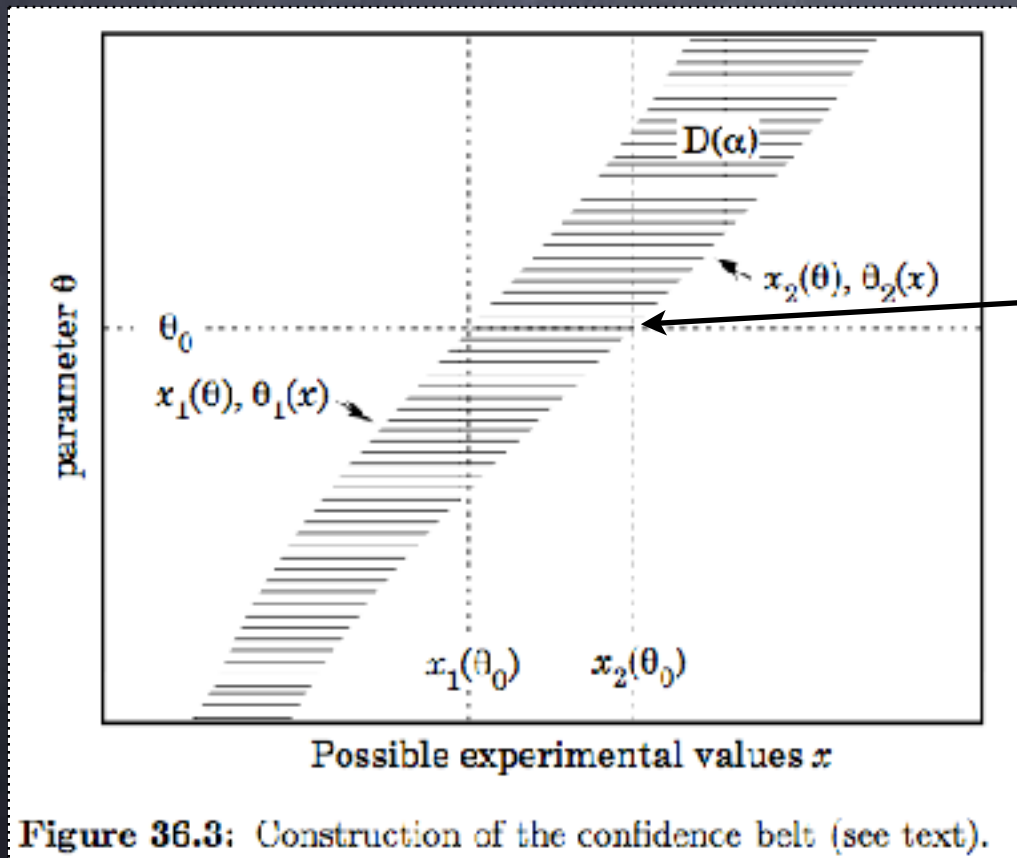
Possible experimental values $x$

**Figure 36.3:** Construction of the confidence belt (see text).

- Need to know the ensemble for every$\theta_0$

- Multi-dimensional space with nuisance parameters more complicated (ugh)

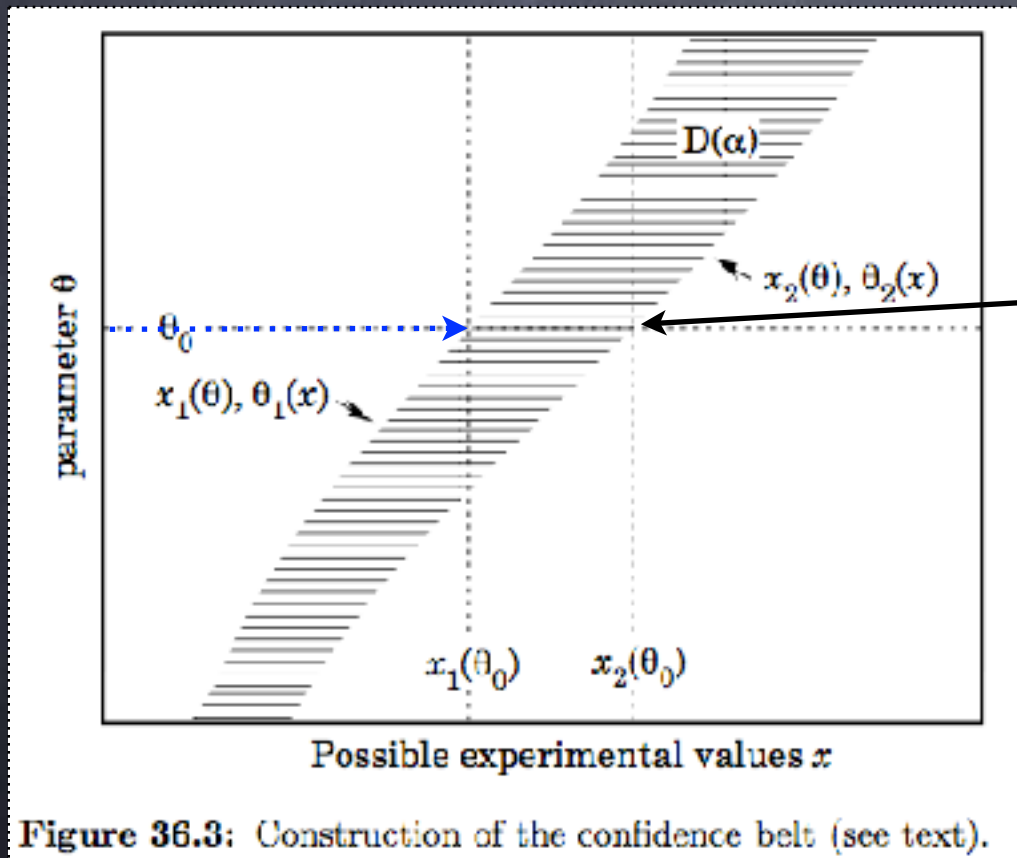# Confidence intervals (Neyman construction)



**Figure 36.3:** Construction of the confidence belt (see text).
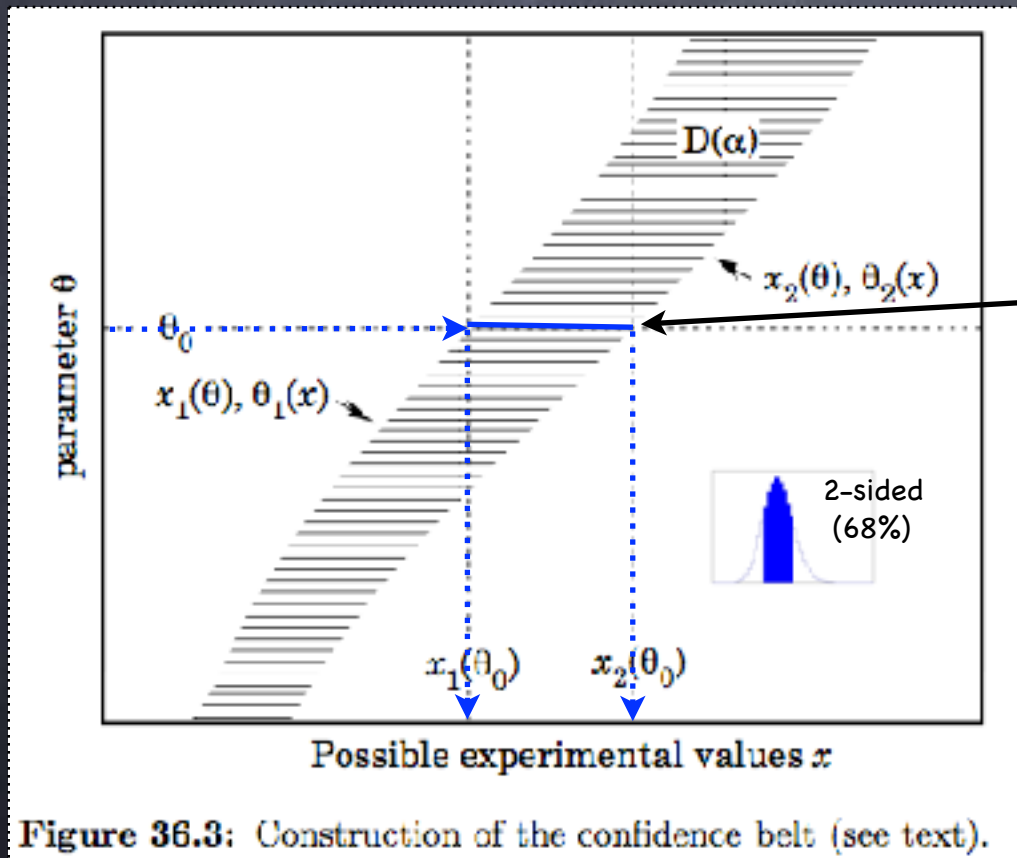
- Need to know the ensemble for every $\theta_0$

- Multi-dimensional space with nuisance parameters more complicated (ugh)

# Confidence intervals (Neyman construction)



**Figure 36.3:** Construction of the confidence belt (see text).
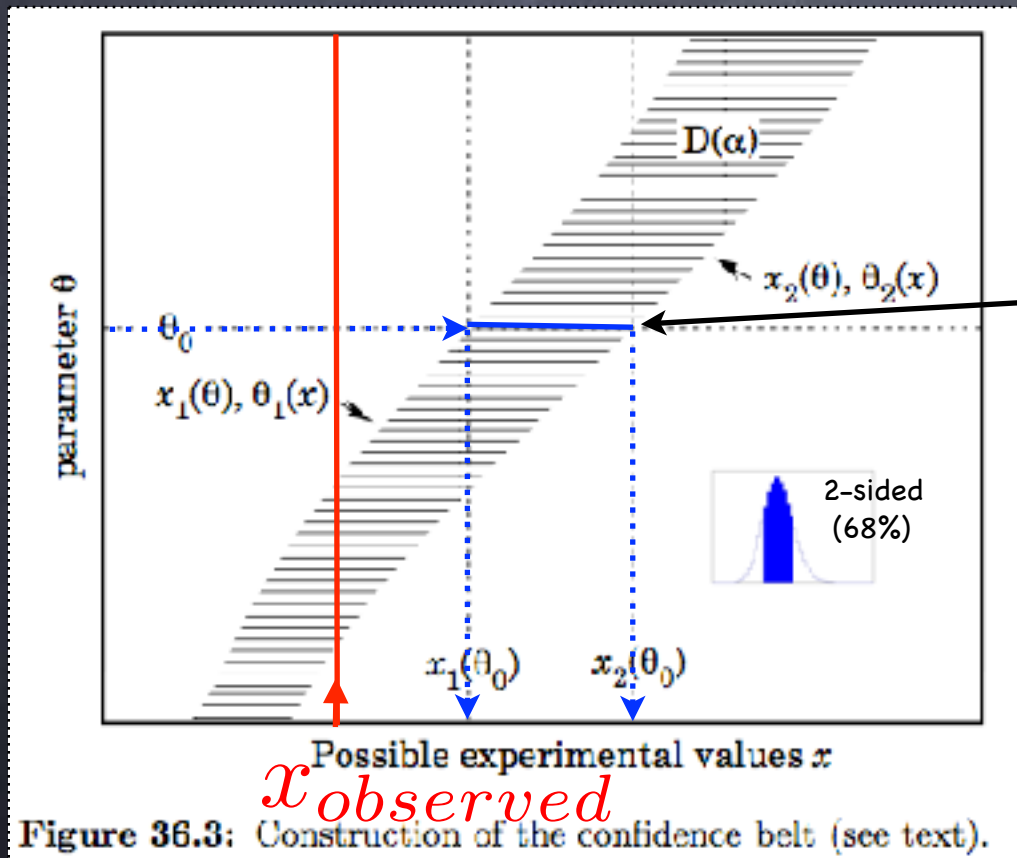
- Need to know the ensemble for every $\theta_0$

- Multi-dimensional space with nuisance parameters more complicated (ugh)

# Confidence intervals (Neyman construction)



**Figure 36.3:** Construction of the confidence belt (see text).
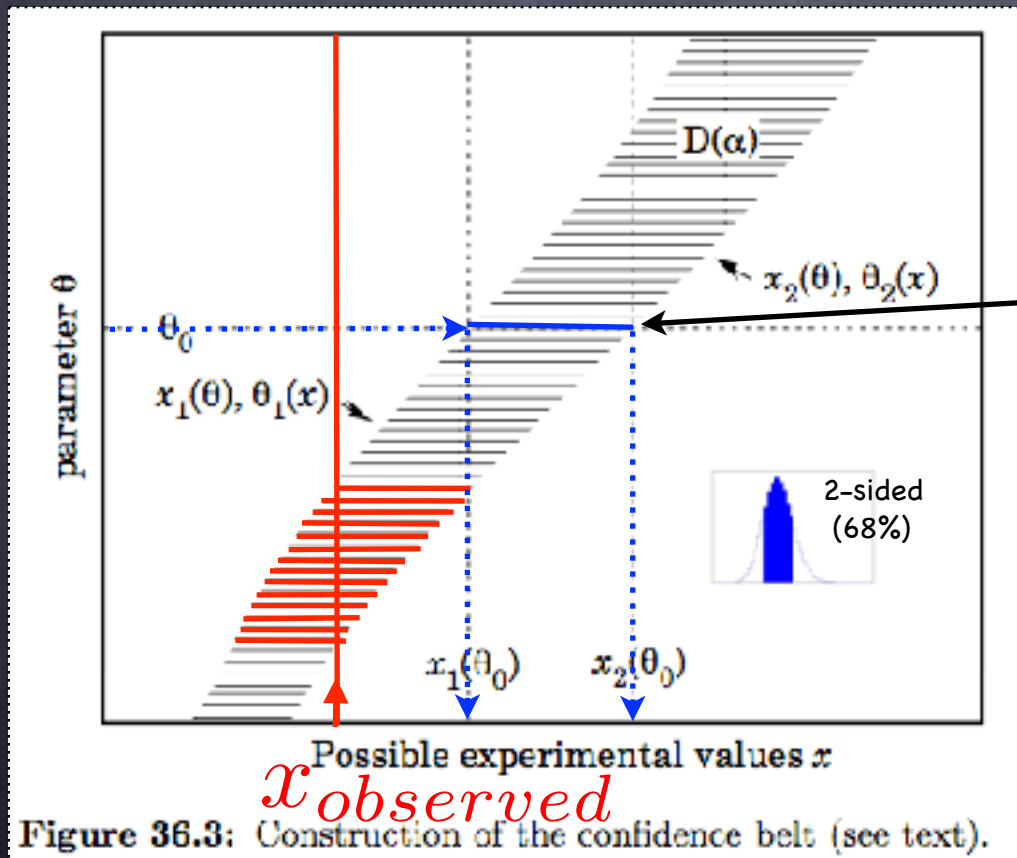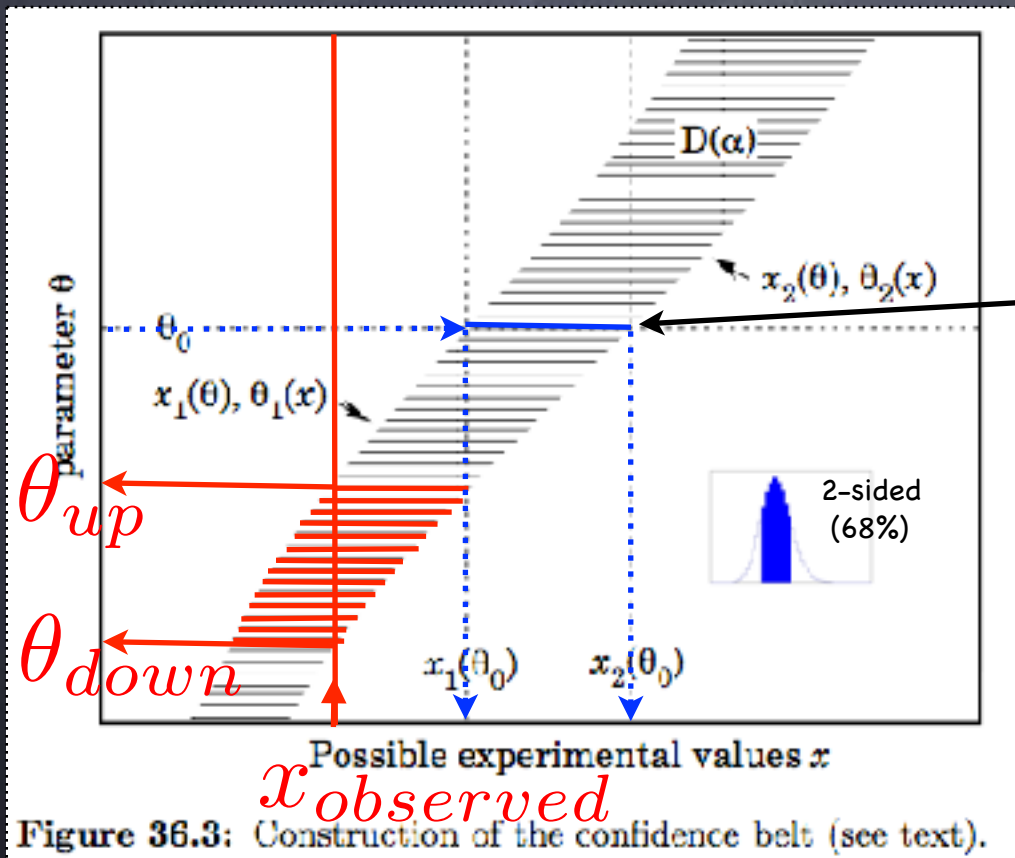
- Need to know the ensemble for every $\theta_0$

- Multi-dimensional space with nuisance parameters more complicated (ugh)

# Confidence intervals (Neyman construction)



**Figure 36.3:** Construction of the confidence belt (see text).

$x_{observed}$

- Need to know the ensemble for every $\theta_0$

- Multi-dimensional space with nuisance parameters more complicated (ugh)

