# Using the Profile Likelihood in Searches for New Physics

Statistical issues relevant to significance of discovery claims
BIRS, Banff, 11-16 July, 2010

Glen Cowan[1], Kyle Cranmer[2], Eilam Gross[3], Ofer Vitells[3]

[1] Physics Department, Royal Holloway, University of London, Egham, TW20 0EX, U.K.
[2] Physics Department, New York University, New York, NY 10003, U.S.A.
[3] Weizmann Institute of Science, Rehovot 76100, Israel

# Outline

Prototype search analysis for LHC

Test statistics based on profile likelihood ratio

Systematics covered via nuisance parameters

Sampling distributions to get significance/sensitivity

Asymptotic formulae from Wilks/Wald

Examples:

$n \sim$ Poisson $(\mu s + b)$, $m \sim$ Poisson$(\tau b)$

Shape analysis

Conclusions

# Prototype search analysis

Search for signal in a region of phase space; result is histogram of some variable $x$ giving numbers:

$$\mathbf{n} = (n_1, \ldots, n_N)$$

Assume the $n_i$ are Poisson distributed with expectation values

$$E[n_i] = \mu s_i + b_i$$

strength parameter

where

$$s_i = s_{\text{tot}} \int_{\text{bin } i} f_s(x; \boldsymbol{\theta}_s)\, dx\,, \quad b_i = b_{\text{tot}} \int_{\text{bin } i} f_b(x; \boldsymbol{\theta}_b)\, dx\,.$$

signal                      background

# Prototype analysis (II)

Often also have a subsidiary measurement that constrains some of the background and/or shape parameters:

$$\mathbf{m} = (m_1, \ldots, m_M)$$

Assume the $m_i$ are Poisson distributed with expectation values

$$E[m_i] = u_i(\boldsymbol{\theta})$$

nuisance parameters $(\boldsymbol{\theta}_\mathrm{s}, \boldsymbol{\theta}_\mathrm{b}, b_\mathrm{tot})$

Likelihood function is

$$L(\mu, \boldsymbol{\theta}) = \prod_{j=1}^{N} \frac{(\mu s_j + b_j)^{n_j}}{n_j!} e^{-(\mu s_j + b_j)} \prod_{k=1}^{M} \frac{u_k^{m_k}}{m_k!} e^{-u_k}$$

# The profile likelihood ratio

Base significance test on the profile likelihood ratio:

maximizes $L$ for specified $\mu$

$$\lambda(\mu) = \frac{L(\mu, \hat{\hat{\boldsymbol{\theta}}})}{L(\hat{\mu}, \hat{\boldsymbol{\theta}})}$$

maximize $L$

The likelihood ratio of point hypotheses gives optimum test (Neyman-Pearson lemma).

The profile LR hould be near-optimal in present analysis with variable $\mu$ and nuisance parameters $\boldsymbol{\theta}$.

# Test statistic for discovery

Try to reject background-only ($\mu = 0$) hypothesis using

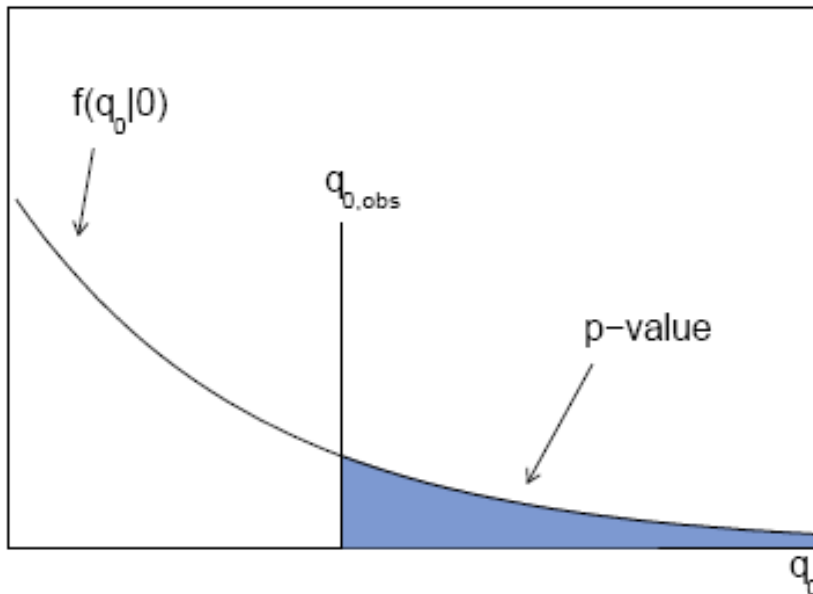$$q_0 = \begin{cases} -2\ln\lambda(0) & \hat{\mu} \geq 0 \\ 0 & \hat{\mu} < 0 \end{cases}$$

i.e. here only regard upward fluctuation of data as evidence against the background-only hypothesis.

# *p*-value for discovery

Large $q_0$ means increasing incompatibility between the data and hypothesis, therefore *p*-value for an observed $q_{0,\text{obs}}$ is

$$p_0 = \int_{q_{0,\text{obs}}}^{\infty} f(q_0|0)\, dq_0$$
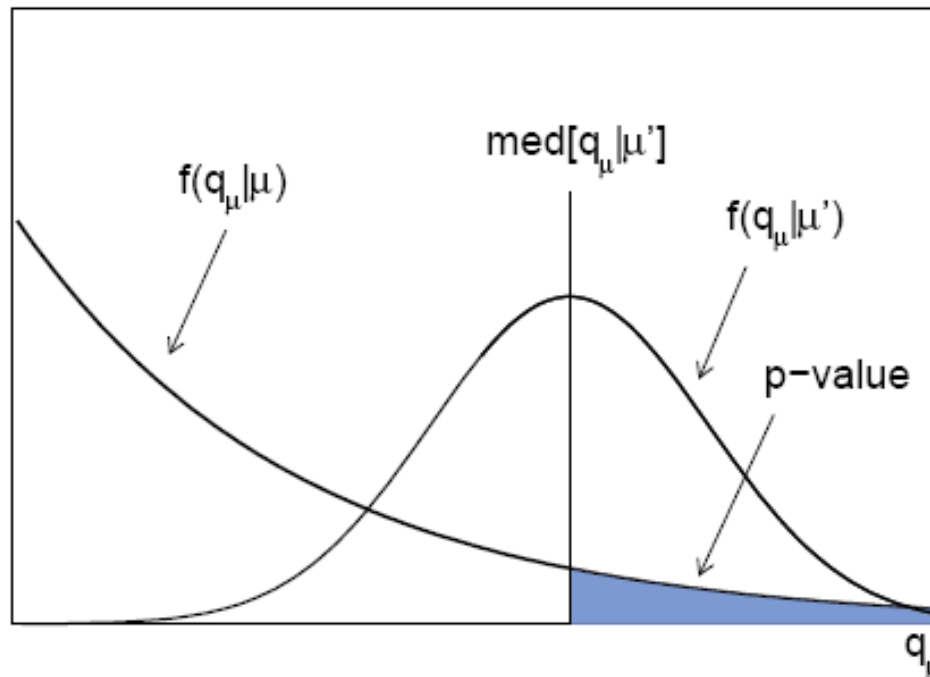
will get formula for this later

From *p*-value get equivalent significance,

$$Z = \Phi^{-1}(1 - p)$$

# Expected (or median) significance / sensitivity

When planning the experiment, we want to quantify how sensitive we are to a potential discovery, e.g., by given median significance assuming some nonzero strength parameter $\mu'$.



So for $p$-value, need $f(q_0|0)$, for sensitivity, will need $f(q_0|\mu')$,

# Wald approximation for profile likelihood ratio

To find $p$-values, we need: $\quad f(q_0|0), \quad f(q_\mu|\mu)$

For median significance under alternative, need: $\quad f(q_\mu|\mu')$

Use approximation due to Wald (1943)

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

sample size

$$\hat{\mu} \sim \text{Gaussian}(\mu', \sigma)$$

i.e., $E[\hat{\mu}] = \mu'$

$\sigma$ from covariance matrix $V$, use, e.g.,

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

# Noncentral chi-square for $-2\ln\lambda(\mu)$

If we can neglect the $O(1/\sqrt{N})$ term, $-2\ln\lambda(\mu)$ follows a noncentral chi-square distribution for one degree of freedom with noncentrality parameter

$$\Lambda = \frac{(\mu - \mu')^2}{\sigma^2}$$

As a special case, if $\mu' = \mu$ then $\Lambda = 0$ and $-2\ln\lambda(\mu)$ follows a chi-square distribution for one degree of freedom (Wilks).

