# Lecture Plan

**Statistics basic concepts** (Monday/Tuesday)

    **Basic ingredients** (PDFs, etc.)

    **Parameter estimation** (maximum likelihood, least-squares, …)

    **Model testing** ($\chi^2$ tests, hypothesis testing, p-values, …)

**These lectures: Computing statistical results**

    **Statistical modeling**

    **Review of model testing**

    **Computing results**

        **Confidence intervals**

        **Discovery**

        **Upper limits**

    **Systematics and profiling**

    **Bayesian techniques**

See also the Hands-on tutorial yesterday covering both sets of lectures.
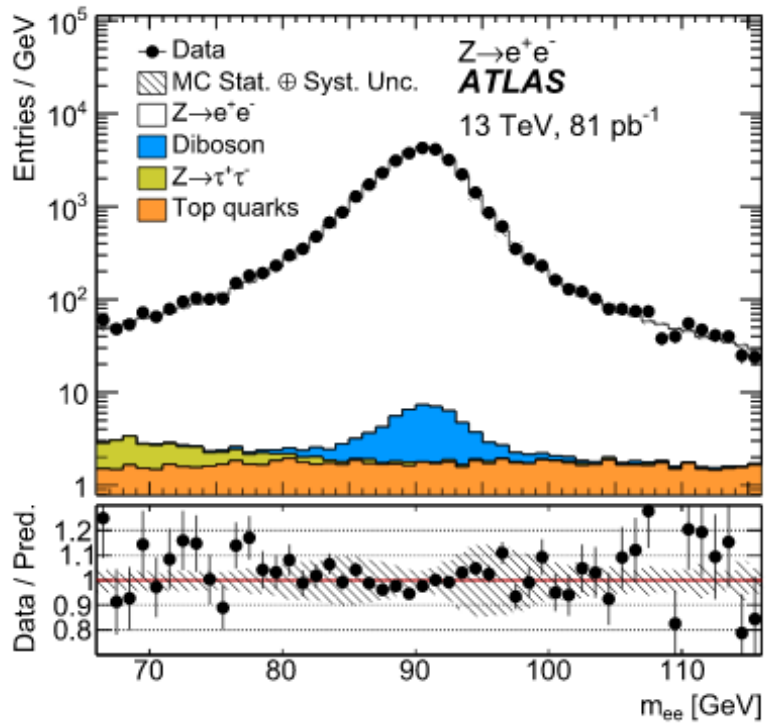
# Statistical Modeling

# Example 1: Z counting

Measure the cross-section (event rate) of the Z→ ee process

**35000 ± 187**

**175 ± 8**

$$\sigma^{fid} = \frac{n_{data} - N_{bkg}}{C_{fid}\, L}$$
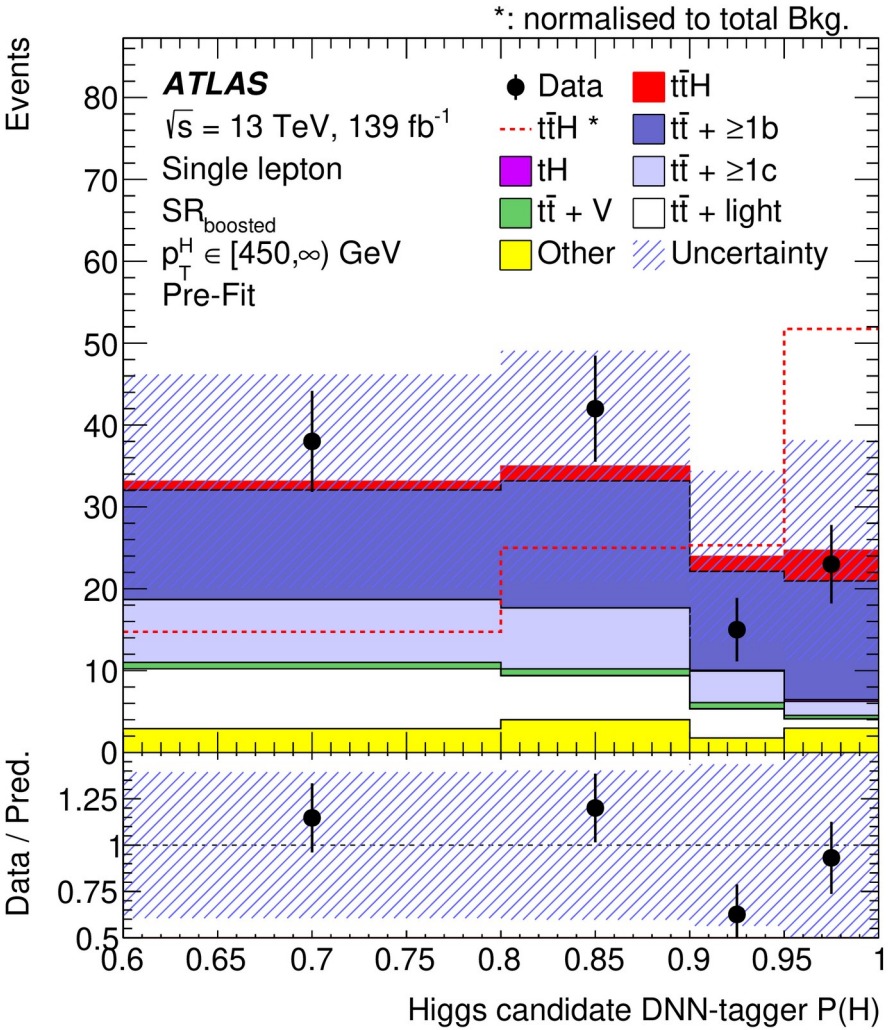
**(81 ± 2) pb$^{-1}$**

**0.552 ± 0.006**



$\sigma^{fid} = 0.781 \ \pm 0.004 \ \text{(stat)} \ \pm 0.018 \ \text{(syst) nb}$

Fluctuations in the data counts

Other uncertainties (assumptions, parameter values)

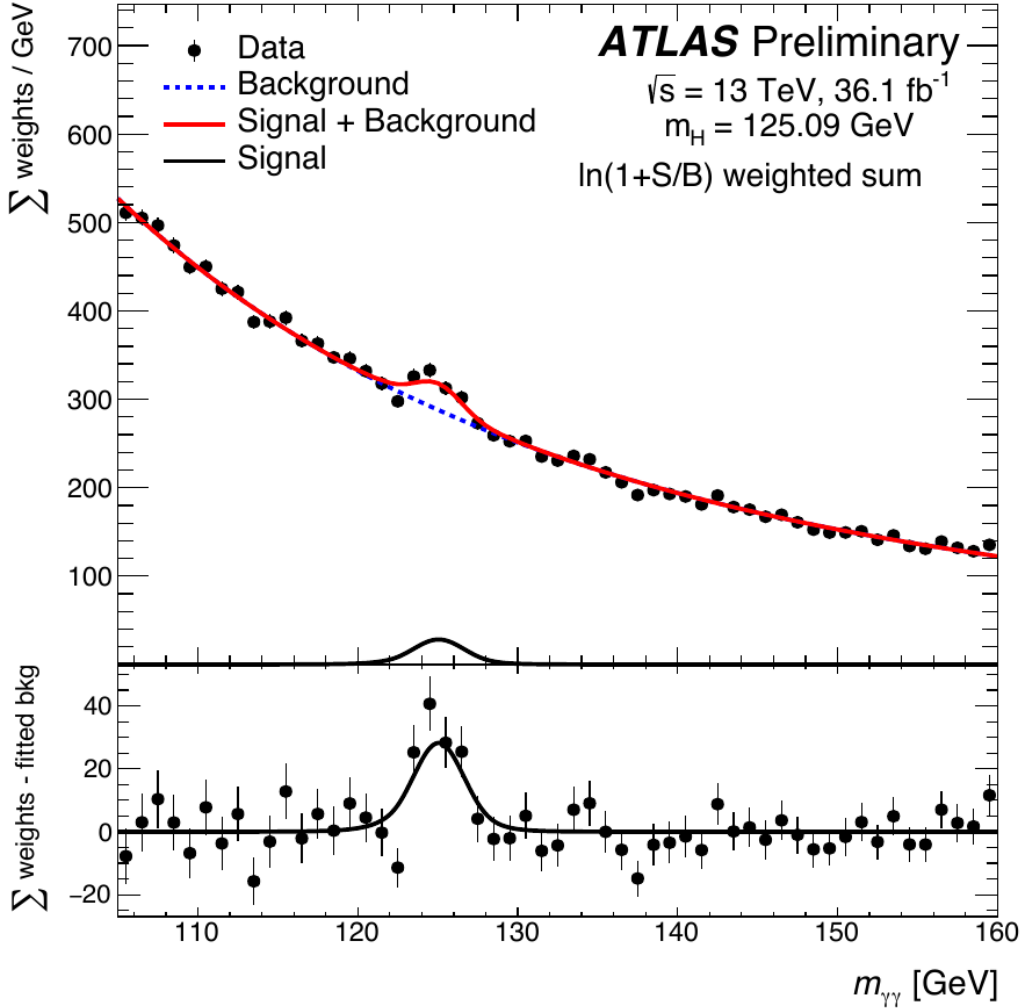**"Single bin counting"** : only data input is $N_{data}$.

# Example 2: ttH→bb

Event counting in different regions:

*Multiple-bin counting*

**Lots of information available**

→ Potentially higher sensitivity

→ How to make optimal use of it ?

# Example 3: unbinned modeling

All modeling done using continuous distributions:

$$P_{\text{total}}(m_{\gamma\gamma}) = \frac{S}{S+B}\, P_{\text{signal}}(m_{\gamma\gamma}; m_H) + \frac{B}{S+B}\, P_{\text{bkg}}(m_{\gamma\gamma})$$

# How to count

Common situation: produce many events N, select a (very) small fraction P

→ In principle, binomial process

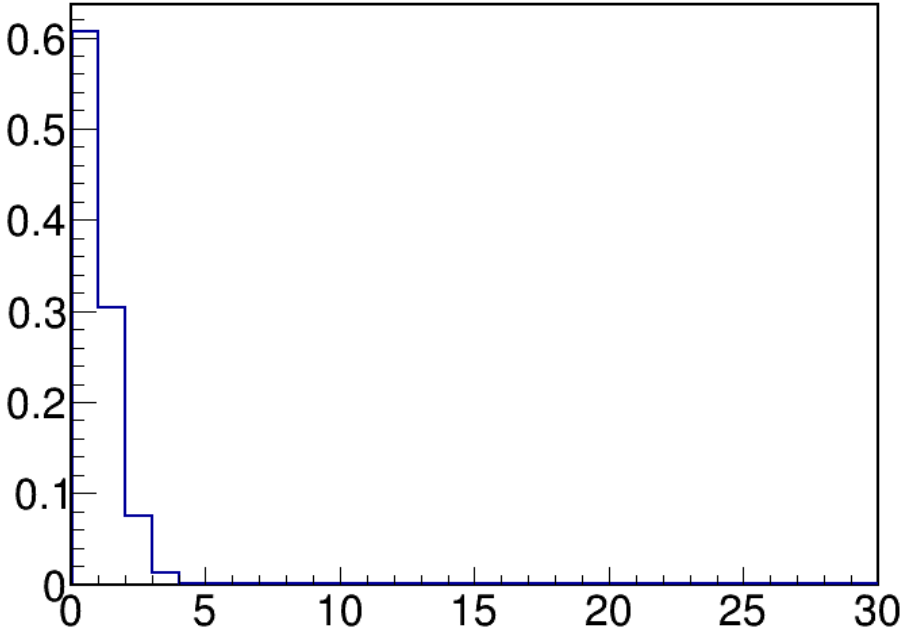→ In practice, $P \ll 1, N \gg 1$, $\Rightarrow$ Poisson approximation.

→ *i.e.* **very rare** process, but **very many trials** so still expect to see good events

**Poisson distribution**

$$P(n;\lambda) = e^{-\lambda} \frac{\lambda^n}{n!}$$

$$(1-P)^{N-n} \overset{n \ll N}{\sim} \left(1 - \frac{\lambda}{N}\right)^N \overset{N \gg 1}{\sim} e^{-\lambda}$$

λ = 0.5



**Mean** = λ

**Variance** = λ

σ = √λ

Central limit theorem :

becomes **Gaussian for large λ** :

$$P(\lambda) \overset{\lambda \to \infty}{\to} G(\lambda, \sqrt{\lambda})$$

# How to count

Common situation: produce many events N, select a (very) small fraction P
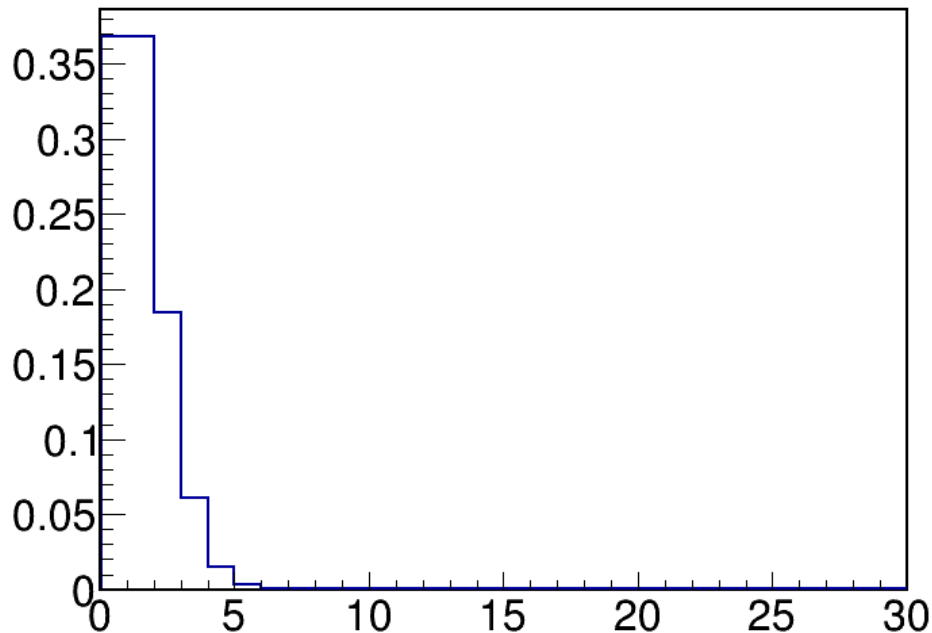
→ In principle, binomial process

→ In practice, **P ≪ 1, N ≫ 1**, ⇒ Poisson approximation.

→ *i.e.* **very rare** process, but **very many trials** so still expect to see good events

**Poisson distribution**

$$P(n;\lambda)=e^{-\lambda}\frac{\lambda^n}{n!}$$

$$(1-P)^{N-n} \overset{n\ll N}{\sim} \left(1-\frac{\lambda}{N}\right)^N \overset{N\gg 1}{\sim} e^{-\lambda}$$



λ = 1

**Mean** = λ

**Variance** = λ

σ = √λ

Central limit theorem :

becomes **Gaussian for large λ** :

$$P(\lambda)\overset{\lambda\to\infty}{\to}G(\lambda,\sqrt{\lambda})$$

# How to count

Common situation: produce many events N, select a (very) small fraction P
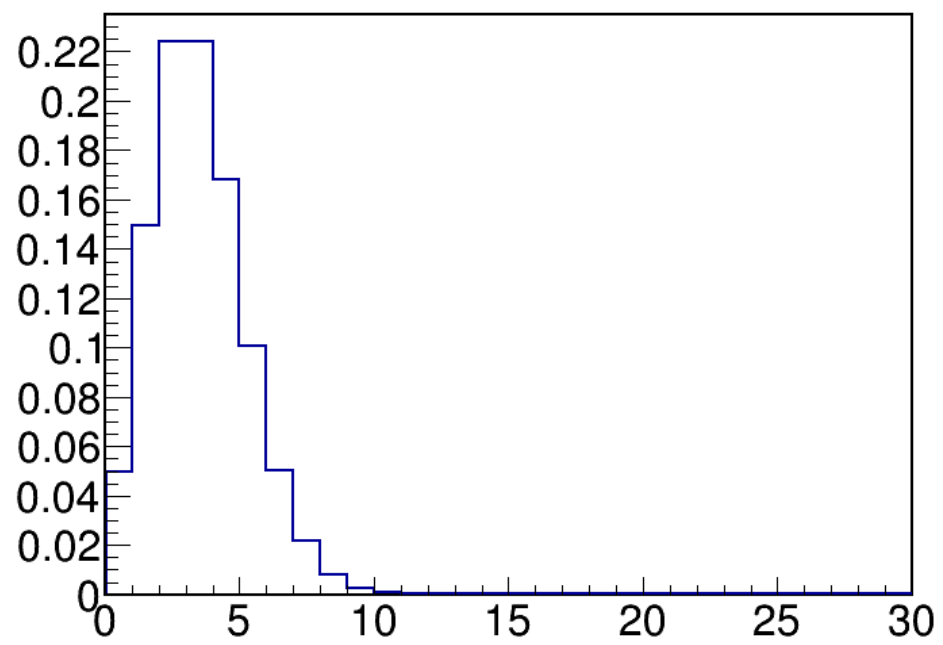
→ In principle, binomial process

→ In practice, **P ≪ 1, N ≫ 1**, ⇒ Poisson approximation.

→ *i.e.* **very rare** process, but **very many trials** so still expect to see good events

**Poisson distribution**
$$P(n;\lambda) = e^{-\lambda}\frac{\lambda^n}{n!}$$

$$(1-P)^{N-n} \overset{n \ll N}{\sim} \left(1-\frac{\lambda}{N}\right)^N \overset{N \gg 1}{\sim} e^{-\lambda}$$

λ = 3

**Mean** = λ

**Variance** = λ

σ = √λ

Central limit theorem :

becomes **Gaussian for large λ** :

$$P(\lambda) \overset{\lambda \to \infty}{\to} G(\lambda, \sqrt{\lambda})$$

# How to count

Common situation: produce many events N, select a (very) small fraction P

→ In principle, binomial process

→ In practice, **P ≪ 1, N ≫ 1**, ⇒ Poisson approximation.

➔ *i.e.* **very rare** process, but **very many trials** so still expect to see good events

**Poisson distribution**

$$P(n;\lambda) = e^{-\lambda} \frac{\lambda^n}{n!}$$

$$(1-P)^{N-n} \overset{n \ll N}{\sim} \left(1 - \frac{\lambda}{N}\right)^N \overset{N \gg 1}{\sim} e^{-\lambda}$$

λ = 5



**Mean** = λ

**Variance** = λ

σ = √λ

Central limit theorem :

becomes **Gaussian for large λ** :

$$P(\lambda) \overset{\lambda \to \infty}{\twoheadrightarrow} G(\lambda, \sqrt{\lambda})$$

# How to count

Common situation: produce many events N, select a (very) small fraction P

→ In principle, binomial process

→ In practice, $P \ll 1$, $N \gg 1$, $\Rightarrow$ Poisson approximation.

→ *i.e.* **very rare** process, but **very many trials** so still expect to see good events

**Poisson distribution**

$$P(n;\lambda) = e^{-\lambda} \frac{\lambda^n}{n!}$$

$$(1-P)^{N-n} \overset{n \ll N}{\sim} \left(1 - \frac{\lambda}{N}\right)^N \overset{N \gg 1}{\sim} e^{-\lambda}$$

λ = 10



**Mean** = λ

**Variance** = λ

σ = √λ

Central limit theorem :

becomes **Gaussian for large λ** :

$$P(\lambda) \overset{\lambda \to \infty}{\to} G(\lambda, \sqrt{\lambda})$$

# How to count

Common situation: produce many events N, select a (very) small fraction P

→ In principle, binomial process

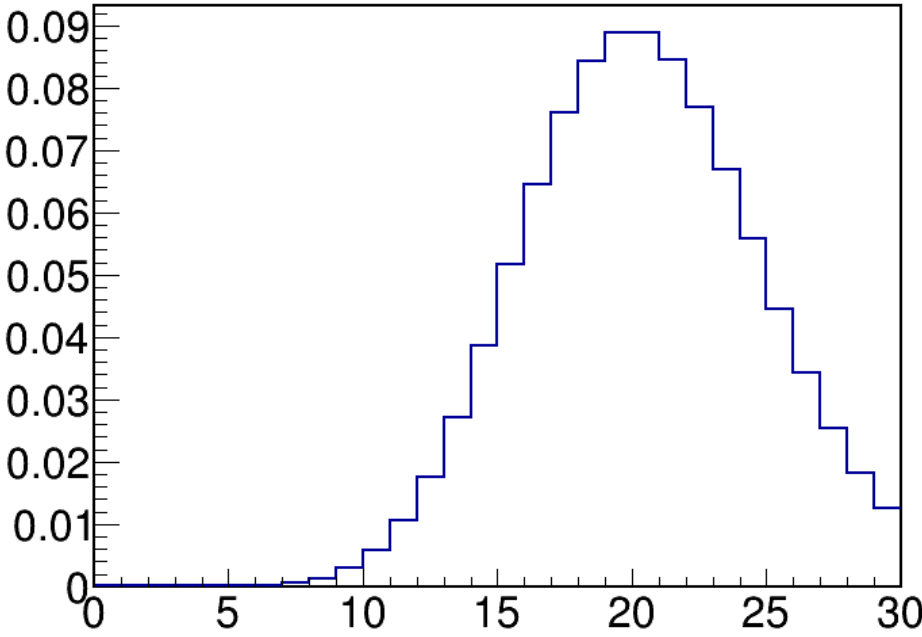→ In practice, **P ≪ 1, N ≫ 1**, ⇒ Poisson approximation.

→ *i.e.* **very rare** process, but **very many trials** so still expect to see good events

**Poisson distribution**

$$P(n;\lambda) = e^{-\lambda}\frac{\lambda^n}{n!}$$

$$(1-P)^{N-n} \overset{n \ll N}{\sim} \left(1-\frac{\lambda}{N}\right)^N \overset{N \gg 1}{\sim} e^{-\lambda}$$

λ = 20



**Mean** = λ

**Variance** = λ

σ = √λ

Central limit theorem :

becomes **Gaussian for large λ** :

$$P(\lambda) \overset{\lambda \to \infty}{\to} G(\lambda, \sqrt{\lambda})$$
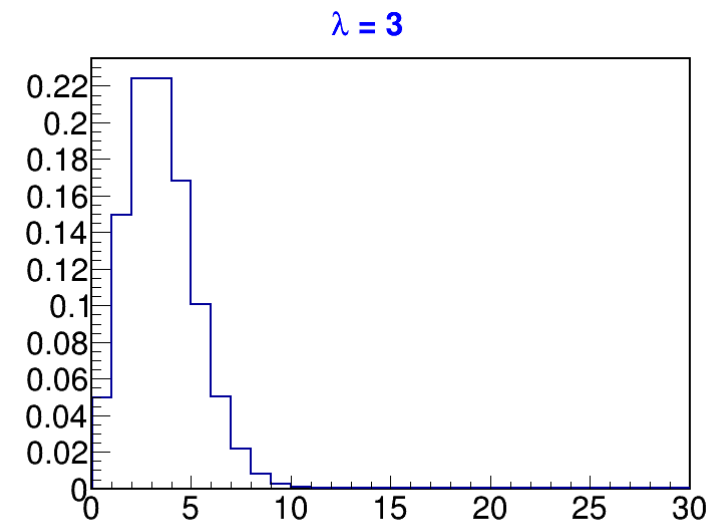
# Statistical Model for Counting

**Observable: number of events n**

Typically both **S**ignal and **B**ackground present:

$$P(n;S,B)=e^{-(S+B)}\frac{(S+B)^n}{n!}$$

**S** : # of events from signal process

**B** : # of events from bkg. process(es)

Model has **parameters S** and **B.**

B can be known a priori or not (S usually not...)