D($\alpha$)

$x_2(\theta), \theta_2(x)$

$x_1(\theta), \theta_1(x)$

$\theta_0$

parameter $\theta$

2-sided (68%)

$x_1(\theta_0)$   $x_2(\theta_0)$

Possible experimental values $x$

# Confidence intervals (Neyman construction)



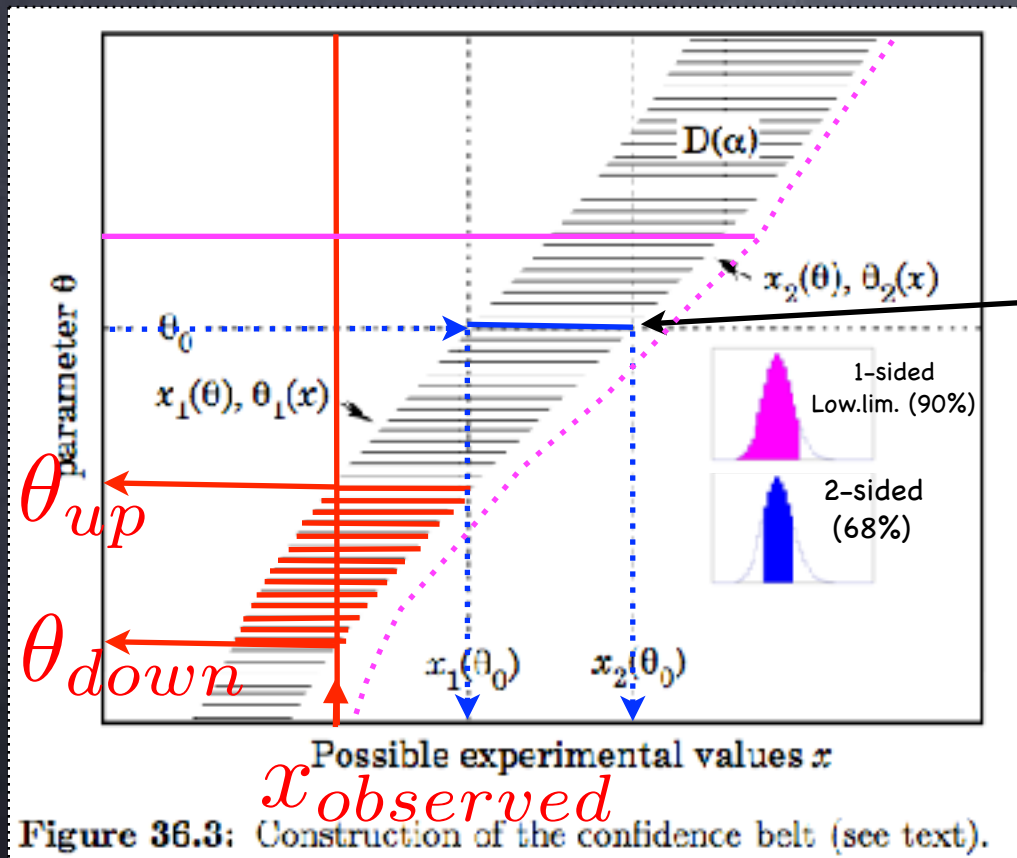Figure 36.3: Construction of the confidence belt (see text).
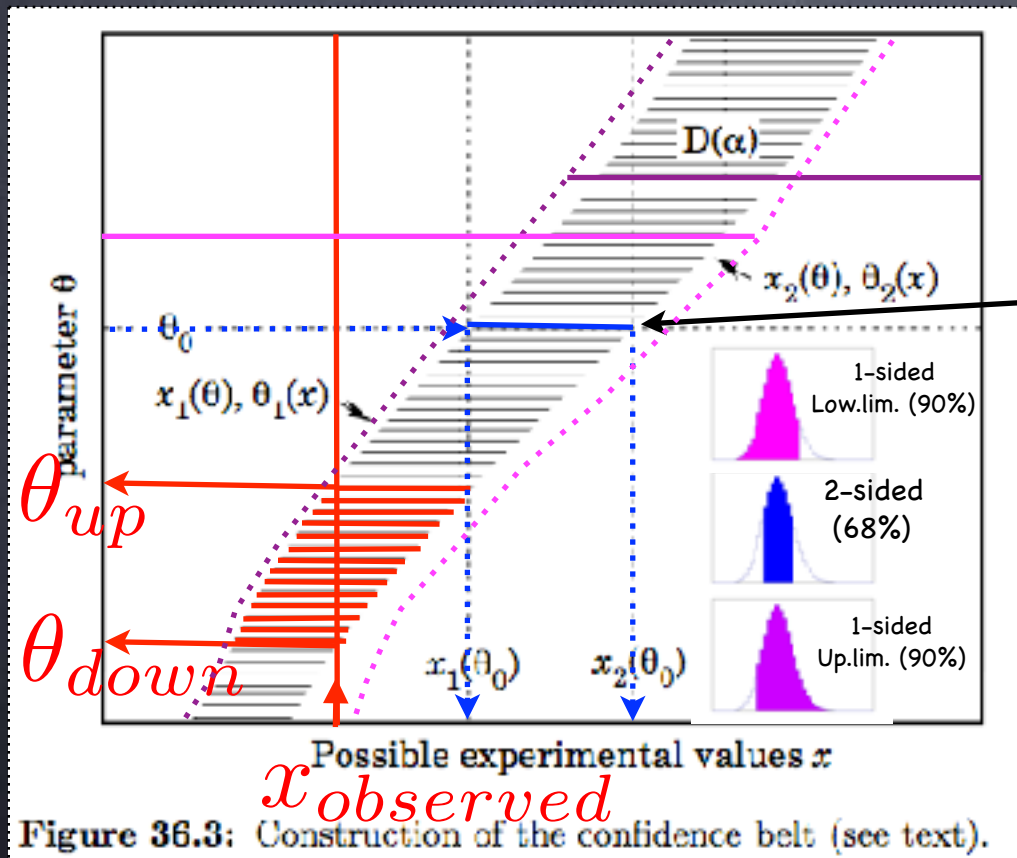
- Need to know the ensemble for every $\theta_0$

- Multi-dimensional space with nuisance parameters more complicated (ugh)

# Confidence intervals (Neyman construction)



Figure 36.3: Construction of the confidence belt (see text).

- Need to know the ensemble for every $\theta_0$

- Multi-dimensional space with nuisance parameters more complicated (ugh)

# Confidence intervals (Neyman construction)



Figure 36.3: Construction of the confidence belt (see text).

- Need to know the ensemble for every $\theta_0$

- Multi-dimensional space with nuisance parameters more complicated (ugh)

# Exam question (Bob Cousins)

For most of this talk[1], I assume familiarity with the 'required reading' for this workshop. But first, let's review the root of the problem as I often explain it to students. (Imagine an oral exam.)

Suppose you have a particle ID detector. You take it to a test beam and measure:

- P(counter says $\pi$ | particle is $\pi$) = 90%
- P(counter says not $\pi$ | particle is $\pi$) = 10%
- P(counter says $\pi$ | particle is not $\pi$) = 1%
- P(counter says not $\pi$ | particle is not $\pi$) = 99%

Then you put the detector in your experiment. You select tracks which the detector says are pions.

Question: What fraction of these tracks are pions?

◎ Related question: What is the probability that a particular track is a pion?

# Bayes vs. freq.

- In many data-dominated situations hardly any difference in reported results, eg. $M_Z = 91.1876 \pm 0.0021$ GeV

  - But interp. not the same!
    Which is B and which is F?
    1) $P(|MZ-91.1876|<0.0021)=68\%$
    2) 68% of such intervals contain the true $M_Z$

- **Small data samples, physical boundries** typically lead to differences

- Doing both analyses and studying the differences can give insights

# Various likelihoods

$$L(n|\mu) = \frac{e^{-\mu}\mu^n}{n!}$$

Poisson, counting (no background)

$$L(n|\mu s + b) = \frac{e^{-(\mu s + b)}(\mu s + b)^n}{n!}$$

Counting, known bkg

$$L(n, m|\mu s + b, \tau) = \frac{e^{-(\mu s+b)}(\mu s+b)^n}{n!}\frac{e^{-\tau b}(\tau b)^m}{m!}$$

Counting "on/off"

$$L(x|x_0, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}}e^{-\frac{(x-x_0)^2}{2\sigma^2}}$$

Gaussian

$$Q = \frac{\prod_{i=1}^{N_{chan}}\frac{e^{-(s_i+b_i)}(s_i+b_i)^{n_i}}{n_i!}}{\prod_{i=1}^{N_{nchan}}\frac{e^{-b_i}b_i^{n_i}}{n_i!}}\frac{\prod_{j=1}^{n_i}\frac{s_i S_i(x_{ij})+b_i B_i(x_{ij})}{s_i+b_i}}{\prod_{j=1}^{n_i}B_i(x_{ij})}$$

Likelihood ratio of marked Poissons in combined channels

# Maximum likelihood

- Ideal estimators of parameters are unbiased and efficient (minimum variance). Not always simultaneously achievable, e.g,

$$s^2 = \frac{1}{N} \sum_i (x_i - \bar{x})^2 \to s^2 = \frac{1}{N-1} \sum_i (x_i - \bar{x})^2$$

- Maximum likelihood (for convenience minimize -ln(L) or even -2ln(L)) is approximately unbiased, efficient for large data samples and widely applicable.

- Wald showed that for a single parameter $\mu$

$$-2 \ln \lambda(\mu) = \left(\frac{\mu - \hat{\mu}}{\sigma}\right)^2 + O(1/\sqrt{N})$$

- Wilks showed that if $\hat{\mu}$ is Gauss-distributed about $\mu$ then

$$-2 \ln \lambda(\mu) \to \chi^2$$

# 1-sided p-values in large-sample limit

```
Double_t Pvalue(Double_t significance) {
  return ROOT::Math::chisquared_cdf_c(pow(significance,2),1)/2;
}
```

| $N_\sigma$ | $\Delta\chi^2$ | $\frac{1}{2}P(\chi^2 > c)$ |
|---|---|---|
| 1 | 1 | 0,159 |
| 2 | 4 | $2.3 \times 10^{-2}$ |
| 3 | 9 | $1.3 \times 10^{-3}$ |
| 4 | 16 | $3.2 \times 10^{-5}$ |
| 5 | 25 | $2.9 \times 10^{-7}$ |



$$q_\mu = -2\ln\frac{L(\mu_{test})}{L(\hat{\mu})} = \chi^2_{test} - \chi^2_{min}$$

# Brief (!) history of limits

- O. Helene (1983) - Bayesian limit with flat prior on signal

- G. Zech (1988) - frequentist interpretation of Helene

- A. Read (1997) - rederived Zech from likelihood ratio and "background conditioning"; $CL_S \approx$ "confidence in the signal-only hypothesis"

- Feldman and Cousins (1998) - auto 2-sided frequentist confidence intervals - "coverage is king" (but tests signal+background hypothesis)

- Birnbaum (1961!!) - support for $CL_S$ in the professional statistics literature - rediscovered by O. Vitells

  - Article links 1, 2, 3

$$CL = \frac{\int_s^\infty \mathcal{L}(s', b)\,ds'}{\int_0^\infty \mathcal{L}(s', b)\,ds'}.$$

$$CL = 1 - \frac{\sum_{n=0}^{n_{obs}} \frac{e^{-(b+s)}(b+s)^n}{n!}}{\sum_{n=0}^{n_{obs}} \frac{e^{-b}b^n}{n!}}.$$

$$CL_s \equiv CL_{s+b}/CL_b.$$



"A concept of statistical evidence is not plausible unless it finds 'strong evidence for H2 as against H1' with small probability (alpha) when H1 is true, and with much larger probability (1 –beta) when H2 is true."

# Origins of CL$_s$

- Almost background-less Higgs searches at LEP1, many different statistical treatments, combination not obvious, LEP2 data was coming

- I proposed simple LR, frequentist approach, CL$_s$ invented to deal robustly with deficits, combination simply adding channels to LR, exclusion with CL$_s$, discovery with CL$_b$, never got to ML for measurement

- Cousins&Highland (hybrid Bayes-frequentist treatment) for (generally small) systematics

# Q$_{LEP}$ and CL$_S$ take hold in DELPHI

# CL$_s$



Counting experiment (Poisson pdfs)



CL$_{s+b}$

1-CL$_b$

$$Q_i = \frac{\dfrac{e^{-(s_i+b_i)}(s_i+b_i)^{n_i^{cand}}}{n_i^{cand}!}}{\dfrac{e^{-b_i}b_i^{n_i^{cand}}}{n_i^{cand}!}}$$

$$-2\ln Q_i = 2\, s_i - 2\, n_i \ln\left(1+\frac{s_i}{b_i}\right)$$

$$CL_{s+b} = P_{s+b}(X \leq X_{obs}),$$

$$P_{s+b}(X \leq X_{obs}) = \int_0^{X_{obs}} \frac{dP_{s+b}}{dX} dX$$

$$CL_b = P_b(X \leq X_{obs}),$$

$$P_b(X \leq X_{obs}) = \int_0^{X_{obs}} \frac{dP_b}{dX} dX$$

$$CL_s \equiv CL_{s+b}/CL_b.$$

$$1 - CL_s \leq CL.$$

# Straightforward LR combination

- Natural combination of channels, extension to discriminant (or counting) per channel

- Learned later <u>Obraztsov (DELPHI 1992)</u>, L3 people proposed similar likelihood but Bayes-like integration of likelihood (implicit uniform prior).

- At LEP eventually 4 experiments, O(10) center of mass energies, O(8) search topologies/channels combined

$$Q = \frac{\prod_{i=1}^{N_{chan}} \frac{e^{-(s_i+b_i)}(s_i+b_i)^{n_i}}{n_i!} \prod_{j=1}^{n_i} \frac{s_i S_i(x_{ij}) + b_i B_i(x_{ij})}{s_i+b_i}}{\prod_{i=1}^{N_{nchan}} \frac{e^{-b_i} b_i^{n_i}}{n_i!} \prod_{j=1}^{n_i} B_i(x_{ij})}$$

# LR from LEP to Tevatron to LHC

| | Test statistic | Nuisance parameters in LR | Randomized in toys | Sampling of test statistic |
|---|---|---|---|---|
| $Q_{\text{LEP}}$ | $-2\ln\dfrac{L(\mu,\tilde{\theta})}{L(0,\tilde{\theta})}$ | Fixed by MC | Nuisance parameters | Hybrid Bayes-frequentist |
| $Q_{\text{Tev}}$ | $-2\ln\dfrac{L(\mu,\hat{\hat{\theta}})}{L(0,\hat{\theta})}$ | Profiled | Nuisance parameters | Hybrid Bayes-frequentist |
| "LHC" $q_\mu$ ($q_0$) | $-2\ln\dfrac{L(\mu(0),\hat{\hat{\theta}})}{L(\hat{\mu},\hat{\theta})}$ | Profiled | External constraints | Frequentist |

# Profile likelihood (MINUIT)

lanl.arXiv.org > physics > arXiv:physics/0403059

Se

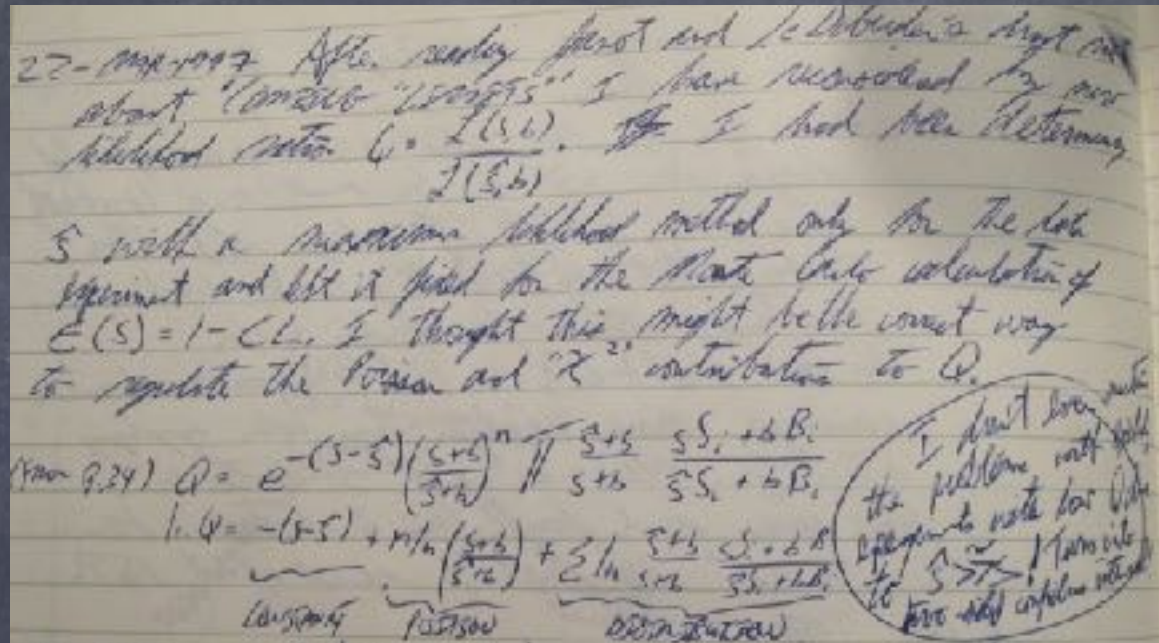**Physics > Data Analysis, Statistics and Probability**

## Limits and Confidence Intervals in the Presence of Nuisance Parameters

Wolfgang A. Rolke, Angel M. Lopez, Jan Conrad

*(Submitted on 9 Mar 2004 (v1), last revised 19 Jan 2009 (this version, v5))*

We study the frequentist properties of confidence intervals computed by the method known to statisticians as the Profile Likelihood. It is seen that the coverage of these intervals is surprisingly good over a wide range of possible parameter values for important classes of problems, in particular whenever there are additional nuisance parameters with statistical or systematic errors. Programs are available for calculating these intervals.