# Distribution of $q_0$

Assuming the Wald approximation, we can write down the full distribution of $q_0$ as

$$f(q_0|\mu') = \left(1 - \Phi\left(\frac{\mu'}{\sigma}\right)\right)\delta(q_0) + \frac{1}{2}\frac{1}{\sqrt{2\pi}}\frac{1}{\sqrt{q_0}}\exp\left[-\frac{1}{2}\left(\sqrt{q_0} - \frac{\mu'}{\sigma}\right)^2\right]$$

The special case $\mu' = 0$ is a "half chi-square" distribution:

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

# Cumulative distribution of $q_0$, significance

From the pdf, the cumulative distribution of $q_0$ is found to be

$$F(q_0|\mu') = \Phi\left(\sqrt{q_0} - \frac{\mu'}{\sigma}\right)$$

The special case $\mu' = 0$ is

$$F(q_0|0) = \Phi\left(\sqrt{q_0}\right)$$

The $p$-value of the $\mu = 0$ hypothesis is

$$p_0 = 1 - F(q_0|0)$$

Therefore the discovery significance $Z$ is simply

$$Z = \Phi^{-1}(1 - p_0) = \sqrt{q_0}$$

# The Asimov data set

To estimate median value of $-2\ln\lambda(\mu)$, consider special data set where all statistical fluctuations suppressed and $n_i$, $m_i$ are replaced by their expectation values (the "Asimov" data set):

$$
\begin{aligned}
n_i &= \mu' s_i + b_i \\
m_i &= u_i
\end{aligned}
$$

$$\longrightarrow \quad \hat{\mu} = \mu' \qquad \hat{\boldsymbol{\theta}} = \boldsymbol{\theta}$$

$$\lambda_A(\mu) = \frac{L_A(\mu, \hat{\hat{\boldsymbol{\theta}}})}{L_A(\hat{\mu}, \hat{\boldsymbol{\theta}})} = \frac{L_A(\mu, \hat{\hat{\boldsymbol{\theta}}})}{L_A(\mu', \boldsymbol{\theta})}$$
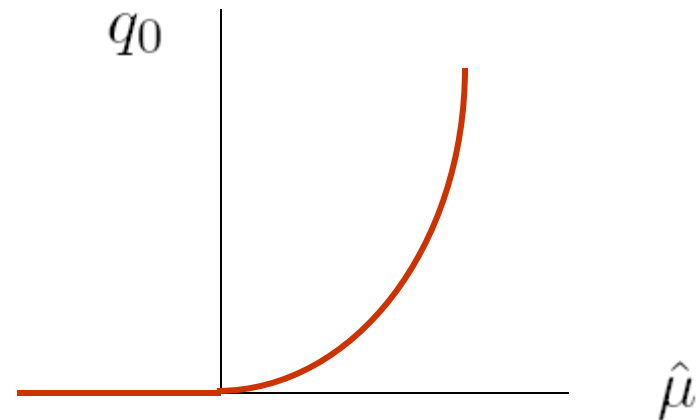
$$-2\ln\lambda_A(\mu) = \frac{(\mu - \mu')^2}{\sigma^2} = \Lambda$$

Asimov value of $-2\ln\lambda(\mu)$ gives non-centrality param. $\Lambda$, or equivalently, $\sigma$

# Relation between test statistics and $\hat{\mu}$

Assuming Wald approximation, the relation between $q_0$ and $\hat{\mu}$ is

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



Monotonic, therefore quantiles of $\hat{\mu}$ map one-to-one onto those of $q_0$, e.g.,

$$\text{med}[q_0] = q_0(\text{med}[\hat{\mu}]) = q_0(\mu') = \frac{\mu'^2}{\sigma^2} = -2\ln\lambda_A(0)$$

$$\text{med}[Z_0] = \sqrt{-2\ln\lambda_A(0)}$$

# Profile likelihood ratio for upper limits

For purposes of setting an upper limit on $\mu$ use

$$q_\mu = \begin{cases} -2\ln\lambda(\mu) & \hat{\mu} \leq \mu \\ 0 & \hat{\mu} > \mu \end{cases} \qquad \text{where} \qquad \lambda(\mu) = \frac{L(\mu, \hat{\hat{\boldsymbol{\theta}}})}{L(\hat{\mu}, \hat{\boldsymbol{\theta}})}$$

Note for purposes of setting an upper limit, one does not regard an upwards fluctuation of the data as representing incompatibility with the hypothesized $\mu$.

Note also here we allow the estimator for $\mu$ be negative (but $\hat{\mu}s_i + b_i$ must be positive).

# Alternative test statistic for upper limits

Assume physical signal model has $\mu > 0$, therefore if estimator for $\mu$ comes out negative, the closest physical model has $\mu = 0$.

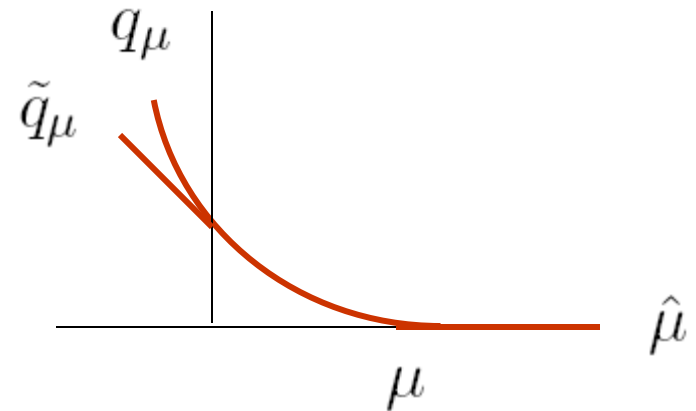Therefore could also measure level of discrepancy between data and hypothesized $\mu$ with

$$\tilde{\lambda}(\mu) = \begin{cases} \dfrac{L(\mu, \hat{\hat{\boldsymbol{\theta}}}(\mu))}{L(\hat{\mu}, \hat{\boldsymbol{\theta}})} & \hat{\mu} \geq 0, \\[2ex] \dfrac{L(\mu, \hat{\hat{\boldsymbol{\theta}}}(\mu))}{L(0, \hat{\hat{\boldsymbol{\theta}}}(0))} & \hat{\mu} < 0 . \end{cases} \qquad \tilde{q}_\mu = \begin{cases} -2 \ln \tilde{\lambda}(\mu) & \hat{\mu} \leq \mu \\[2ex] 0 & \hat{\mu} > \mu \end{cases}$$

Performance not identical to but very close to $q_\mu$ (of previous slide). $q_\mu$ is simpler in important ways.

# Relation between test statistics and $\hat{\mu}$

Assuming the Wald approximation for $-2\ln\lambda(\mu)$, $q_\mu$ and $\tilde{q}_\mu$ both have monotonic relation with $\mu$.

$$q_\mu = \begin{cases} \frac{(\mu-\hat{\mu})^2}{\sigma^2} & \hat{\mu} < \mu \\ 0 & \hat{\mu} > \mu \end{cases}$$



$$\tilde{q}_\mu = \begin{cases} \frac{\mu^2}{\sigma^2} - \frac{2\mu\hat{\mu}}{\sigma^2} & \hat{\mu} < 0 \\ \frac{(\mu-\hat{\mu})^2}{\sigma^2} & 0 \leq \hat{\mu} \leq \mu \\ 0 & \hat{\mu} > \mu \, , \end{cases}$$

And therefore quantiles of $q_\mu$, $\tilde{q}_\mu$ can be obtained directly from those of $\hat{\mu}$ (which is Gaussian).