→ Example: **assume B is known,** use **measured n** to find out about **S.**
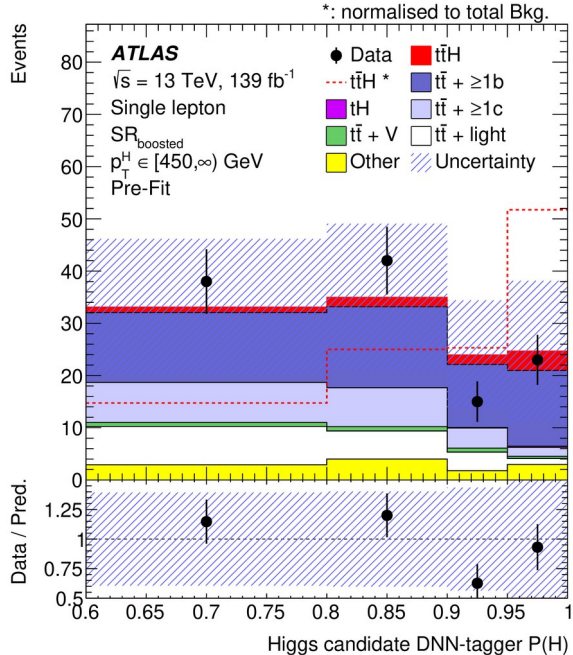
# Multiple counting bins

Count in bins of a variable $\Rightarrow$ *histogram* $n_1 \ldots n_N$.
(N : number of bins)

**Per-bin fractions (=shapes) of Signal and Background**

$$P(\{n_i\}; S, B) = \prod_{i=1}^{N} e^{-(Sf_{S,i} + Bf_{B,i})} \frac{(Sf_{S,i} + Bf_{B,i})^{n_i}}{n_i!}$$

**Poisson distribution in each bin**



**Shapes f** typically obtained from simulated events (*Monte Carlo*)

$\rightarrow$ HEP: typically excellent modeling from simulation, although some uncertainties need to be accounted for.

However not always possible to generate sufficiently large MC samples
**MC stat fluctuations** can create artefacts, especially for $S \ll B$.

# Model Parameters

Model typically includes:

- **Parameters of interest** (POIs) : what we want to measure

    $\rightarrow$ **S**, **m$_W$**, ...

- **Nuisance parameters** (NPs) : other parameters needed to define the model

    $\rightarrow$ Background levels (**B**)

    $\rightarrow$ For binned data, **f$^{sig}_i$** , **f$^{bkg}_i$**

NPs must be either:

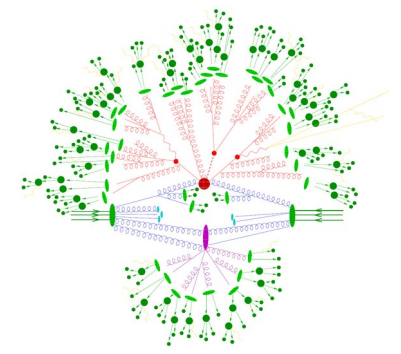$\rightarrow$ **Known a priori** (within uncertainties) or

$\rightarrow$ **Constrained by the data**

# Takeaways

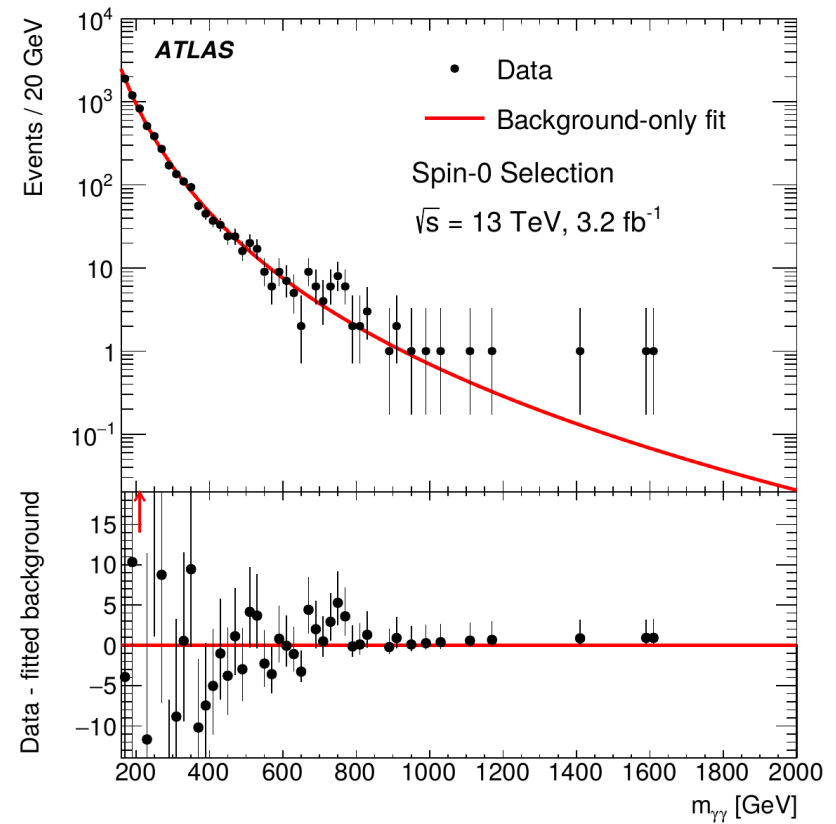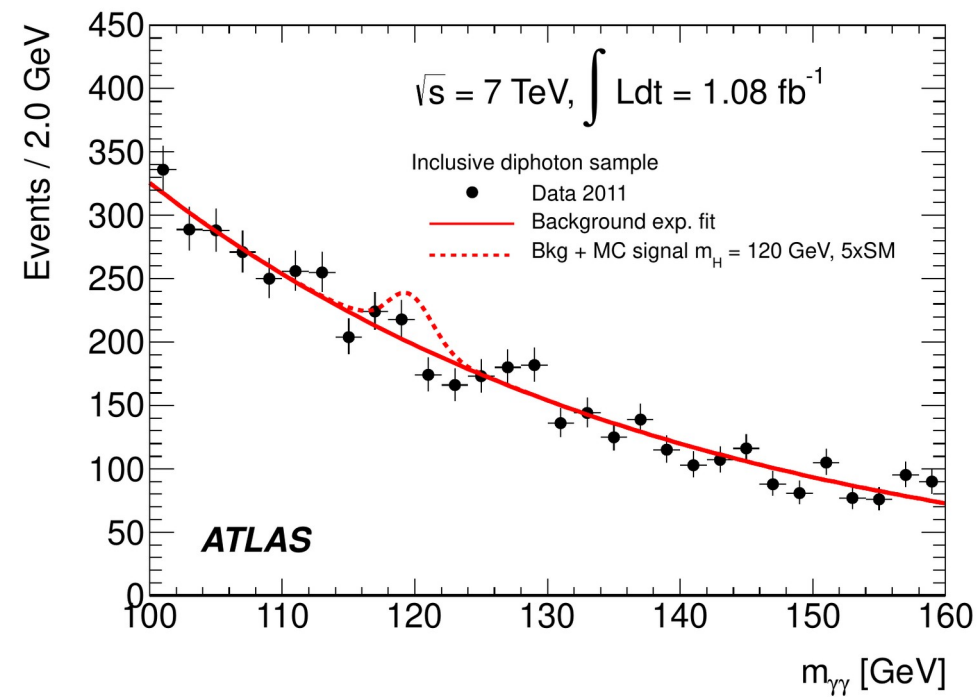Random data must be described using a statistical model:

| Description | Observable | Likelihood |
|---|---|---|
| Counting | n | **Poisson** $$P(n;S,B)=e^{-(S+B)}\frac{(S+B)^n}{n!}$$ |
| Binned shape analysis | $n_i$, i = 1 .. $N_{bins}$ | **Poisson product** $$P(n_i;S,B)=\prod_{i=1}^{n_{bins}} e^{-(S f_i^{sig} + B f_i^{bkg})}\frac{(S f_i^{sig} + B f_i^{bkg})^{n_i}}{n_i!}$$ |
| Unbinned shape analysis | $m_i$, i = 1 .. $n_{evts}$ | **Extended Unbinned Likelihood** $$P(m_i;S,B)=\frac{e^{-(S+B)}}{n_{evts}!}\prod_{i=1}^{n_{evts}} S\, P_{sig}(m_i)+B\, P_{bkg}(m_i)$$ |

Model can include multiple **categories**, each with a separate description

Includes **parameters of interest** (POIs) but also **nuisance parameters** (NPs)
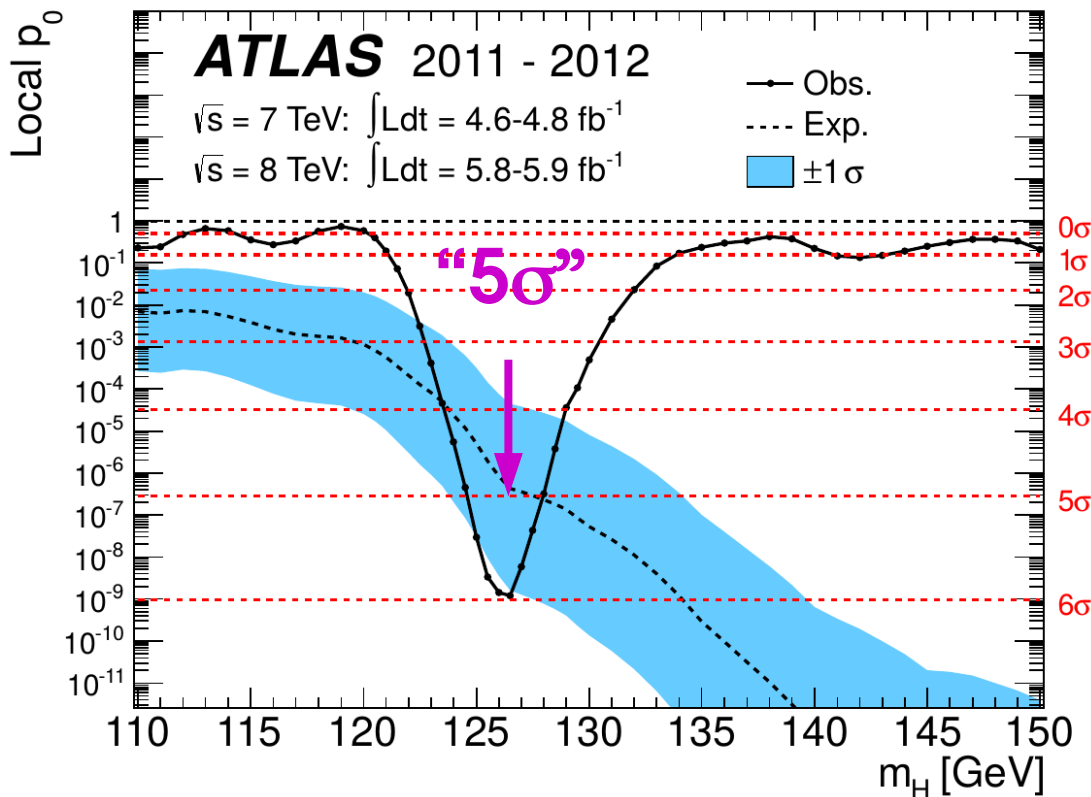
**Next step**: use the model to obtain information on the POIs
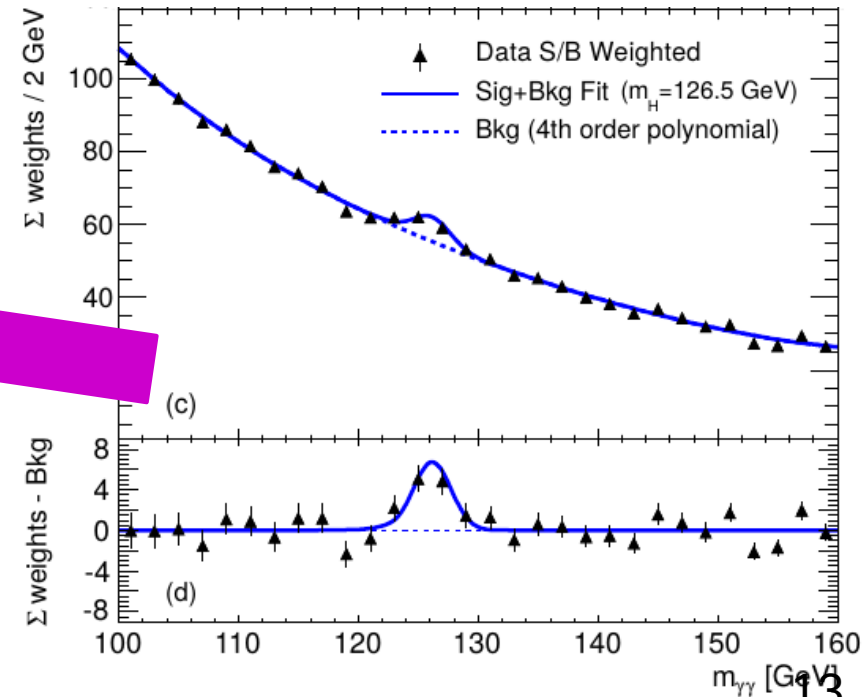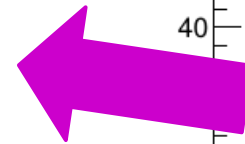
# Hypothesis Testing and discovery

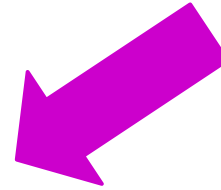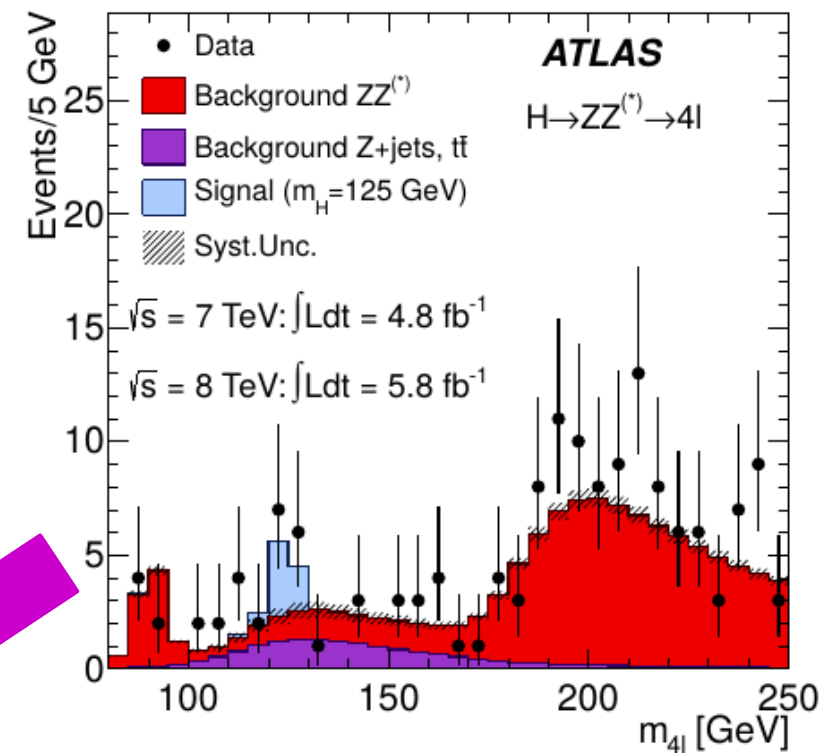# Discovery Testing

We see an unexpected feature in our data, is it a signal for new physics or a fluctuation ?
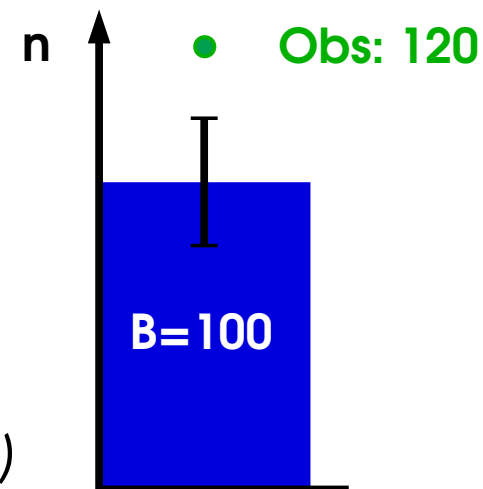
e.g. Higgs discovery :  **"We have 5σ" !**



Phys. Lett. B 716 (2012) 1-29

# Discovery Testing

Say we have a Gaussian measurement with a background **B=100**, and we measure **n=120**

Did we just discover something ? *Maybe :-) (but not very likely)*


Obs: 120
B=100

The measured signal is **S = 20**.     $S = n_{obs} - B$

Uncertainty on B is $\sqrt{B} = 10$

⇒ Significance Z = 2

⇒ we are **~2σ** away from S=0.

$$Z = \frac{S}{\sqrt{B}}$$


Obs: 120
$\sqrt{B}=10$
B=100     n

**Gaussian quantiles** :

Z = 2 happens **$p_0$ ~ 2.3%** of the time if S=0

**P-value:**     $$p_0 = 1 - \Phi(Z)$$

$$\Phi(Z) = \int_{-\infty}^{Z} G(u;0,1) \, du$$

⇒ Rare, but not exceptional

# Discovery Testing



| $n_{obs}$ | S | Z | $p_0$ | |
|---|---|---|---|---|
| 105 | 5 | 0.5σ | 31% | |
| 110 | 10 | 1σ | 16% | |
| 120 | 20 | 2σ | 2.3% | |
| 130 | 30 | 3σ | 0.1% | Evidence |
| 150 | 50 | 5σ | 3 10⁻⁷ | Discovery |

Straightforward in this Gaussian case

Need to be able to do the same in more complex cases:

- **Determine S**
- **Compute Z and $p_0$**

# What is PDF is for

Model describes the distribution of the observable: **P(<span style="color:red">data</span>; <span style="color:blue">parameters</span>)**

⇒ Possible outcomes of the experiment, for given parameter values

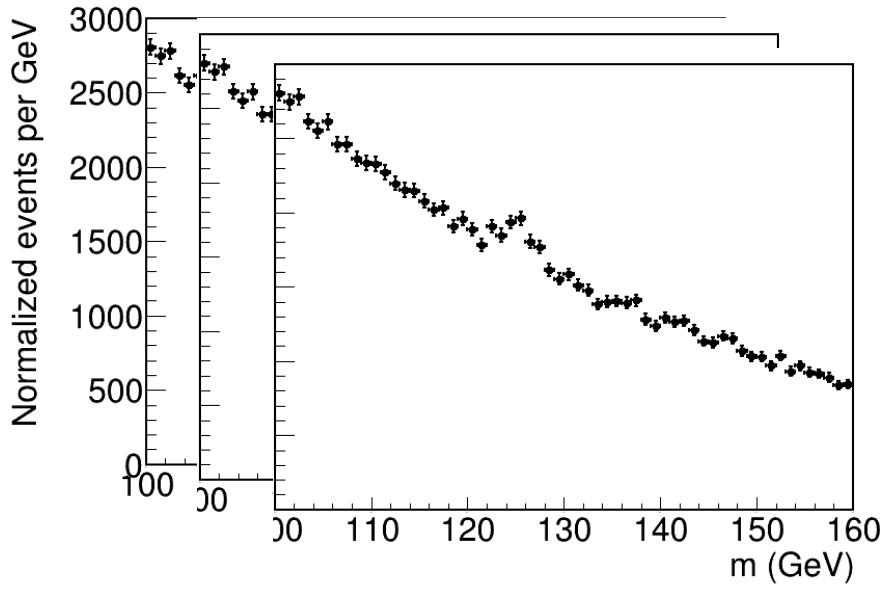Can draw random events according to PDF : **generate** *pseudo-data*

$$P(\lambda = 5)$$   ⟹   **2, 5, 3, 7, 4, 9, ....**

Each entry = separate "experiment"

**Generate**



**Unbinned**

# What is PDF is also for: Likelihood

Model describes the distribution of the observable: **P(data; parameters)**

⇒ Possible outcomes of the experiment, for given parameter values

We want the **other** direction: **use data to get information on parameters**

$$P(\lambda = ?)$$

**2**

**Estimate**



**Likelihood**: L(parameters) = P(data; parameters)

→ **same as the PDF,** but seen as function of the parameters

# Maximum Likelihood Estimation

To estimate a parameter μ, find the **value μ̂ that maximizes** L(μ)

**Maximum Likelihood Estimator (MLE) μ̂:**

$$\hat{\mu} = arg\,max\,L(\mu)$$



**MLE**: the value of μ for which **this data** was *most likely to occur*

**The MLE is a function of the data** – itself an **observable**

*No guarantee* it is the true value (data may be "unlikely") but sensible estimate

# Gaussian case

**data**



**Best-fit** of Gaussian PDF mean to observed data

# Gaussian case



Best-fit of Gaussian PDF mean to observed data

# Gaussian case



**Best-fit** of Gaussian PDF mean to observed data

# Multiple Gaussian bins



-2 log Likelihood:

$$\lambda(\mu) = -2\log L(\mu) = \sum_{i=1}^{N_{\text{bins}}} \left(\frac{n_i - \mu_i}{\sigma_i}\right)^2$$

**Maximum likelihood** $\Leftrightarrow$ Minimum $\chi^2$

$\Leftrightarrow$ Least-squares

minimization

However typically need to perform non-linear minimization.

HEP practice:

- **MINUIT** (C++ library within ROOT, numerical gradient descent)
- **scipy.minimize** – using NumPy/TensorFlow/PyTorch/... backends
  - $\rightarrow$ Usual methods – gradient-based, etc.

# Multiple Gaussian bins



-2 log Likelihood:

$$\lambda(\mu) = -2\log L(\mu) = \sum_{i=1}^{N_{\text{bins}}} \left(\frac{n_i - \mu_i}{\sigma_i}\right)^2$$

**Maximum likelihood** $\Leftrightarrow$ Minimum $\chi^2$

$\Leftrightarrow$ Least-squares

minimization

However typically need to perform non-linear minimization.

HEP practice:

- **MINUIT** (C++ library within ROOT, numerical gradient descent)
- **scipy.minimize** – using NumPy/TensorFlow/PyTorch/… backends
  $\rightarrow$ Usual methods – gradient-based, etc.

# Multiple Gaussian bins



-2 log Likelihood:

$$\lambda(\mu) = -2 \log L(\mu) = \sum_{i=1}^{N_{bins}} \left(\frac{n_i - \mu_i}{\sigma_i}\right)^2$$

**Maximum likelihood** $\Leftrightarrow$ Minimum $\chi^2$

$\Leftrightarrow$ Least-squares

minimization

However typically need to perform non-linear minimization.

HEP practice:
- **MINUIT** (C++ library within ROOT, numerical gradient descent)
- **scipy.minimize** – using NumPy/TensorFlow/PyTorch/… backends

    → Usual methods – gradient-based, etc.