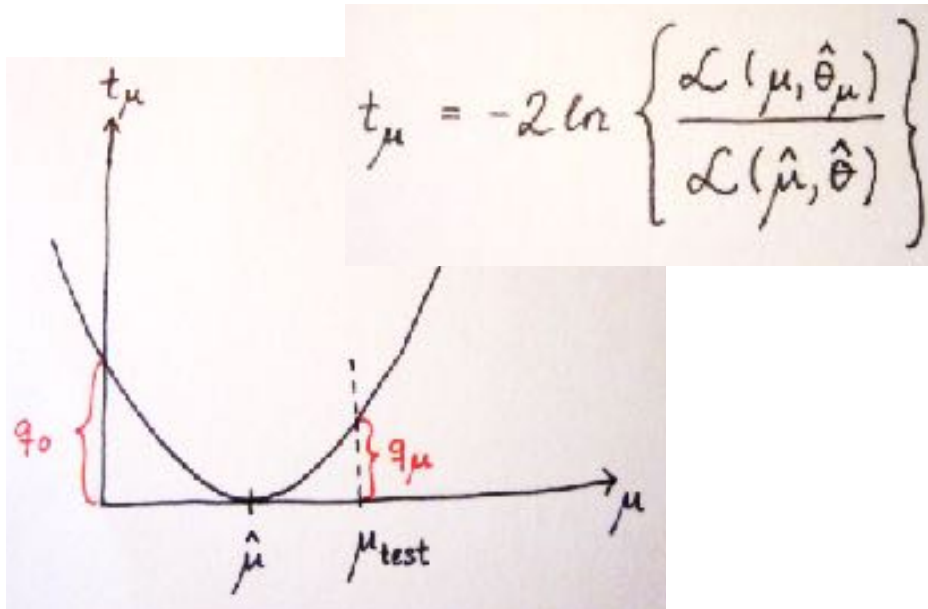# Curiousity: PL considered at LEP times



- I abandoned it to avoid 2-sided intervals (Feldman&Cousins!) - don't want to exclude if there is a nice fat excess!

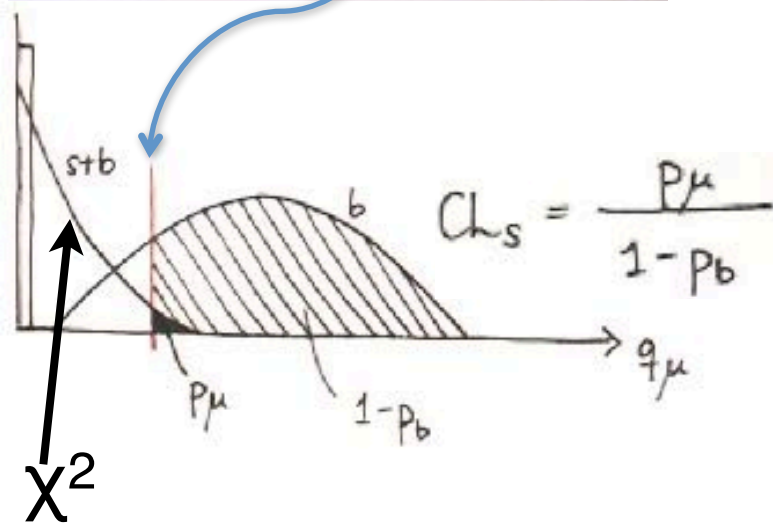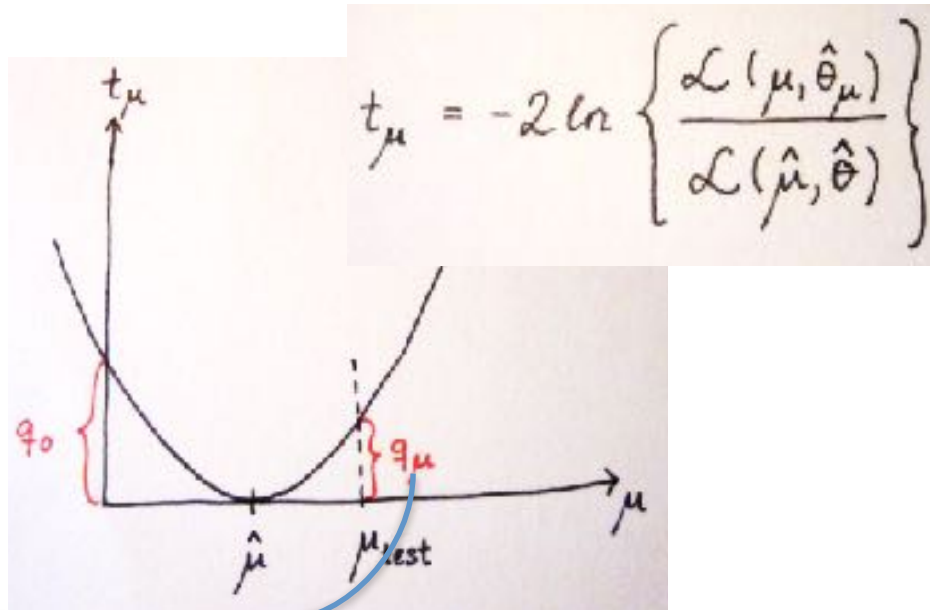- ~10 years later CCGV elegant solution:

$$q_\mu = \begin{cases} -2\ln\lambda(\mu) & \hat{\mu} \le \mu \\ 0 & \hat{\mu} > \mu \end{cases}$$

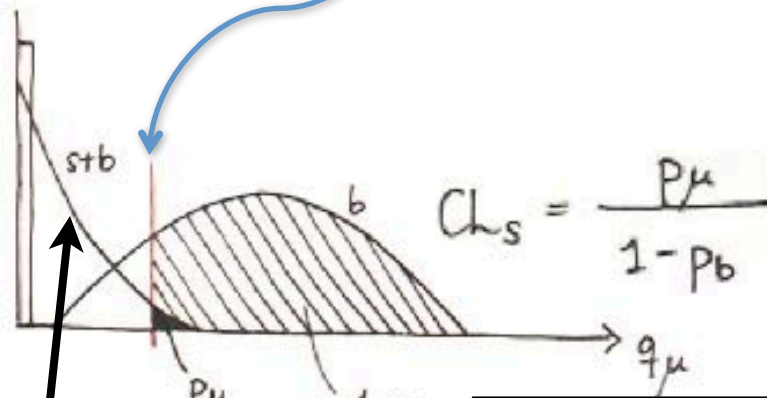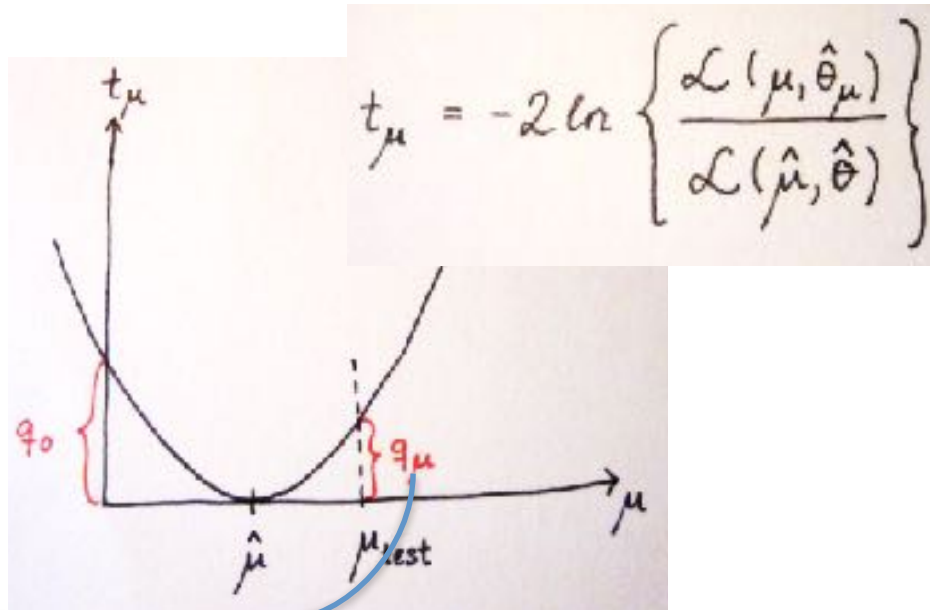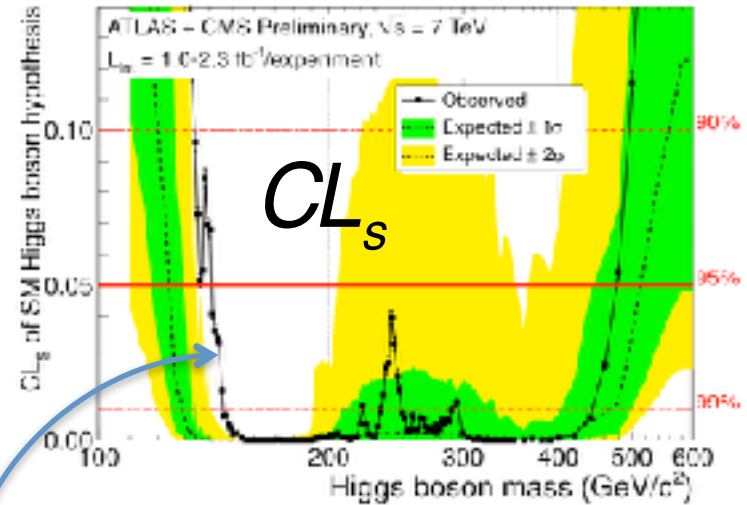# Profile likelihood ratio: CL$_s$ and $\mu_{95}^{up}$



$$t_\mu = -2 \ln \left\{ \frac{\mathcal{L}(\mu, \hat{\theta}_\mu)}{\mathcal{L}(\hat{\mu}, \hat{\theta})} \right\}$$

$\chi^2$

# Profile likelihood ratio: $CL_s$ and $\mu_{95}^{up}$

$$t_\mu = -2 \ln \left\{ \frac{\mathcal{L}(\mu, \hat{\hat{\theta}}_\mu)}{\mathcal{L}(\hat{\mu}, \hat{\theta})} \right\}$$

$$CL_s = \frac{p_\mu}{1 - p_b}$$

$\chi^2$
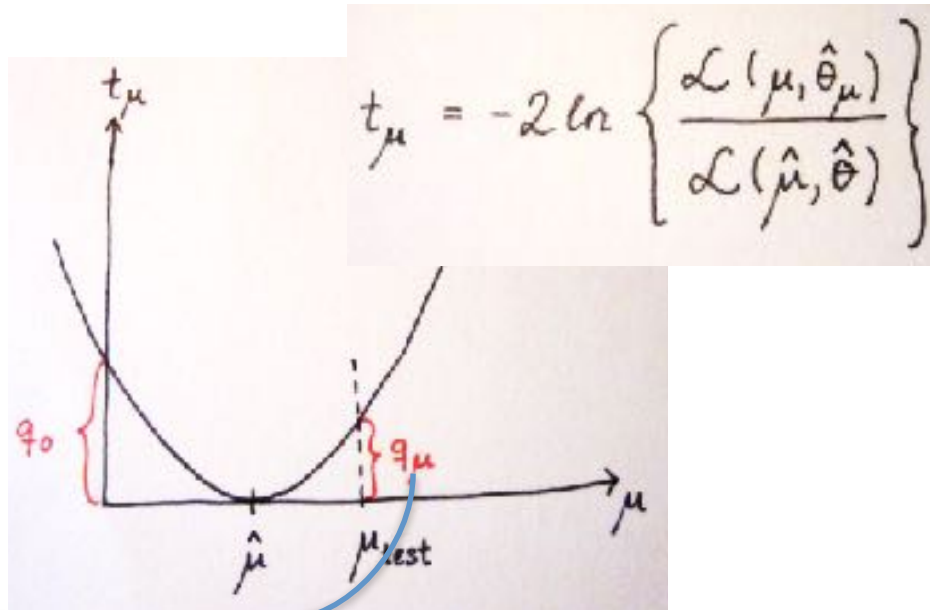
# Profile likelihood ratio: $CL_s$ and $\mu_{95}^{up}$



$$t_\mu = -2 \ln \left\{ \frac{\mathcal{L}(\mu, \hat{\hat{\theta}}_\mu)}{\mathcal{L}(\hat{\mu}, \hat{\theta})} \right\}$$

$CL_s$

$$CL_s = \frac{p_\mu}{1-p_b}$$

$\chi^2$

$p_\mu$ : test signal+background

$CL_s$ : ~test signal

# Profile likelihood ratio: $CL_s$ and $\mu_{95}^{up}$

$$t_\mu = -2 \ln \left\{ \frac{\mathcal{L}(\mu, \hat{\hat{\theta}}_\mu)}{\mathcal{L}(\hat{\mu}, \hat{\theta})} \right\}$$

$CL_s$

$$\mu_{95}^{up} = \mu(CL_s = 0.05)$$

$$CL_s = \frac{p_\mu}{1 - p_b}$$

$\chi^2$

p$_\mu$ : test signal+background

CL$_s$ : ~test signal

ATLAS Higgs - Aspen 2012 - A. L. Read

# Profile likelihood ratio: $CL_s$ and $\mu_{95}^{up}$

$$t_\mu = -2 \ln \left\{ \frac{\mathcal{L}(\mu, \hat{\hat{\theta}}_\mu)}{\mathcal{L}(\hat{\mu}, \hat{\theta})} \right\}$$



$CL_s$

$$\mu_{95}^{up} = \mu(CL_s = 0.05)$$

$$CL_s = \frac{p_\mu}{1 - p_b}$$

$\chi^2$

$\mu_{95}^{up}$

$p_\mu$ : test signal+background
$CL_s$ : ~test signal

# Profile likelihood ratio: p₀ and μ̂

[LHCHCG Combination Procedures](#)

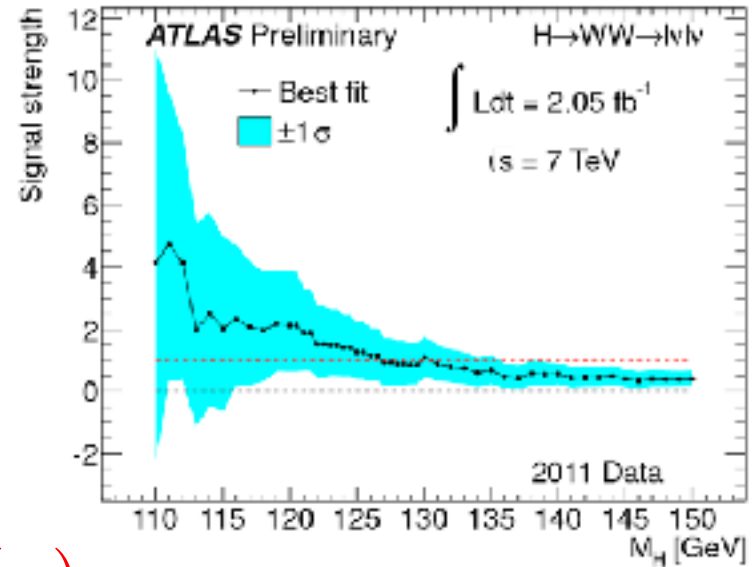$$= -2 \ln \left\{ \frac{\mathcal{L}(\mu, \hat{\hat{\theta}}_\mu)}{\mathcal{L}(\hat{\mu}, \hat{\theta})} \right\}$$



$$\text{P.S.} \, q_{LEP}(\mu) = q_\mu - q_0$$

$\chi^2$

# Profile likelihood ratio: $p_0$ and $\hat{\mu}$

$$t_\mu = -2 \ln \left\{ \frac{\mathcal{L}(\mu, \hat{\hat{\theta}}_\mu)}{\mathcal{L}(\hat{\mu}, \hat{\theta})} \right\}$$



$$\text{P.S.} \, q_{LEP}(\mu) = q_\mu - q_0$$

$\chi^2$

# Profile likelihood ratio: $p_0$ and $\hat{\mu}$

LHCHCG Combination Procedures

$$t_\mu = -2 \ln \left\{ \frac{\mathcal{L}(\mu, \hat{\hat{\theta}}_\mu)}{\mathcal{L}(\hat{\mu}, \hat{\theta})} \right\}$$
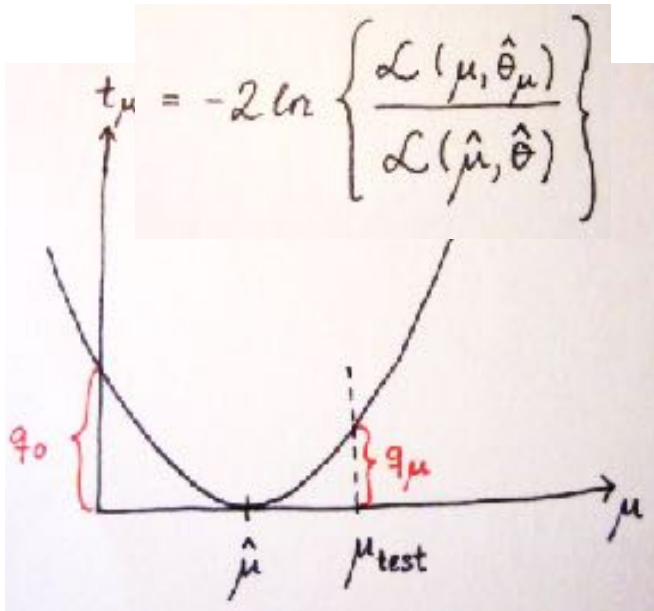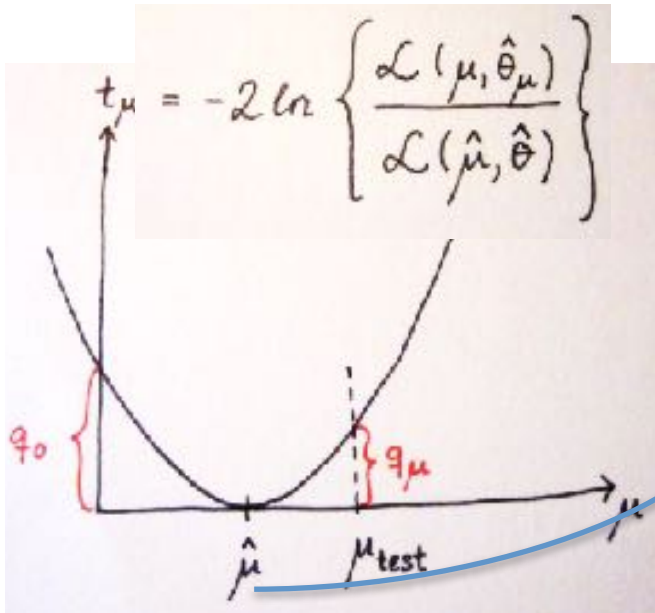
$\hat{\mu}$ to estimate signal strength