# Distribution of $q_\mu$

Similar results for $q_\mu$

$$f(q_\mu|\mu') = \Phi\left(\frac{\mu' - \mu}{\sigma}\right)\delta(q_\mu) + \frac{1}{2}\frac{1}{\sqrt{2\pi}}\frac{1}{\sqrt{q_\mu}}\exp\left[-\frac{1}{2}\left(\sqrt{q_\mu} - \frac{(\mu - \mu')}{\sigma}\right)^2\right]$$

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

$$F(q_\mu|\mu') = \Phi\left(\sqrt{q_\mu} - \frac{(\mu - \mu')}{\sigma}\right)$$

$$p_\mu = 1 - F(q_\mu|\mu) = 1 - \Phi\left(\sqrt{q_\mu}\right)$$

# Distribution of $\tilde{q}_\mu$

Similar results for $\tilde{q}_\mu$

$$f(\tilde{q}_\mu|\mu') = \Phi\left(\frac{\mu'-\mu}{\sigma}\right)\delta(\tilde{q}_\mu)$$

$$+ \begin{cases} \frac{1}{2}\frac{1}{\sqrt{2\pi}}\frac{1}{\sqrt{\tilde{q}_\mu}}\exp\left[-\frac{1}{2}\left(\sqrt{\tilde{q}_\mu}-\frac{(\mu-\mu')}{\sigma}\right)^2\right] & 0 < \tilde{q}_\mu \leq \mu^2/\sigma^2 \\ \frac{1}{\sqrt{2\pi}\sigma}\exp\left[-\frac{1}{2}\frac{(\tilde{q}_\mu-(\mu^2-2\mu\mu')/\sigma^2)^2}{(2\mu/\sigma)^2}\right] & \tilde{q}_\mu > \mu^2/\sigma^2 \end{cases}$$

$$F(\tilde{q}_\mu|\mu') = \begin{cases} \Phi\left(\sqrt{\tilde{q}_\mu}-\frac{(\mu-\mu')}{\sigma}\right) & 0 < \tilde{q}_\mu \leq \mu^2/\sigma^2 \,, \\ \Phi\left(\frac{\tilde{q}_\mu-(\mu^2-2\mu\mu')/\sigma^2}{2\mu/\sigma}\right) & \tilde{q}_\mu > \mu^2/\sigma^2 \,. \end{cases}$$

# Monte Carlo test of asymptotic formula

$$n \sim \text{Poisson}(\mu s + b)$$

$$m \sim \text{Poisson}(\tau b)$$

Here take $\tau = 1$.

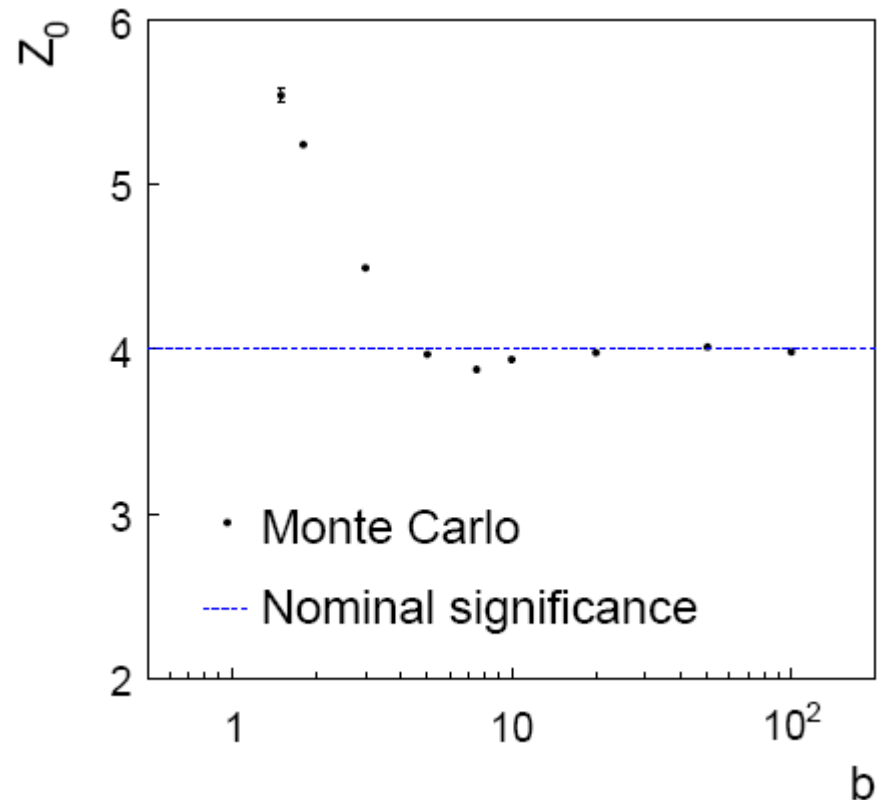Asymptotic formula is good approximation to $5\sigma$ level ($q_0 = 25$) already for $b \sim 20$.

# Monte Carlo test of asymptotic formulae

Significance from asymptotic formula, here $Z_0 = \sqrt{q_0} = 4$, compared to MC (true) value.

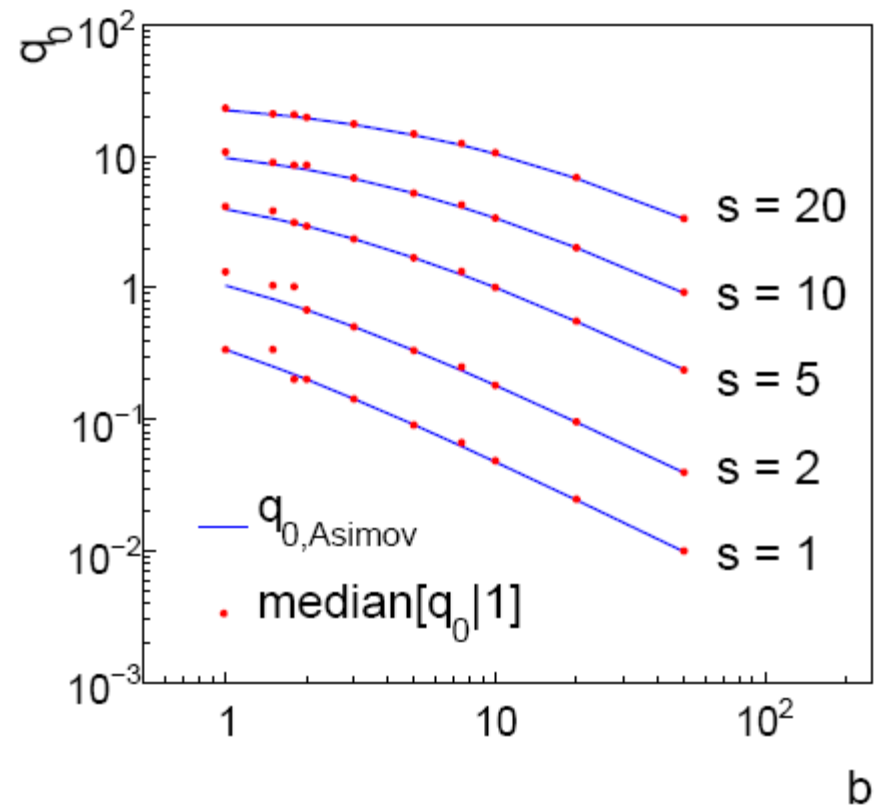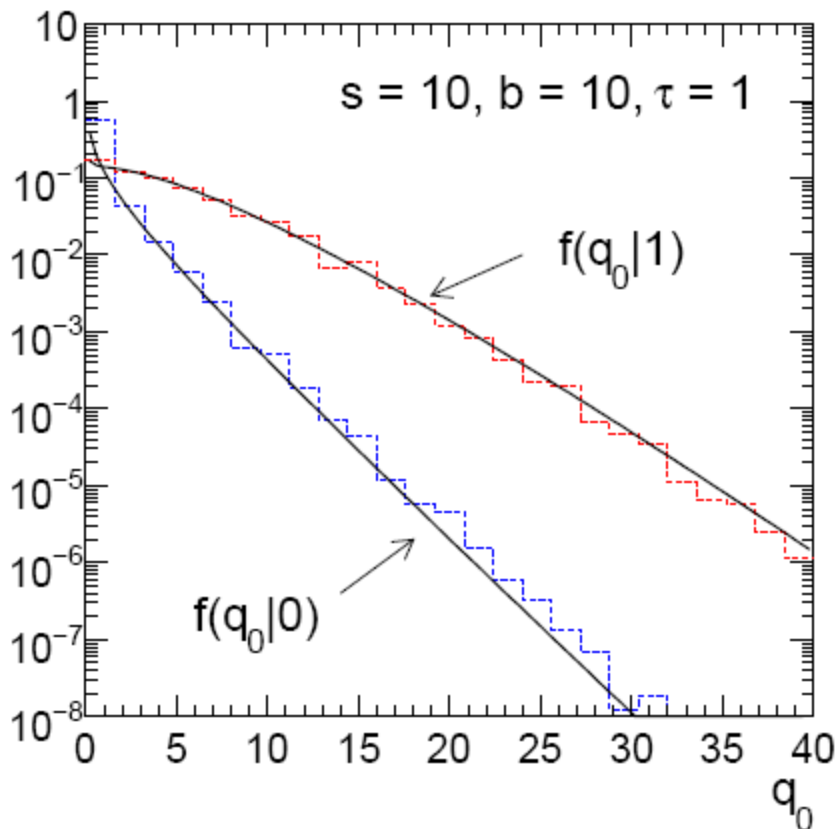For very low $b$, asymptotic formula underestimates $Z_0$.

Then slight overshoot before rapidly converging to MC value.