# Multiple Gaussian bins



-2 log Likelihood:

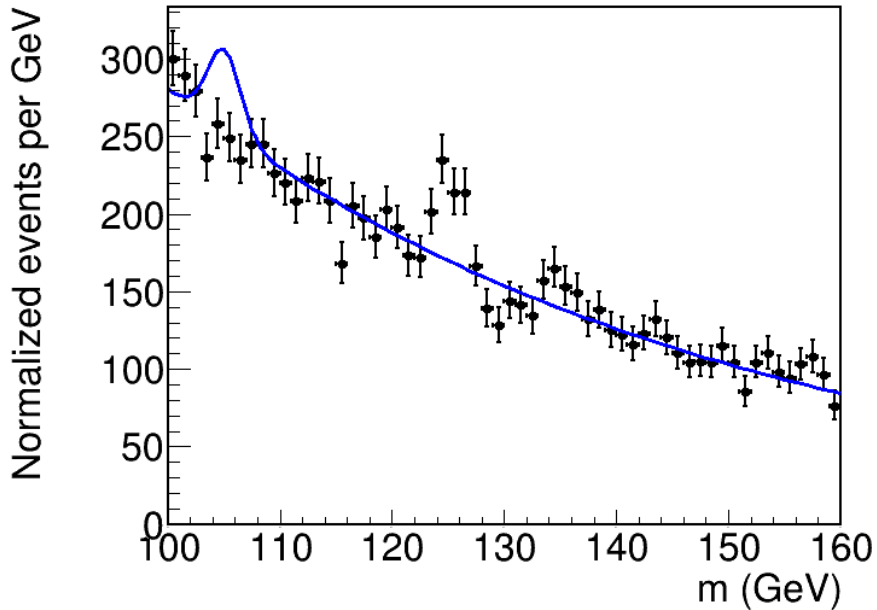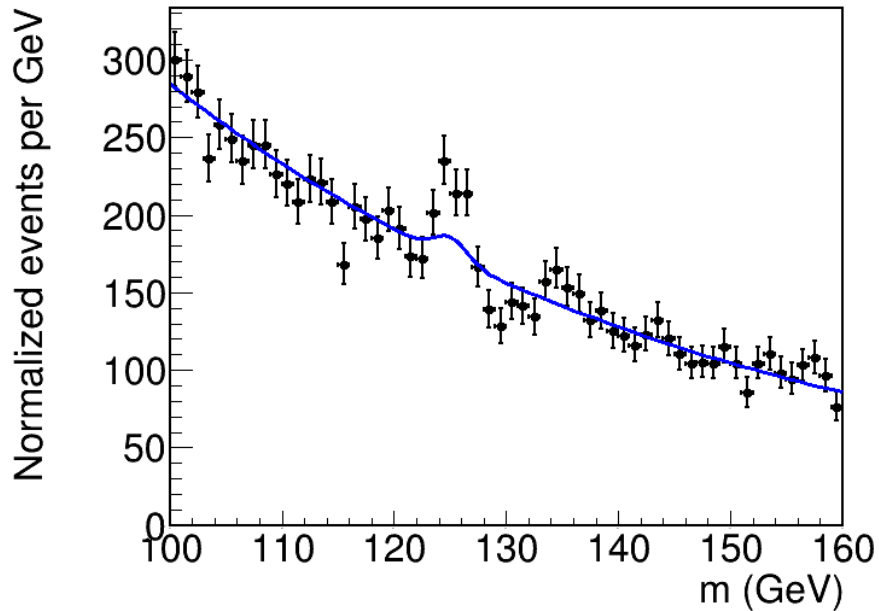$$\lambda(\mu) = -2\log L(\mu) = \sum_{i=1}^{N_{bins}} \left(\frac{n_i - \mu_i}{\sigma_i}\right)^2$$
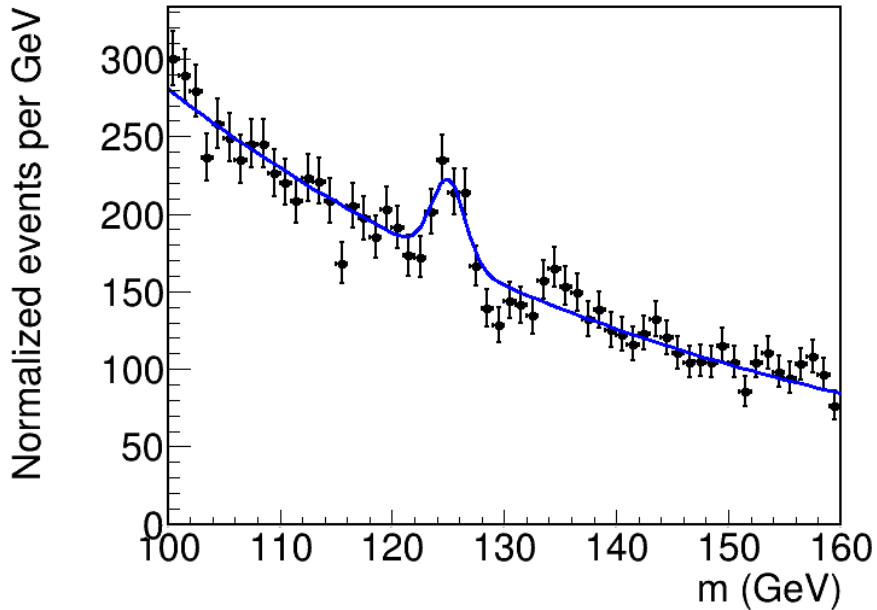
**Maximum likelihood** $\Leftrightarrow$ Minimum $\chi^2$

$\Leftrightarrow$ Least-squares

minimization

However typically need to perform non-linear minimization.

HEP practice:
- **MINUIT** (C++ library within ROOT, numerical gradient descent)
- **scipy.minimize** – using NumPy/TensorFlow/PyTorch/... backends
  $\rightarrow$ Usual methods – gradient-based, etc.

# Hypothesis Testing

**Null Hypothesis**: assumption on POIs, say value of S (e.g. $H_0 : S=0$)

→ **Goal** : decide if $H_0$ is favored or disfavored using a test based on the data

| Possible outcomes: | Data disfavors $H_0$ (Discovery claim) | Data favors $H_0$ (Nothing found) |
|---|---|---|
| $H_0$ is false (New physics!) | Discovery! | Missed discovery |
| $H_0$ is true (Nothing new) | False discovery | No new physics, None found |

*"... the null hypothesis is never proved or established, but is possibly disproved, in the course of experimentation. Every experiment may be said to exist only to give the facts a chance of disproving the null hypothesis." – R. A. Fisher*

# Hypothesis Testing

**Hypothesis**: assumption on model parameters, say value of S (e.g. $H_0 : S=0$)

| | Data disfavors $H_0$ (Discovery claim) | Data favors $H_0$ (Nothing found) |
|---|---|---|
| $H_0$ is false (New physics!) | Discovery! | Type-II error (Missed discovery) |
| $H_0$ is true (Nothing new) | Type-I error (False discovery) | No new physics, none found |

**p-value, significance**

**Lower Type-I errors** ⇔ **Higher Type-II errors** and vice versa: cannot have everything!

→ **Goal**: test that minimizes Type-II errors **for a given level of Type-I error**.



22 /

# ROC Curves

**"Receiver operating characteristic" (ROC) Curve:**

→ Shows Type-I vs Type-II rates for different selections

→ All curves monotonically decrease from (0,1) to (1,0)

→ Better discriminators more bent towards (1,1)

Better

Better

No discrimination

$1 - \varepsilon_{\text{Type-I}} \ (= 1 - \varepsilon_B)$

$1 - \varepsilon_{\text{Type-II}} \ (= \varepsilon_S)$

0

1

1

→ **Goal**: test that minimizes Type-II errors **for given level of Type-I error**.

→ Usually set predefined level of **acceptable Type-I error** (e.g. "5σ")

S = 0

BSM

Type-II Error

Type-I error p-value

Discriminant observable

23 /

# ROC Curves

**"Receiver operating characteristic" (ROC) Curve:**

→ Shows Type-I vs Type-II rates for different selections

→ All curves monotonically decrease from (0,1) to (1,0)

→ Better discriminators more bent towards (1,1)

Better

No discrimination

Better

$1 - \varepsilon_{Type-I}$ $(= 1 - \varepsilon_B)$

$1 - \varepsilon_{Type-II}$ $(= \varepsilon_S)$

0

1

1

→ **Goal**: test that minimizes Type-II errors **for given level of Type-I error.**

→ Usually set predefined level of **acceptable Type-I error** (e.g. "5σ")

S = 0

BSM

Type-II Error

Type-I error p-value

Discriminant observable

23 /

# ROC Curves

**"Receiver operating characteristic" (ROC) Curve:**

→ Shows Type-I vs Type-II rates for different selections
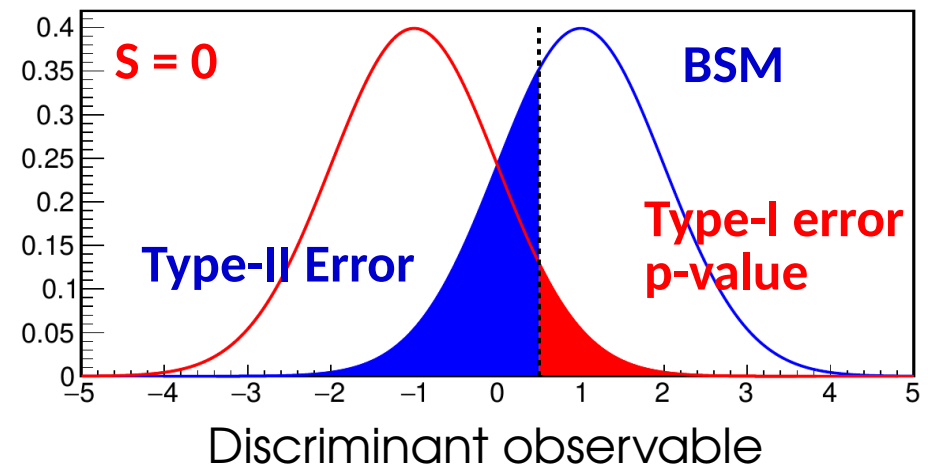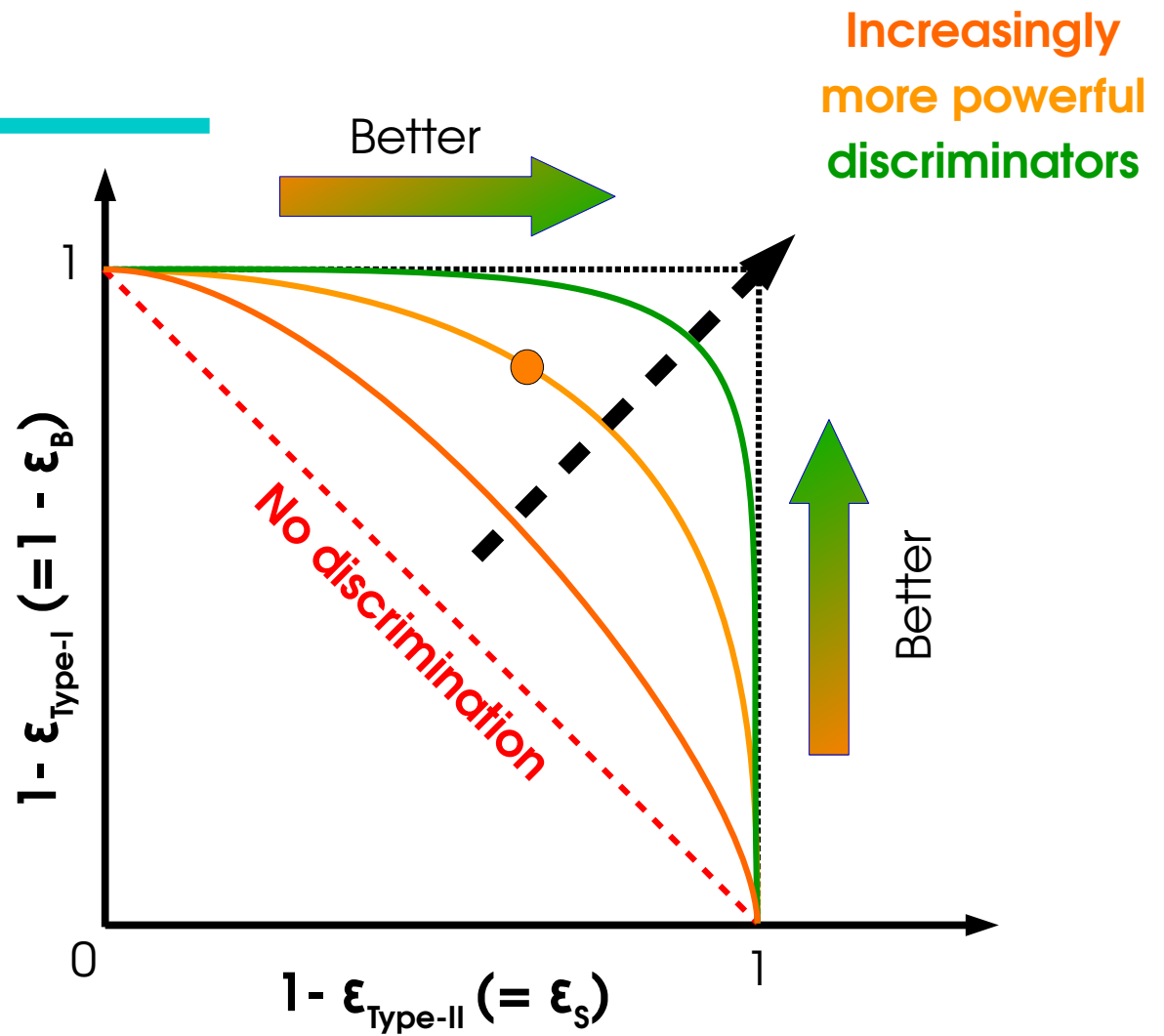
→ All curves monotonically decrease from (0,1) to (1,0)

→ Better discriminators more bent towards (1,1)

Better

$1 - \varepsilon_{\text{Type-I}}$ ($= 1 - \varepsilon_B$)

No discrimination

Better

$1 - \varepsilon_{\text{Type-II}}$ ($= \varepsilon_S$)

0                                                                1

→ **Goal**: test that minimizes Type-II errors **for given level of Type-I error**.

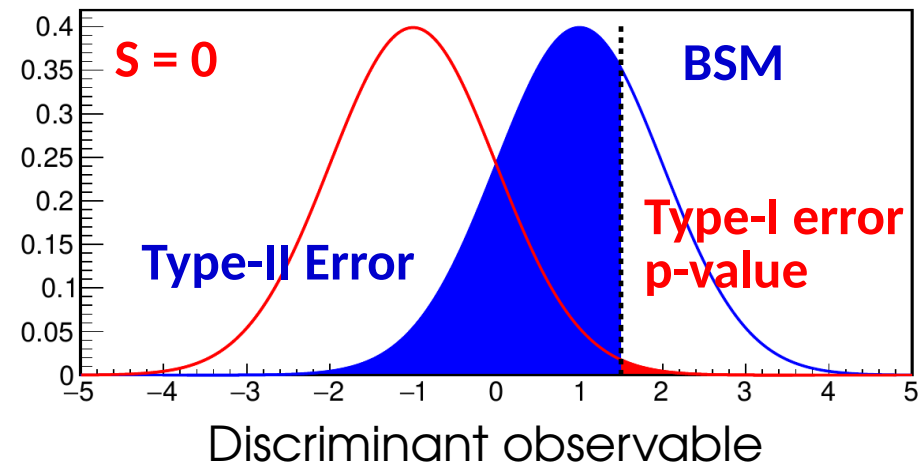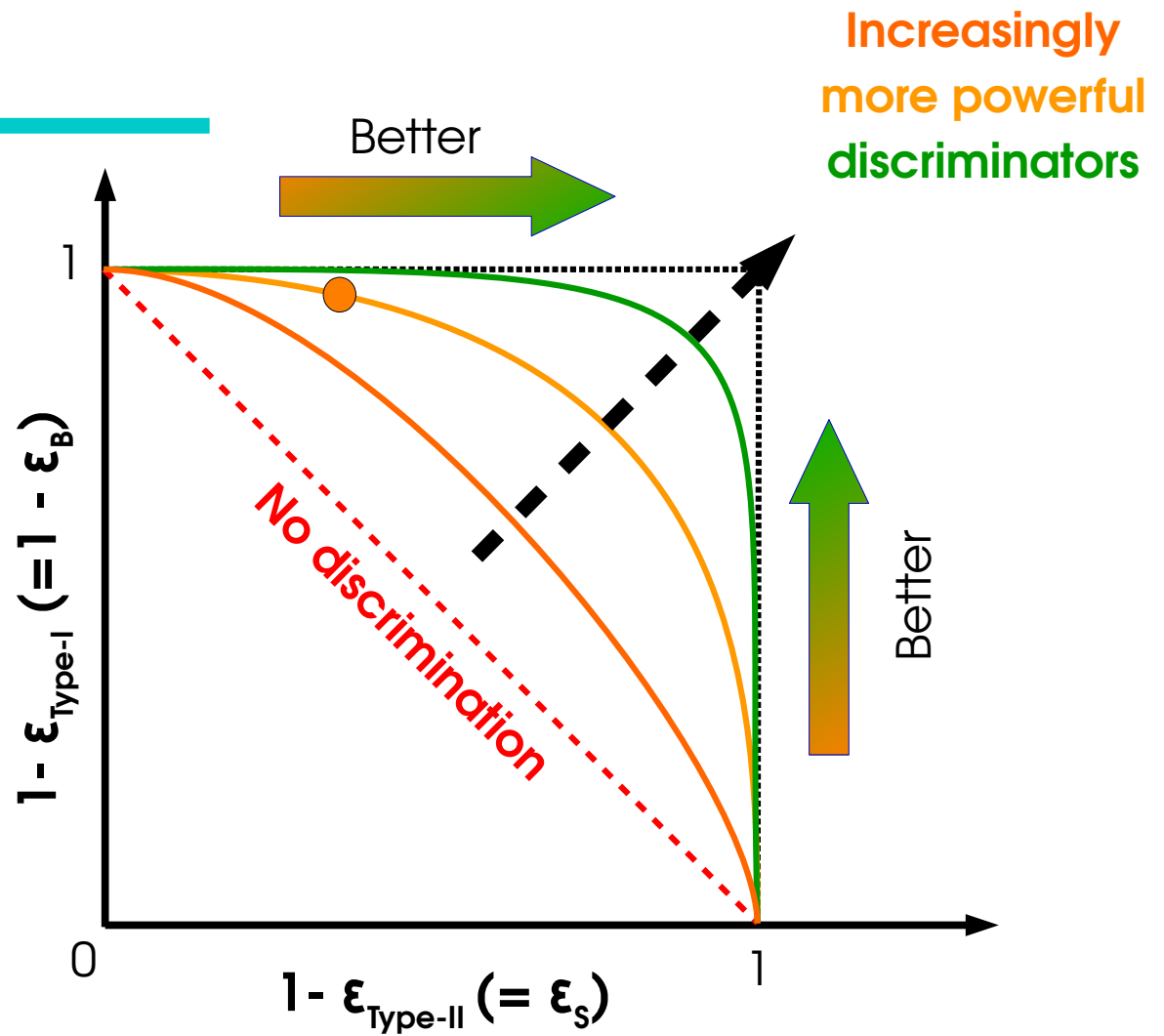→ Usually set predefined level of **acceptable Type-I error** (e.g. "5σ")

S = 0

BSM

Type-II Error

Type-I error p-value

Discriminant observable

−5  −4  −3  −2  −1  0  1  2  3  4  5

0.4  0.35  0.3  0.25  0.2  0.15  0.1  0.05  0

23 /

# Hypothesis Testing with Likelihoods

## Neyman-Pearson Lemma

*When comparing two hypotheses $H_0$ and $H_1$, the*

*optimal discriminator is the **Likelihood ratio** (LR)*

$$\frac{L(H_1 ; data)}{L(H_0 ; data)}$$

e.g. $\dfrac{L(S = 5 ; data)}{L(S = 0 ; data)}$

**Caveat**: Strictly true only for *simple hypotheses* (no free parameters)

As for MLE, choose the hypothesis that is more likely **given the data we have**.

→ **Minimizes Type-II uncertainties** for given level of Type-I uncertainties

→ Always need an **alternate hypothesis** to test against the **null**.

→ **In the following**: all tests based on LR, will focus on p-values (Type-I errors), trusting that Type-II errors are anyway as small as they can be...

# Discovery: Test Statistic

**Discovery :**

- $H_0$ : **background only (S = 0)** against

- $H_1$: presence of a signal (**S > 0**)

$S=0$

$H_0 \longrightarrow H_1$

$\rightarrow$ For $H_1$, any S > 0 is possible, which to use ? **The one preferred by the data, Ŝ.**

$\Rightarrow$ Use Likelihood ratio: $\quad \dfrac{L(S=0)}{L(\hat{S})}$

$\rightarrow$ In fact use the **test statistic** $\quad q_0 \;=\; -2\log\dfrac{L(S=0)}{L(\hat{S})}$

**Note**: for Ŝ < 0, set $q_0$=0 to reject negative signals (*"one-sided* test statistic")

# Discovery p-value

Large values of $-2\log\dfrac{L(S=0)}{L(\hat{S})}$ if:

$\Rightarrow$ observed $\hat{S}$ is far from 0

$\Rightarrow$ **H₀(S=0)** *disfavored* compared to **H₁(S≠0)**.

**How large q₀ before we can exclude H₀ ?**

(and **claim a discovery!**)

$\rightarrow$ Need small Type-I rate (falsely rejecting H₀)

$\rightarrow$ Type-I error rate, a.k.a. the *p-value* : $\boldsymbol{p_0 = \displaystyle\int_{q_0^{obs}}^{\infty} f(q_0|S=0)\,dq_0}$

= *Fraction of outcomes that are*

**At least as extreme** *(signal-like)* **as data**, *when H₀ is true (no signal).*

data prefer S = 0 ⟷ data prefer S > 0



$\hat{S} \leq 0$    $f(q_0|S=0)$

Observed value $q_0^{obs}$

large $\hat{S}$

$q_0$

# Asymptotic distribution of q₀

**Gaussian regime for Ŝ** (e.g. large $n_{evts}$, Central-limit theorem) :

*Wilk's Theorem*: $q_0$ distributed as $\chi^2$ ($n_{par}$) for S = 0

$\Rightarrow n_{par} = 1$ : $\sqrt{q_0}$ is distributed as a Gaussian

$\Rightarrow$ Can compute p-values from Gaussian quantiles

$$p_0 = 1 - \Phi\left(\sqrt{q_0}\right)$$

$\Rightarrow$ Even more simply, the significance is:

$$Z = \sqrt{q_0}$$

Typically works well already for for event counts of O(5) and above $\Rightarrow$ Widely applicable

(*) 1-line "proof" : asymptotically L and S are Gaussian, so

$$L(S) = \exp\left[-\frac{1}{2}\left(\frac{S-\hat{S}}{\sigma}\right)^2\right] \Rightarrow q_0 = \left(\frac{\hat{S}}{\sigma}\right)^2 \Rightarrow \sqrt{q_0} = \frac{\hat{S}}{\sigma} \sim G(0,1) \Rightarrow q_0 \sim \chi^2(n_{dof}=1)$$

# Homework 1: Gaussian Counting

**Count number of events n in data**

→ Assume n large enough so process is Gaussian

→ Assume B is known, and we measure S

**Likelihood :**
$$L(S; n_{obs}) = e^{-\frac{1}{2}\left(\frac{n_{obs}-(S+B)}{\sqrt{S+B}}\right)^2}$$



$\sqrt{(S+B)}$

$n_{obs}$

S+B

→ Find the best-fit value (MLE) Ŝ for the signal

  (can use λ = -2 log L instead of L for simplicity)

→ Find the expression of $q_0$ for Ŝ > 0.