ATLAS Preliminary    H→WW→lνlν
Best fit    $\int L dt = 2.05$ fb$^{-1}$
±1σ    $\sqrt{s} = 7$ TeV
2011 Data

$\hat{\mu}$

P.S. $q_{LEP}(\mu) = q_\mu - q_0$

$\chi^2$
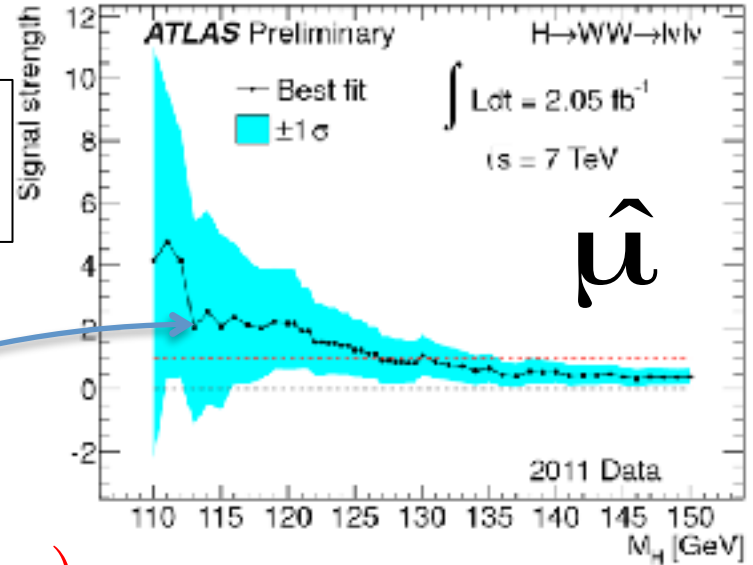
# Profile likelihood ratio: $p_0$ and $\hat{\mu}$

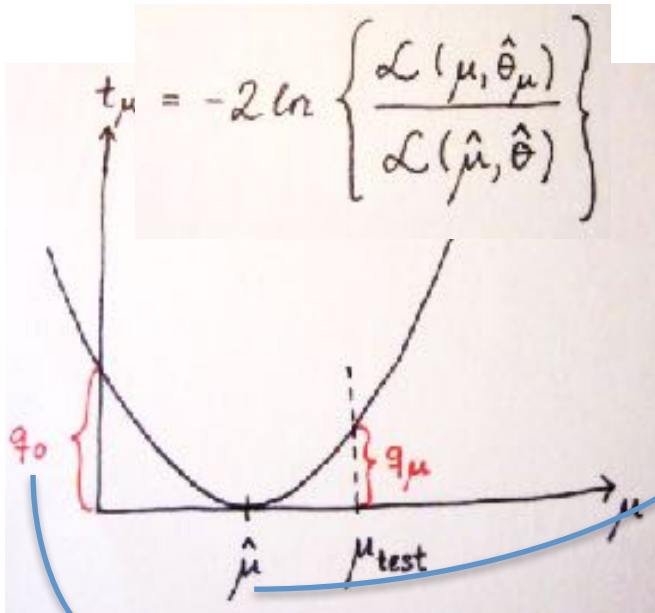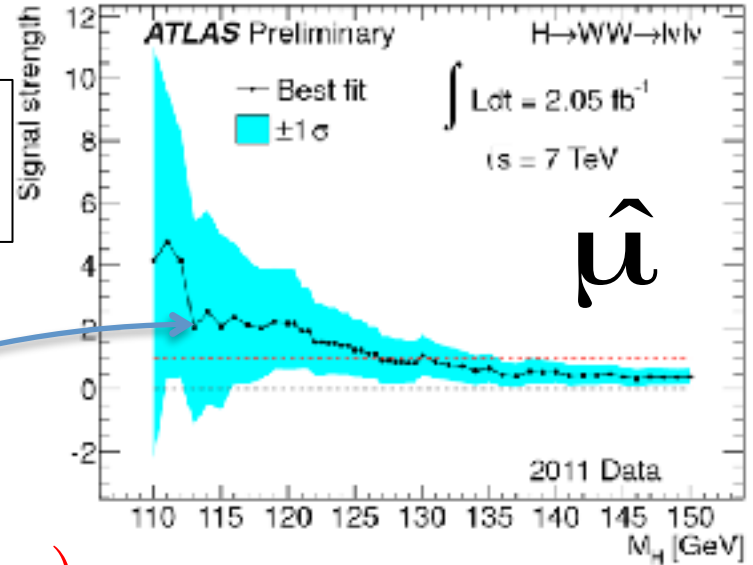LHCHCG Combination Procedures



$\hat{\mu}$ to estimate signal strength

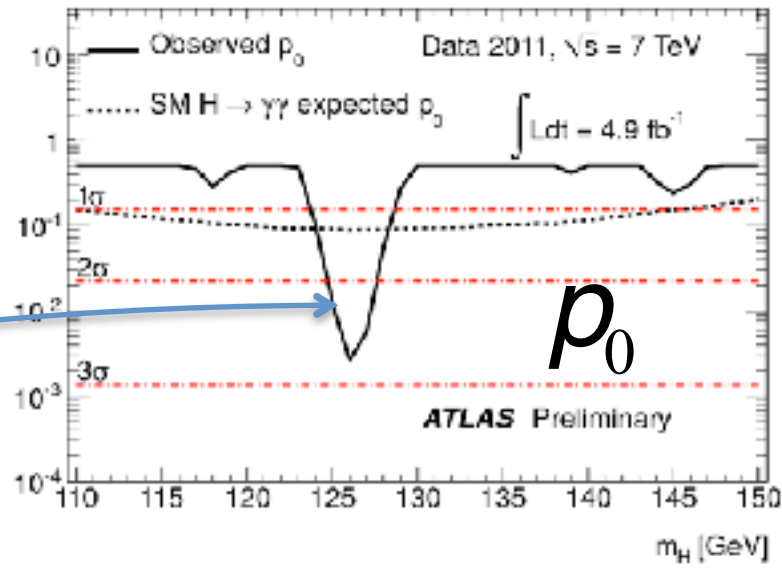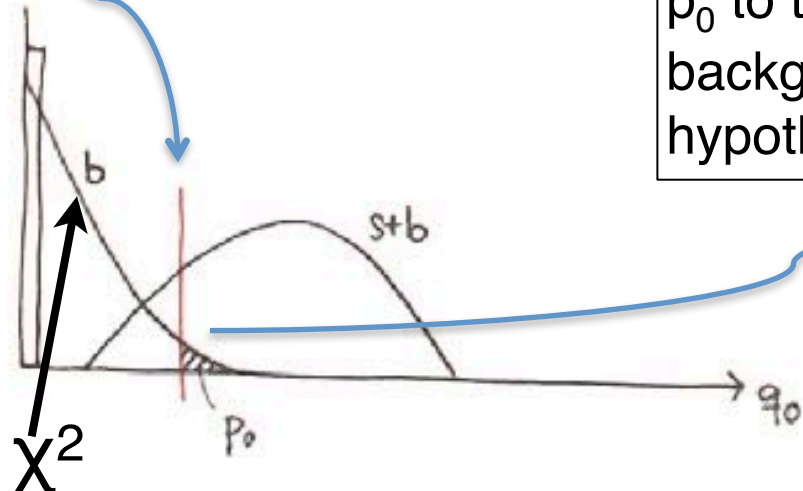$$t_\mu = -2\ln\left\{\frac{\mathcal{L}(\mu, \hat{\hat{\theta}}_\mu)}{\mathcal{L}(\hat{\mu}, \hat{\theta})}\right\}$$

$\hat{\mu}$

P.S. $q_{LEP}(\mu) = q_\mu - q_0$

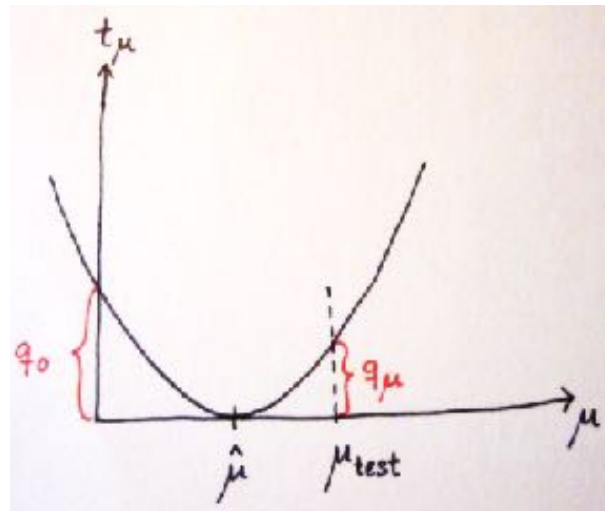$p_0$ to test background hypothesis

$p_0$

# Combined Results

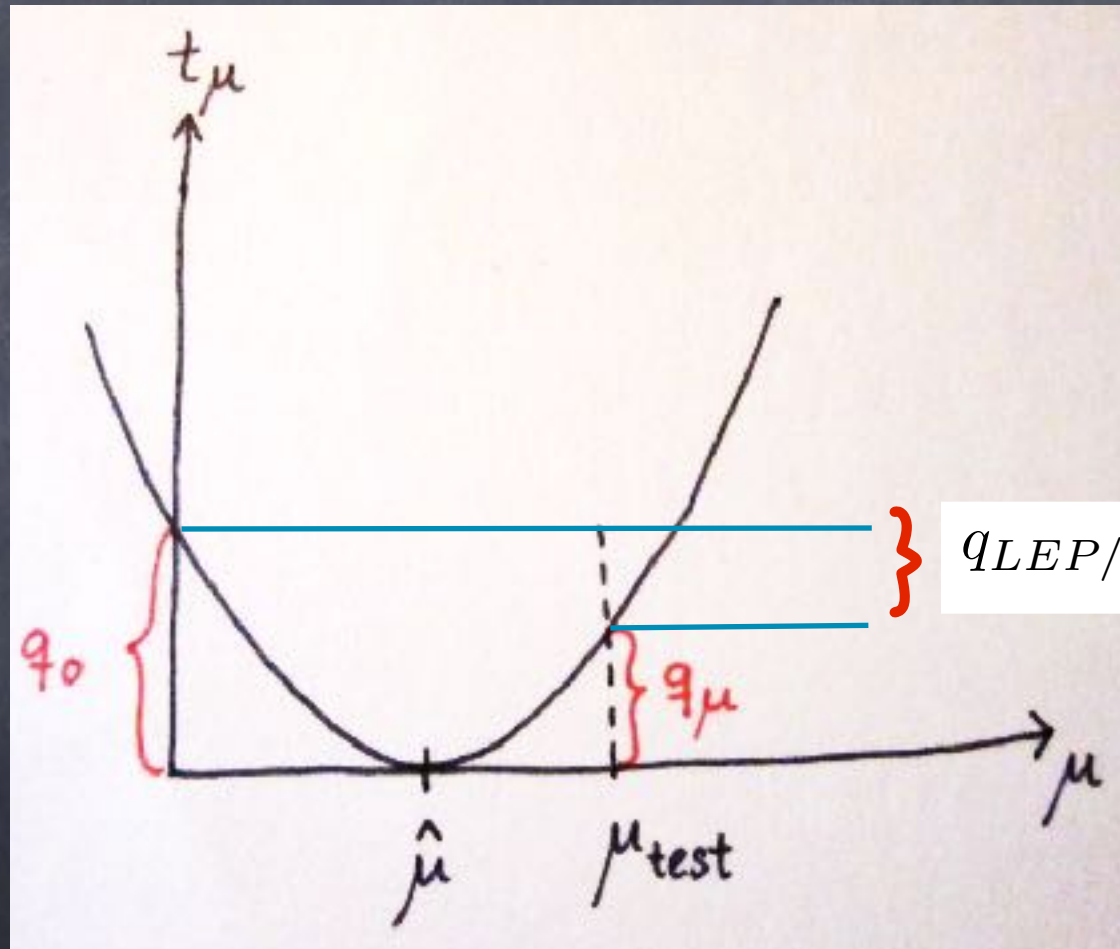$$L(m_H, \mu, \vec{\vartheta}) = \prod_i L_i(m_H, \mu, \vec{\vartheta}_i)$$

$$\bar{x} = \frac{\sum\limits_{i=1}^{n} x_i/\sigma_i^2}{1/\sigma^2}$$

$$\frac{1}{\sigma^2} = \frac{1}{\sum\limits_{i=1}^{n} 1/\sigma_i^2}$$

$$t_\mu = -2 \ln \left\{ \frac{\mathcal{L}(\mu, \hat{\theta}_\mu)}{\mathcal{L}(\hat{\mu}, \hat{\theta})} \right\}$$

# Q<sub>LEP</sub> (Q<sub>TeV</sub> w/o nuisances)



$$q_{LEP/TeV} = q_\mu - q_0$$

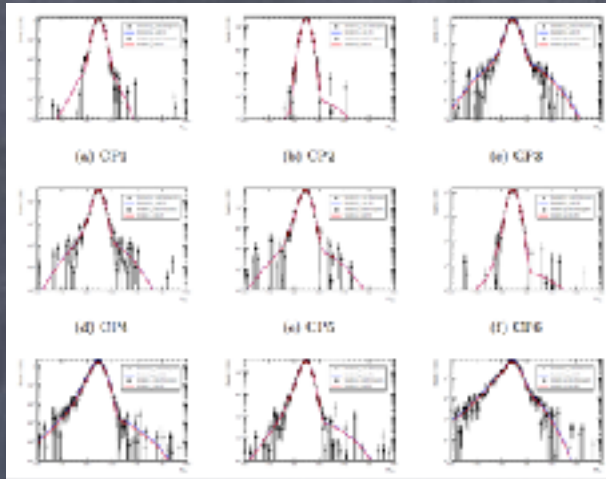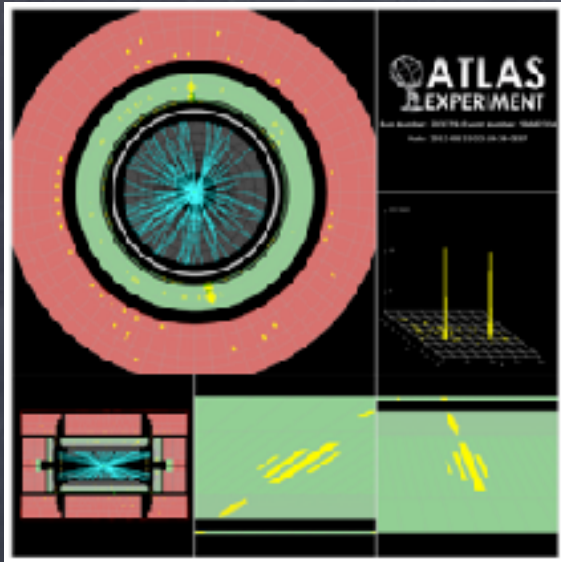# Importance of nuisance parameters



- **background**, uncertainty, uncertainties among most frequent words in ATLAS Higgs boson discovery paper

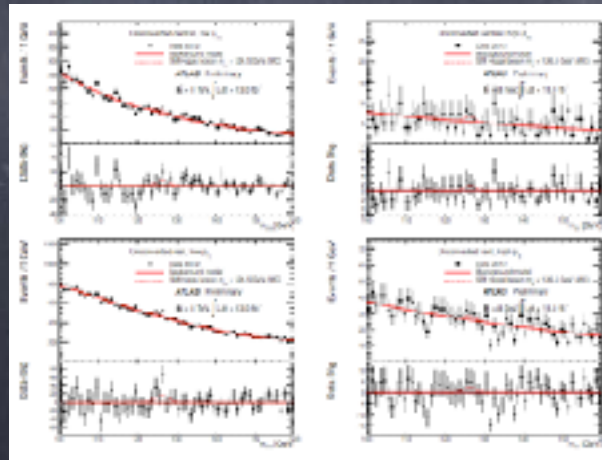Parameterized signal and/or background
models
e.g. ATLAS H->$\gamma\gamma$ search

# 9 categories of unbinned likelihood



Parameterized signal model from fits to MC

Background model: selected functions with unconstrained nuisance parameters

| Name | Criteria | | |
|------|----------|---|---|
| CP1 | unconverted | central | low $p_{Tt}$ |
| CP2 | unconverted | central | high $p_{Tt}$ |
| CP3 | unconverted | non-central | low $p_{Tt}$ |
| CP4 | unconverted | non-central | high $p_{Tt}$ |
| CP5 | converted | central | low $p_{Tt}$ |
| CP6 | converted | central | high $p_{Tt}$ |
| CP7 | converted | non-central | low $p_{Tt}$ |
| CP8 | converted | non-central | high $p_{Tt}$ |
| CP9 | converted | transition | |