# Monte Carlo test of asymptotic formulae

Asymptotic $f(q_0|1)$ good already for fairly small samples.

Median$[q_0|1]$ from Asimov data set; good agreement with MC.
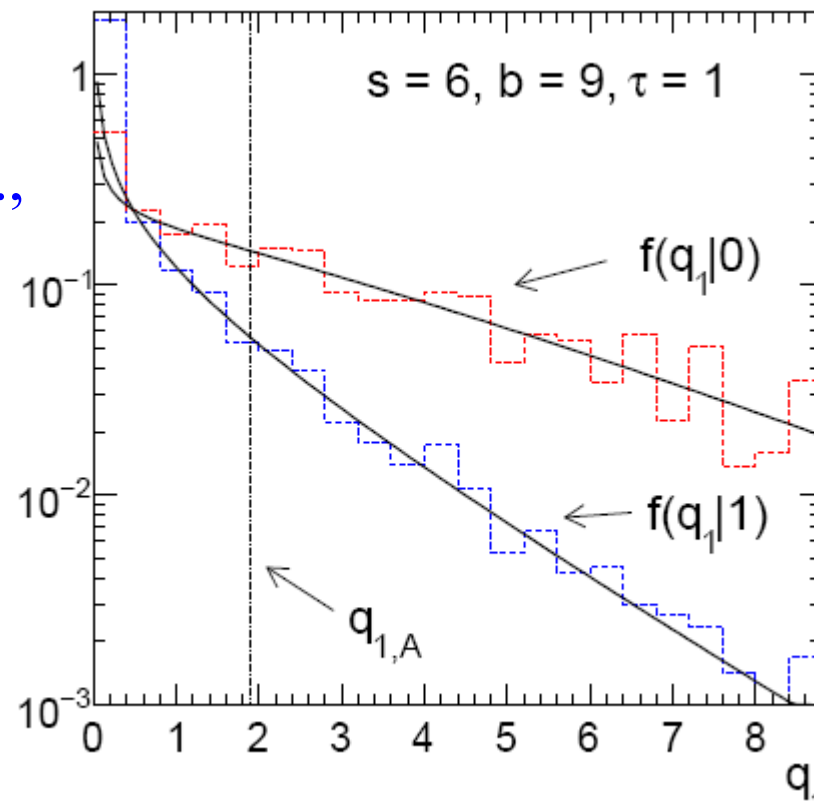
# Monte Carlo test of asymptotic formulae

Consider again $n \sim$ Poisson $(\mu s + b)$, $m \sim$ Poisson$(\tau b)$
Use $q_\mu$ to find $p$-value of hypothesized $\mu$ values.

E.g. $f(q_1|1)$ for $p$-value of $\mu = 1$.

Typically interested in 95% CL, i.e.,
$p$-value threshold = 0.05, i.e.,
$q_1 = 2.69$ or $Z_1 = \sqrt{q_1} = 1.64$.
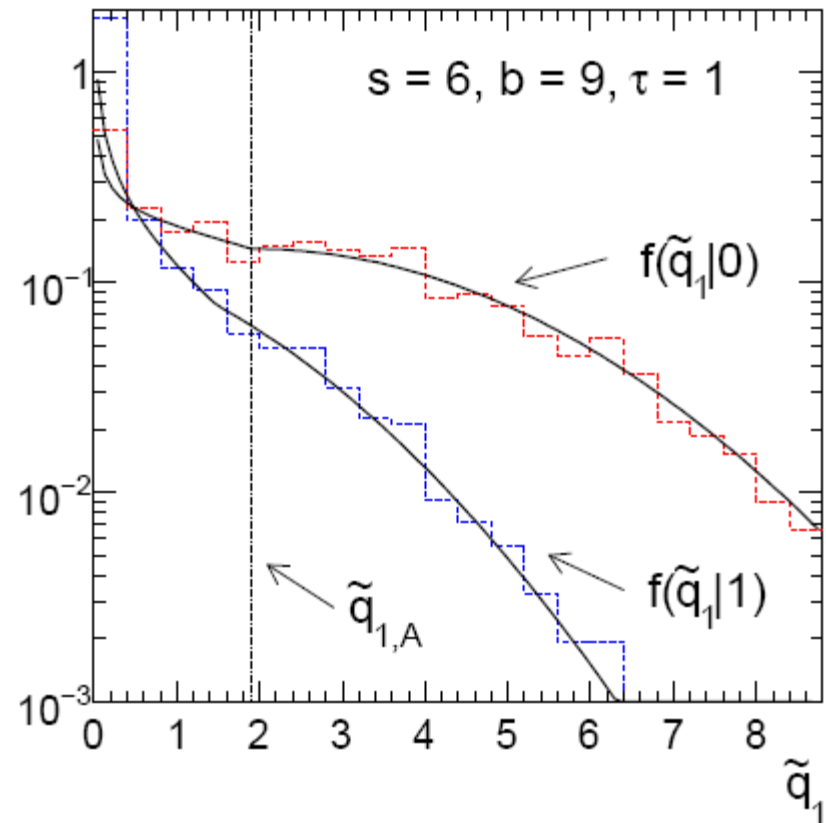
Median$[q_1|0]$ gives "exclusion sensitivity".

Here asymptotic formulae good for $s = 6$, $b = 9$.



$s = 6, b = 9, \tau = 1$

$f(q_1|0)$

$f(q_1|1)$

$q_{1,A}$
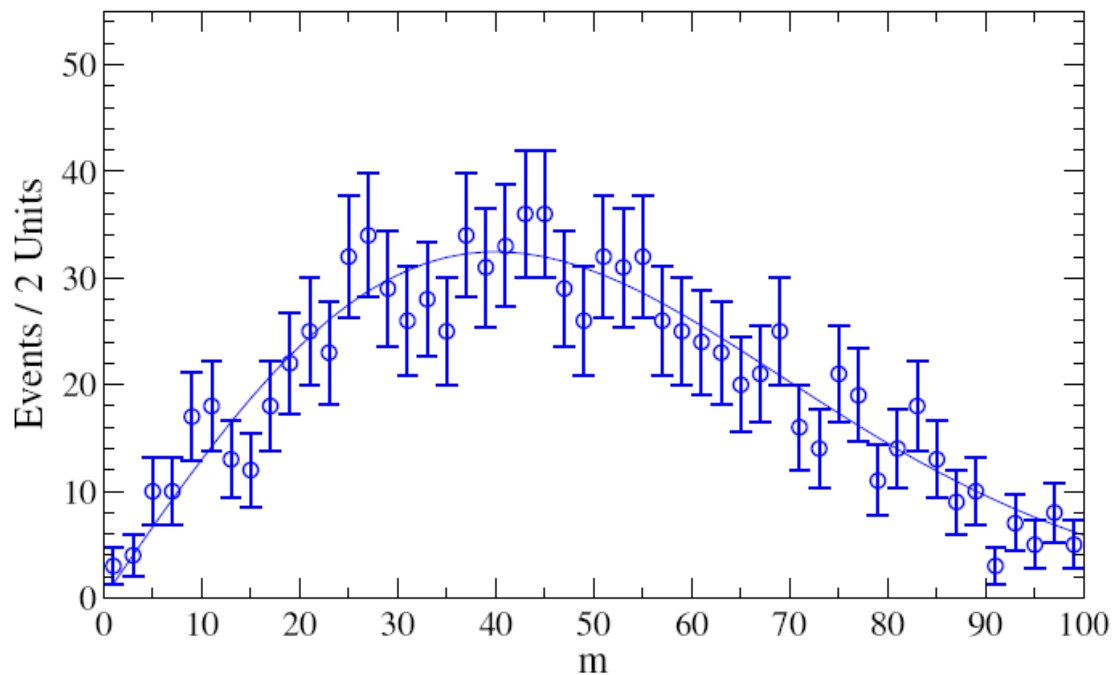
# Monte Carlo test of asymptotic formulae

Same message for test based on $\tilde{q}_\mu$.