→ Find the expression for the significance

$$Z = \frac{\hat{S}}{\sqrt{B}}$$

# Homework 2: Poisson Counting

Same problem but now **not** assuming Gaussian behavior:

$$L(S;n) = e^{-(S+B)}(S+B)^n$$

→ As before, compute $\hat{S}$, and $q_0$

(Can remove the n! constant since we're only dealing with L ratios)

→ Compute $Z = \sqrt{q_0}$, assuming asymptotic behavior

**Solution:**

$$Z = \sqrt{2\left[(\hat{S}+B)\log\left(1+\frac{\hat{S}}{B}\right) - \hat{S}\right]}$$

Exact result can be obtained using

pseudo-experiments → close to $\sqrt{q_0}$ result

**Asymptotic formulas justified by Gaussian regime, but remain valid even for small values of S+B (down to 5 events!)**



Eur.Phys.J.C71:1554,2011

See G. Cowan's slides for the case with B uncertainty

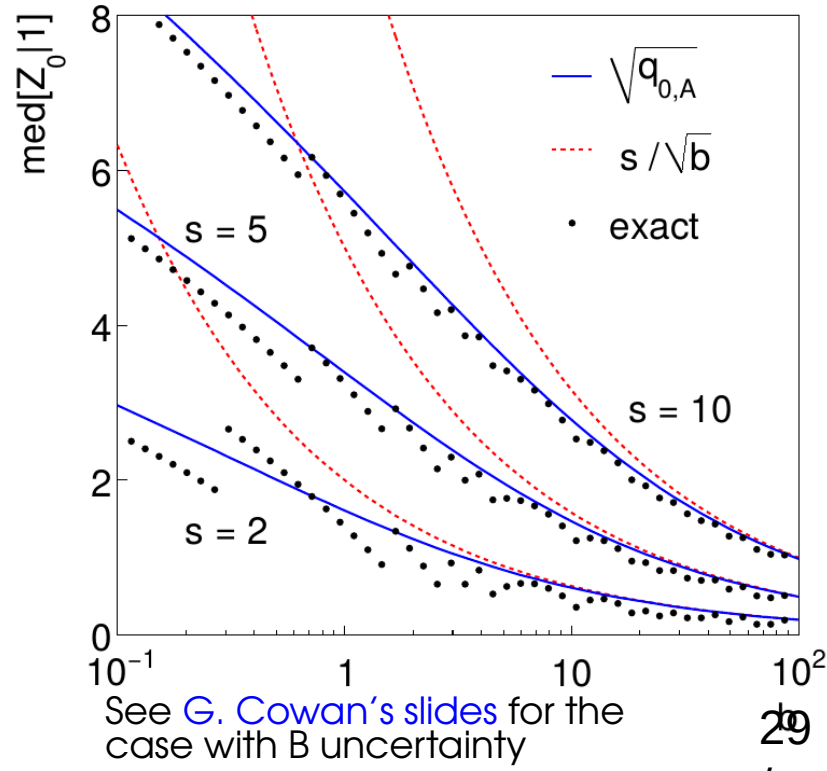# Discovery Thresholds

**Evidence** : $3\sigma \Leftrightarrow p_0 = 0.3\% \Leftrightarrow 1$ chance in 300

**Discovery**: $5\sigma \Leftrightarrow p_0 = 3 \cdot 10^{-7} \Leftrightarrow 1$ chance in 3.5M

**Why so high thresholds ?** (from Louis Lyons):

- **Look-elsewhere effect**: searches typically cover
  multiple independent regions $\Rightarrow$ Higher chance
  to have a fluctuation "somewhere"

  $N_{trials} \sim 1000$ : **local $5\sigma \Leftrightarrow O(10^{-4})$** more reasonable

- **Mismodeled systematics**: factor 2 error in
  syst-dominated analysis $\Rightarrow$ factor 2 error on Z...

- **History**: $3\sigma$ and $4\sigma$ excesses do occur regularly, for the reasons above

*Extraordinary claims require extraordinary evidence!*



30 /

# Takeaways

Given a statistical model P(data; μ), define likelihood **L(μ) = P(data; μ)**

**To estimate a parameter**, use the value **μ̂** that maximizes L(μ) → best-fit value

**To decide between hypotheses** $H_0$ and $H_1$, use the **likelihood ratio** $\dfrac{L(H_0)}{L(H_1)}$

To test for **discovery**, use $\quad q_0 = -2\log \dfrac{L(S=0)}{L(\hat{S})} \quad \hat{S} \geq 0$

For large enough datasets (n >~ 5), $\quad Z = \sqrt{q_0}$

For a **Gaussian** measurement, $\quad Z = \dfrac{\hat{S}}{\sqrt{B}}$

For a **Poisson** measurement, $\quad Z = \sqrt{2\left[(\hat{S}+B)\log\left(1+\dfrac{\hat{S}}{B}\right) - \hat{S}\right]}$

# Confidence Intervals

# Confidence Intervals

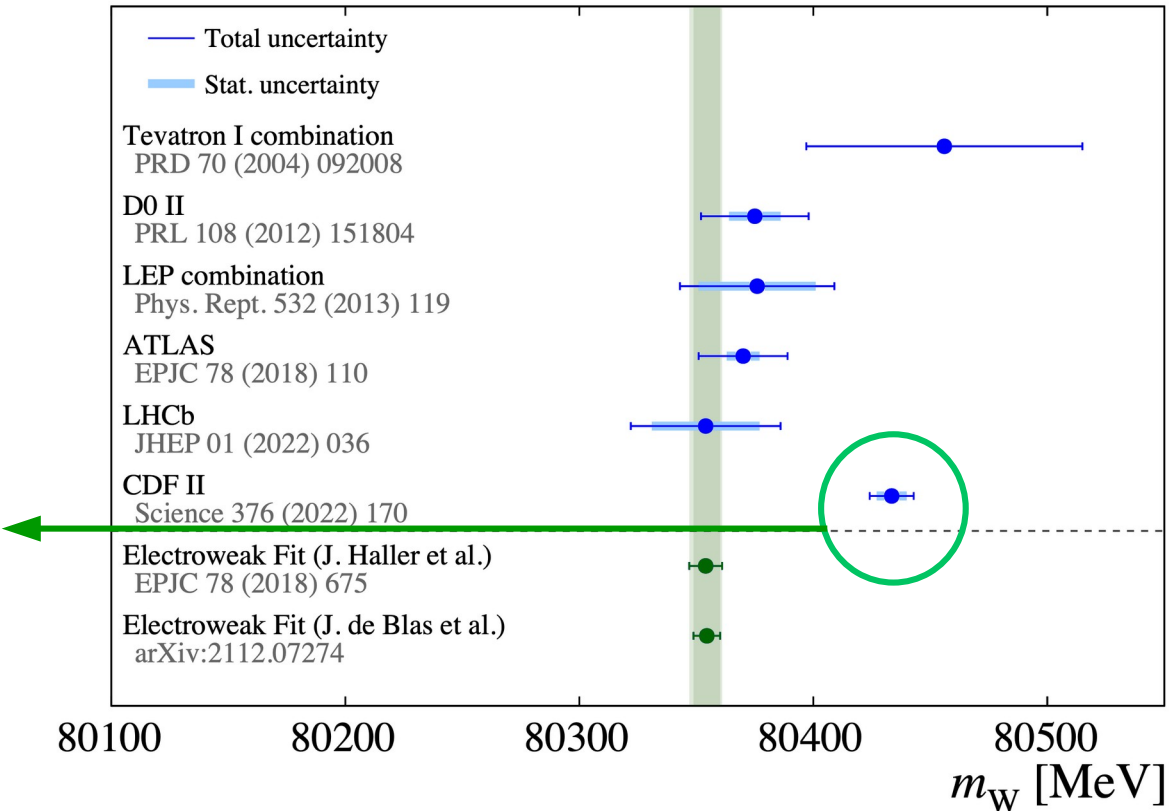Last lecture we saw how to estimate (=compute) the value of a parameter

**Maximum Likelihood Estimator (MLE) $\hat{\mu}$:**

$$\hat{\mu} = arg\,max\,L(\mu)$$

However we also need to estimate the associated uncertainty.

**What is the meaning of an uncertainty ?**

We don't know what the true value is, but **there is a 68% chance that it is within the error bar**

# Gaussian confidence intervals

Consider a Gaussian likelihood:

$$L(\mu) = \exp\left[-\frac{1}{2}\left(\frac{n-\mu}{\sigma}\right)^2\right]$$

$$P(\mu - \sigma < n < \mu + \sigma) \geq 68.3\%$$

$$P(n - \sigma < \mu < n + \sigma) \geq 68.3\%$$

**Still a statement on n!**

$\mu = n \pm \sigma$ at 68% CL ("1σ")

**The reported interval n ± σ will contain the true value of μ 68.3% of the time**

# Gaussian confidence intervals



Consider a Gaussian likelihood:

$$L(\mu) = \exp\left[-\frac{1}{2}\left(\frac{n-\mu}{\sigma}\right)^2\right]$$

$$P(\mu - \sigma < n < \mu + \sigma) \geq 68.3\,\%$$

$$P(n - \sigma < \mu < n + \sigma) \geq 68.3\,\%$$

**Still a statement on n!**

**μ = n ± σ at 68% CL ("1σ")**

**The reported interval n ± σ will contain the true value of μ 68.3% of the time**

# Gaussian confidence intervals



Consider a Gaussian likelihood:

$$L(\mu) = \exp\left[-\frac{1}{2}\left(\frac{n-\mu}{\sigma}\right)^2\right]$$
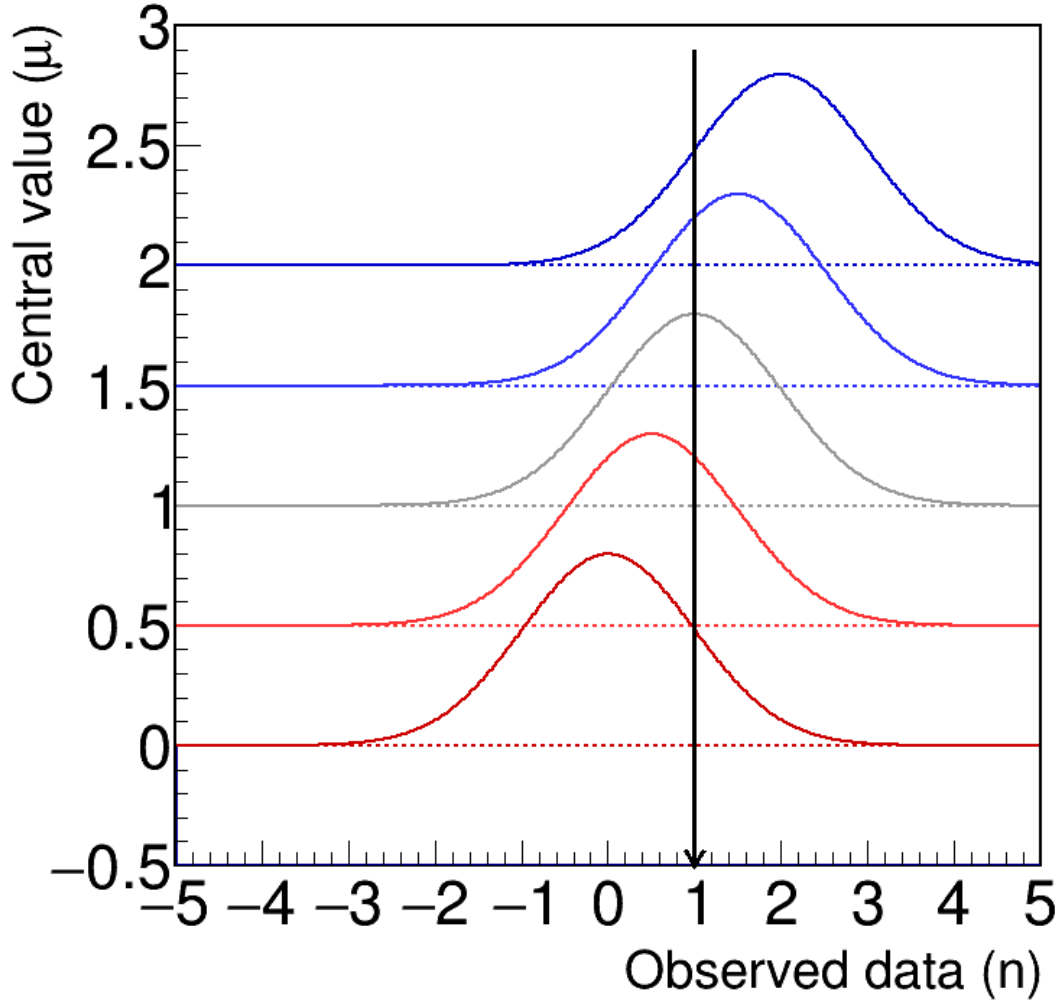
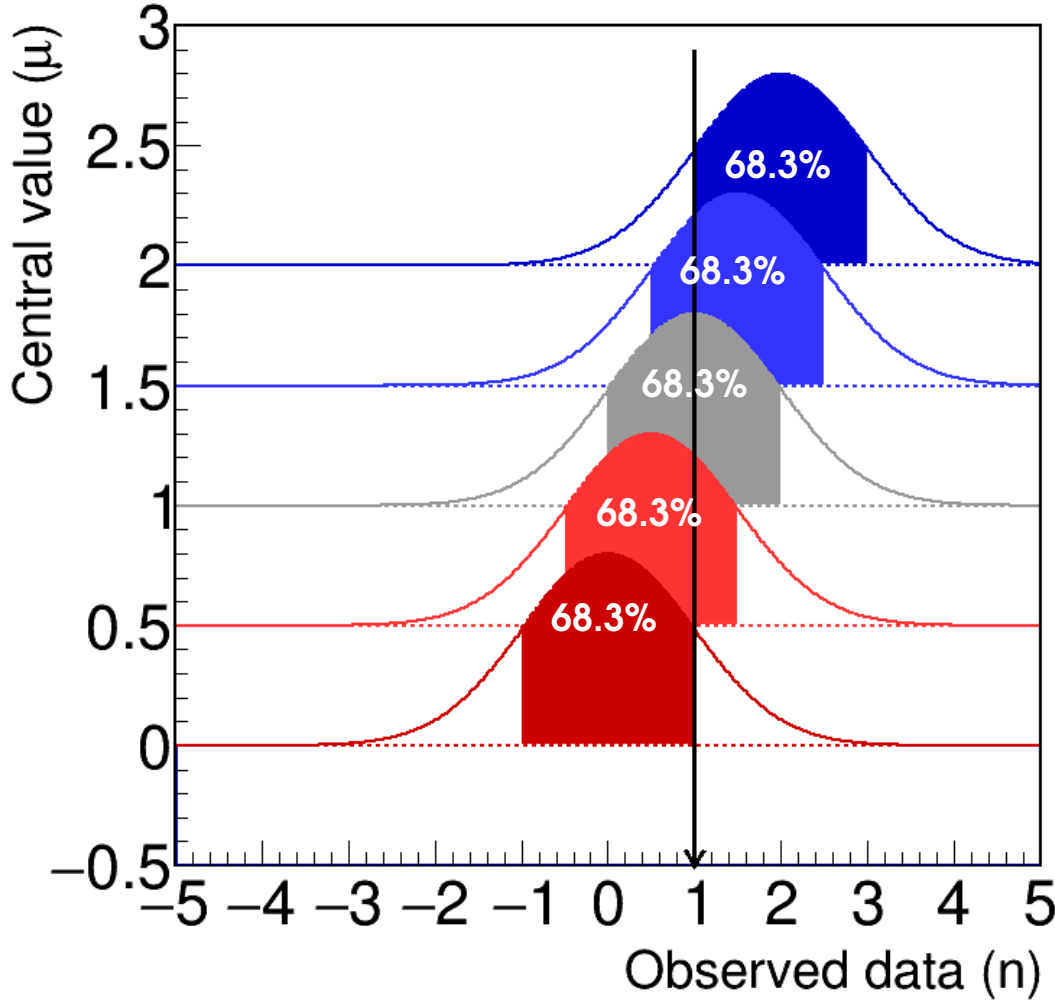$$P(\mu - \sigma < n < \mu + \sigma) \geq 68.3\,\%$$

$$P(n - \sigma < \mu < n + \sigma) \geq 68.3\,\%$$

**Still a statement on n!**

**μ = n ± σ at 68% CL ("1σ")**

**The reported interval n ± σ will contain the true value of μ 68.3% of the time**

# Gaussian confidence intervals

**Frequentist interpretation**

If we would repeat the same experiment multiple times, with true value μ*, then 68.3% of the 1σ intervals would contain μ*.

→ **Crucially, this works even if we do not know μ* !**

For each experiment, get the interval

**μ = n ± σ at 68% CL ("1σ")**

**The reported interval n ± σ will contain the true value of μ 68.3% of the time**

# Neyman Construction

**General case:** build 1σ intervals of observed values for each true value

⇒ *Confidence belt*



Peak Position

$P(\mu; \mu^*)$

**68% intervals for $\hat{\mu}$**

True value $\mu^*$

**Observed value $\hat{\mu}$**

# Neyman Construction

**General case:** build 1σ intervals of observed values for each true value

*⇒ Confidence belt*



Peak Position

$P(\mu; \mu^*)$

68% intervals for $\hat{\mu}$

Observed value $\hat{\mu}$

True value $\mu^*$

# Inversion using the Confidence Belt

**General case:** Intersect belt with given $\hat{\boldsymbol{\mu}}$ , get $\quad P\left(\hat{\mu} - \sigma_\mu^- < \mu^* < \hat{\mu} + \sigma_\mu^+\right) = 68\,\%$

→ Same as before for Gaussian, works also when $P(\mu^{obs}|\mu)$ varies with μ.

# Inversion using the Confidence Belt

**General case:** Intersect belt with given $\hat{\boldsymbol{\mu}}$ , get $\quad P\left(\hat{\mu} - \sigma_\mu^- < \mu^* < \hat{\mu} + \sigma_\mu^+\right) = 68\,\%$

$\rightarrow$ Same as before for Gaussian, works also when $P(\mu^{obs}|\mu)$ varies with $\mu$.