4/9 categories

26

# Various terms in L

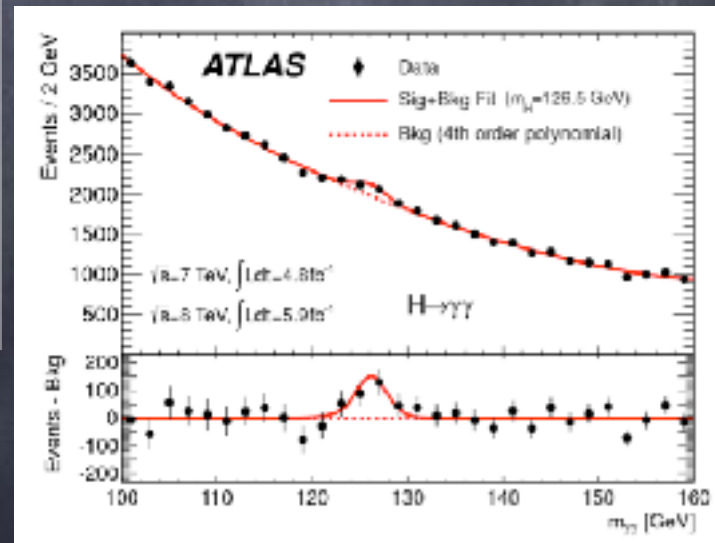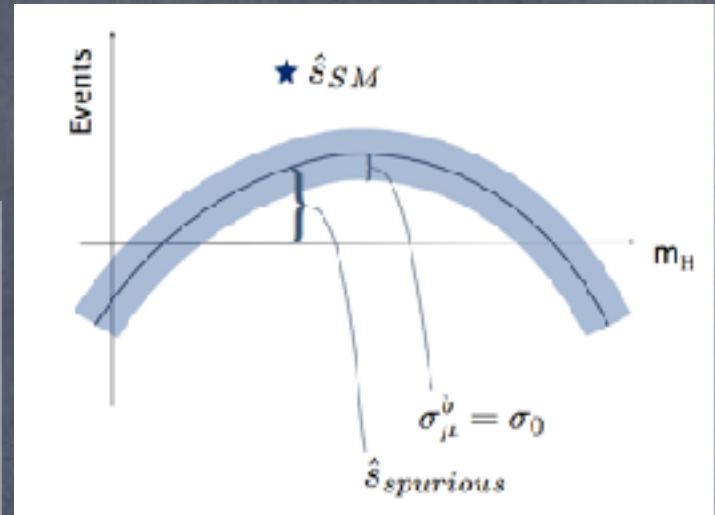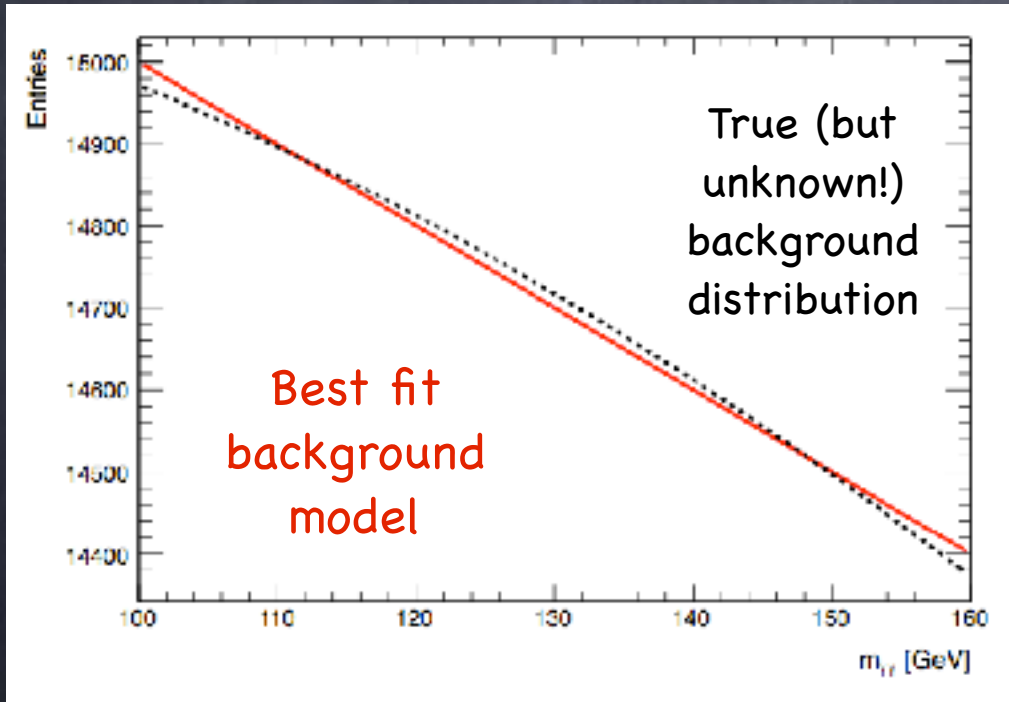$$\mathcal{L}_c(\mu, \boldsymbol{\theta}_c) = e^{-N_c} \prod_{n=1}^{N_c} \mathcal{L}_{c,n}(m_{\gamma\gamma}(n); \mu, \boldsymbol{\theta}_c)$$

L per event in a category

$$\mathcal{L}_{c,n}(m_{\gamma\gamma}(n); \mu, \boldsymbol{\theta}_c) = N_{s,c}(\mu, \boldsymbol{\theta}_c^{norm}) f_{s,c}(m_{\gamma\gamma}; \boldsymbol{\theta}_c^{shape})$$

Mass distribution $+ N_{bkg,c} f_{bkg,c}(m_{\gamma\gamma}; \boldsymbol{\theta}_c^{bkg})$ ,

$$
\begin{aligned}
N_{s,c}(\mu, \boldsymbol{\theta}_c^{norm}) = {} & \mu [N_c^{ggH,SM}(\boldsymbol{\theta}_c^{ggH}) + N_c^{VBF,SM}(\boldsymbol{\theta}_c^{VBF}) \\
& + N_c^{WH,SM}(\boldsymbol{\theta}_c^{WH}) + N_c^{ZH,SM}(\boldsymbol{\theta}_c^{ZH}) + N_c^{ttH,SM}(\boldsymbol{\theta}_c^{ttH})] \\
& \cdot K_{BR}(\theta_{BR}) K_{lumi}(\theta_{lumi}) K_{eff}(\theta_{eff}) K_{isol}(\theta_{isol}) \\
& K_{pile-up}(\theta_{pile-up}) K_{EScale}(\theta_{EScale}) \\
& K_{pile-up,c}(\theta_{pile-up,c}) K_{mat,c}(\theta_{mat}) \\
& + \sigma_{spurious,c} \theta_{spurious,c} .
\end{aligned}
$$

(8.12)

Signal normalization

# Distinguish signal from spurious signal



True (but unknown!) background distribution

Best fit background model

$\star \; \hat{s}_{SM}$

$\sigma^{\hat{\nu}}_{\mu} = \sigma_0$

$\hat{s}_{spurious}$

# Model tests (on MC)



- 9 categories

- No CPU time for full simulation

- 3 MC generators, don't expect them to perfectly reproduce the background data

- Select parameterizations which can incorporate shape uncertainty in unconstrained nuisance parameters without producing false signals

# BG model selection



| Category | Function | Max $|S_{sp}|$ (in GeV) | $\%\sqrt{B}$ ($N_S$) | $S_{sp}$ ($m_y$) | $\sigma^{fit}$ | $\sigma^{syst}$ | Pass | Pass-all |
|----------|----------|-------------------------|----------------------|-------------------|----------------|-----------------|------|----------|
| CP1 | Exp | -4.7 (125) | -48 (11) | -35 (14) | 0.79 | -0.35 | | |
| CP1 | Epoly2 | 2.1 (117) | 18 (12) | 13 (16) | 0.70 | 0.13 | ✓ | ✓ |
| CP2 | Exp | -0.93 (110) | -18 (11.8) | -6.4 (3.5) | 0.43 | -0.064 | ✓ | ✓ |
| CP3 | Exp | 13 (117) | 50 (23) | 35 (33) | 0.71 | 0.35 | | |
| CP3 | Epoly2 | 9.2 (112) | 41 (23) | 26 (30) | 0.64 | 0.26 | | |
| CP3 | Epoly3 | 3.4 (111) | 15 (22) | 8.8 (38) | 0.59 | 0.088 | ✓ | ✓ |
| CP3 | Bern3 | 5.8 (111) | 26 (22) | 18 (36) | 0.62 | 0.16 | | |
| CP3 | Bern4 | 2.9 (111) | 13 (22) | 7.1 (40) | 0.56 | 0.071 | ✓ | ✓ |
| CP4 | Exp | 3.5 (132) | 19 (2.0) | 7.2 (0.9) | 0.39 | 0.072 | ✓ | ✓ |
| CP5 | Exp | -4.4 (125) | -64 (5.8) | -34 (13) | 0.64 | -0.34 | | |
| CP5 | Epoly2 | 1.8 (117) | 22 (7.4) | 10 (16) | 0.47 | 0.10 | ✓ | ✓ |
| CP6 | Exp | 0.57 (110) | 27 (1.48) | 8.6 (2.4) | 0.28 | 0.060 | ✓ | ✓ |
| CP7 | Exp | 5.5 (122) | 29 (23) | 18 (37) | 0.60 | 0.17 | | |
| CP7 | Epoly2 | 5.8 (122) | 26 (22) | 14 (40) | 0.56 | 0.14 | | |
| CP7 | Epoly3 | -6.1 (110) | -29 (22) | -13 (46) | 0.46 | -0.13 | ✓ | ✓ |
| CP7 | Bern3 | -6.3 (110) | -29 (22) | -14 (46) | 0.47 | -0.14 | ✓ | ✓ |
| CP7 | Bern4 | -4.6 (110) | -41 (24) | -2.6 (50) | 0.43 | -0.068 | ✓ | ✓ |
| CP8 | Exp | 0.45 (134) | 18 (9.2) | 5.7 (7.9) | 0.19 | 0.057 | ✓ | ✓ |
| CP9 | Exp | -16 (130) | -179 (9.1) | -69 (28) | 0.33 | -0.39 | | |
| CP9 | Epoly2 | -3.2 (110) | -33 (10.9) | -6.3 (30) | 0.26 | -0.063 | ✓ | ✓ |

- the exponential function

$$Ne^{-\vartheta m_{\gamma\gamma}},\qquad (3.25)$$

where $N$ and $\vartheta$ were the fitted parameters – the normalization and slope of the exponential, respectively;

- the exponential polynomial of order $n$ (orders 2 and 3 were used)

$$e^{\sum_{i=0}^{n}\delta_i m_{\gamma\gamma}^i},\qquad (3.26)$$

where $\delta_i$ were the fitted parameters. Note that the latter $i$ is not an index, but the power $m_{\gamma\gamma}$ is raised to. The normalization, $N$, is described by the first term, $e^{\delta_0}$;
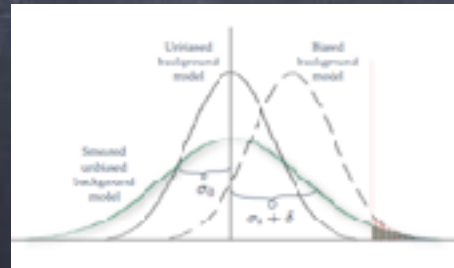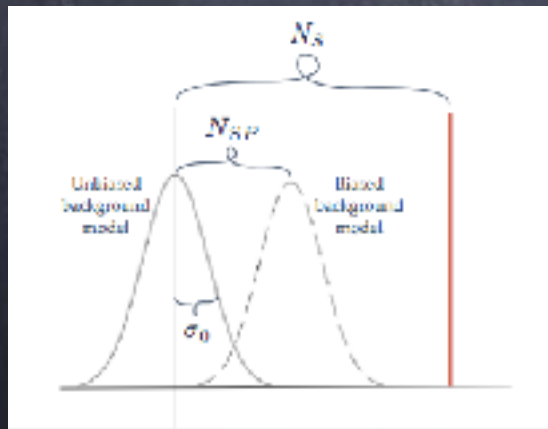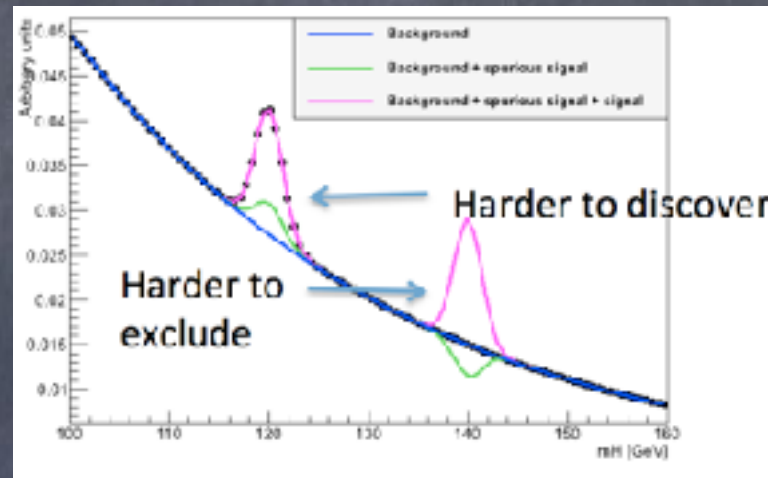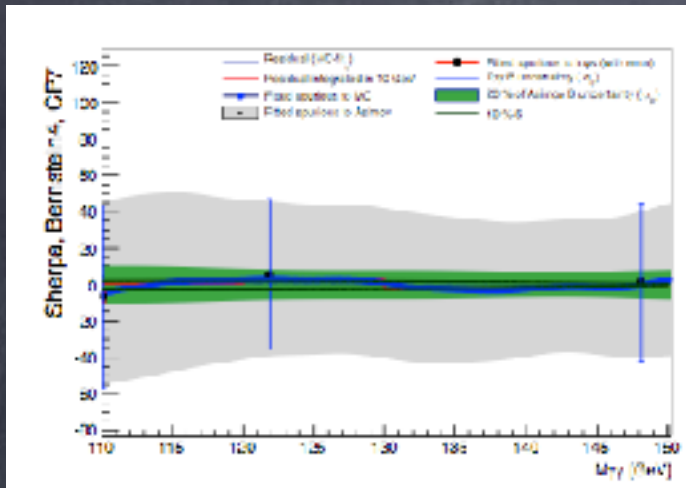
- the Bernstein polynomial of order $n$ (orders 3–7 were used)

$$b_n(t) = \sum_{i=0}^{n} \beta_i \binom{n}{i} t^i (1-t)^{n-i},\qquad (3.27)$$

where $t = \frac{m_{\gamma\gamma}[GeV]-110}{60}$, and where $\beta_i$ were fitted parameters.

Maximum spurious
signal amplitude

# Residual (unknown!) bias: Spurious signal term in likelihood





$$\chi^2 = \frac{(n - (\mu + \delta))^2}{\sigma^2} + \frac{\delta^2}{\sigma_s^2}$$

$$\hat{\mu} = n, \delta = 0$$

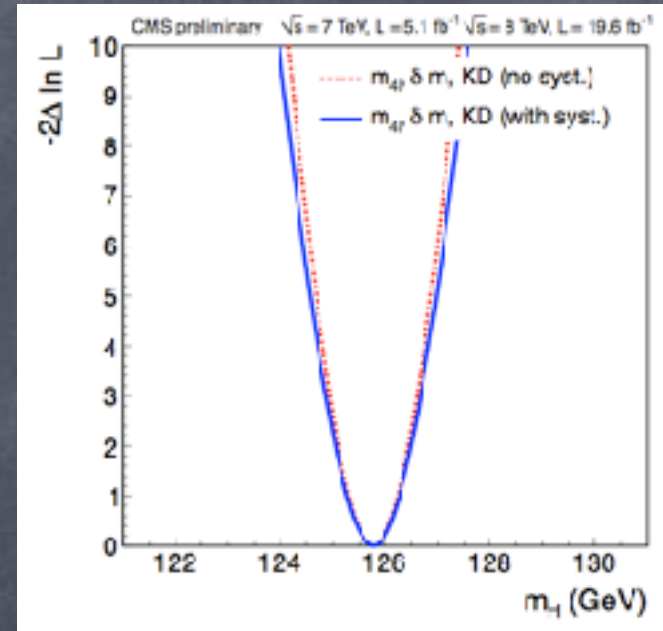$$\sigma_\mu = \sqrt{\sigma^2 + \sigma_s^2}$$