$q_\mu$ and $\tilde{q}_\mu$ give similar tests to the extent that asymptotic formulae are valid.
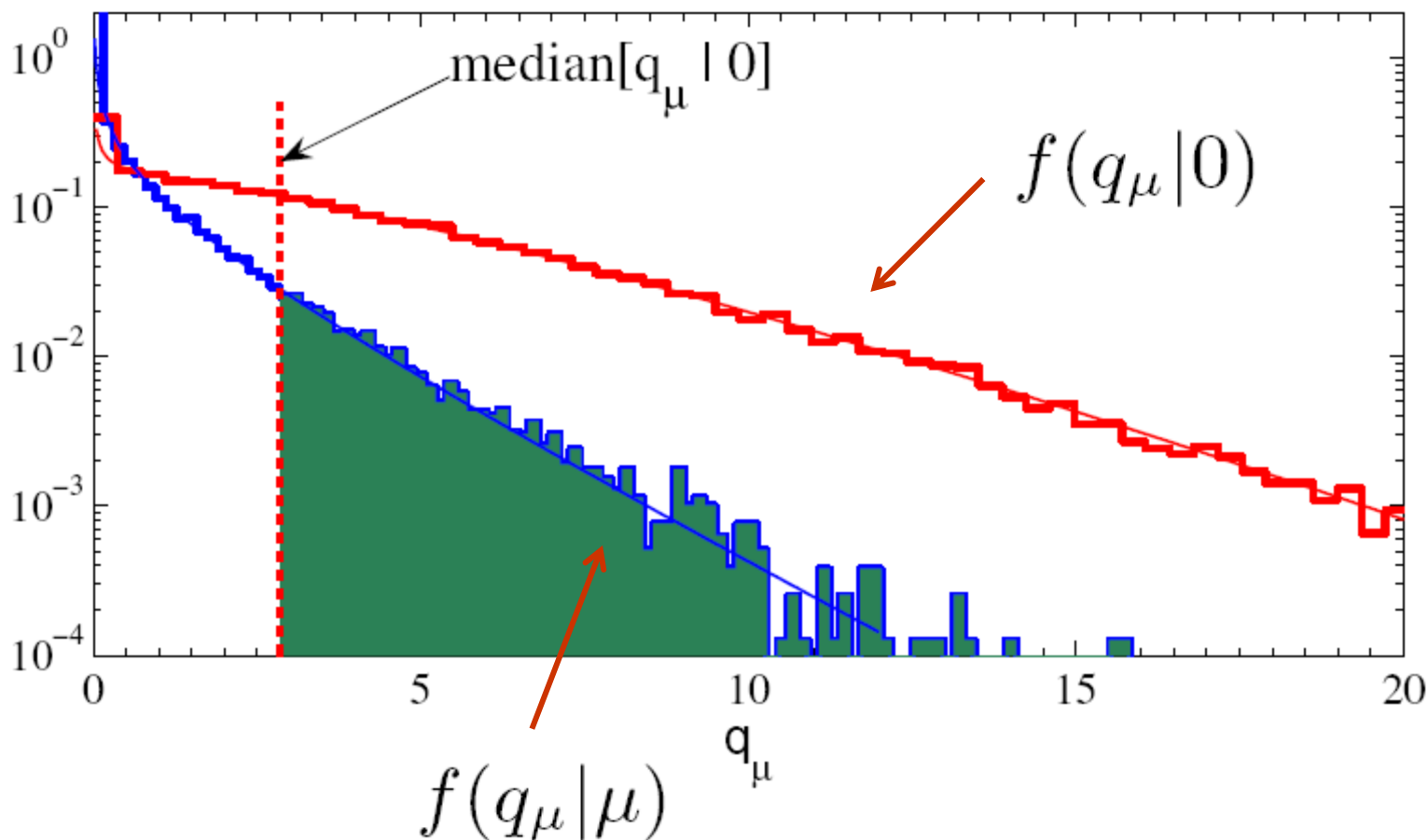
# Example 2: Shape analysis

Look for a Gaussian bump sitting on top of:



$$L(\mu, \theta) = \prod_{i=1}^{N} \frac{(\mu s_i + \theta f_{\mathrm{b},i})^{n_i}}{n_i!} e^{-(\mu s_i + \theta f_{\mathrm{b},i})}$$

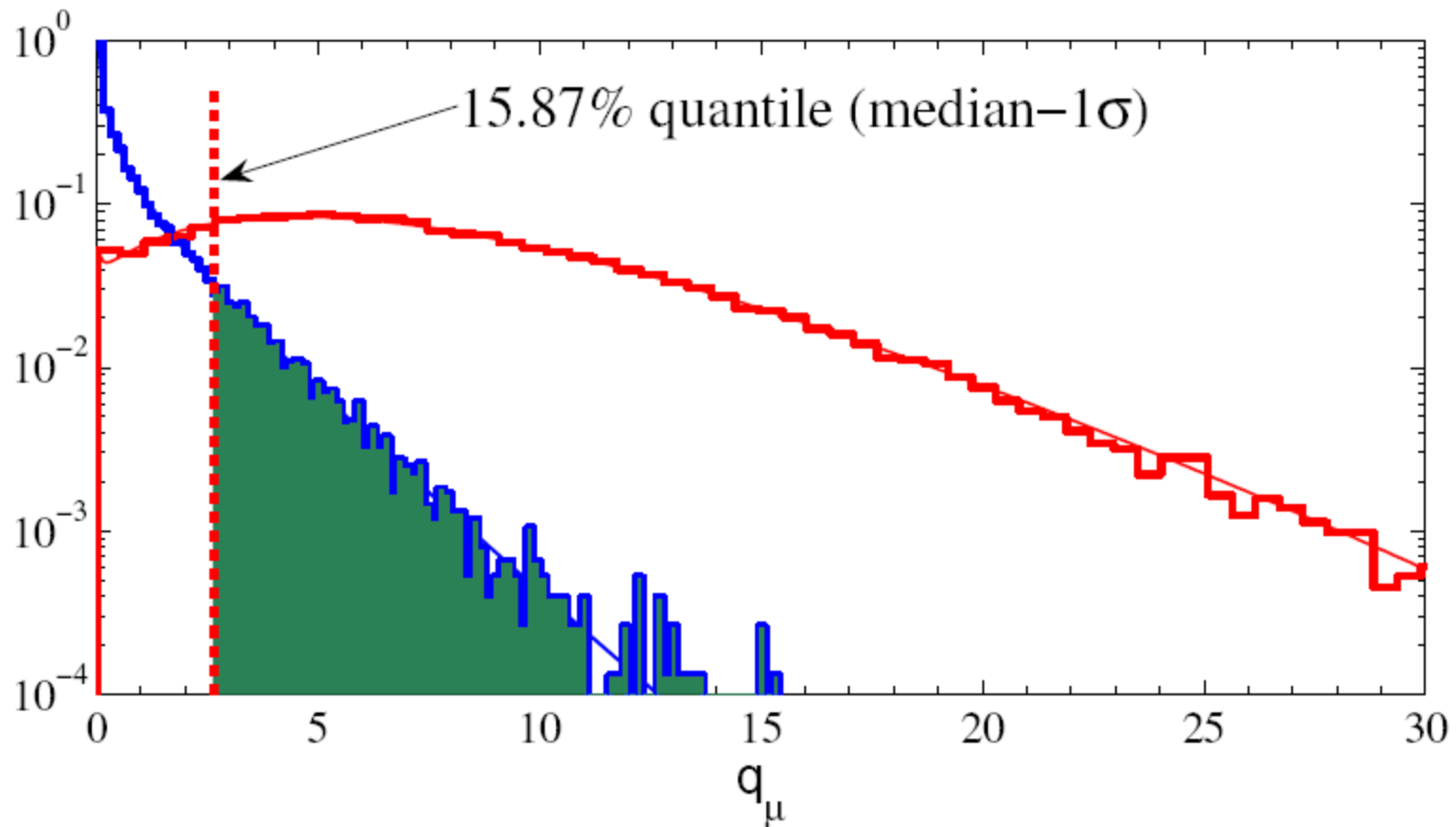# Monte Carlo test of asymptotic formulae

Distributions of $q_\mu$ here for $\mu$ that gave $p_\mu = 0.05$.