**True value $\mu^*$**

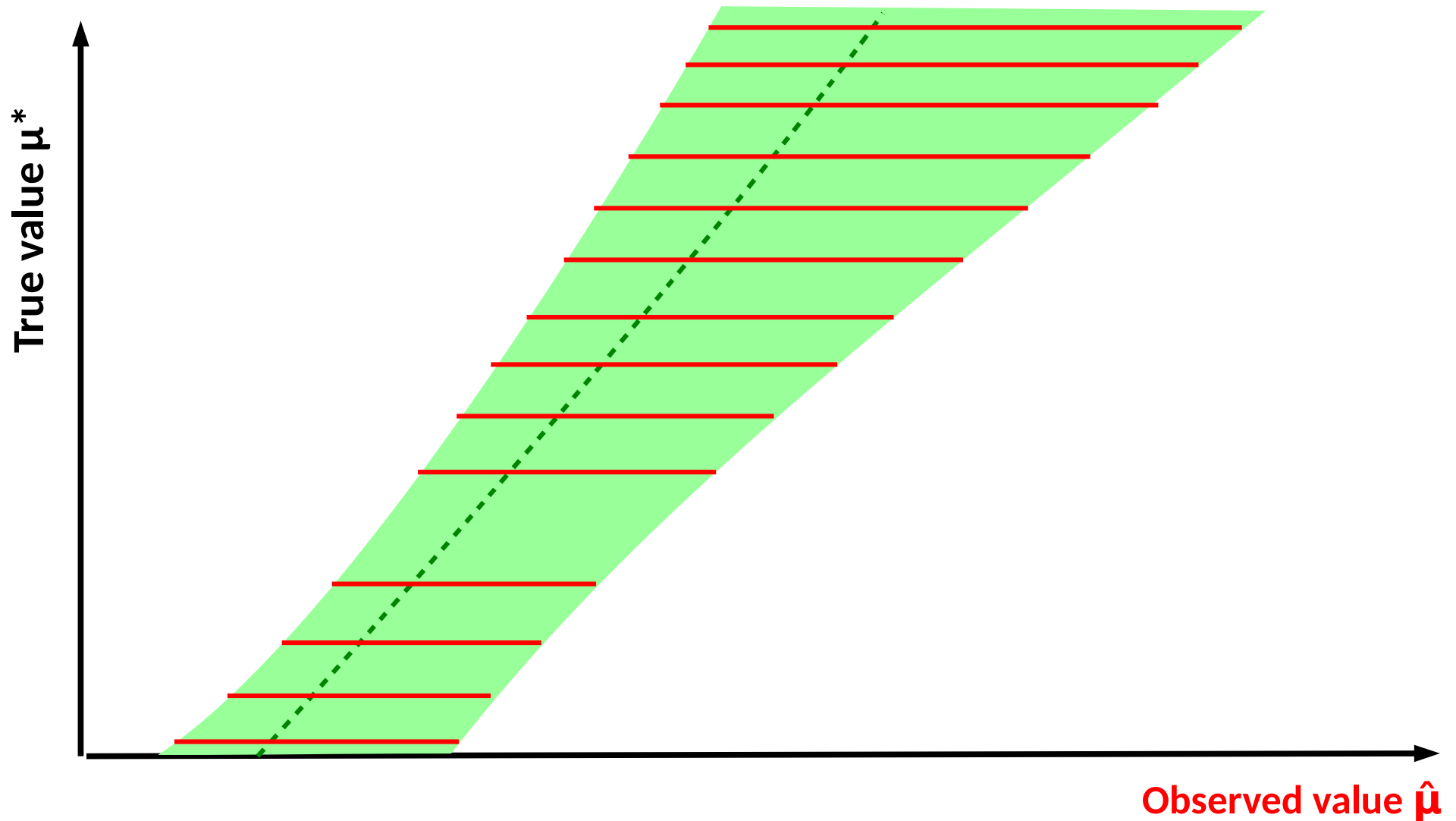$\hat{\boldsymbol{\mu}}$

**Observed value $\hat{\boldsymbol{\mu}}$**

# Inversion using the Confidence Belt

**General case:** Intersect belt with given $\hat{\mu}$ , get $\quad P\left(\hat{\mu} - \sigma_{\mu}^{-} < \mu^{*} < \hat{\mu} + \sigma_{\mu}^{+}\right) = 68\,\%$
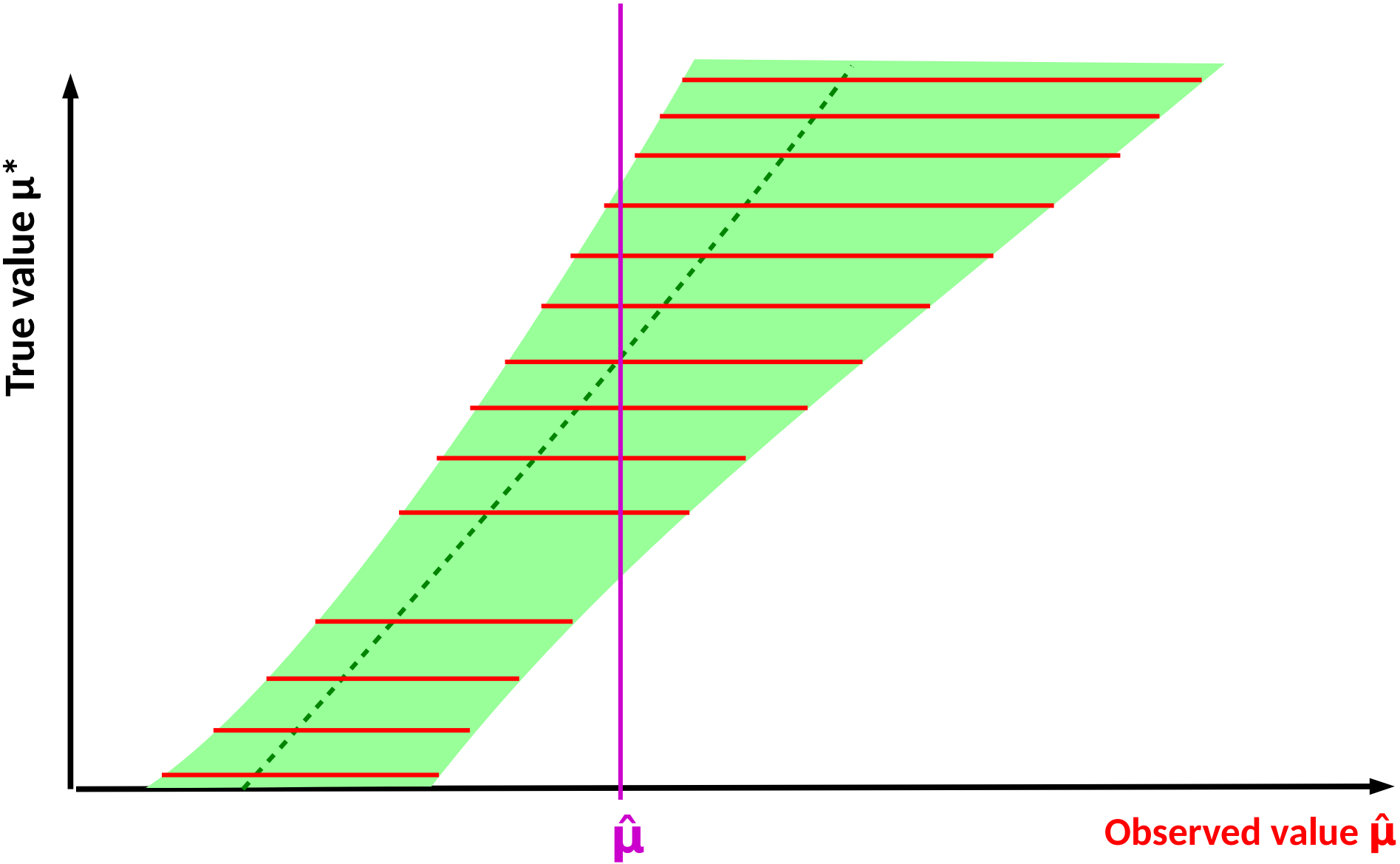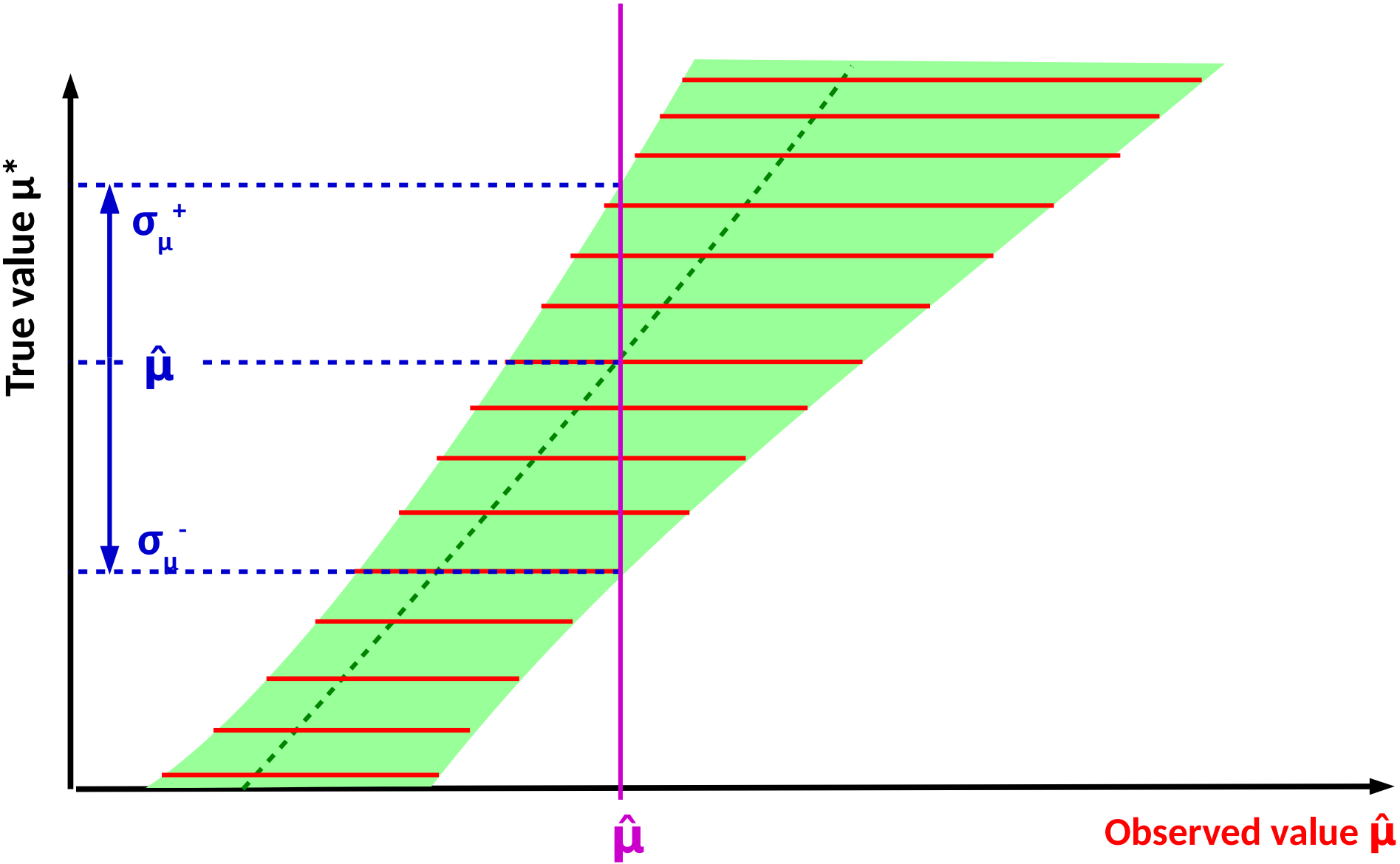
→ Same as before for Gaussian, works also when $P(\mu^{obs}|\mu)$ varies with μ.

# Inversion using the Confidence Belt

**General case:** Intersect belt with given $\hat{\mu}$, get $\quad P\left(\hat{\mu} - \sigma_{\mu}^{-} < \mu^{*} < \hat{\mu} + \sigma_{\mu}^{+}\right) = 68\,\%$

$\rightarrow$ Same as before for Gaussian, works also when $P(\mu^{obs}|\mu)$ varies with $\mu$.



$\sigma_{\mu}$ comes from the model,

**not the data**

$\rightarrow$ data only provides $\hat{\mu}$.

$\sigma_{\mu}^{+}$ from *negative* side of $\hat{\mu}$ intervals

$\sigma_{\mu}^{-}$ from *positive* side of $\hat{\mu}$ intervals

**Problem: Doesn't generalize well to many parameters in realistic models**

# General case: Likelihood Intervals

**Confidence intervals from L(µ):**

- Test various values µ using the **Profile Likelihood Ratio t(µ)**

- Minimum (=0) for µ=µ̂, rises away from µ̂.

- Good properties thanks to the Neyman-Pearson lemma.

$$t(\mu) = -2\log\frac{L(\mu)}{L(\hat{\mu})}$$

**Probability to observe the data for best-fit µ̂.**



ATLAS-CONF-2017-047

ATLAS Preliminary
$\sqrt{s}$ = 13 TeV, 36.1 fb$^{-1}$
$H\to\gamma\gamma$ and $H\to ZZ^*\to 4l$
$m_H$ = 125.09 GeV

— Combination
········· $H\to\gamma\gamma$
—·—·— $H\to ZZ^*\to 4l$

$\mu = 1.09 \pm 0.12$

**Gaussian L(µ):**

$$L(\mu) = \exp\left[-\frac{1}{2}\left(\frac{n-\mu}{\sigma}\right)^2\right]$$

$$t(\mu) = \left(\frac{n-\mu}{\sigma}\right)^2$$

- $t(\mu)$ is parabolic, distributed as a $\chi^2$
- Minimum occurs at $\mu = \hat{\mu}$
- 1σ interval $[\mu_-, \mu_+]$ given by $t(\mu_\pm)= 1$

38
/

# General case: Likelihood Intervals

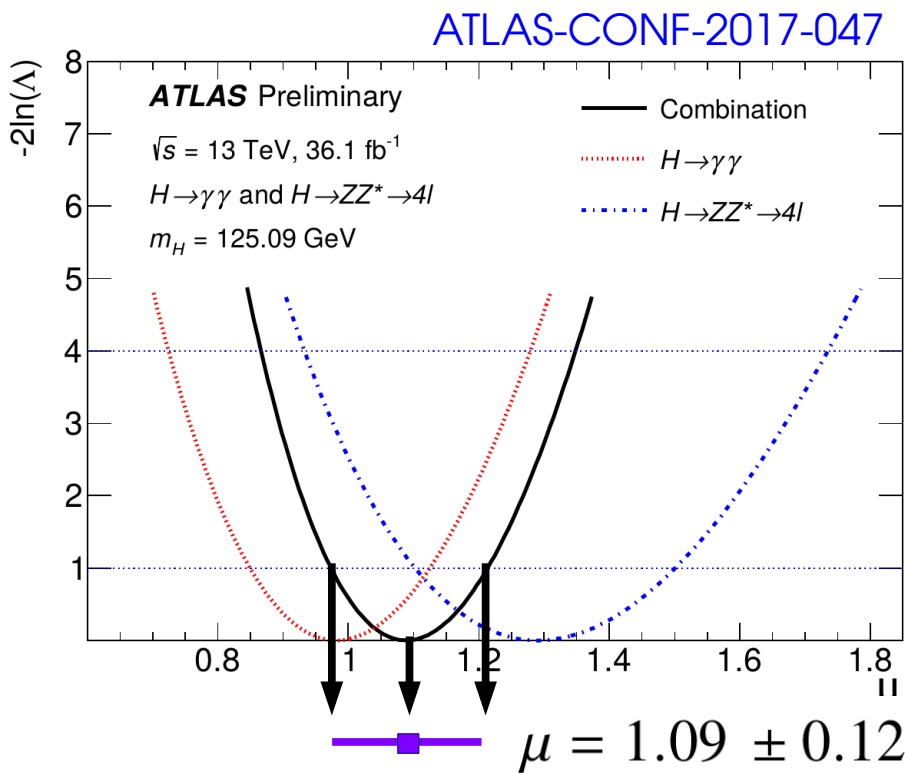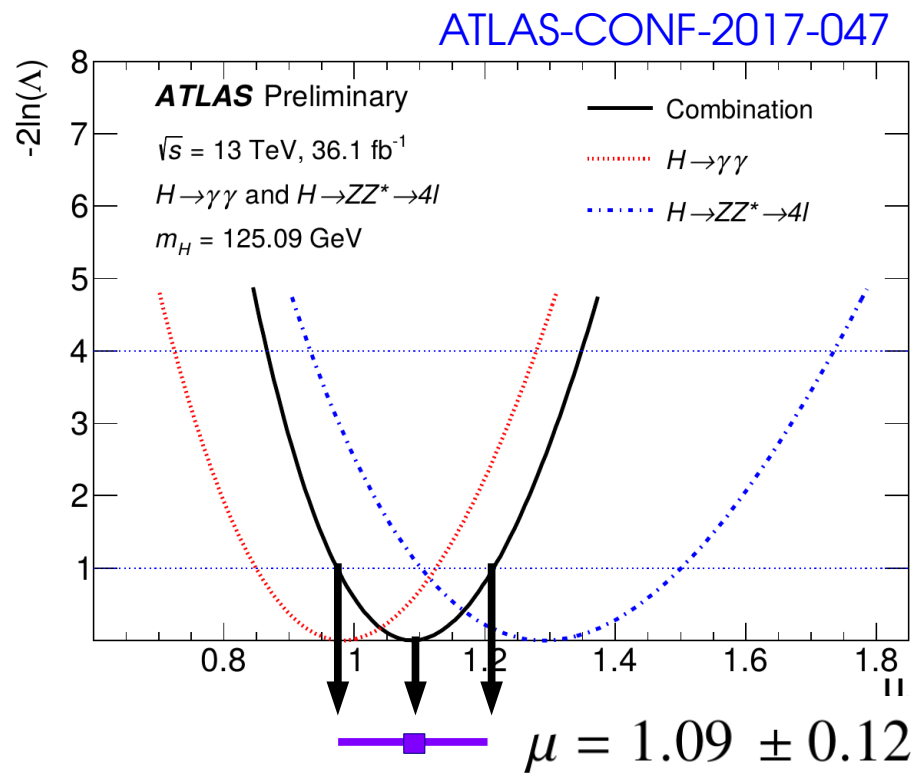**Confidence intervals from L(μ):**

- Test various values μ using the **Profile Likelihood Ratio t(μ)**

- Minimum (=0) for μ=μ̂, rises away from μ̂.

- Good properties thanks to the Neyman-Pearson lemma.

$$t(\mu) = -2\log\frac{L(\mu)}{L(\hat{\mu})}$$

**General case:**

- Generally not a perfect parabola
- Minimum still at **μ = μ̂**

**Asymptotic approximation**

- Compute t(μ) using the exact L(μ)
- Assume t(μ) ~ χ² as for Gaussian (*"Wilks' Theorem"*)

**1σ interval [μ_-, μ_+] given by t(μ_±)= 1**



ATLAS-CONF-2017-047

**ATLAS** Preliminary
$\sqrt{s}$ = 13 TeV, 36.1 fb$^{-1}$
$H \to \gamma\gamma$ and $H \to ZZ^* \to 4l$
$m_H$ = 125.09 GeV

— Combination
······· $H \to \gamma\gamma$
–·–·– $H \to ZZ^* \to 4l$

$\mu = 1.09 \pm 0.12$
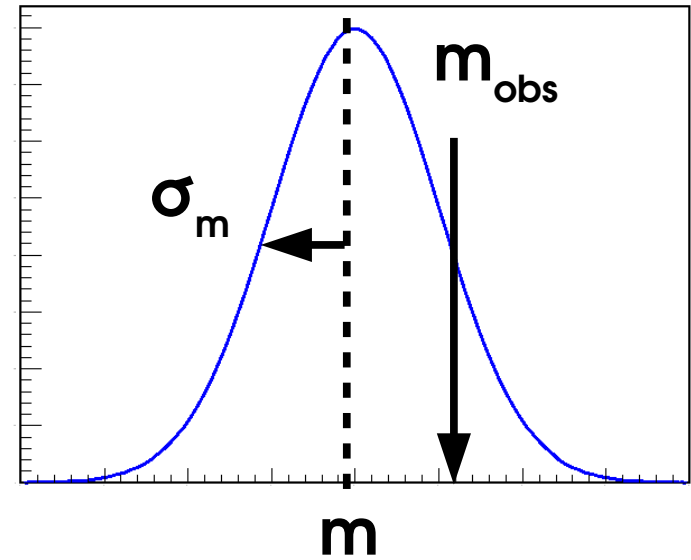
# Homework 3: Gaussian Case

Consider a parameter m (e.g. Higgs boson mass) whose measurement is Gaussian with known width $\sigma_m$, and we measure $m_{obs}$:

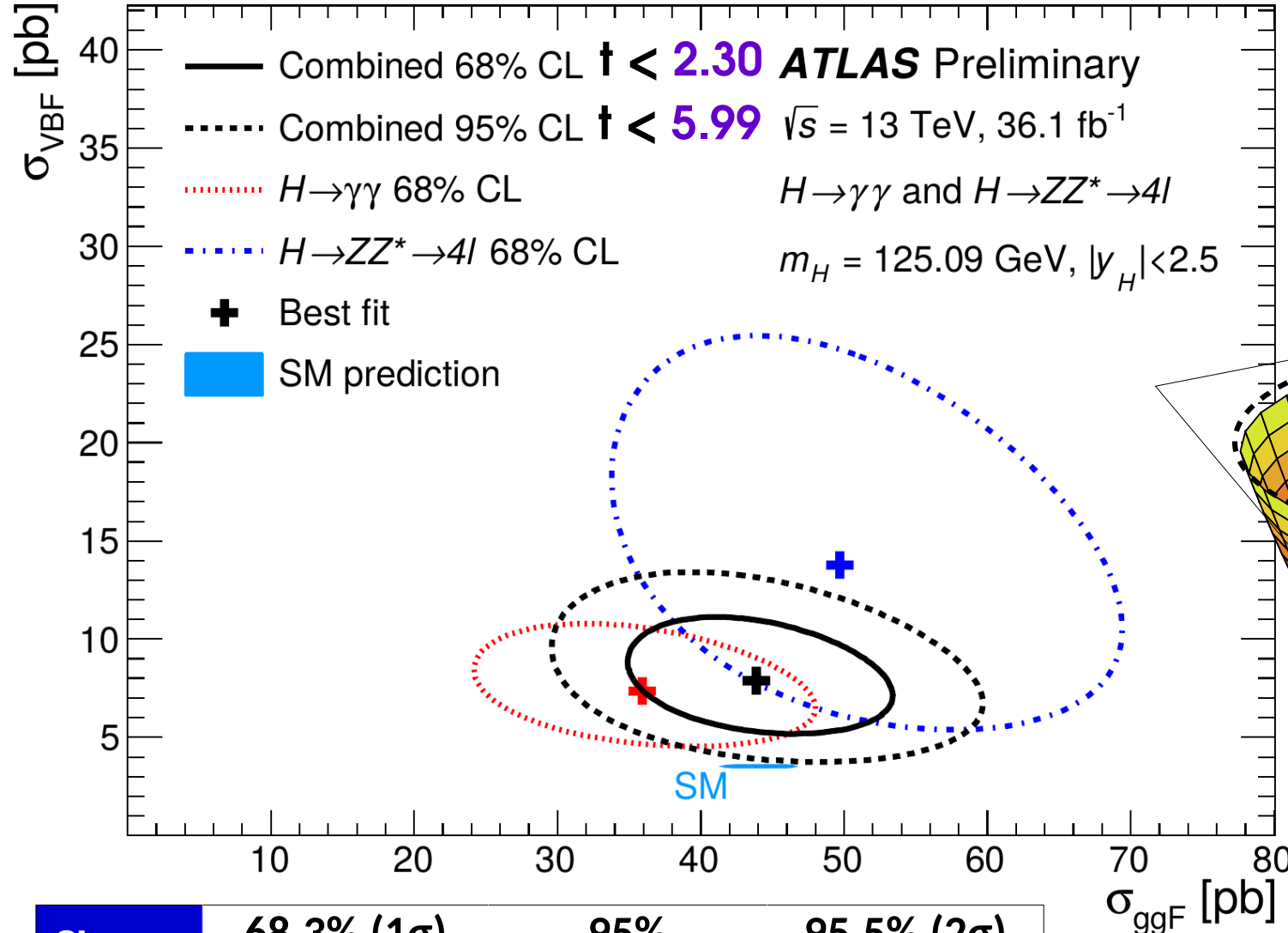$$L(m\,;m_{obs}) = e^{-\frac{1}{2}\left(\frac{m-m_{obs}}{\sigma_m}\right)^2}$$

$\rightarrow$ Compute the best-fit value (MLE) $\hat{m}$

$\rightarrow$ Compute $t_m$

$\rightarrow$ Compute the 1-$\sigma$ (Z=1, ~68% CL) interval on m

**Solution:** $\quad m = m_{obs} \pm \sigma_m$

$\rightarrow$ As expected!

$\rightarrow$ General method can be applied in the same way to more complex cases

# 2D Example: Higgs $\sigma_{VBF}$ vs. $\sigma_{ggF}$

$$t = -2\log\frac{L(X_0, Y_0)}{L(\hat{X}, \hat{Y})}$$

$$\sim \chi^2(N_{dof} = 2)$$

Legend for plot:
- Combined 68% CL $t < 2.30$ **ATLAS** Preliminary
- Combined 95% CL $t < 5.99$ $\sqrt{s}$ = 13 TeV, 36.1 fb$^{-1}$
- $H \to \gamma\gamma$ 68% CL
- $H \to ZZ^* \to 4l$ 68% CL
- Best fit
- SM prediction

$H \to \gamma\gamma$ and $H \to ZZ^* \to 4l$

$m_H$ = 125.09 GeV, $|y_H| < 2.5$

Krishnavedala - Own work, CC BY-SA 3.0, https://commons.wikimedia.org/w/index.php?curid=15278826

| CL | 68.3% (1σ) | 95% | 95.5% (2σ) |
|---|---|---|---|
| 1D Z² | 1.00 | 3.84 | 4.00 |
| 2D Z² | 2.30 | 5.99 | 6.18 |

**Gaussian case**: elliptic paraboloid surface

41/

# Reparameterization
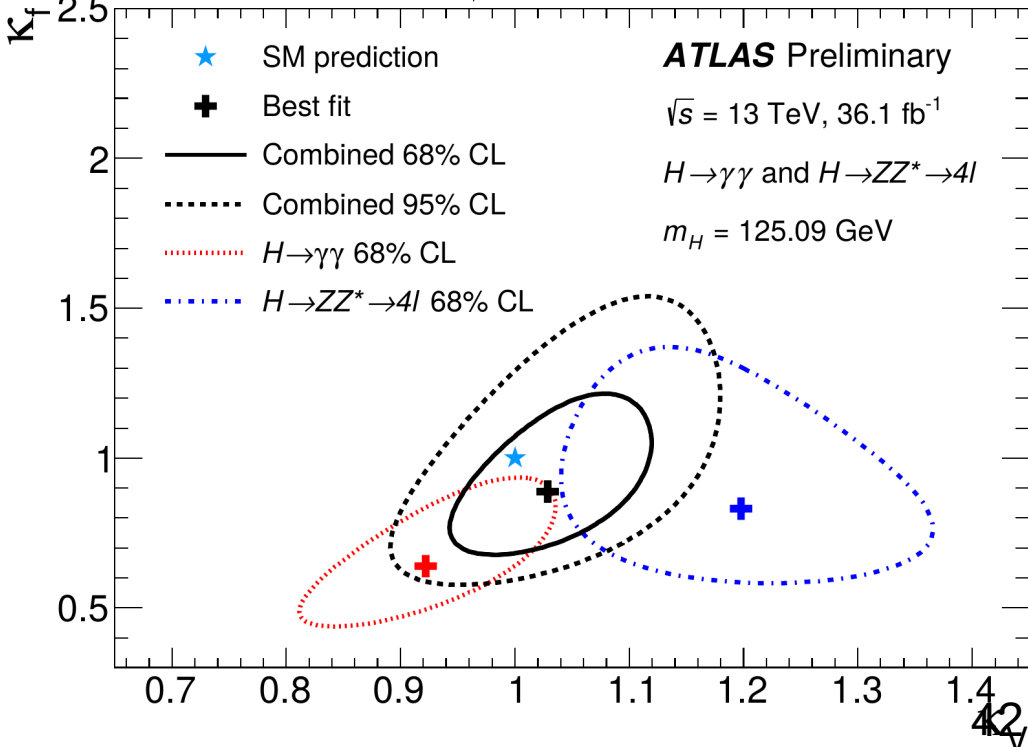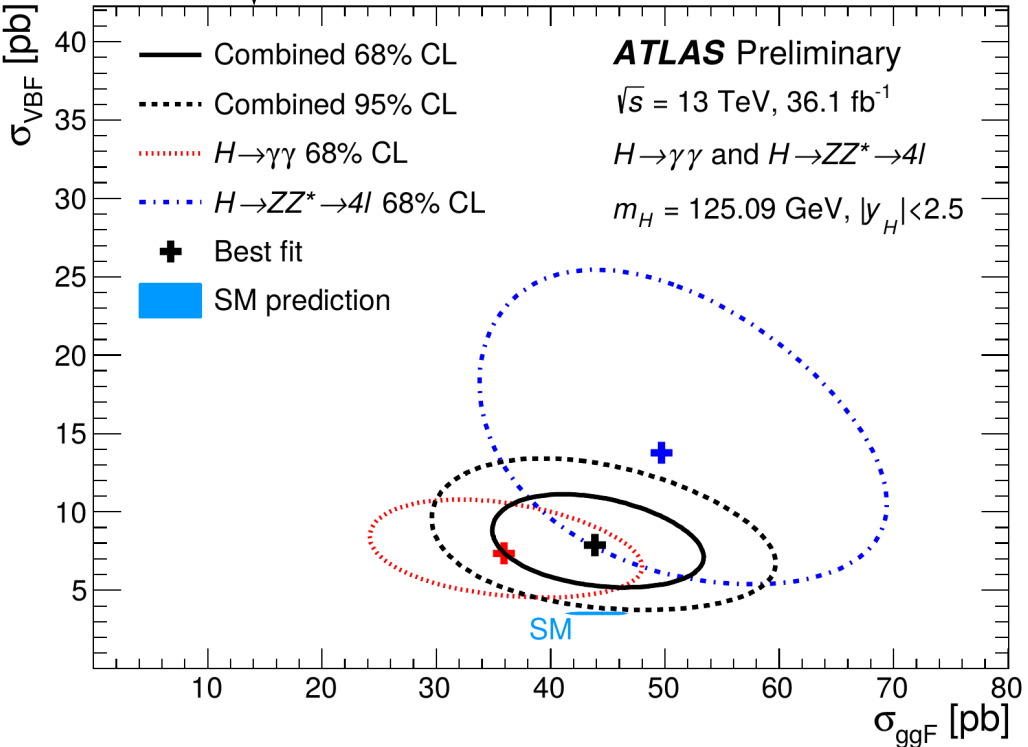
Start with basic measurement in terms of e.g. **σ×B**