# Nuisance parameters

- NP's broaden the likelihood profile for the parameter of interest



$$\frac{\partial \chi^2}{\partial \delta} = 0$$

$$\frac{\partial \chi^2}{\partial \mu} = 0$$

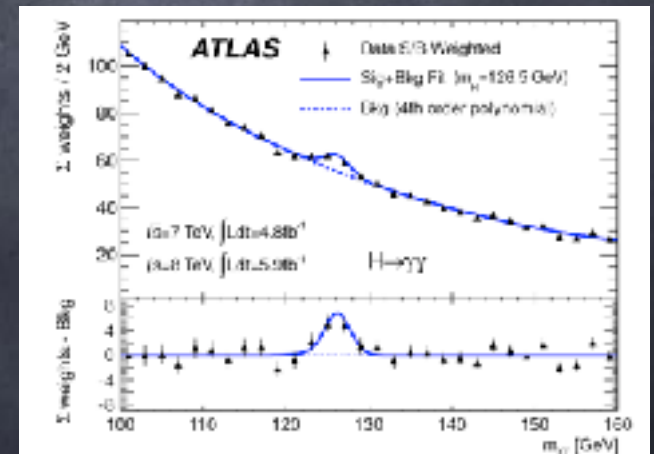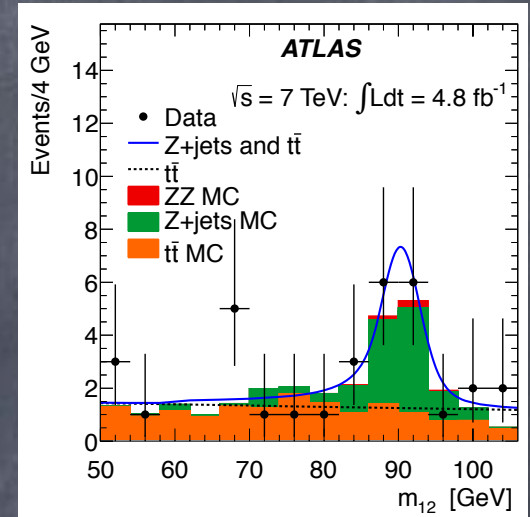$$\frac{1}{\sigma_\mu^2} = \frac{1}{2}\frac{\partial^2 \chi^2}{\partial \mu^2}$$

$$\chi^2 = \frac{(n - (\mu + \delta))^2}{\sigma^2} + \frac{\delta^2}{\sigma_s^2}$$

$$\hat{\mu} = n, \hat{\delta} = 0$$

$$\sigma_\mu = \sqrt{\sigma^2 + \sigma_s^2}$$

# Nuisance parameters

- Parameters fitted directly to the data but no real interest

  - E.g. parametric background; both shape and normalization uncertainty

- Parameters from external estimates that incorporate systematic uncertainty

  - E.g. luminosity, signal theory, mass resolution, electron, muon and jet energy scales
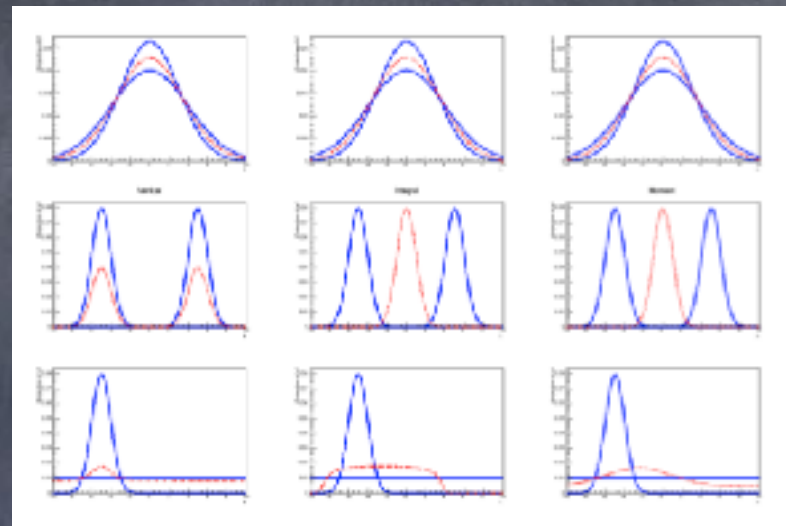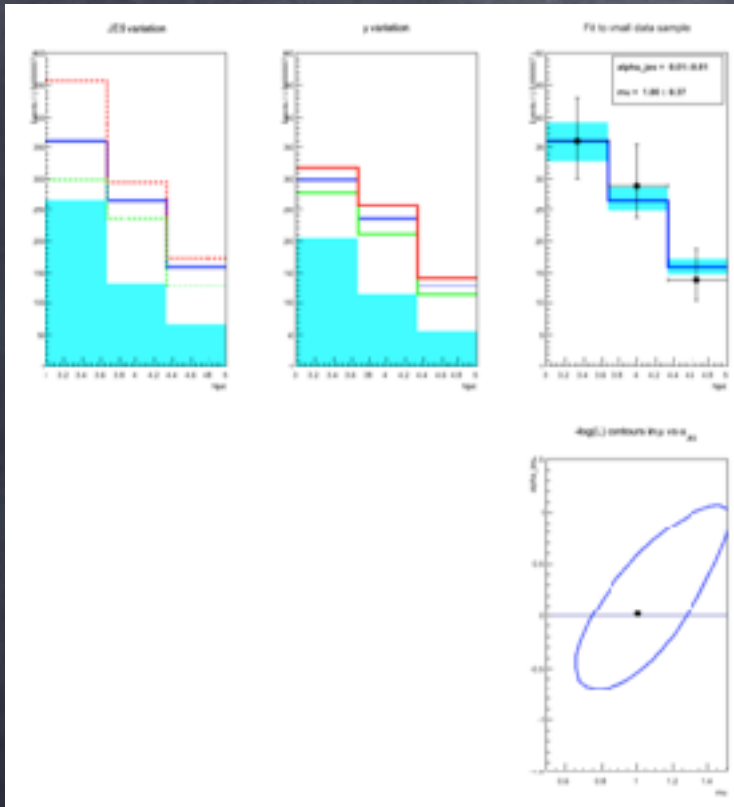
# Constraining nuisance parameters with data

- In the profile likelihood priors are implemented as constraints with external pseudo-measurements (which in many but not all cases are real measurements).

$$L(data \mid \mu, \theta) = Poisson(data \mid \mu s(\theta) + b(\theta)) \times p(\theta \mid \tilde{\theta})$$

Signal region main measurement      Control region auxiliary measurement

- If signal and background are ambiguous (e.g. counting events) the constraints (e.g. prior on the background) may break the ambiguity but uncertainty is governed by the constraint/prior.

- If there is a contraint/prior but the signal and background are NOT ambiguous (e.g. there is a mass or ionization distribution which partially discriminates between them) the uncertainty is reduced by the information (via the fit) in the data.

# Shape systematics





- Don't always have parameterized shape

- Interpolate between templates (interpolation distance is nuisance parameter)

- Various interpolation strategies in ROOT, tradeoff between speed and accuracy (and sometimes unintended consequences)

Linear interpolation of histograms

A.L. Read[1]

University of Oslo, Department of Physics, P.O. Box 1048, Blindern

http://dx.doi.org/10.1016/S0168-9002(98)01517-3, How to Cite or Li

NUCLEAR
INSTRUMENTS
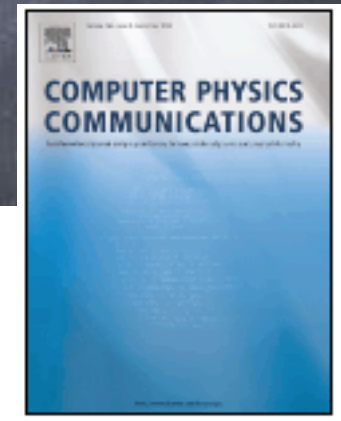& METHODS
IN
PHYSICS
RESEARCH

# MC statisics

- In HEP the simulations tend to be computationally expensive – limited MC statistics is sometimes a real issue.

- Put a Poisson term (nuisance) on each bin. The higher the MC stats the more this will constrain the shape to the predicted shape. If the statistics are poor the data will constrain the background shape at the cost of reduced sensitivity to the signal (i.e. higher uncertainty).

- Usually based on Barlow-Beeston:

## Fitting using finite Monte Carlo samples

Roger Barlow , Christine Beeston

Department of Physics, Manchester University, Manchester M13 9PL, UK

http://dx.doi.org/10.1016/0010-4655(93)90005-W, How to Cite or Link Using DOI

COMPUTER PHYSICS COMMUNICATIONS

# AA - Asymptotics and Asimov dataset



arXiv.org > physics > arXiv:1007.1727

**Physics > Data Analysis, Statistics and Probability**

## Asymptotic formulae for likelihood-based tests of new physics

Glen Cowan, Kyle Cranmer, Eilam Gross, Ofer Vitells

(Submitted on 10 Jul 2010 (v1), last revised 3 Oct 2010 (this version, v2))

We describe likelihood-based statistical tests for use in high energy physics for the discovery of new phenomena and for construction of confidence intervals on model parameters. We focus on the properties of the test procedures that allow one to account for systematic uncertainties. Explicit formulae for the asymptotic distributions of test statistics are derived using results of Wilks and Wald. We motivate and justify the use of a representative data set, called the "Asimov data set", which provides a simple method to obtain the median experimental sensitivity of a search or measurement as well as fluctuations about this expectation.

# AA - Asymptotics and Asimov dataset

$$-2\ln\lambda(\mu) = \frac{(\mu-\hat{\mu})^2}{\sigma^2} + \mathcal{O}(1/\sqrt{N})$$

$$V_{ij}^{-1} = -E\left[\frac{\partial^2\ln L}{\partial\theta_i\partial\theta_j}\right]$$

$$\sigma^2 = V_{00}$$

$$q_0 = \begin{cases} \hat{\mu}^2/\sigma^2 & \hat{\mu}\geq 0, \\ 0 & \hat{\mu} < 0, \end{cases}$$

$$f(q_0|0) = \frac{1}{2}\delta(q_0) + \frac{1}{2}\frac{1}{\sqrt{2\pi}}\frac{1}{\sqrt{q_0}}e^{-q_0/2}$$

$$V_{jk}^{-1} = -E\left[\frac{\partial^2\ln L}{\partial\theta_j\partial\theta_k}\right] = -\frac{\partial^2\ln L_{\mathrm{A}}}{\partial\theta_j\partial\theta_k}$$

$$n_{i,\mathrm{A}} = E[n_i] = \nu_i = \mu' s_i(\boldsymbol{\theta}) + b_i(\boldsymbol{\theta}),$$

$$m_{i,\mathrm{A}} = E[m_i] = u_i(\boldsymbol{\theta}).$$

Compact formulae for both observed results and expectations (including fluctuation bands)

38

A. Read, U. Oslo

# Curiosity: Precursor to Asimov dataset in LEP (DELPHI) Higgs combination code

```
      SUBROUTINE explnQnom(s)
*.--------------------------------------------------------------
*   Compute the expected Likelihood Ratio for the combined counting and
*   invariant mass (or other discriminating variable) measurement experiment in
*   multiple channels. This only works for combinations where for each
*   channels the number of background and signal bins is indentical. This
*   is fast and simple to compute and can serve as a precise check
*   of Monte Carlo and semi-analytic computations.
*
*   The expected -2lnQ (Q is likelihood ratio) is computed both for
*   background-only and signal+background hypotheses.
*
*   10.12.99 Add the RMS of the distributions of -2lnQ for signal+background
*            and background-only experiments.
*.--------------------------------------------------------------
```

```
lrwt = log(1. + si*bkgprdx(i)/bi/sigprdx(i))
lnqisb = -si + (si+bi)*lrwt
lnqib  = -si +         bi*lrwt
avg2lnqsb8 = avg2lnqsb8 + lnqisb
avg2lnqb8  = avg2lnqb8  + lnq1b

r2lnqisb = 4.*(si+bi)*lrwt**2
r2lnqib  = 4.*(   bi)*lrwt**2
avgr2lnqsb8 = avgr2lnqsb8 + r2lnqisb
avgr2lnqb8  = avgr2lnqb8  + r2lnqib
```

- But unlike CCGV not possible to treat nuisance parameters

## Combination Details

- At first (one or two combinations), ATLAS results were fully based on toys

- As model grew, these became impractical
  - ~570 nuisance parameters at time of discovery
  - ~310 of these are due to MC stats, treated Barlow-Beeston style

- ~10-30 minutes per fit  ›20-60 minutes per toy
  - $\mathcal{O}$(millions) CPU hours to produce full result

# Data-driven methods

- HEP depends heavily on Monte Carlo calculations of physics processes and detector response for both signals (known and hypothetical) and backgrounds.

- Sometime we just don't know and/or have reason not to trust the MC results.

- Various data-driven methods used to estimate background in signal region.

  - Fits (unbinned, many bins) with sidebands

  - Variations of "on-off": ABCD, Matrix method, fit to shapes derived from well-understood (signal-free) control regions, …

# Other data-driven methods (ABCD)
# (Variations of on-off and sideband fits)

- **K**nown small backgrounds (e.g. electroweak processes):

$$\mu^K_{A,B,C,D}$$

- Poorly known ("**U**nknown") backgrounds:



Signal region

A: $\mu^U$

B: $\mu^U \tau_B$

C: $\mu^U \tau_C$

D: $\mu^U \tau_B \tau_C$

- Naively:

$$\mu^U = N_c \frac{N_B}{N_D}$$

- Correlations should be accounted for as well...

"Let's write down the likelihood function"

$$L(n_A, n_B, n_C, n_D | \mu, \theta_{\omega}) = \Pi_{i=A,B,C,D} \frac{e^{-\mu_i} \mu_i^{n_i}}{n_i!}$$

$$\mu_A = \mu + \mu_A^K + \mu^U$$

$$\mu_B = b\mu + \mu_B^K + \mu^U \tau_B$$

$$\mu_C = c\mu + \mu_C^K + \mu^U \tau_C$$

$$\mu_D = d\mu + \mu_D^K + \mu^U \tau_B \tau_C$$

# "Matrix method"
## (Variation of Bob Cousins' homework problem)

◉ Suppose you know you have only two particle species in your sample and know how to tag them but don't know the mixture (e.g. pions and electrons).

◉ N = true number, C = Counted by experiment

$$\pi^\pm \to \pi^\pm \qquad r \qquad (\text{"real"})$$
$$\pi^\pm \to e^\pm \qquad 1 - r$$
$$e^\pm \to \pi^\pm \qquad f \qquad (\text{"fake"})$$
$$e^\pm \to e^\pm \qquad 1 - f$$

$$\begin{pmatrix} C_\pi \\ C_e \end{pmatrix} = \begin{pmatrix} r & f \\ 1-r & 1-f \end{pmatrix} \begin{pmatrix} N_\pi \\ N_e \end{pmatrix}$$

$$\begin{pmatrix} N_\pi \\ N_e \end{pmatrix} = \begin{pmatrix} r & f \\ 1-r & 1-f \end{pmatrix}^{-1} \begin{pmatrix} C_\pi \\ C_e \end{pmatrix}$$

◉ Homework: reformulate as 2-bin maximum likelihood (note: Nr. parameters=Nr. measurements - why is this "bad"?)

# Event selection
and
# Multi-variate analysis
(often ML these days)

# Selecting events

Suppose we have a data sample with two kinds of events, corresponding to hypotheses $H_0$ and $H_1$ and we want to select those of type $H_0$.

Each event is a point in $\vec{x}$ space. What decision boundary should we use to accept/reject events as belonging to event type $H_0$?

Probably start with cuts:

$$x_i < c_i$$

$$x_j < c_j$$

# Other ways to select events

Or maybe use some other sort of decision boundary:

linear

or nonlinear



How can we do this in an 'optimal' way?

**Likelihood ratio**
$$Q=L(H0)/L(H1)$$
**or approximation in case of complexity**

# Neyman-Pearson lemma

# Machine learning (wikipedia)

**4 Approaches**

- Worth understanding and learning to use Boosted Decision Tree (BDT) – frequently used in HEP, relatively fast and effective

# Look-elsewhere effect (LEE)

SPECIAL ARTICLE - TOOLS FOR EXPERIMENT AND THEORY

### Trial factors for the look elsewhere effect in high energy physics

Eilam Gross and Ofer Vitells

$$TF = \frac{p_0^{global}}{p_0^{local}}$$

- ⊙ Rule of thumb for trials factor used before LHC

$$TF \sim \frac{\Delta m}{\sigma_m} \; \textbf{?}$$

- ⊙ Eilam and Ofer discovered that trials factor grows with significance Z (ROT ~OK for Z=3)