# Using $f(q_\mu|0)$ to get error bands

We are not only interested in the median[qm|0]; we want to know how much statistical variation to expect from a real data set.
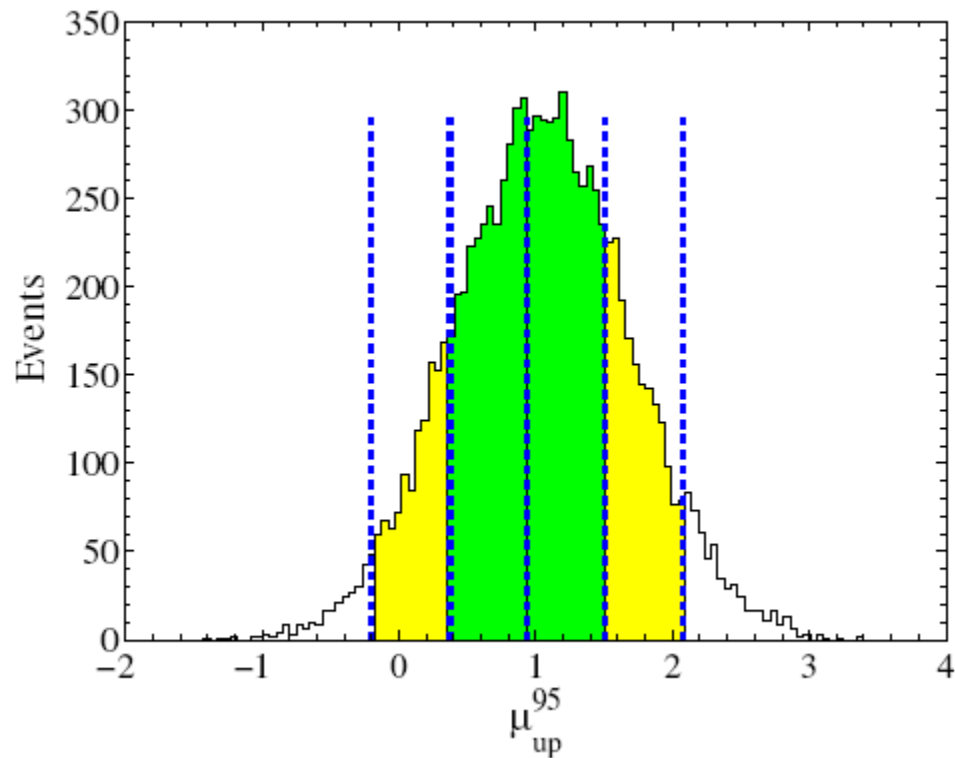
But we have full $f(q_\mu|0)$; we can get any desired quantiles.



15.87% quantile (median$-1\sigma$)

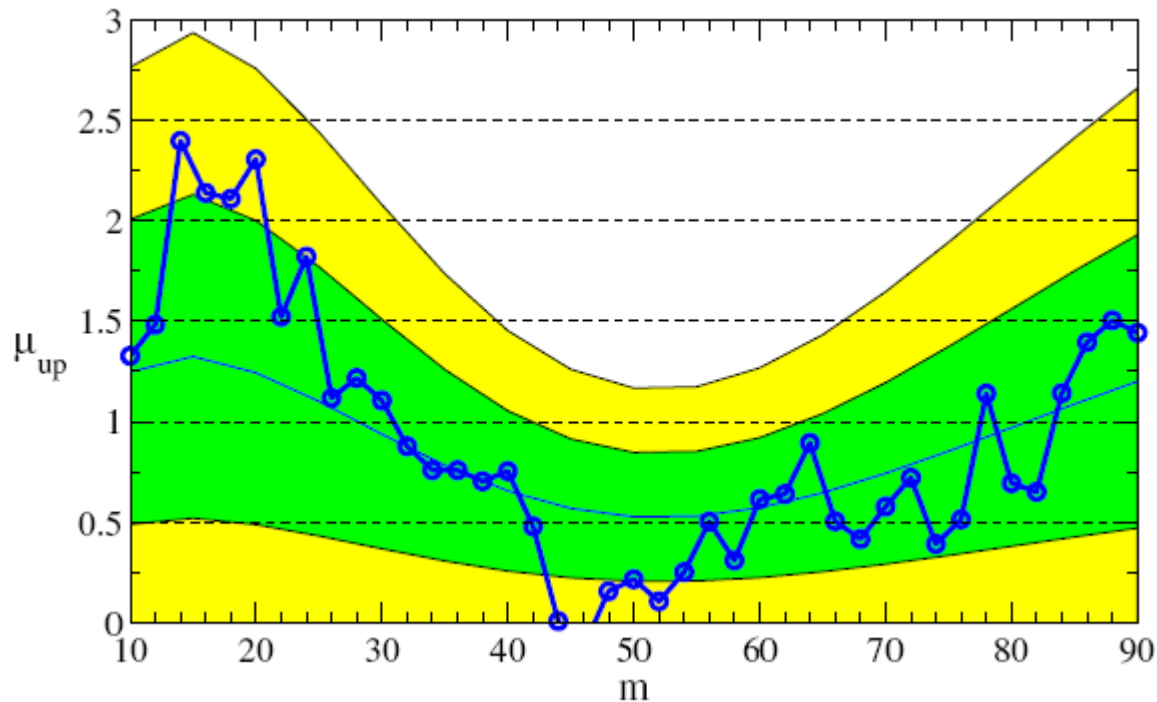# Distribution of upper limit on $\mu$

$\pm 1\sigma$ (green) and $\pm 2\sigma$ (yellow) bands from MC;

Vertical lines from asymptotic formulae

# Limit on $\mu$ versus peak position (mass)

$\pm 1\sigma$ (green) and $\pm 2\sigma$ (yellow) bands from asymptotic formulae;

Points are from a single arbitrary data set.

# Summary

Asymptotic distributions of profile LR applied to an LHC search.

Wilks: $f(q_\mu|\mu)$ for $p$-value of $\mu$.

Wald approximation for $f(q_\mu|\mu')$.

"Asimov" data set used to estimate median $q_\mu$ for sensitivity.

Gives $\sigma$ of distribution of estimator for $\mu$.

Asymptotic formulae especially useful for estimating sensitivity in high-dimensional parameter space.

Can always check with MC for very low data samples and/or when precision crucial.

Implementation in RooStats (ongoing)

# Extra slides

# Profile likelihood ratio for unified interval

We can also use directly

$$t_\mu = -2\ln\lambda(\mu) \qquad \text{where} \qquad \lambda(\mu) = \frac{L(\mu, \hat{\hat{\boldsymbol{\theta}}})}{L(\hat{\mu}, \hat{\boldsymbol{\theta}})}$$

as a test statistic for a hypothesized $\mu$.

Large discrepancy between data and hypothesis can correspond either to the estimate for $\mu$ being observed high or low relative to $\mu$.

This is essentially the statistic used for Feldman-Cousins intervals (here also treats nuisance parameters).

# Distribution of $t_\mu$

Using Wald approximation, $f(t_\mu|\mu')$ is noncentral chi-square for one degree of freedom:

$$f(t_\mu|\mu') = \frac{1}{2\sqrt{t_\mu}}\frac{1}{\sqrt{2\pi}}\left[\exp\left(-\frac{1}{2}\left(\sqrt{t_\mu}+\frac{\mu-\mu'}{\sigma}\right)^2\right)+\exp\left(-\frac{1}{2}\left(\sqrt{t_\mu}-\frac{\mu-\mu'}{\sigma}\right)^2\right)\right]$$

Special case of $\mu = \mu'$ is chi-square for one d.o.f. (Wilks).

The $p$-value for an observed value of $t_\mu$ is

$$p_\mu = 1 - F(t_\mu|\mu) = 2\left(1-\Phi\left(\sqrt{t_\mu}\right)\right)$$

and the corresponding significance is

$$Z_\mu = \Phi^{-1}(1-p_\mu) = \Phi^{-1}\left(2\Phi\left(\sqrt{t_\mu}\right)-1\right)$$

# Combination of channels

For a set of independent decay channels, full likelihood function is product of the individual ones:

$$L(\mu, \boldsymbol{\theta}) = \prod_i L_i(\mu, \boldsymbol{\theta}_i)$$

For combination need to form the full function and maximize to find estimators of $\mu$, $\boldsymbol{\theta}$.

$\rightarrow$ ongoing ATLAS/CMS effort with RooStats framework

Trick for median significance: estimator for $\mu$ is equal to the Asimov value $\mu'$ for all channels separately, so for combination,

$$\lambda_{\mathrm{A}}(\mu) = \prod_i \lambda_{\mathrm{A},i}(\mu) \qquad \text{where} \qquad \lambda_{\mathrm{A},i}(\mu) = \frac{L_i(\mu, \hat{\hat{\boldsymbol{\theta}}})}{L_i(\mu', \boldsymbol{\theta})}$$

# Discovery significance for $n \sim$ Poisson($s + b$)

Consider again the case where we observe $n$ events , model as following Poisson distribution with mean $s + b$ (assume $b$ is known).

1) For an observed $n$, what is the significance $Z_0$ with which we would reject the $s = 0$ hypothesis?

2) What is the expected (or more precisely, median ) $Z_0$ if the true value of the signal rate is $s$?

# Gaussian approximation for Poisson significance

For large $s + b$, $n \to x \sim \text{Gaussian}(\mu, \sigma)$, $\mu = s + b$, $\sigma = \sqrt{(s + b)}$.