→ How to measure derived quantities (couplings, parameters in some theory model, etc.) ?
→ **just reparameterize the likelihood:**

*e.g.* Higgs couplings: $\sigma_{ggF}$, $\sigma_{VBF}$ sensitive to Higgs coupling modifiers $\kappa_V$, $\kappa_F$.

$$\sigma_{ggF} \rightarrow \sigma_{ggF}(\kappa_V, \kappa_F)$$

$$L(\sigma_{ggF}, \sigma_{VBF}) \xrightarrow{\hspace{3cm}} L(\sigma_{ggF}(\kappa_V, \kappa_F), \sigma_{VBF}(\kappa_V, \kappa_F)) \equiv L'(\kappa_V, \kappa_F)$$

$$\sigma_{VBF} \rightarrow \sigma_{VBF}(\kappa_V, \kappa_F)$$
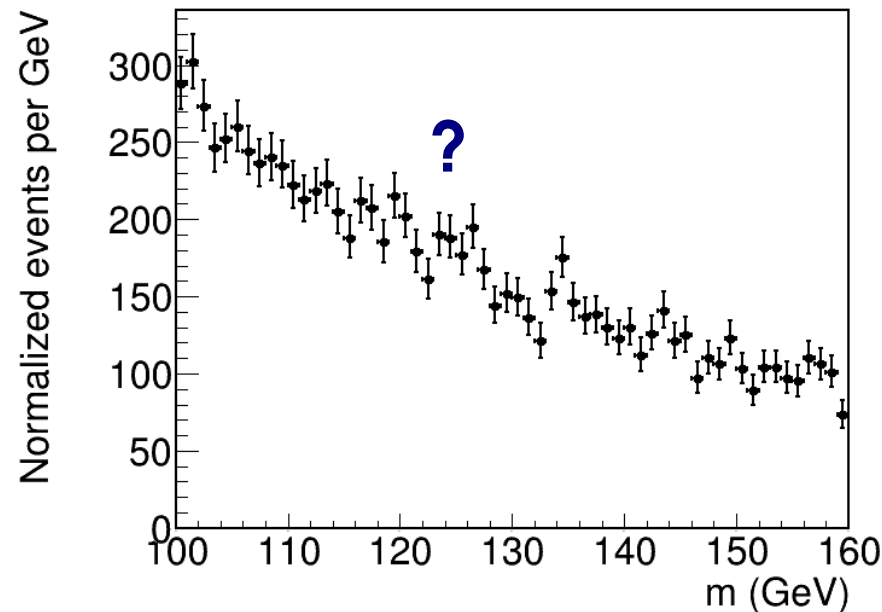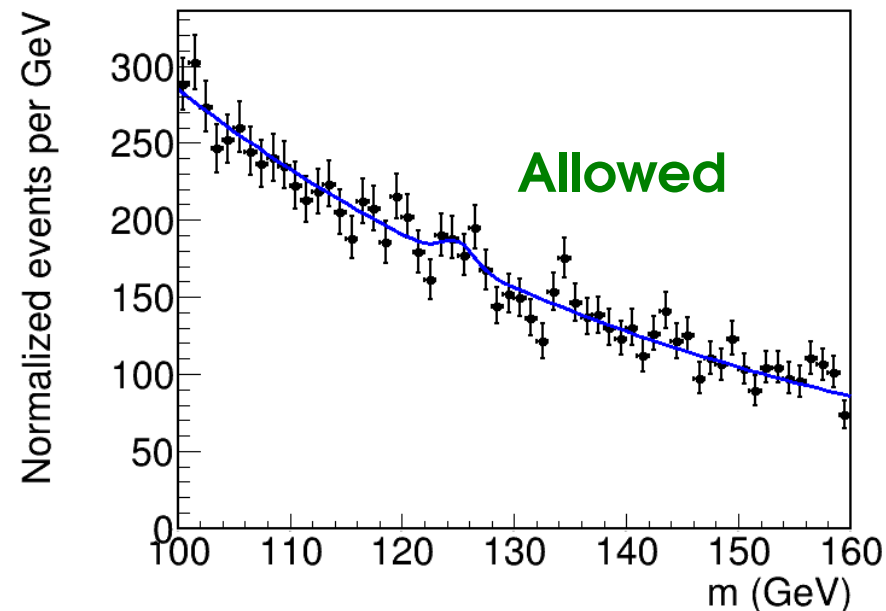
# Upper Limits

# Hypothesis tests for Limits

If no signal in data, testing for discovery not very relevant (report 0.2σ excess ?)

→ More interesting to **exclude large signals**

⇒ **Upper limits on signal yield**

→ Typically report **95% CL** upper limit (p-value = 5%) : "S < $S_0$ @ 95% CL"

# Hypothesis tests for Limits

If no signal in data, testing for discovery not very relevant (report 0.2σ excess ?)

→ More interesting to **exclude large signals**

⇒ **Upper limits on signal yield**

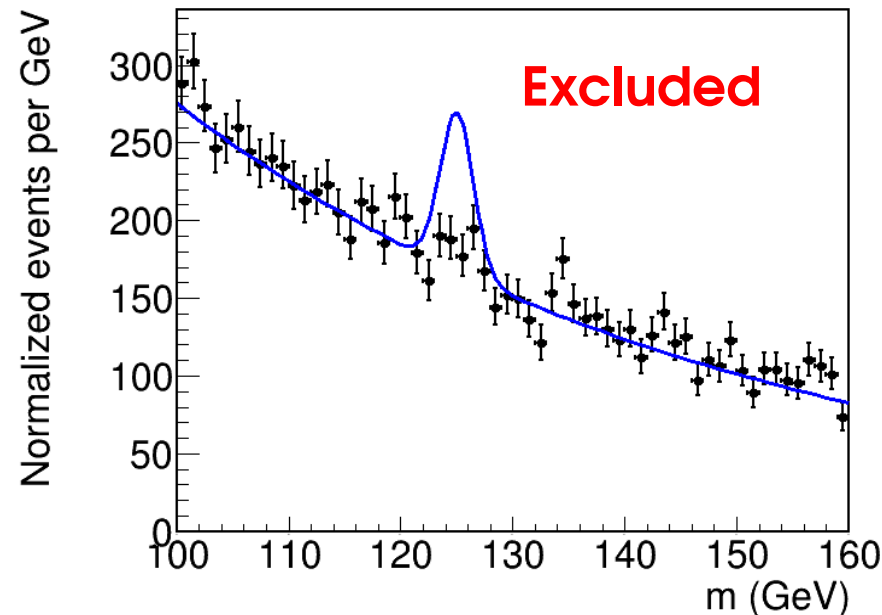→ Typically report **95% CL** upper limit (p-value = 5%) : "S < $S_0$ @ 95% CL"
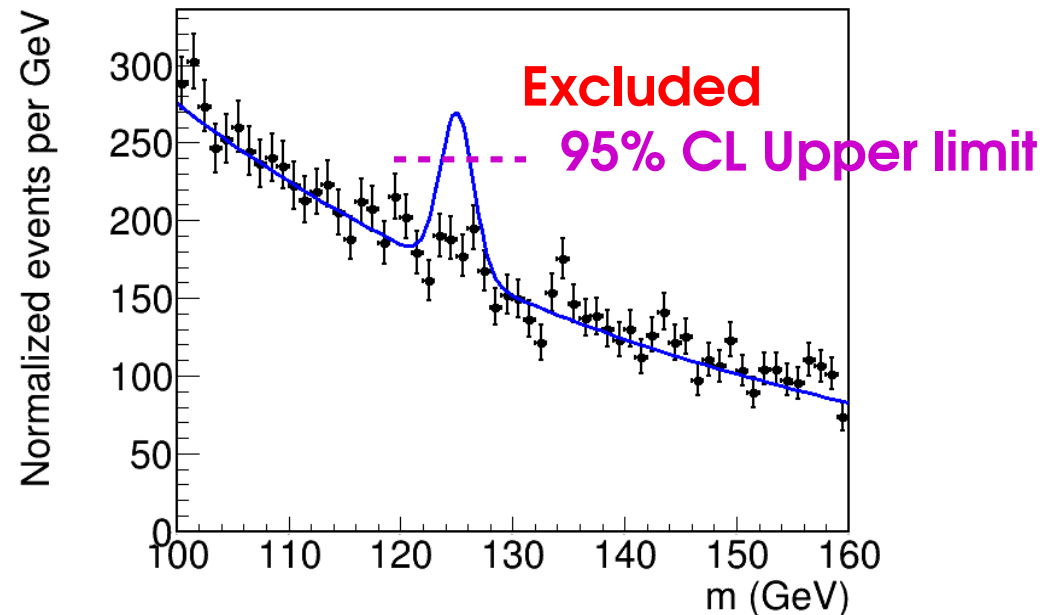
# Hypothesis tests for Limits

If no signal in data, testing for discovery not very relevant (report 0.2σ excess ?)

→ More interesting to **exclude large signals**

⇒ **Upper limits on signal yield**

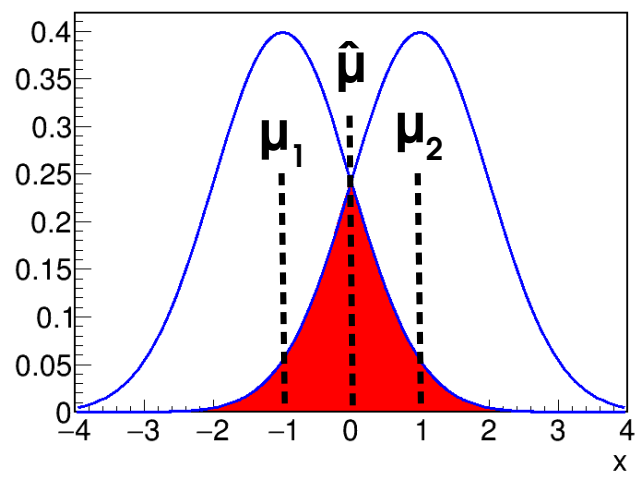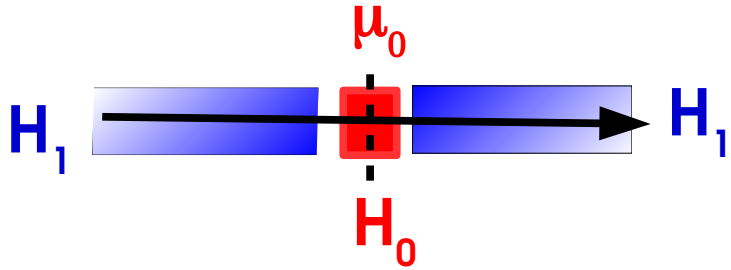→ Typically report **95% CL** upper limit (p-value = 5%) : "S < $S_0$ @ 95% CL"

# Hypothesis tests for Limits

If no signal in data, testing for discovery not very relevant (report $0.2\sigma$ excess ?)

$\rightarrow$ More interesting to **exclude large signals**

$\Rightarrow$ **Upper limits on signal yield**

$\rightarrow$ Typically report **95% CL** upper limit (p-value = 5%) : "$S < S_0$ @ 95% CL"

# Test Statistics for Limit-Setting

**Interval :**

$H_0 : \mu = \mu_0$

$H_1 : \mu \neq \mu_0$

$\mu_0$

$H_1$ — $H_0$ → $H_1$

Try to exclude μ values away from $\hat{\mu}$.

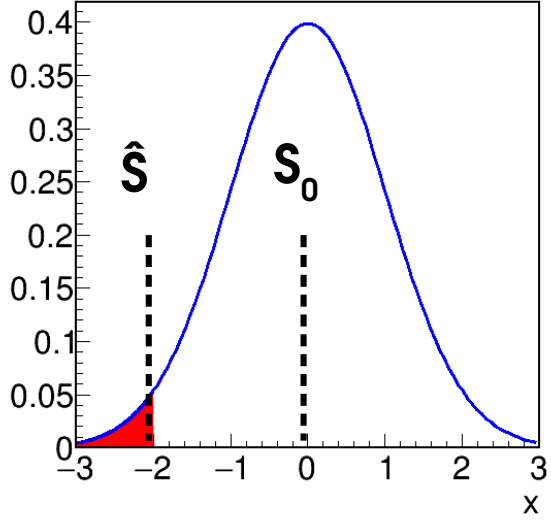$$t(\mu_0) = -2\log \frac{L(\mu = \mu_0)}{L(\hat{\mu})}$$

"**Two-sided**" test

**Limit-setting**

$H_0 : S = S_0$

$H_1 : S < S_0$

$S_0$

$H_1$ — → $H_0$

$$q(S_0) = \begin{cases} -2\log \dfrac{L(S = S_0)}{L(\hat{S})} & S_0 > \hat{S} \\[2em] 0 & S_0 \leq \hat{S} \end{cases}$$

Try to exclude values of S that are above $\hat{S}$.

⇒ "**One-sided**" test : only interested in excluding above

Discovery is also one-sided, for S>0 !

# *Inversion* : Getting the limit for a given CL

**Procedure:**

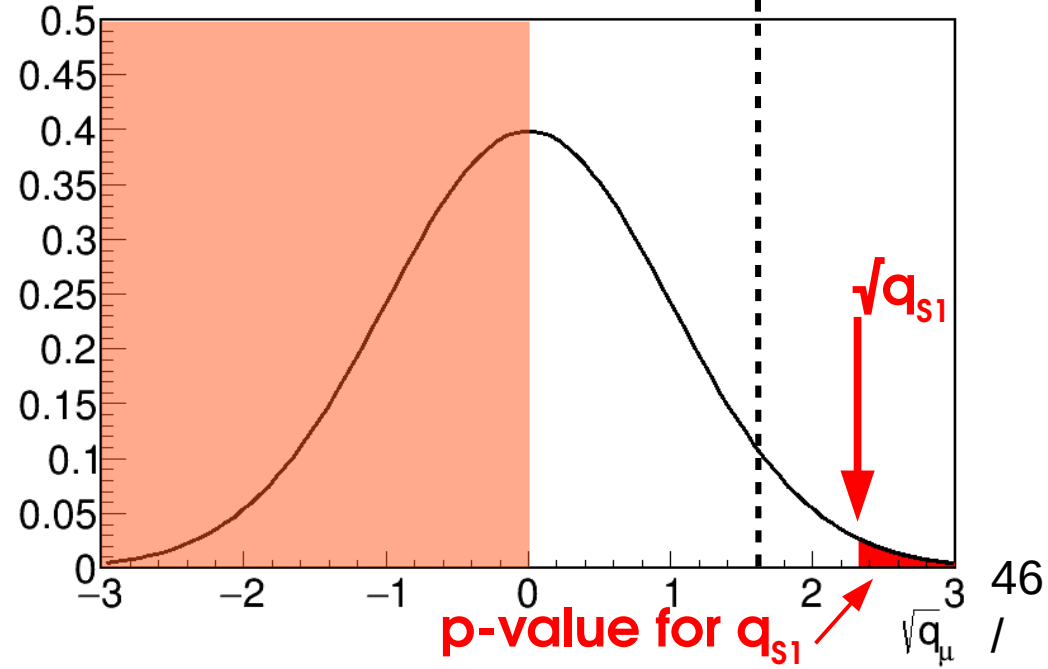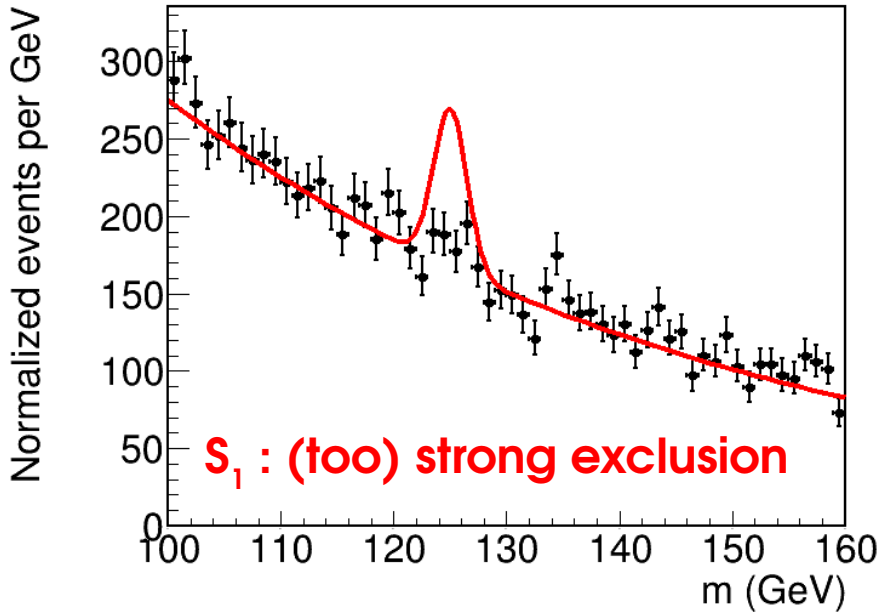→ Compute $q(S_0)$ for some $S_0$,
   get the **exclusion p-value $p(S_0)$**.

**Asymptotics:**

$$p(S_0) = 1 - \Phi\left(\sqrt{q(S_0)}\right)$$

| CL | p | Region |
|----|-----|----------------------|
| **90%** | **10%** | $\sqrt{q(S)} > 1.28$ |
| **95%** | **5%** | $\sqrt{q(S)} > 1.64$ |
| **99%** | **1%** | $\sqrt{q(S)} > 2.33$ |

→ **Adjust $S_0$** to get the desired exclusion
   **Asymptotics**: need $\sqrt{q(S_{95})} = 1.64$ for **95% CL**



$S_1$ : (too) strong exclusion

$\sqrt{q(S)} = 1.64$
$(p = 5\%)$

$\sqrt{q_{S1}}$

p-value for $q_{S1}$

$\sqrt{q_\mu}$

46 /

# *Inversion* : Getting the limit for a given CL

**Procedure:**

$\rightarrow$ Compute q($S_0$) for some $S_0$, get the **exclusion p-value p($S_0$).**

**Asymptotics:** $$p(S_0) = 1 - \Phi\left(\sqrt{q(S_0)}\right)$$

| CL | p | Region |
|----|----|--------|
| **90%** | **10%** | $\sqrt{q(S)} > 1.28$ |
| **95%** | **5%** | $\sqrt{q(S)} > 1.64$ |
| **99%** | **1%** | $\sqrt{q(S)} > 2.33$ |

$\rightarrow$ **Adjust $S_0$** to get the desired exclusion

**Asymptotics**: need $\sqrt{q(S_{95})} = 1.64$ for **95% CL**



$\sqrt{q(S)} = 1.64$
(p = 5%)

$\sqrt{q_{S2}}$

$\sqrt{q_{S1}}$

$S_2$ : no exclusion

46

# *Inversion* : Getting the limit for a given CL

**Procedure:**

→ Compute $q(S_0)$ for some $S_0$,

   get the **exclusion p-value $p(S_0)$**.

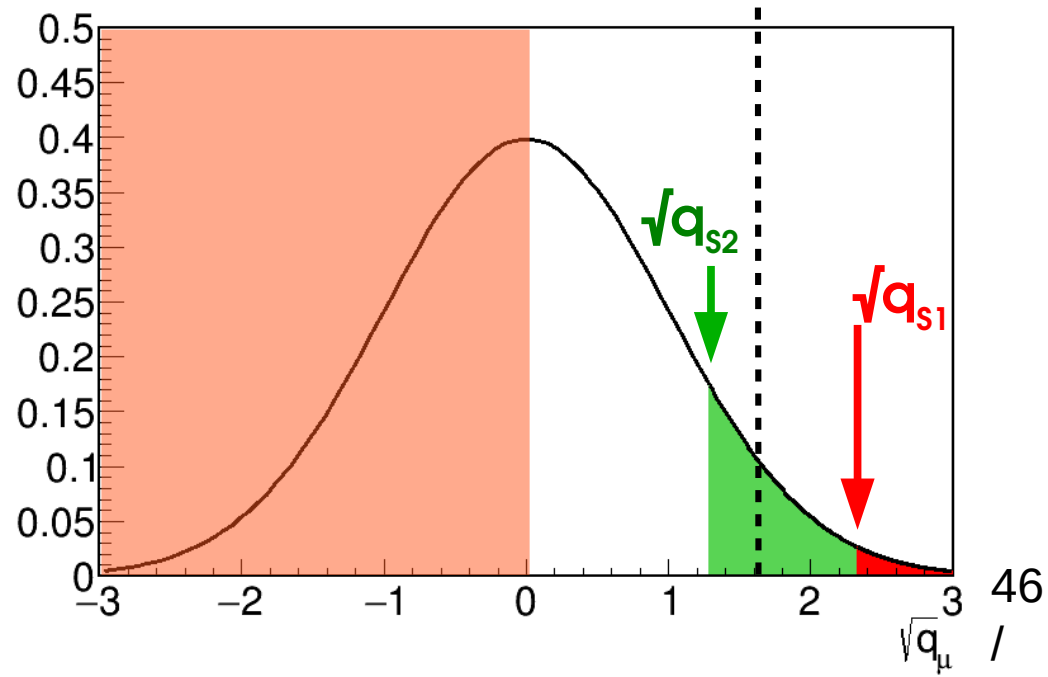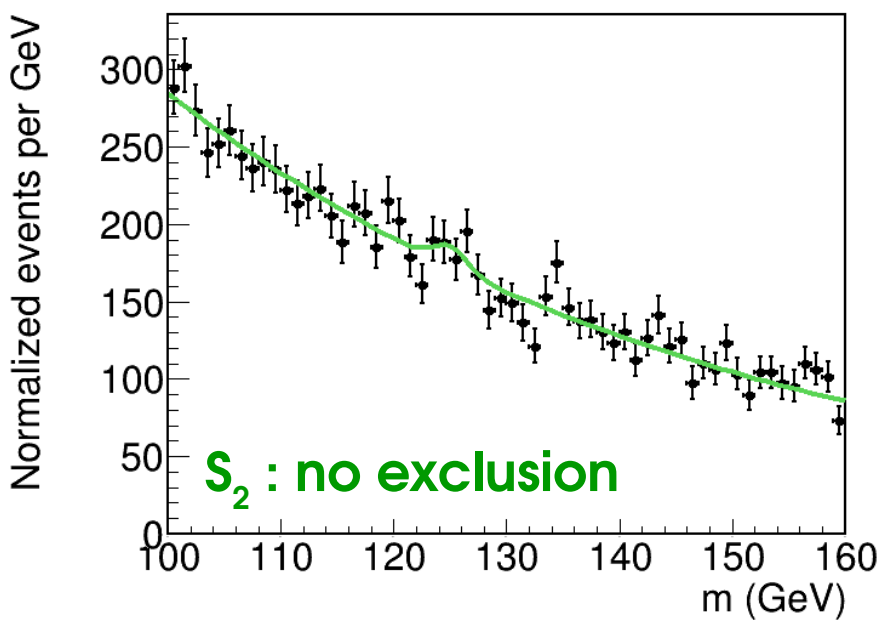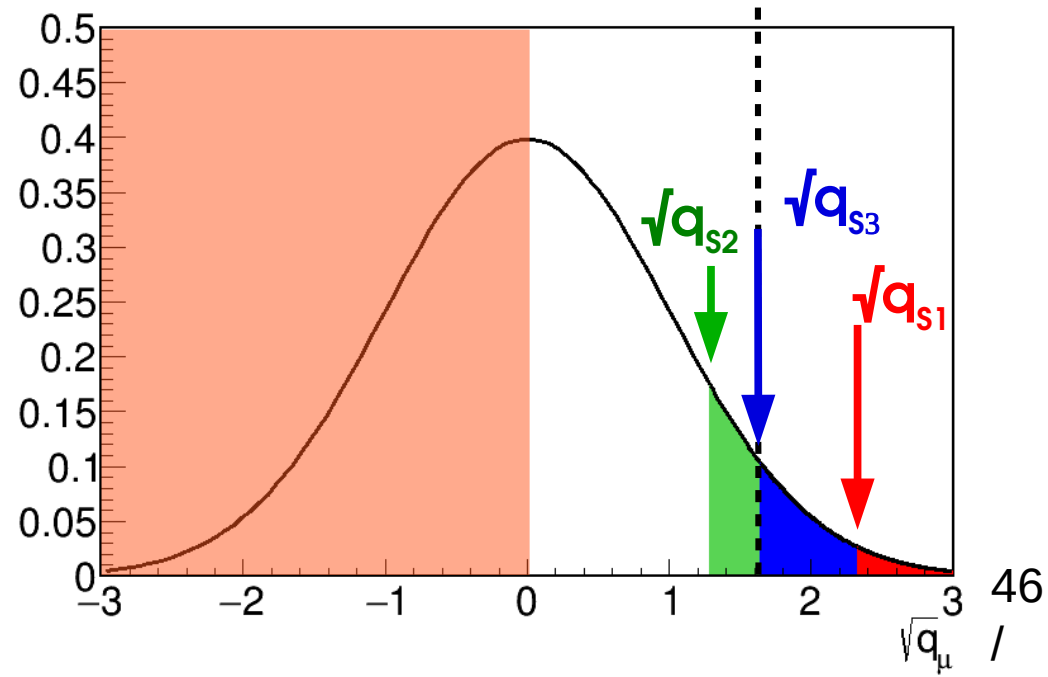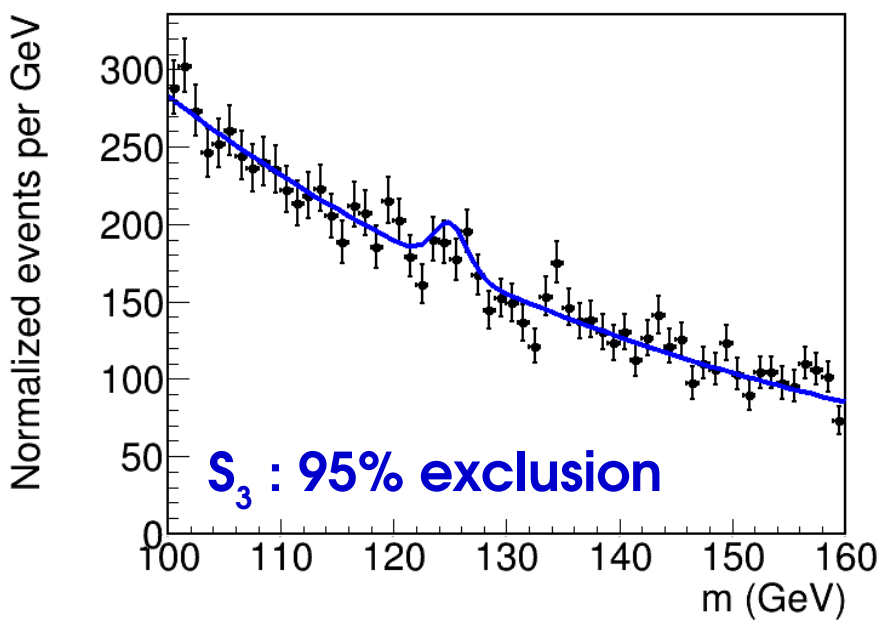   **Asymptotics:** $\quad p(S_0) = 1 - \Phi\!\left(\sqrt{q(S_0)}\right)$

| CL | p | Region |
|-----|------|-----------------|
| 90% | 10% | $\sqrt{q(S)} > 1.28$ |
| 95% | 5% | $\sqrt{q(S)} > 1.64$ |
| 99% | 1% | $\sqrt{q(S)} > 2.33$ |