$$TF \simeq 1 + \sqrt{\frac{\pi}{2}} \mathcal{N} Z$$

# Look-elsewhere effect (LEE)



3 crossings

$$p_0^{\text{global}} \simeq p_0^{\text{local}} + \; < N(q_{\text{ref}}) > e^{-(q - q_{\text{ref}})/2}$$
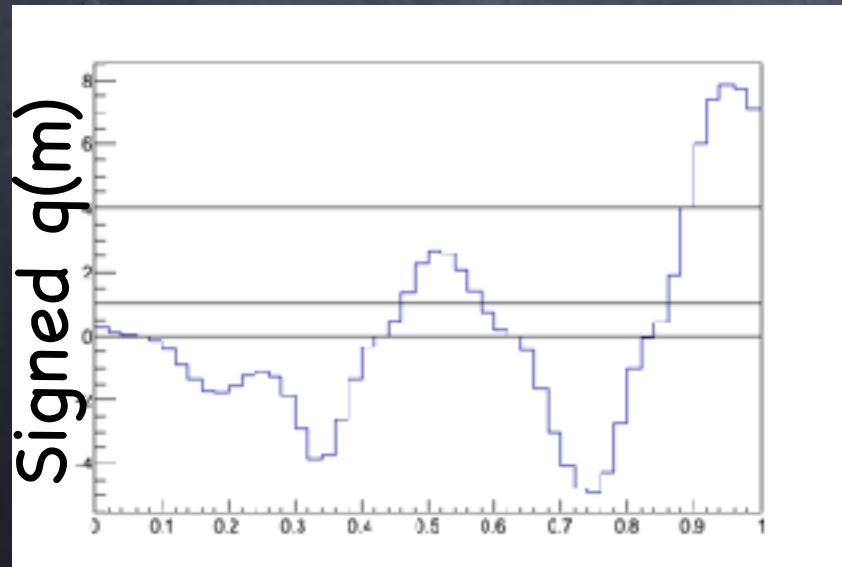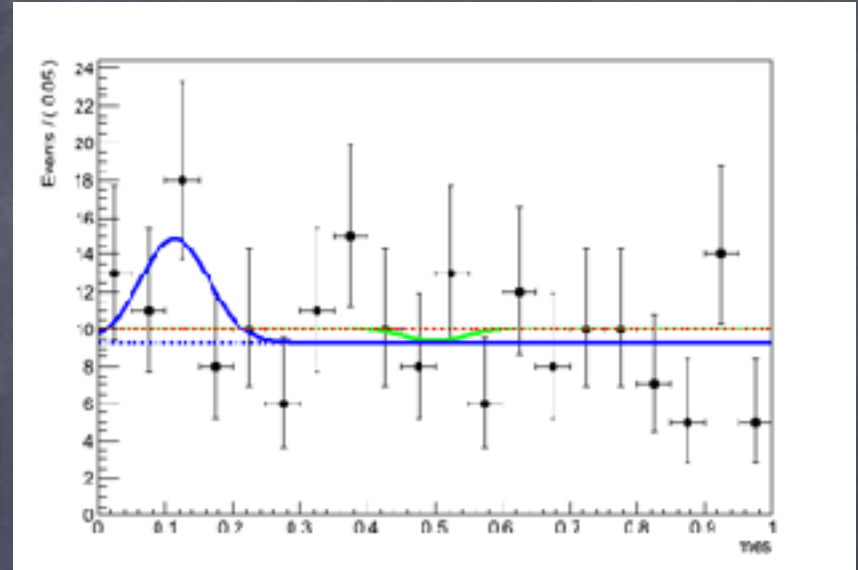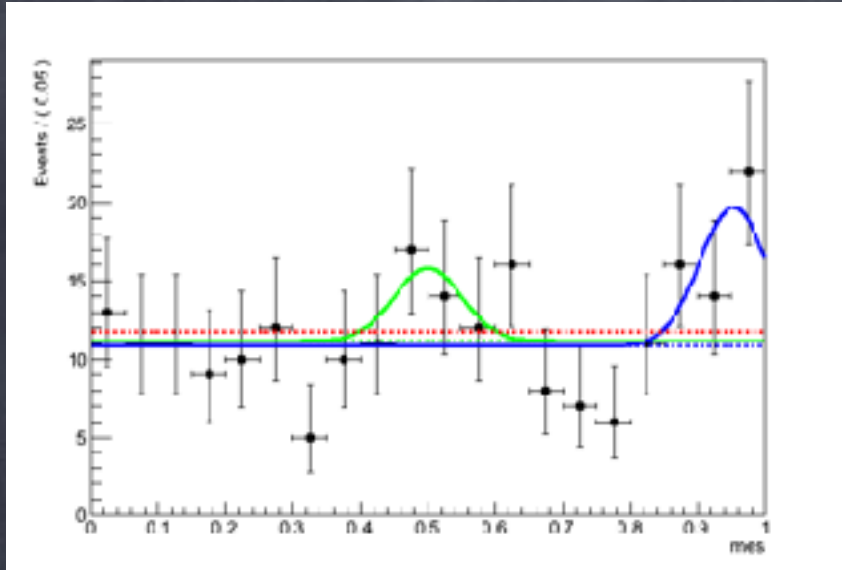
$$TF = \frac{P(q(\hat{m}) > Z^2)}{P(q(m) > Z^2)} \simeq 1 + \mathcal{N} \frac{P(\chi_2^2 > Z^2)}{P(\chi_1^2 > Z^2)}$$



1M fits

# Fit to background toy
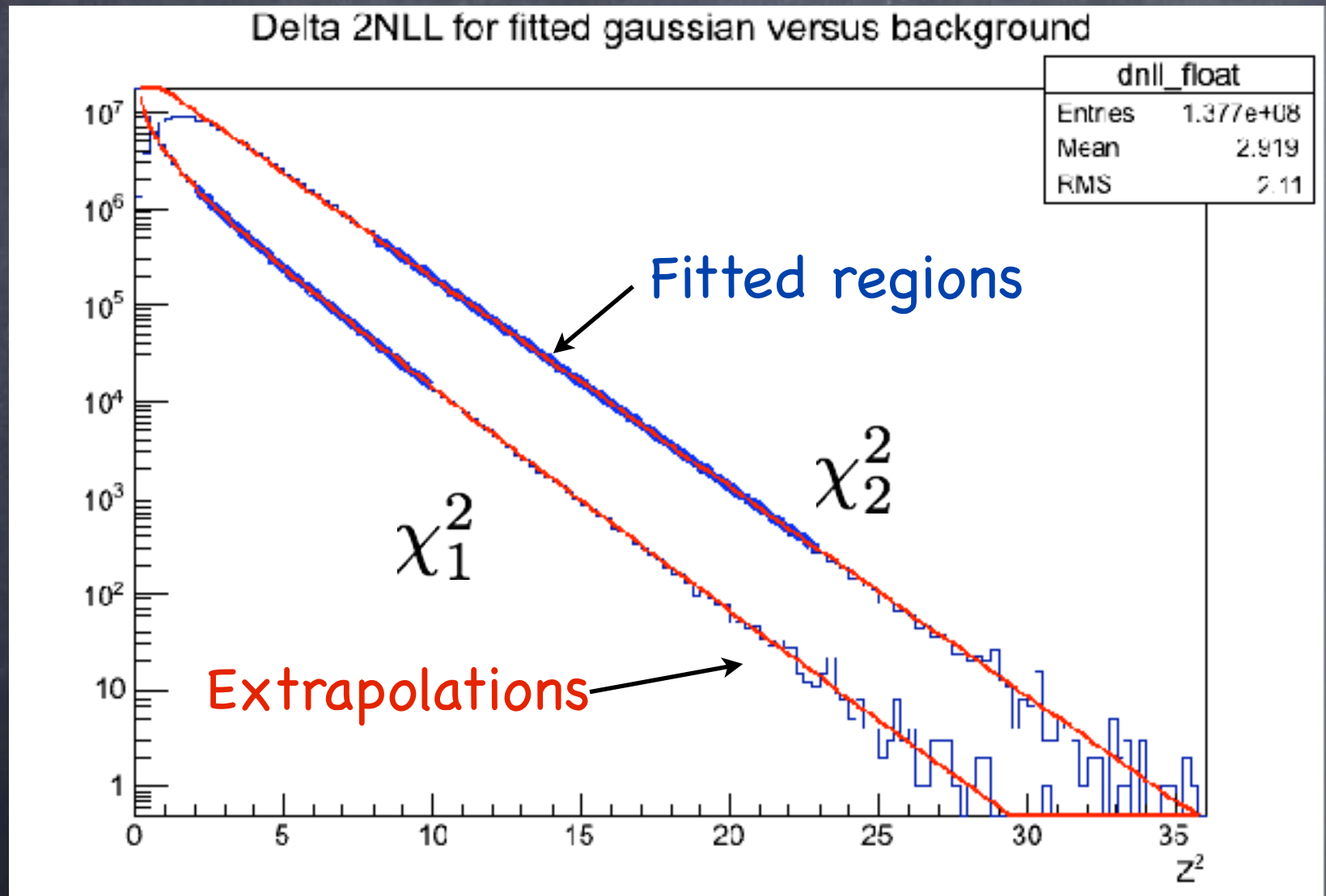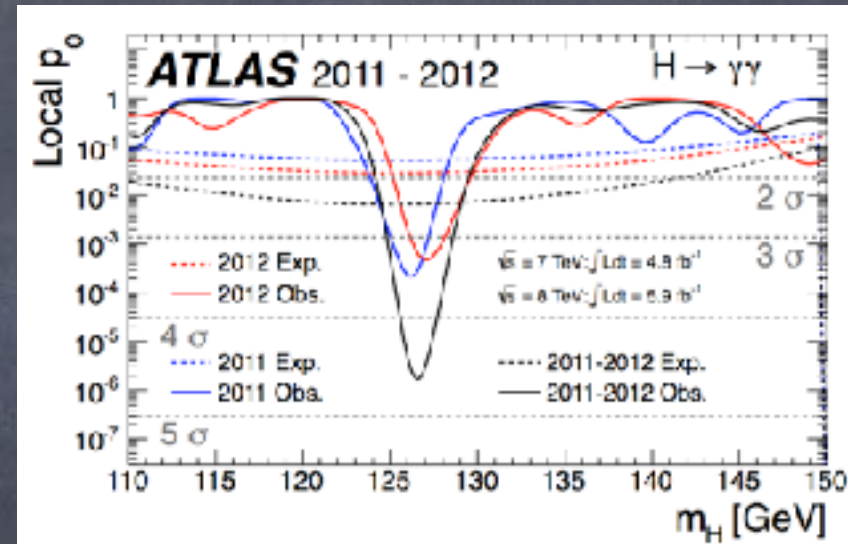
# 2 examples toy fits

52

# 138 Mfits - it all checks out



Delta 2NLL for fitted gaussian versus background

| dnll_float | |
|---|---|
| Entries | 1.377e+08 |
| Mean | 2.919 |
| RMS | 2.11 |

Fitted regions

$\chi_1^2$

$\chi_2^2$
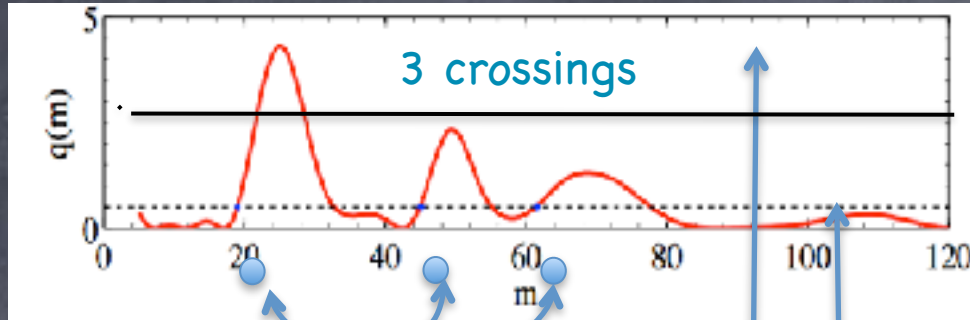
Extrapolations

$z^2$

# Energy (mass) scale systematic uncertainties

- Local look-elsewhere effect when combining channels with different energy (mass) scales, e.g. electrons, photons, muons, jets

- Not accounted for in asymptotic expressions, nor in the classical look elsewhere effect
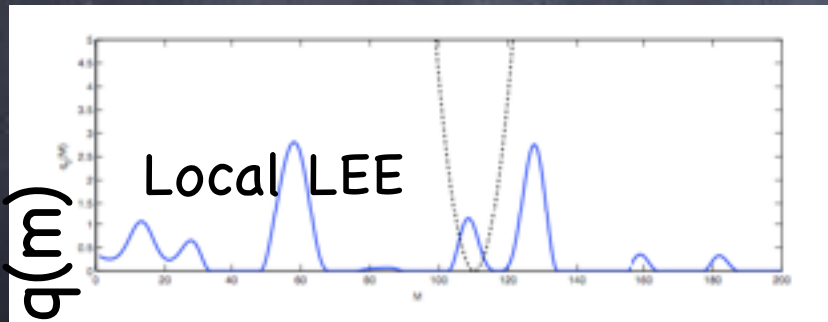


- Illustration: Imagine we had (illegally!) aligned the red and blue curves by hand before combining...

   - i.e. we don't yet use the Higgs boson for detector calibration!!

# Example: 1 uncertain mass scale



3 crossings

Usual LEE

$$p_0^{\text{global}} \simeq p_0^{\text{local}} + < N(q_{\text{ref}}) > e^{-(q-q_{\text{ref}})/2}$$

Local LEE

q(m)

m

$\Delta$ - mass internal

$\sigma_m$ - mass resolution

(u=q=-2ln$\Delta$L)

$$\mathbb{E}[N_u'] \leq \frac{1}{2}\mathbb{P}(\chi^2 > u) + \mathcal{N}_1 e^{-u/2}\frac{\sqrt{2\pi}\sigma_M}{|\Delta|}$$

Leadbetter (1965),
O. Vitells (2012)

P.S. Ofer, please publish your work!!

# Extrapolate ESS correction



q0 distribution

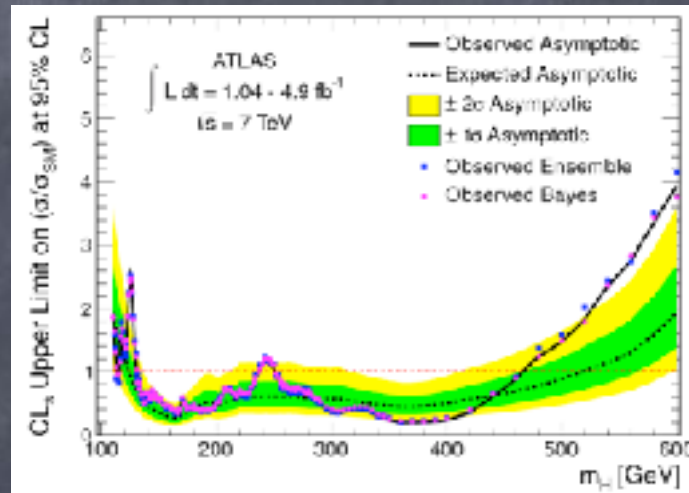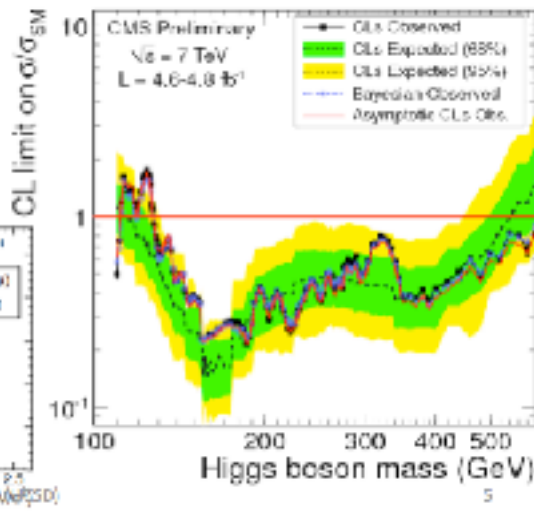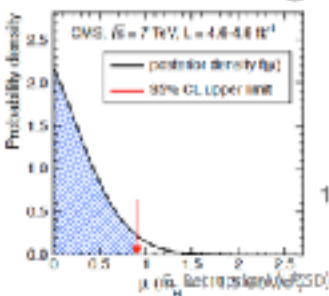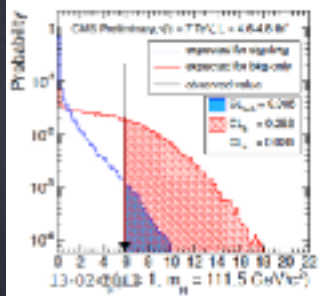| hq0 | |
|---|---|
| Entries | 182798 |
| Mean | 0.8642 |
| RMS | 1.455 |

$$p_0 = (1 - \epsilon)\frac{1}{2}P(\chi^2 > q_0) + \frac{\epsilon}{2}e^{-q_0/2}$$

- In practice, several energy scales

- Don't need $O(10/p_0)$ fits to MC toys to estimate tiny effect!

- Several nuisance parameters, perform empirical fit

- $O(0.1\sigma)$ effect around $5\sigma$ for ATLAS

# What about Bayesian methodology in LHC Higgs boson searches?



- Up to Moriond 2012, CMS produced limits with three prescriptions, to check robustness.
  - CLs using Toy MC
  - CLs using asymptotics
  - Bayesian w/ flat prior

$$L(\mu) = \frac{1}{C} \int_\theta p(\text{data}|\mu s + b) \; \rho_\theta(\theta) \; \pi_\mu(\mu) \; d\theta.$$

$$\int_0^{\mu_{95\% \, CL}} L(\mu) \; d\mu = 0.95.$$

- Limits, with flat prior, very consistent with CLs limits derived in frequentist framework

- No serious attempt (yet) to quantify excess at 125/6 GeV with Bayes factors

# What about Bayesian methodology in LHC Higgs boson searches?

$$L(\mu) = \frac{1}{C} \int_\theta p(\text{data}|\mu s + b) \, \rho_0(\theta) \, \pi_\mu(\mu) \, d\theta.$$

$$\int_0^{\mu_{95\%CL}} L(\mu) \, d\mu = 0.95.$$

- ◉ Louis Lyons and David van Dyk (Statistics, Imperial College) want to analyse Higgs boson discovery in Bayesian framework

  $$B_{12} = \frac{\text{pr}(\mathbf{D}|H_1)}{\text{pr}(\mathbf{D}|H_2)}$$

  - ◉ "Bayesian wear their priors on their sleeves"

  $$\text{pr}(\mathbf{D}|H_k) = \int \text{pr}(\mathbf{D}|\theta_k, H_k)\pi(\theta_k|H_k) \, d\theta_k$$

  - ◉ However, statistical procedures applied to Higgs boson discovery "among the most rigorous of complex scientific data today"
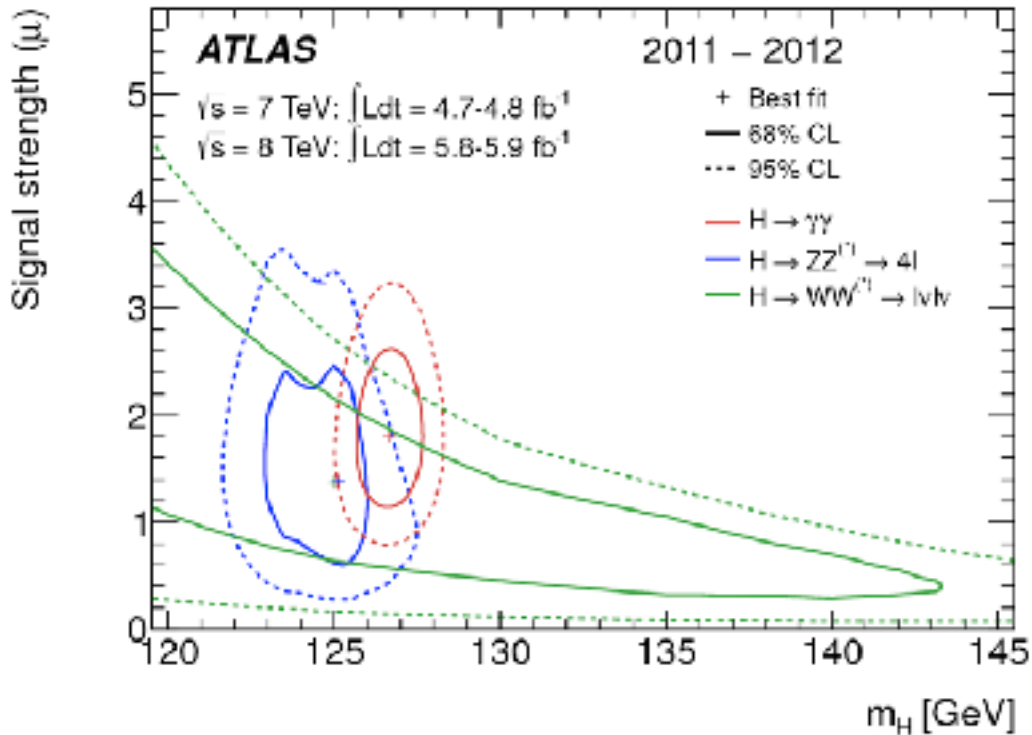
# Limits interpretation

- freq: upper limit on $\mu$ at 95% CL does NOT mean that $P(\mu < \mu_{up}) = 5\%$ ! Only conclusion is that we didn't see anything in the data consistent with $\mu \geq \mu_{up}$ (with a method that is guarantied to be wrong 5% of the time).