For observed value $x_{\text{obs}}$, $p$-value of $s = 0$ is $\text{Prob}(x > x_{\text{obs}} \mid s = 0)$,:

$$p_0 = 1 - \Phi\left(\frac{x_{\text{obs}} - b}{\sqrt{b}}\right)$$

Significance for rejecting $s = 0$ is therefore

$$Z_0 = \Phi^{-1}(1 - p_0) = \frac{x_{\text{obs}} - b}{\sqrt{b}}$$

Expected (median) significance assuming signal rate $s$ is

$$\text{median}[Z_0 | s + b] = \frac{s}{\sqrt{b}}$$

# Better approximation for Poisson significance

Likelihood function for parameter $s$ is

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

or equivalently the log-likelihood is

$$\ln L(s) = n \ln(s+b) - (s+b) - \ln n!$$

Find the maximum by setting $\qquad \frac{\partial \ln L}{\partial s} = 0$

gives the estimator for $s$: $\qquad \hat{s} = n - b$

# Approximate Poisson significance (continued)

The likelihood ratio statistic for testing $s = 0$ is

$$q_0 = -2 \ln \frac{L(0)}{L(\hat{s})} = 2 \left( n \ln \frac{n}{b} + b - n \right) \quad \text{for } n > b, \; 0 \text{ otherwise}$$

For sufficiently large $s + b$, (use Wilks' theorem),

$$Z_0 \approx \sqrt{q_0} = \sqrt{2 \left( n \ln \frac{n}{b} + b - n \right)} \quad \text{for } n > b, \; 0 \text{ otherwise}$$

To find median$[Z_0|s+b]$, let $n \rightarrow s + b$,

$$\text{median}[Z_0|s + b] \approx \sqrt{2 \left( (s + b) \ln(1 + s/b) - s \right)}$$

This reduces to $s/\sqrt{b}$ for $s \ll b$.

# Higgs search with profile likelihood

Combination of Higgs boson search channels (ATLAS)

*Expected Performance of the ATLAS Experiment:  Detector, Trigger and Physics*, arXiv:0901.0512, CERN-OPEN-2008-20.

Standard Model Higgs channels considered (more to be used later):

$H \rightarrow \gamma\gamma$

$H \rightarrow WW^{(*)} \rightarrow e\nu\mu\nu$

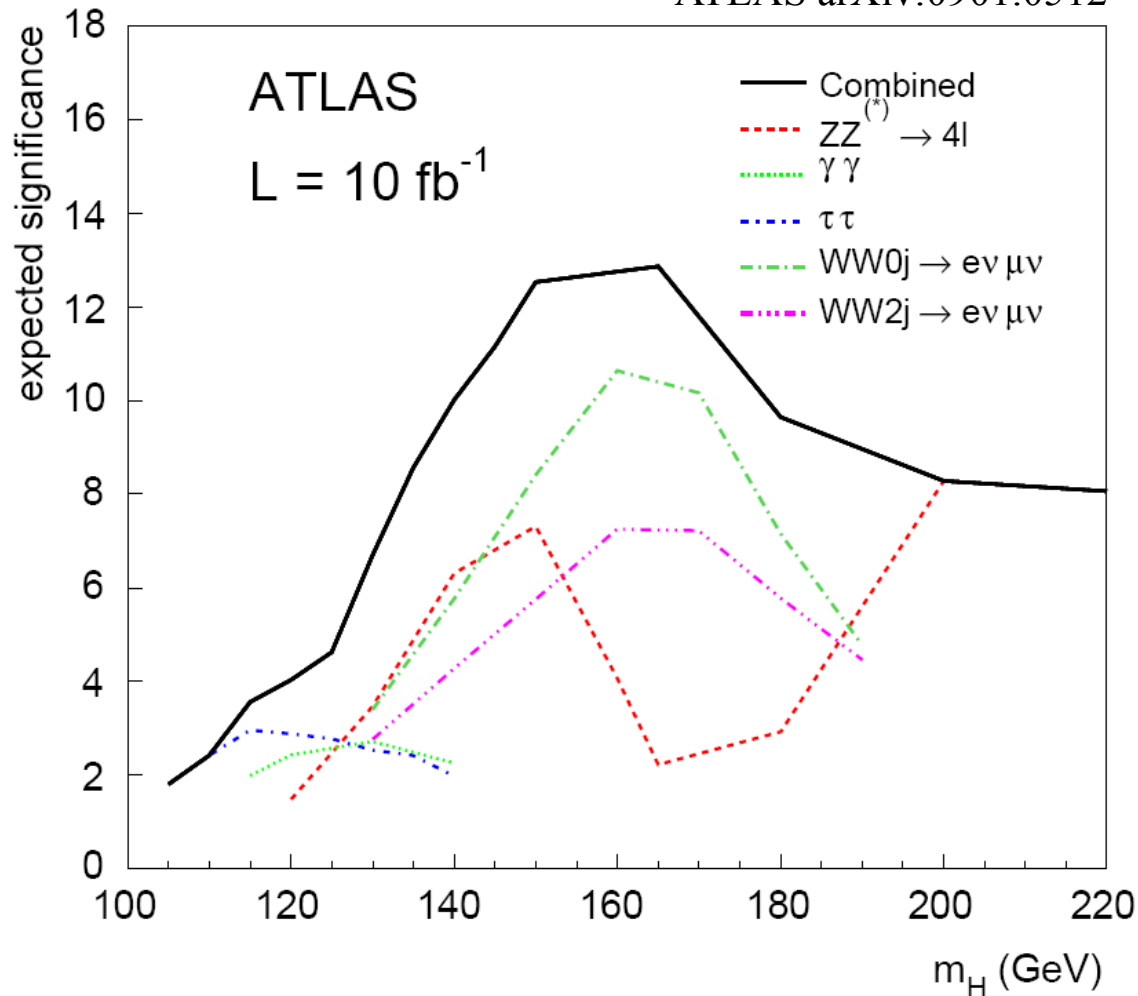$H \rightarrow ZZ^{(*)} \rightarrow 4l \ \ (l = e, \mu)$

$H \rightarrow \tau^{+}\tau^{-} \rightarrow ll, lh$

Used profile likelihood method for systematic uncertainties:
background rates, signal & background shapes.

# Combined median significance



ATLAS arXiv:0901.0512

N.B. illustrates statistical method, but study did not include all usable Higgs channels.

# An example: ATLAS Higgs search

(ATLAS Collab., CERN-OPEN-2008-020)