→ **Adjust $S_0$** to get the desired exclusion

   **Asymptotics**: need $\sqrt{q(S_{95})} = 1.64$ for **95% CL**



$\sqrt{q(S)} = 1.64$
$(p = 5\%)$

$S_3$ : 95% exclusion

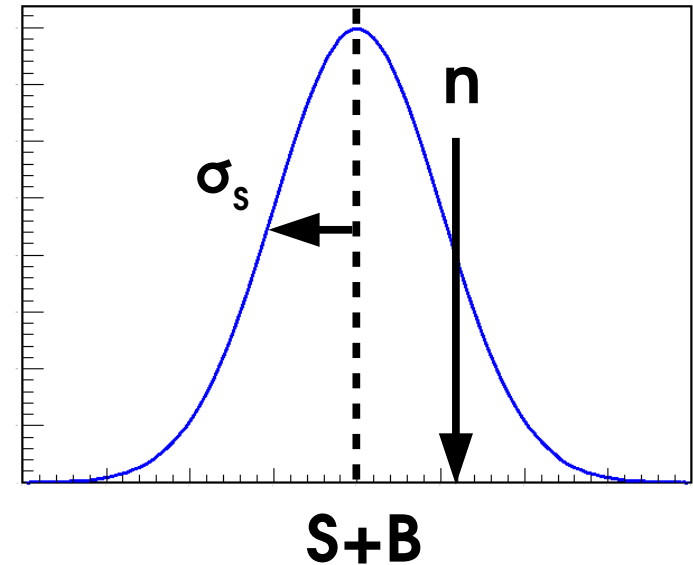# Homework 4: Gaussian Example

Usual Gaussian counting example with known B:

$$L(S;n) = e^{-\frac{1}{2}\left(\frac{n-(S+B)}{\sigma_s}\right)^2}$$

$\sigma_s \sim \sqrt{B}$ for small S

**Reminder:** Significance: $Z = \hat{S}/\sigma_s$

$\rightarrow$ Compute $q_{S0}$

$\rightarrow$ Compute the 95% CL upper limit on S, $S_{up}$, by solving $\sqrt{q_{S0}} = 1.64$.

**Solution:** $\quad S_{up} = \hat{S} + 1.64\,\sigma_S \ \ \text{at 95 \% CL}$

# CL$_s$

Upper limits sometimes take negative values (exclude all S>0 !)

Known feature – to avoid, usual solution in HEP is to use **CL$_s$** "modified p-value"
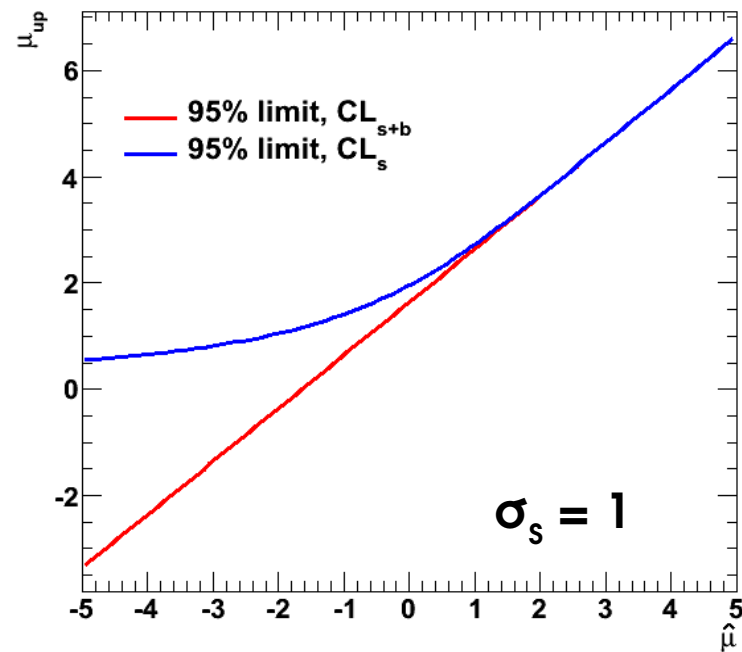
$$p_{CL_s} = \frac{p(S_0)}{p_B}$$

**Usual p-value for S=S$_0$**

**P-value for S=0**

⇒ Compute exclusion relative to that of S=0

→ Somewhat ad-hoc, but good properties...

**Ŝ ∼ 0 ⇒ p$_B$ ∼ O(1), p$_{CLs}$ ∼ p(S$_0$) no change**

**Ŝ ≪ 0 ⇒ p$_B$ ≪ 1, p$_{CLs}$ ≫ p(S$_0$) no exclusion at S=0**



95% limit, CL$_{s+b}$
95% limit, CL$_s$

$\sigma_s = 1$

**Drawback**: *overcoverage*

→ limit is claimed to be 95% CL, but actually >95% CL for small p$_B$.

48 /

# Homework 5: CL$_s$ : Gaussian Case

Usual Gaussian counting example with known B:

$$L(S;n) = e^{-\frac{1}{2}\left(\frac{n-(S+B)}{\sigma_S}\right)^2}$$

$\sigma_s \sim \sqrt{B}$ for small S



**Reminder**

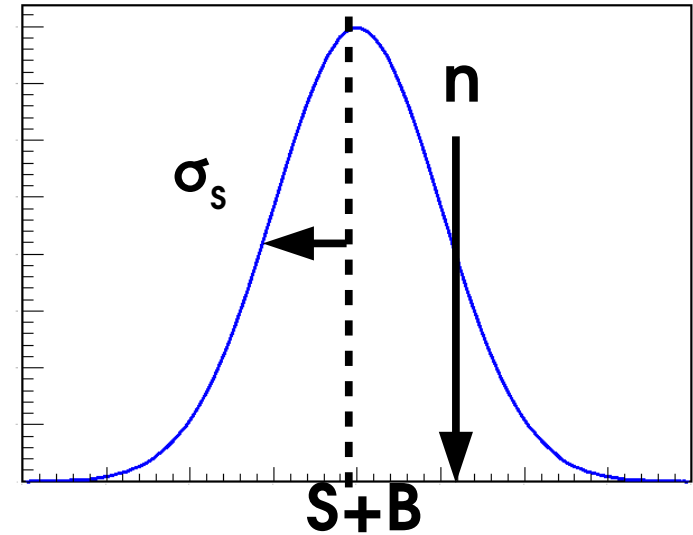CL$_{s+b}$ limit:    $S_{up} = \hat{S} + 1.64\,\sigma_S$ **at 95 % CL**

**CL$_s$ upper limit** :

→ Compute p$_{s0}$ (same as for CLs+b)

→ Compute 1-p$_B$ (hard!)

**Solution:**    $S_{up} = \hat{S} + \left[\Phi^{-1}\left(1 - 0.05\,\Phi\left(\hat{S}/\sigma_s\right)\right)\right]\sigma_S$ at 95 % CL

for $\hat{S} \sim 0$,    $S_{up} = \hat{S} + 1.96\,\sigma_S$ at 95 % CL

# Homework 6: CL$_S$ Rule of Thumb for $n_{obs}$=0

Same exercise, for the Poisson case with $n_{obs}$ = 0. Perform an exact computation of the 95% CLs upper limit based on the definition of the p-value:

**p-value** : *sum probabilities of cases at least as extreme as the data*

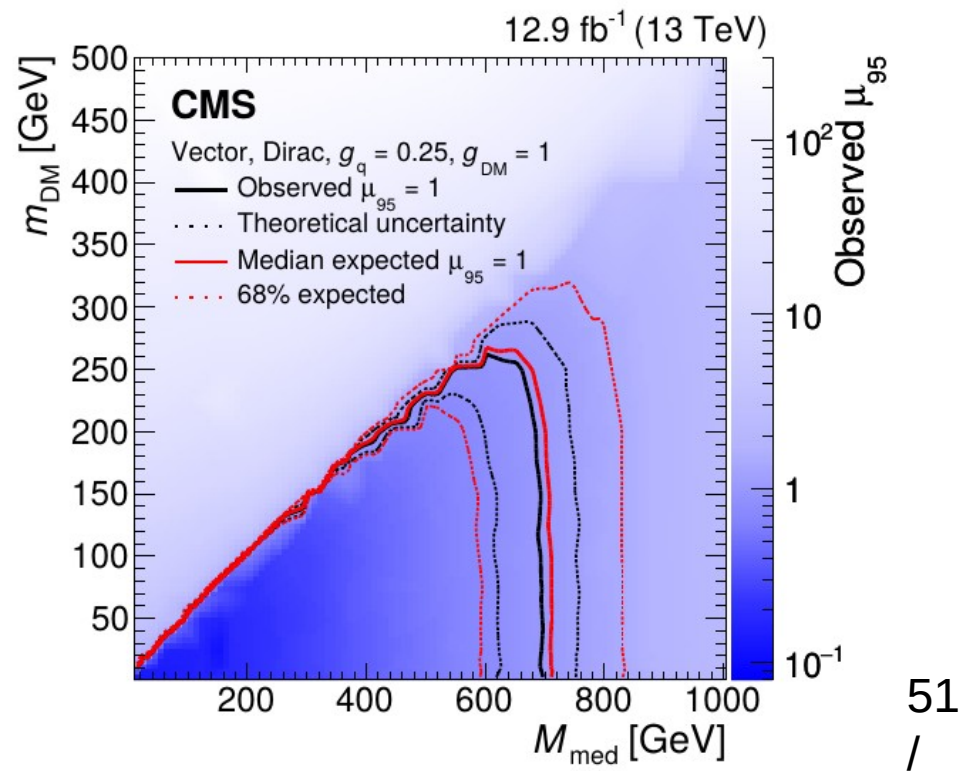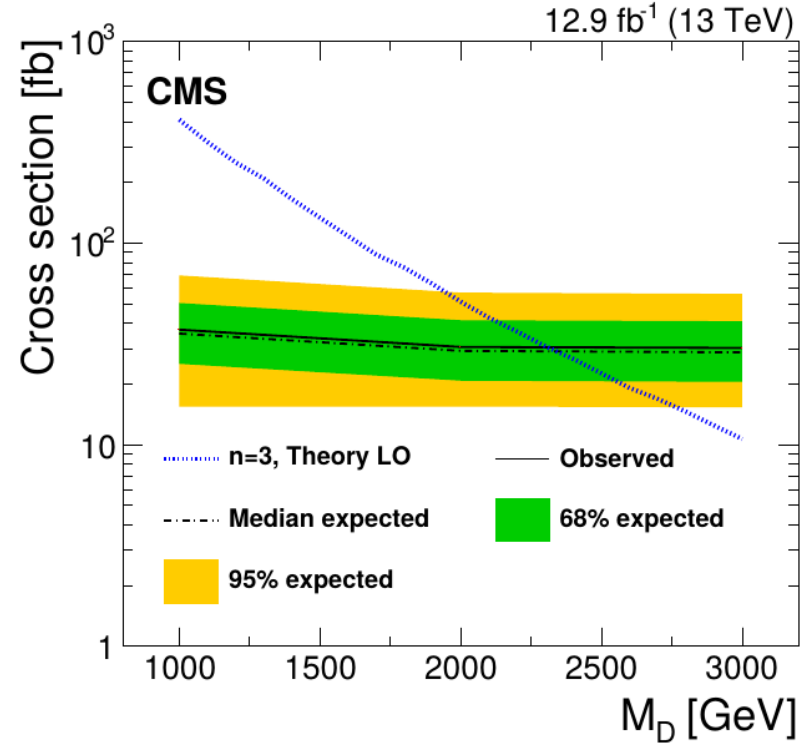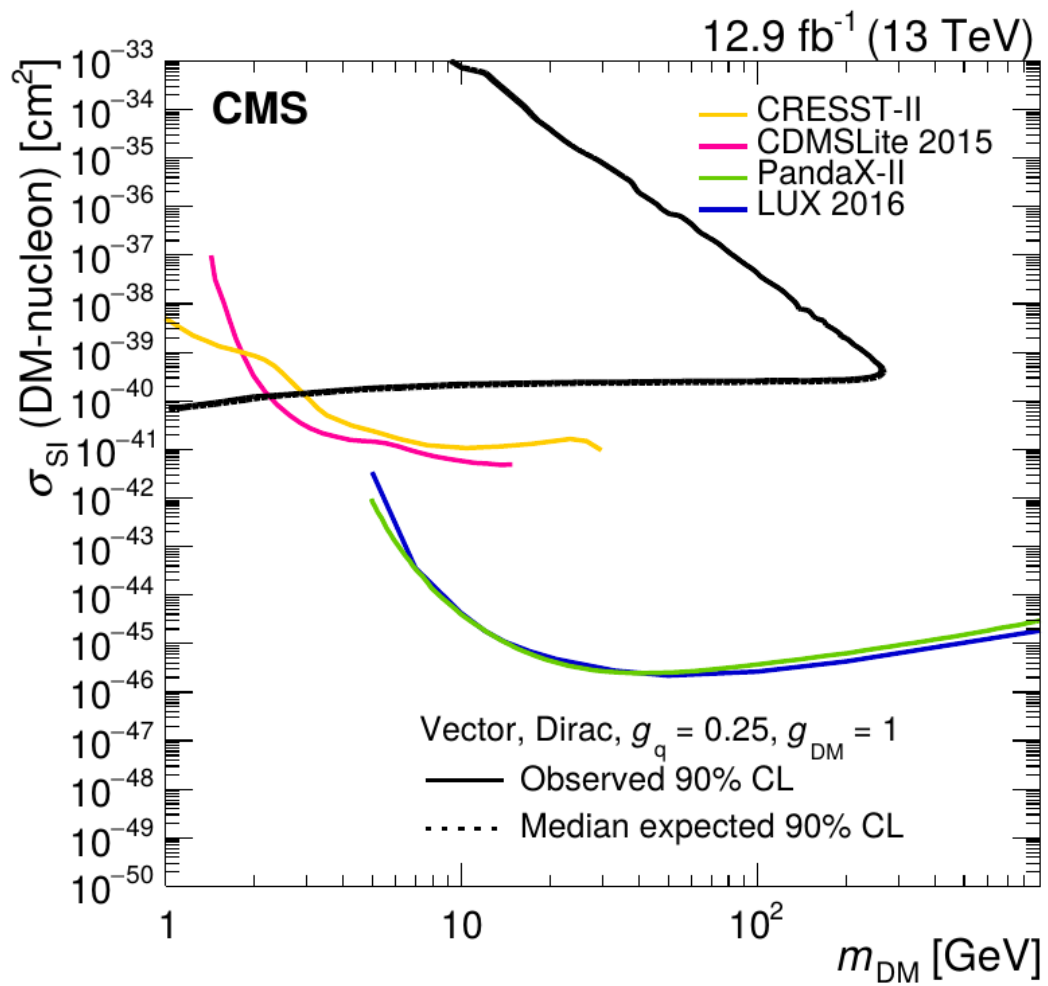**Hint**: for $n_{obs}$=0, there are no "more extreme" cases (cannot have n<0 !), so

$p_{s0}$ = Poisson(n=0 | $S_0$+B) and 1 - $p_B$ = Poisson(n=0 | B)

**Solution:** $$S_{up}(n_{obs}=0) = \log(20) = 2.996 \approx 3$$

$\Rightarrow$ **Rule of thumb: when $n_{obs}$ = 0, the 95% CL$_s$ limit is 3 events (for any B)**

# Reparameterization: Limits

CMS Run 2 Monophoton Search: measured $N_s$ in a counting experiment reparameterized according to various DM models

# Generating Pseudo-data

Model describes the distribution of the observable: **P(data; parameters)**

⇒ Possible outcomes of the experiment, for given parameter values

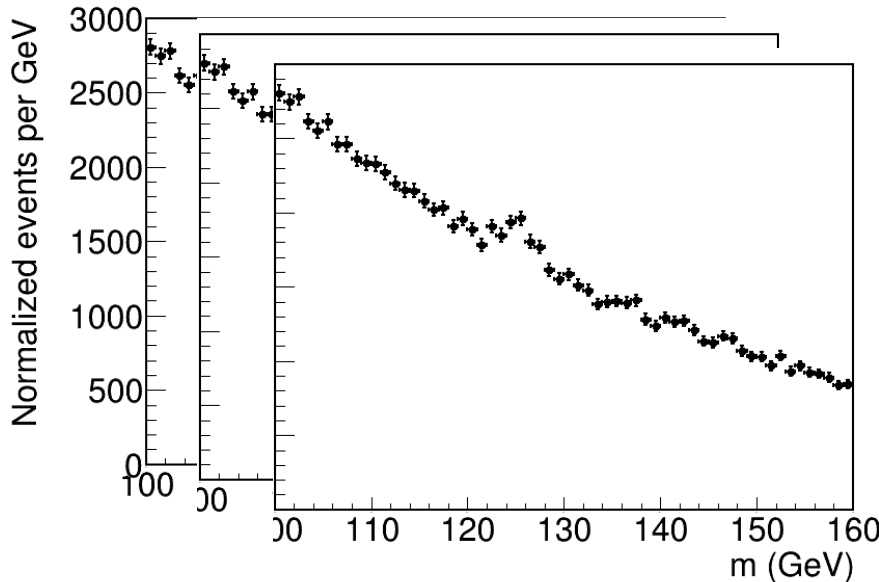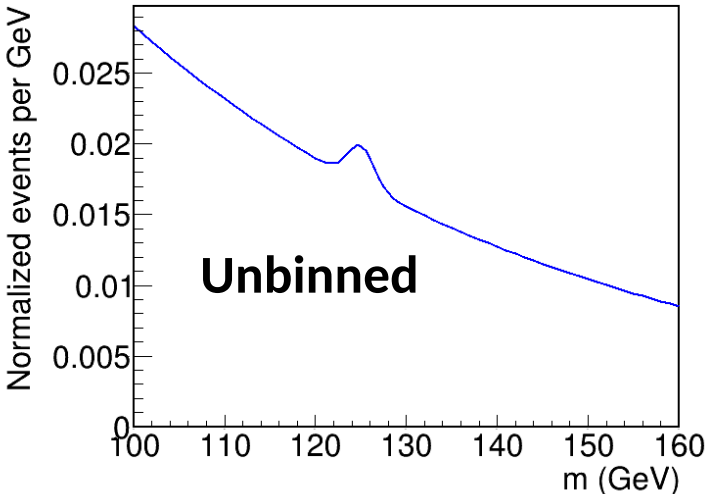Can draw random events according to PDF : **generate** *pseudo-data*

$$P(\lambda = 5)$$

➡ **2, 5, 3, 7, 4, 9, ....**

Each entry = separate "experiment"

**Generate**



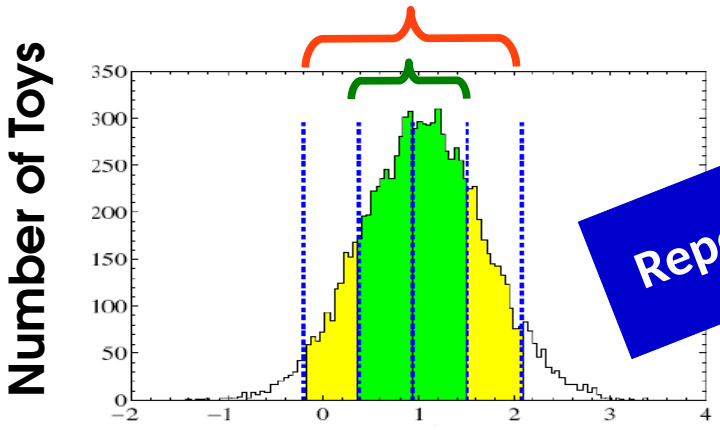Unbinned

# Expected Limits: Toys

*Expected* **results**: median outcome under a given hypothesis

→ usually B-only for searches, but other choices possible.