- Bayes: upper limit on $\mu$ at 5% (1-95%) of posterior density DOES mean $P(\mu < \mu_{up}) = 5\%$ BUT there is a prior that the physics of $\mu$ exists.

- $CL_s$ has similar interpretation as freq but protected against obvious wrong freq limits for insensitive experiments

    - Cost of robustness is overcoverage (e.g. wrong less than 5% of time for 95% CL)

    - Otherwise many same features as Bayes limits

        - "Lucky" background fluctuations don't give obviously optimistic limits

        - Increased uncertainty doesn't improve a search

        - Adding a low-sensitivity channel hardly improves the search

# From exclusion to discovery to measurement

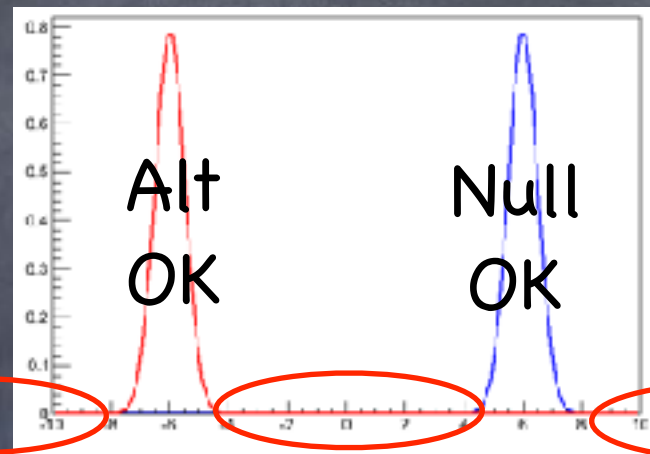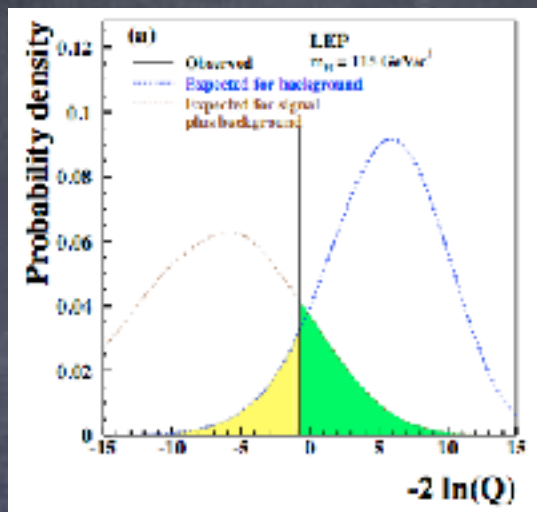Release 1 by 1 the model assumptions in the statistical model used in the search, e.g.

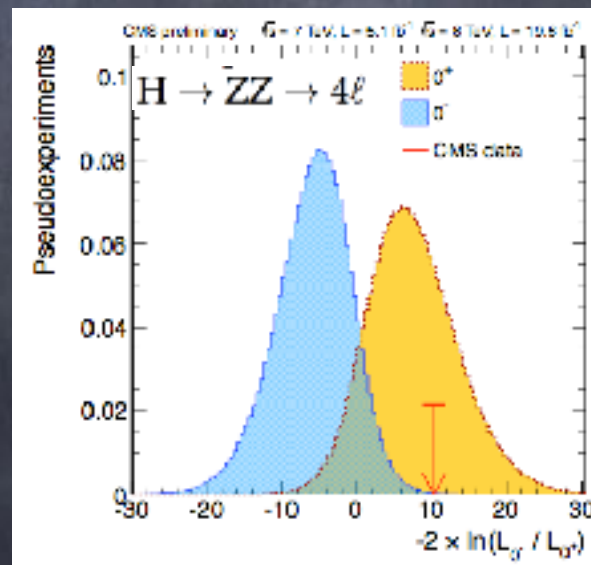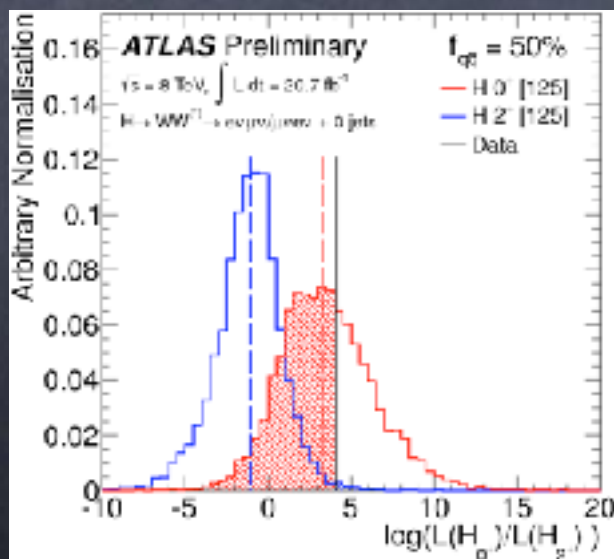| Background (scan $m_H$) | $\lambda(\mu = 0, m_H) = \dfrac{L(\mu = 0, m_H, \hat{\hat{\theta}})}{L(\hat{\mu}, m_H, \hat{\theta})}$ |
|---|---|
| Signal (scan $m_H$) | $\lambda(\mu, m_H) = \dfrac{L(\mu, m_H, \hat{\hat{\theta}})}{L(\hat{\mu}, m_H, \hat{\theta})}$ |
| Mass consistency | $\lambda(m_H) = \dfrac{L(m_H, \hat{\hat{\mu}}_1, \hat{\hat{\mu}}_2, \hat{\hat{\theta}})}{L(m_{\hat{1}H}, m_{\hat{2}H}, \hat{\mu}_1, \hat{\mu}_2, \hat{\theta})}$ |
| Mass | $\lambda(m_H) = \dfrac{L(m_H, \hat{\hat{\mu}}_1, \hat{\hat{\mu}}_2, \hat{\hat{\theta}})}{L(\hat{m_H}, \hat{\mu}_1, \hat{\mu}_2, \hat{\theta})}$ |
| Signal and mass | $\lambda(\mu, m_H) = \dfrac{L(\mu, m_H, \hat{\hat{\theta}}_\mu)}{L(\hat{\mu}, \hat{m_H}, \hat{\theta}_\mu)}$ |

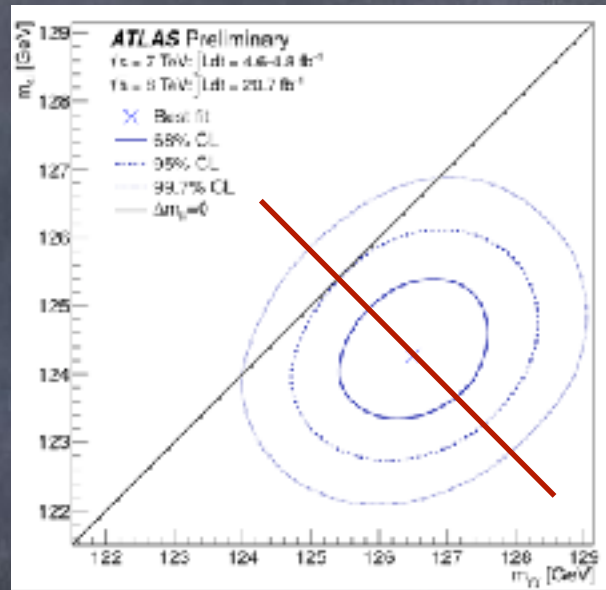# Signal strength vs mass



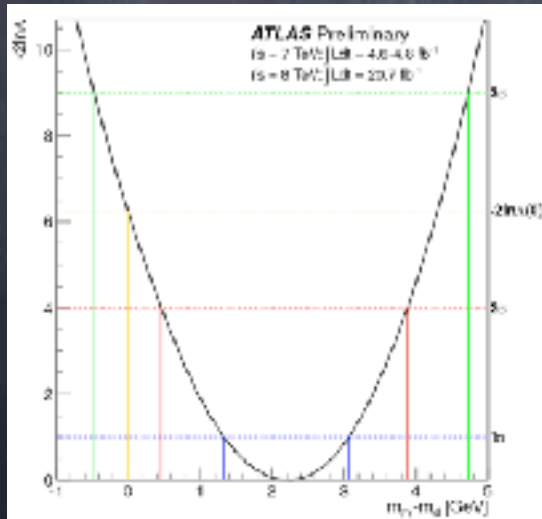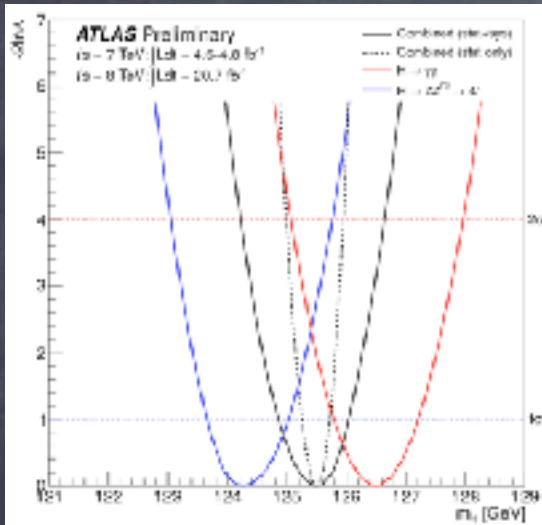- Contours not shown for μ→0...

# Testing J^P - 2 point hyp. test



Alt
OK

Null
OK

!!?!!

# Mass measurements



Compatibility, combination

# Implementation

K. Cranmer, G. Lewis, L. Moneta, A. Shibata, and W. Verkerke, *HistFactory: A tool for creating statistical models for use with RooFit and RooStats*, CERN-OPEN-2012-016 (2012). http://cdsweb.cern.ch/record/1456844.

L. Moneta, K. Belasco, K. S. Cranmer, S. Kreiss, A. Lazzaro, et al., *The RooStats Project*, PoS ACAT2010 (2010) 057, arXiv:1009.1003 [physics.data-an].

W. Verkerke and D. Kirkby, *The RooFit toolkit for data modeling*, Tech. Rep. physics/0306116, SLAC, Stanford, CA, Jun, 2003. arXiv:physics/0306116 [physics.data-an].

CERN Program Library Long Writeup D506

## MINUIT

Function Minimization and Error Analysis

Reference Manual

Version 94.1

F. James

# Summary

- Statistical practices in HEP evolved during the Higgs boson searches from LEP to Tevatron to LHC

  - Profile likelihood (ratio) used for searches as well as measurements (MINUIT fits at the base)

  - The full chain from exclusion to measurements via discovery carried out in a common framework

- Bayes and non-standard treatment of limits (CL$_s$) widely used in HEP

*improbabile rerum cotidie fieri*

## The Unconditional Ensemble

$$L(data \mid \mu, \theta) = Poisson(data \mid \mu s(\theta) + b(\theta)) \times p(\theta \mid \tilde{\theta})$$

Signal region main measurement

Control region auxiliary measurement

The nuisance parameters represent uncertain aspects of the model (background normalization and shape, systematic uncertainties) :

- Initial measured value of the parameter $\tilde{\theta}_0$  $G(\theta \mid \tilde{\theta}_0, \sigma)$

Marumi Kado

- First fit the data (typically under the alternate hypothesis) $\hat{\hat{\theta}}_A$

- The nuisance parameter $\theta$ is fixed for generation to default measured value $\hat{\hat{\theta}}_A$

- The auxiliary measurement $\tilde{\theta}$ is randomized  $G(\tilde{\theta} \mid \hat{\hat{\theta}}_A, \sigma)$

- Fit $\hat{\theta}, \hat{\tilde{\theta}}$ in toys  $G(\theta \mid \tilde{\theta}, \sigma)$
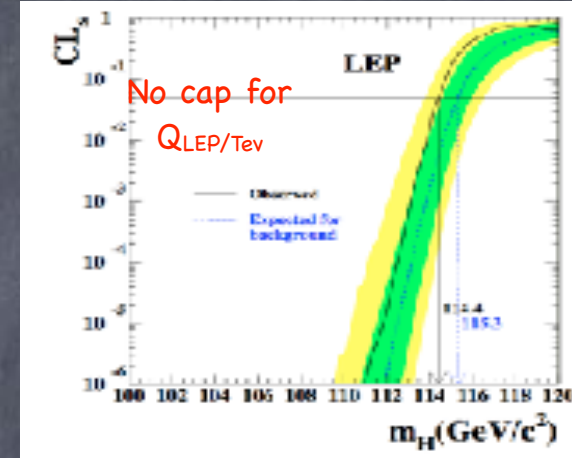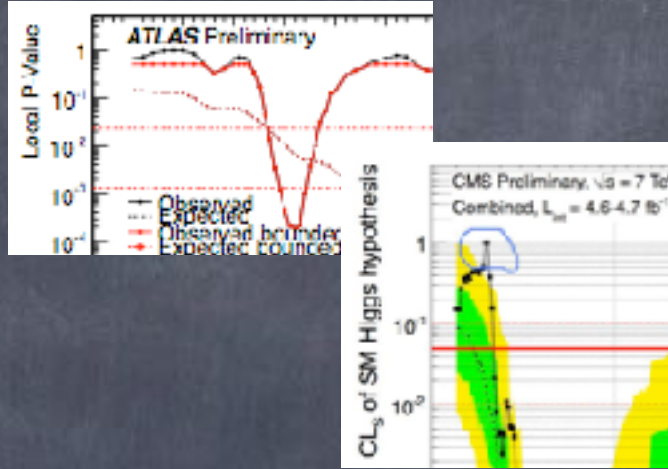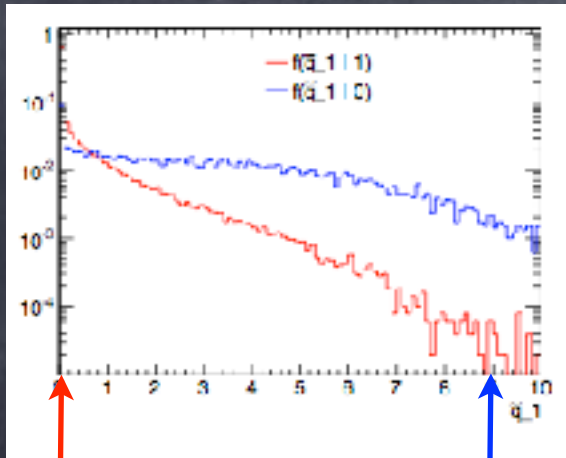
- Build a combined likelihood $\mathcal{L}(\mu, \theta) = (\prod_{i}^{N} \mathcal{L}_i^0(\mu, \theta_i)) \times (\prod_{j}^{M} \mathcal{A}(\theta^j))$
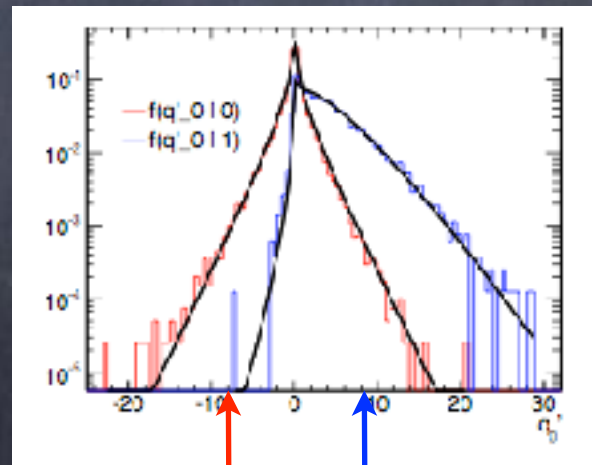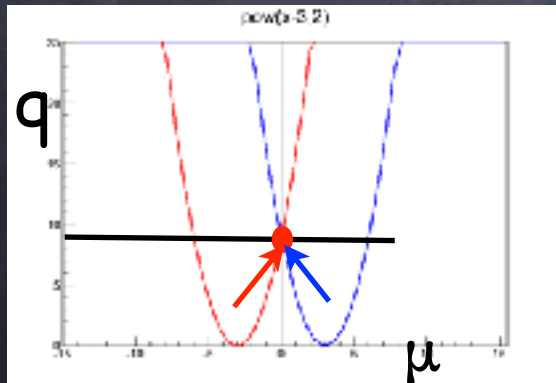
  - $\theta$ is now the set of all *unique* nuisance parameters
  - Some $\theta_i^j$ are shared between channels. This must be recognized to ensure proper correlation.

A. Armbruster

# Uncapping (open issue)



$$\hat{\mu} < 0 \rightarrow q_0 = 0$$

- No change in interpretation of limit or significance

- Visualize deficits for $p_0$ and excesses for $CL_s$

- Need to convince CMS colleagues... :-)