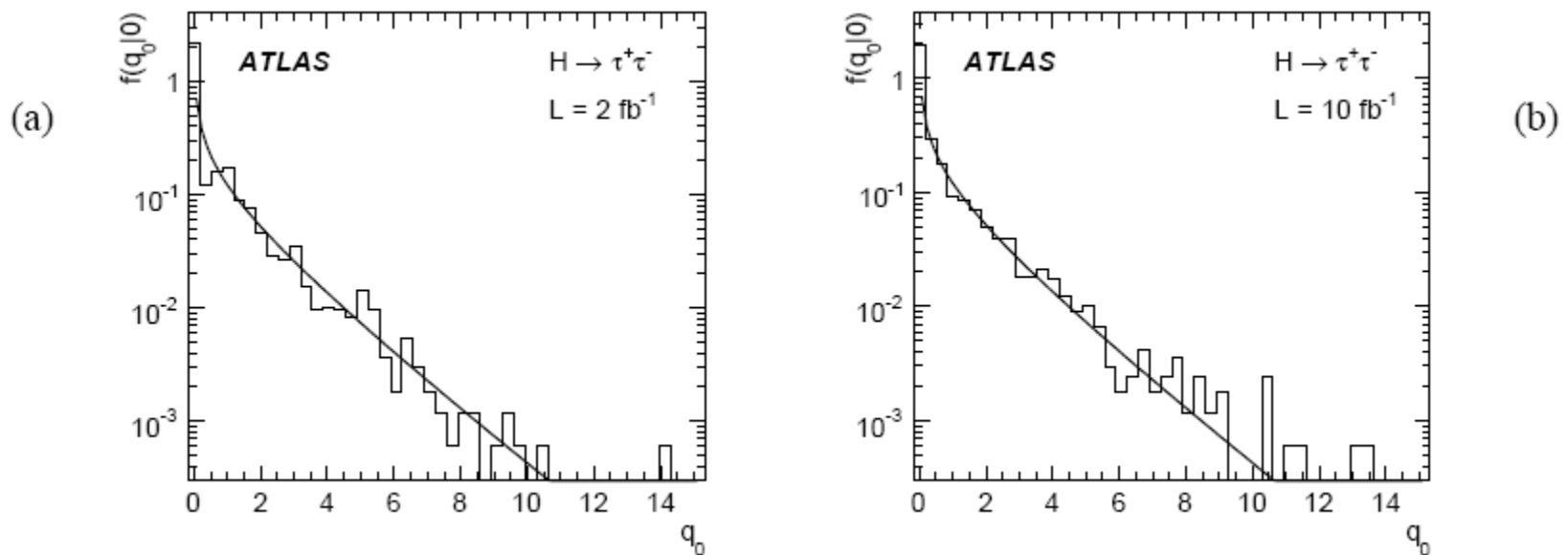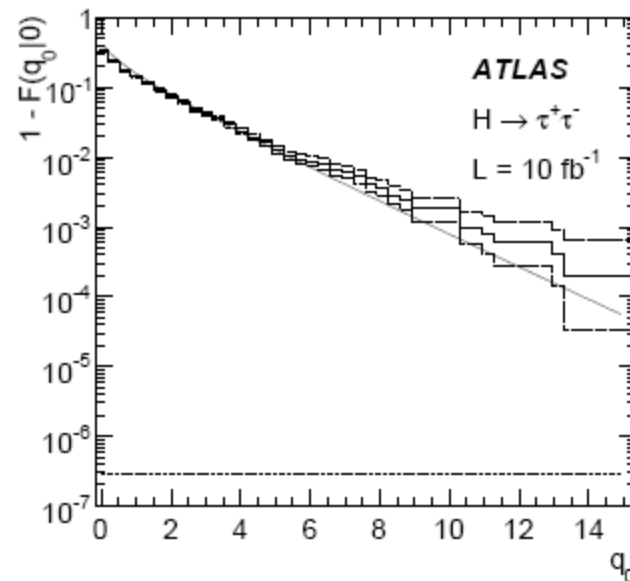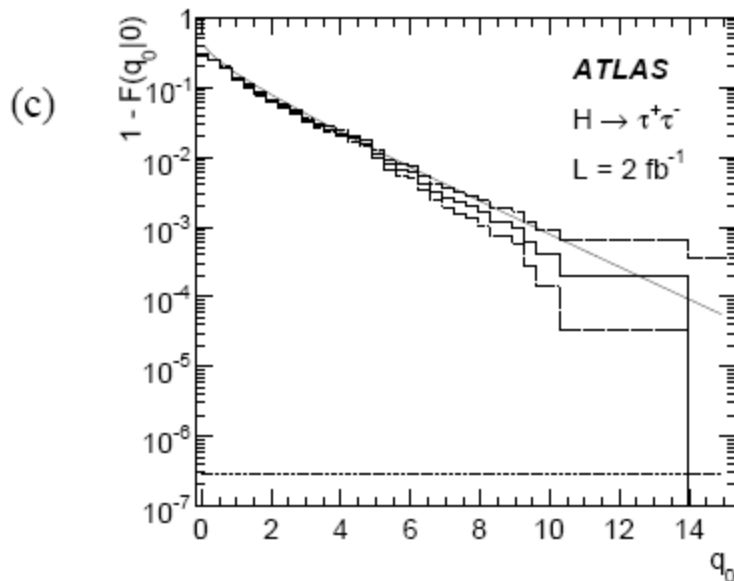**Statistical Combination of Several Important Standard Model Higgs Boson Search Channels.**



Figure 12: The distribution of the test statistic $q_0$ for $H \rightarrow \tau^+\tau^-$ under the null background-only hypothesis, for $m_H = 130\,\text{GeV}$ with an integrated luminosity of 2 (a) and 10 (b) $\text{fb}^{-1}$. A $\frac{1}{2}\chi_1^2$ distribution is superimposed. Figures (c) and (d) show $1 - F(q_0)$ where $F(q_0)$ is the corresponding cumulative distribution. The small excess of events at high $q_0$ is statistically compatible with the expected curves, as can be seen by comparison with the dotted histograms that show the 68.3% central confidence intervals for $p = 1 - F(q_0|0)$. The lower dotted line at $2.87 \times 10^{-7}$ shows the $5\sigma$ discovery threshold.
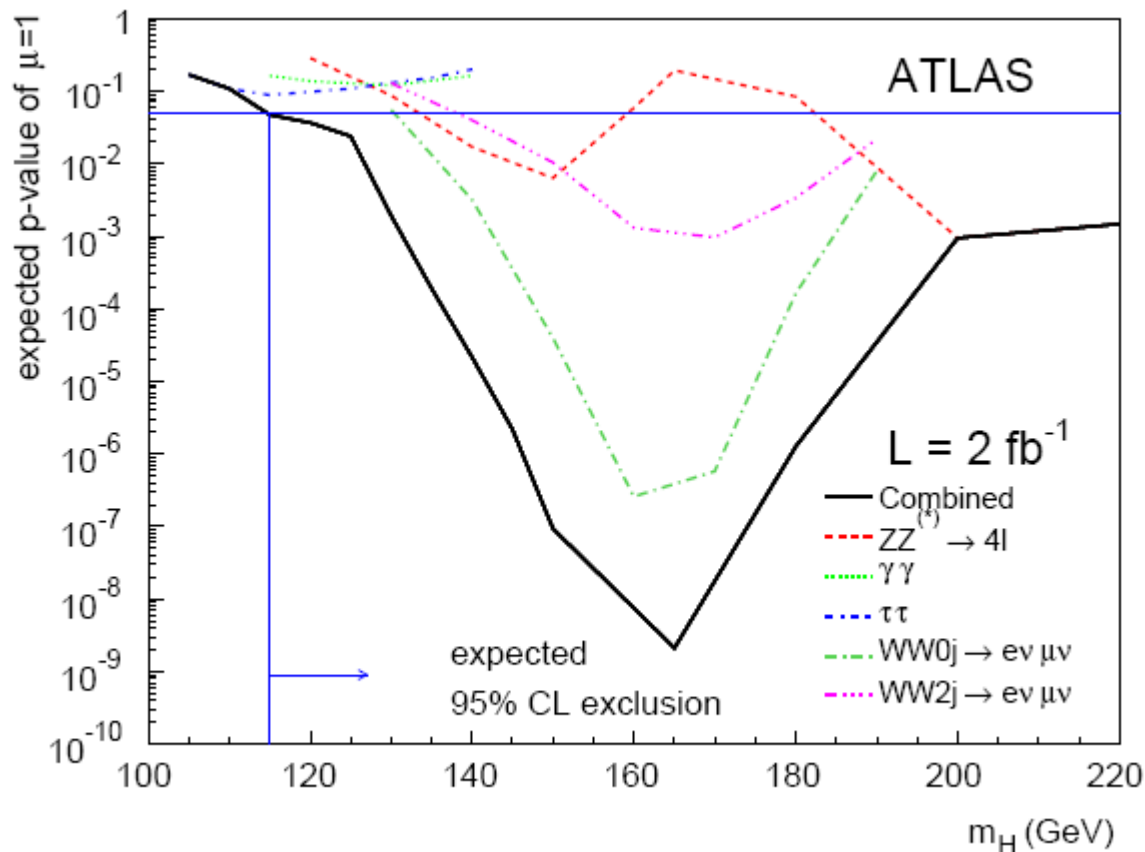
# Cumulative distributions of $q_0$

To validate to $5\sigma$ level, need distribution out to $q_0 = 25$, i.e., around $10^8$ simulated experiments.

Will do this if we really see something like a discovery.

# Example: exclusion sensitivity

Median *p*-value of $\mu = 1$ hypothesis versus Higgs mass assuming background-only data (ATLAS, arXiv:0901.0512).

# Confidence intervals by inverting a test

Confidence intervals for a parameter $\theta$ can be found by defining a test of the hypothesized value $\theta$ (do this for all $\theta$):

Specify values of the data that are 'disfavoured' by $\theta$ (critical region) such that $P$(data in critical region) $\leq \gamma$ for a prespecified $\gamma$, e.g., 0.05 or 0.1.

If data observed in the critical region, reject the value $\theta$.

Now invert the test to define a confidence interval as:

set of $\theta$ values that would not be rejected in a test of size $\gamma$ (confidence level is $1 - \gamma$).

The interval will cover the true value of $\theta$ with probability $\geq 1 - \gamma$.

Equivalent to confidence belt construction; confidence belt is acceptance region of a test.

# Relation between confidence interval and *p*-value

Equivalently we can consider a significance test for each hypothesized value of $\theta$, resulting in a *p*-value, $p_\theta$.

If $p_\theta < \gamma$, then we reject $\theta$.

The confidence interval at CL $= 1 - \gamma$ consists of those values of $\theta$ that are not rejected.

E.g. an upper limit on $\theta$ is the greatest value for which $p_\theta \geq \gamma$.

In practice find by setting $p_\theta = \gamma$ and solve for $\theta$.