Two main ways to compute:
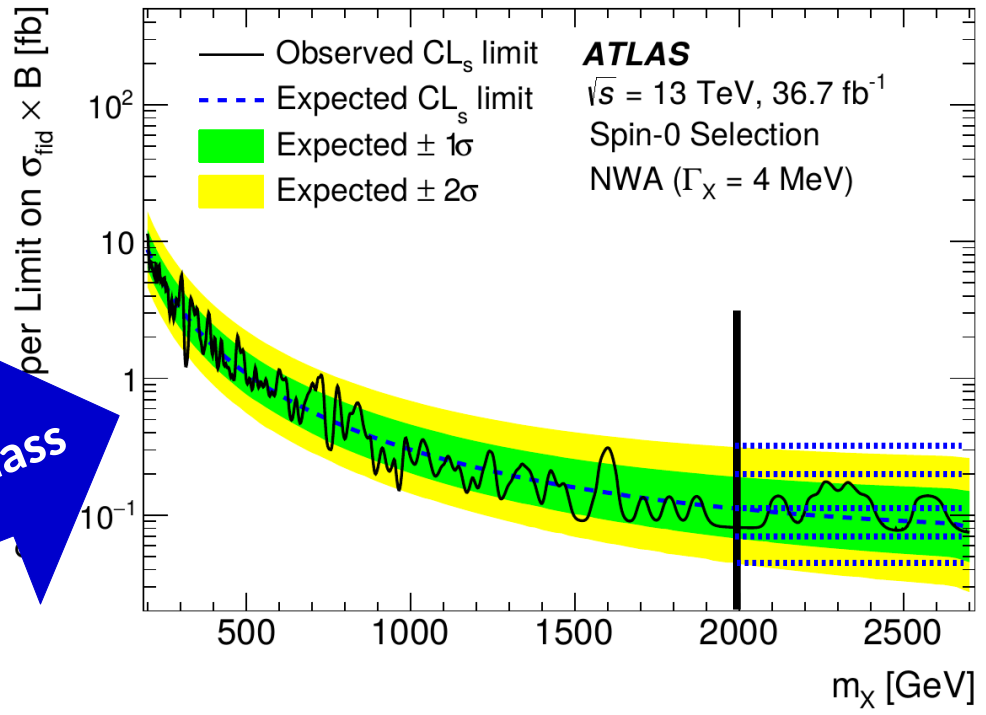
→ **Pseudo-experiments (*toys*):**

- Generate a pseudo-dataset in B-only hypothesis
- Compute limit
- Repeat and histogram the results
- Central value = median, bands based on quantiles

**68% of toys**    **95% of toys**

Phys. Lett. B 775 (2017) 105



**Repeat for each mass**

Eur.Phys.J.C71:1554,2011

**Number of Toys**

**Computed limit**

# Expected Limits: Asimov Datasets

*Expected* **results**: median outcome under a given hypothesis

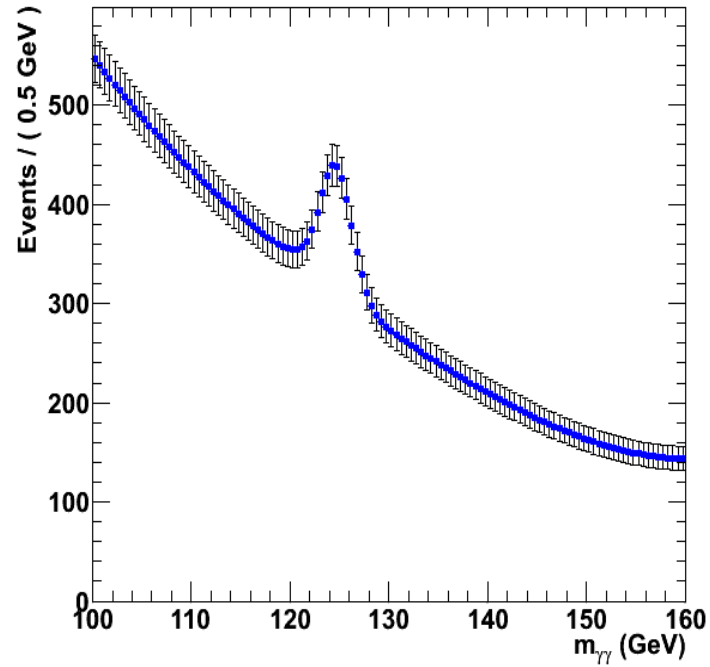→ usually B-only for searches, but other choices possible.

Two main ways to compute:

Strictly speaking, Asimov dataset if

$$\hat{X} = X_0 \text{ for all parameters X,}$$

where $X_0$ is the generation value

→ *Asimov Datasets*

- Generate a "perfect dataset" – *e.g.* for binned data, set bin contents carefully, no fluctuations.

- Gives the median result immediately:

  **median(toy results)  ↔  result(median dataset)**

- Get bands from asymptotic formulas:
  Band width

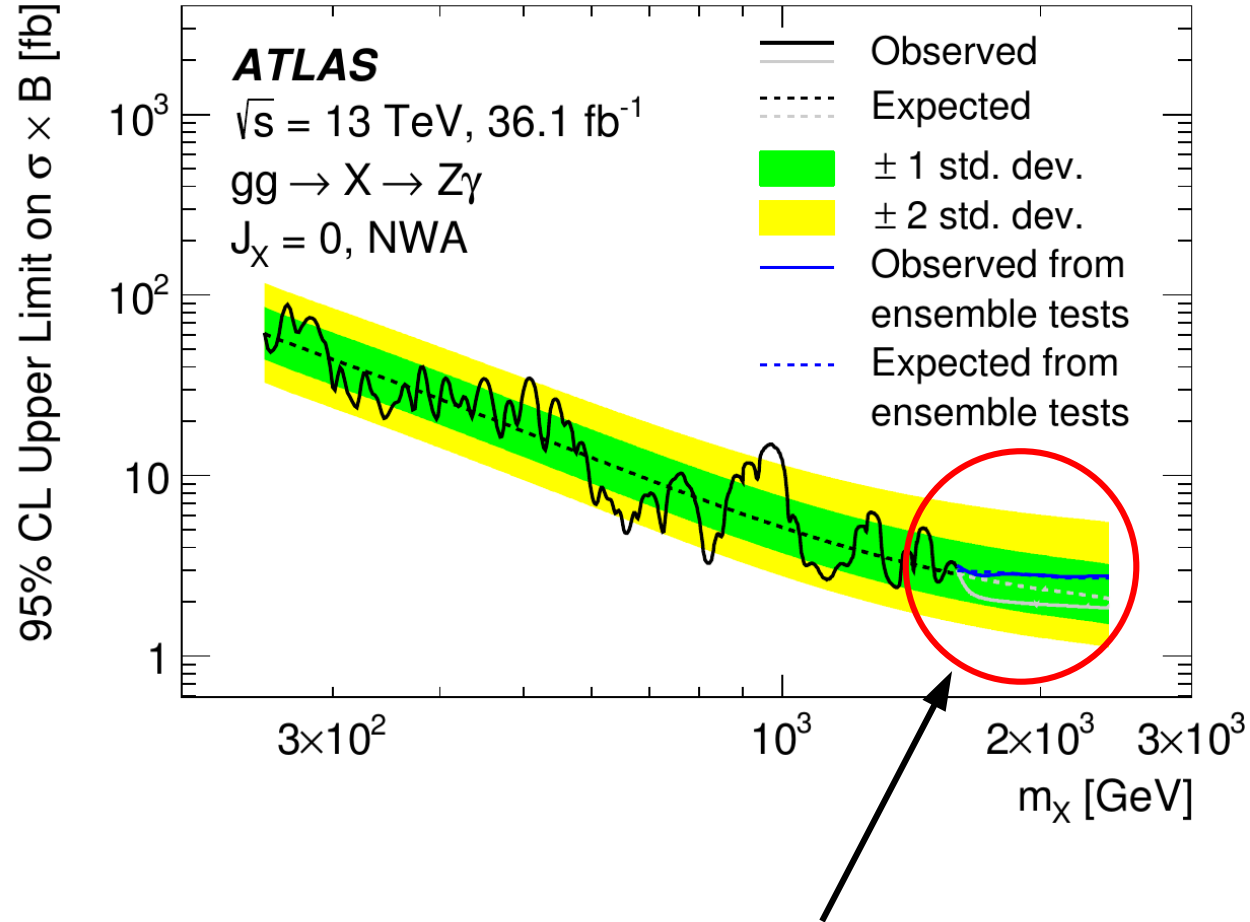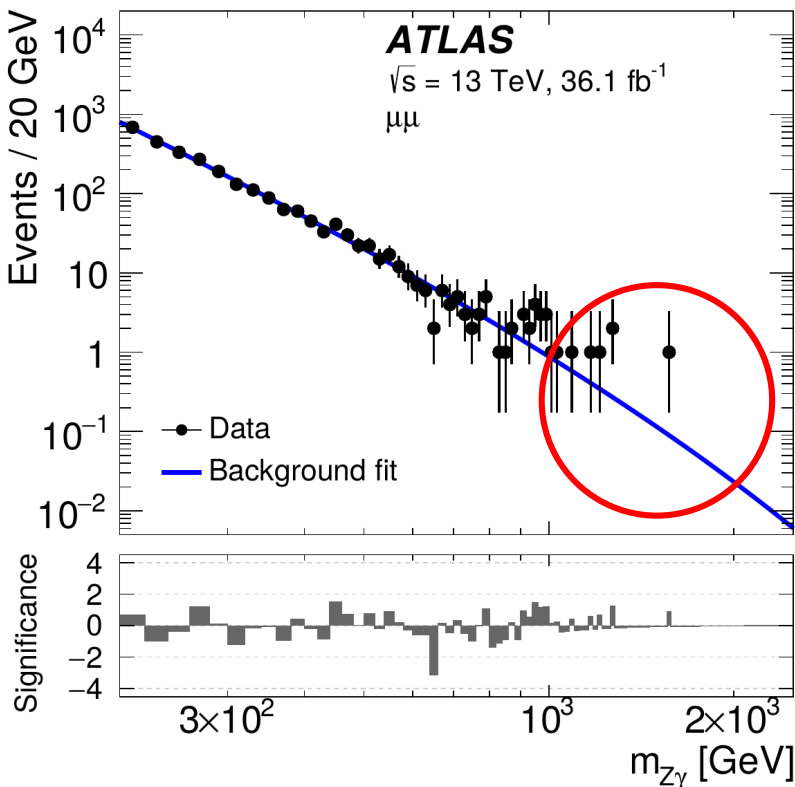$$\sigma^2_{S_0, A} = \frac{S_0^2}{q_{S_0}(\mathbf{Asimov})}$$

⊕ **Much faster (1 "toy")**

⊖ **Relies on Gaussian approximation**

ATLAS X→Zγ Search: covers 200 GeV < $m_X$ < 2.5 TeV
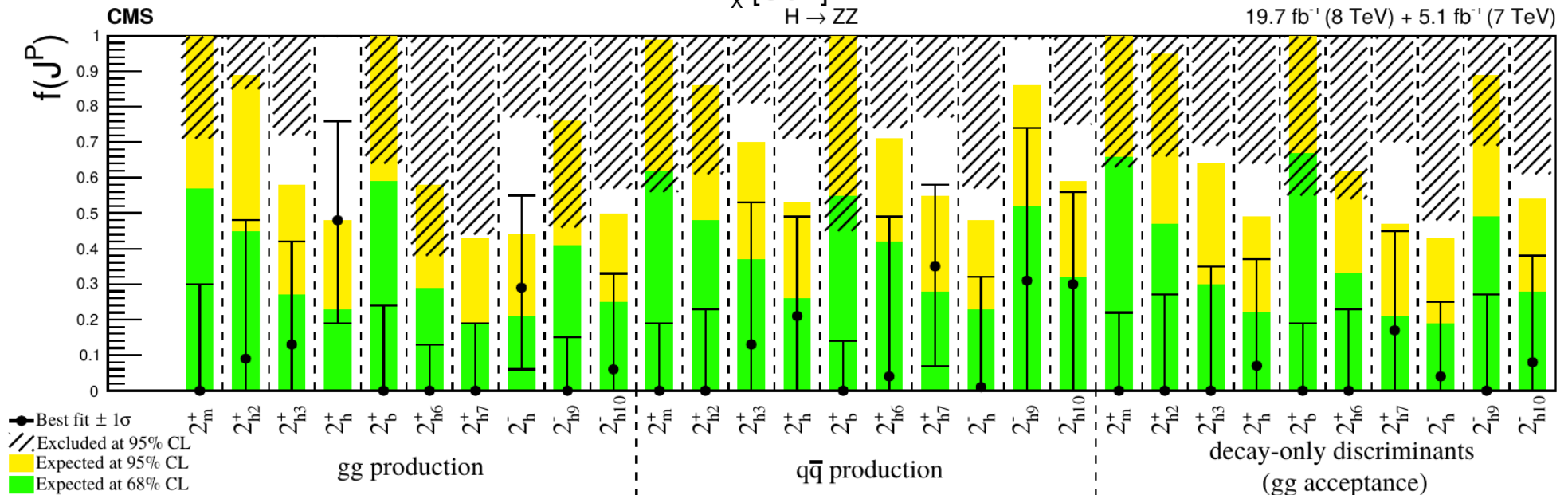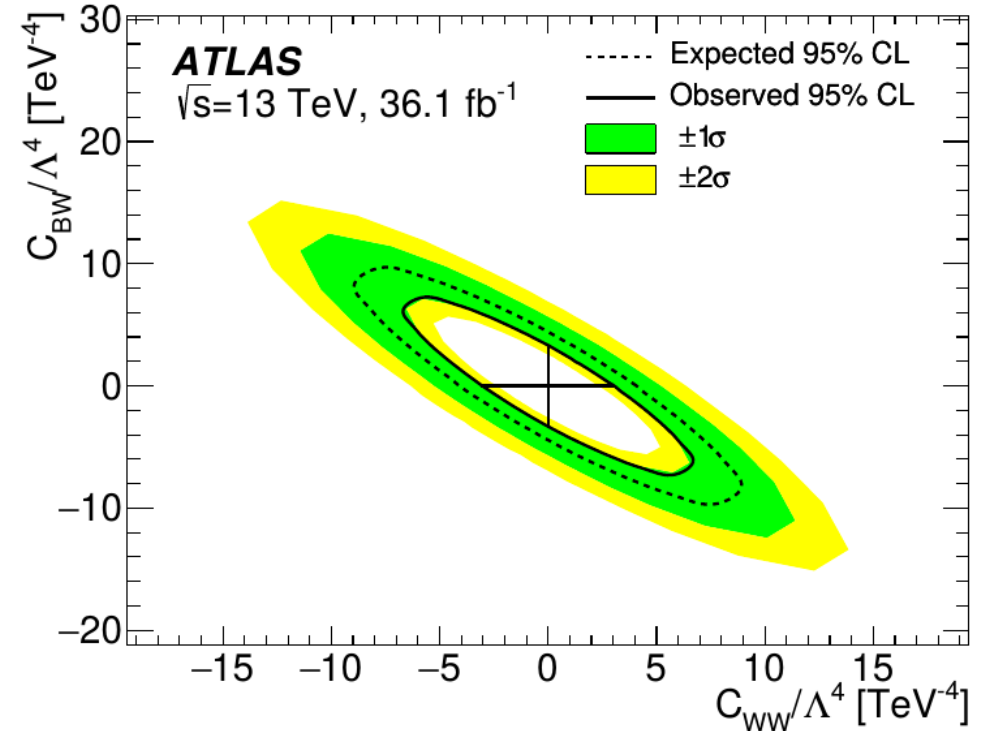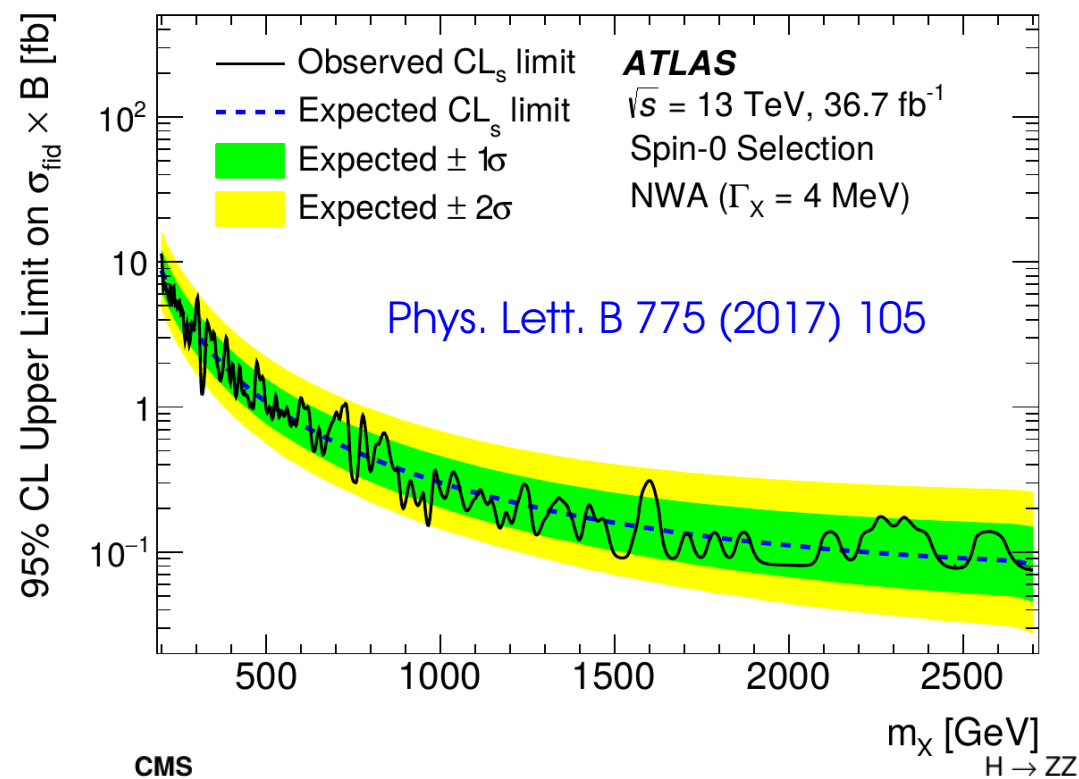
For $m_X$ > 1.6 TeV, low event counts ⇒ derive results from toys



Asimov results (**in gray**) give optimistic result compared to toys (**in blue**)

# Upper Limit Examples



ATLAS 2015-2016 4l aTGC Search

Phys. Lett. B 775 (2017) 105

Phys. Rev. D 92 (2015) 012004

56 /

# Takeaways

**Confidence intervals**: use $t_{\mu_0} = -2\log\dfrac{L(\mu = \mu_0)}{L(\hat{\mu})}$

→ Crossings with $t_{\mu_0} = Z^2$ for $\pm Z\sigma$ intervals (in 1D)

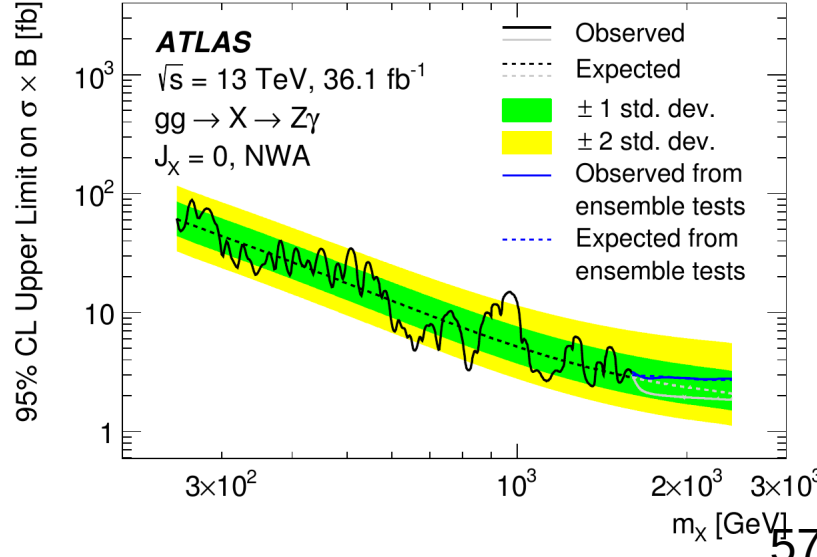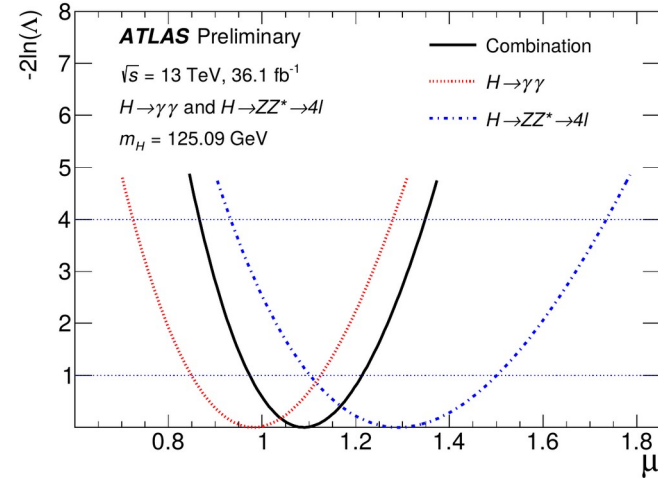**Gaussian regime**: $\mu = \hat{\mu} \pm \sigma_\mu$ (1σ interval)

**Limits** : use LR-based test statistic:

$$q_{S_0} = -2\log\frac{L(S = S_0)}{L(\hat{S})} \qquad S_0 \geq \hat{S}$$

→ Use **CL$_s$ procedure** to avoid negative limits

**Gaussian regime**, n~0: **S < Ŝ + 1.96σ at 95% CL**

**Poisson regime**, n=0 : **S$_{up}$ = 3 events at 95% CL**





57
/
/

# Extra Slides

# Rare Processes ?

**HEP** : almost always use Poisson distributions. Why ?

**ATLAS** :

- **Event rate ~ 1 GHz**

  ($L \sim 10^{34}$ cm$^{-2}$s$^{-1}$~10 nb$^{-1}$/s, $\sigma_{tot}$~$10^8$ nb, )

- **Trigger rate ~ 1 kHz**

  (Higgs rate **~ 0.1 Hz**)

$\Rightarrow$ **p ~ $10^{-6}$ ≪ 1** ($p_{H \to \gamma\gamma}$ ~ **$10^{-13}$**)

A day of data: **N ~ $10^{14}$ ≫ 1**

$\Rightarrow$ **Poisson regime! Similarly true in many other physics situations.**

(Large N = design requirement, to get not-too-small $\lambda = Np$...)

## proton - (anti)proton cross sections



W.J. Stirling, private communication

WJS2012

# Unbinned Shape Analysis

**Observable**: set of values $m_1 \ldots m_n$, one per event

$\rightarrow$ Describe shape of the **distribution of m**

$\rightarrow$ Deduce the **probability to observe $m_1 \ldots m_n$**

**H$\rightarrow\gamma\gamma$-inspired example:**

- **Gaussian signal** $\qquad P_{\text{signal}}(m) = G(m; m_H, \sigma)$
- **Exponential bkg** $\qquad P_{\text{bkg}}(m) = \alpha \, e^{-\alpha m}$
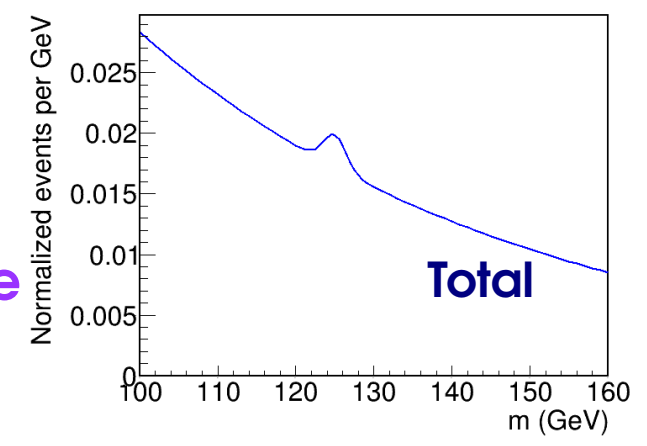
Expected yields : **S, B**

$\Rightarrow$ **Total PDF for a single event**:

$$P_{\text{total}}(m) = \frac{S}{S+B} G(m; m_H, \sigma) + \frac{B}{S+B} \alpha \, e^{-\alpha m}$$

$\Rightarrow$ **Total PDF for a dataset**

**Probability to observe the value $m_i$**

**Probability to observe n events**

$$P\left(\{m_i\}_{i=1\ldots n}\right) = e^{-(S+B)} \frac{(S+B)^n}{n!} \prod_{i=1}^{n} \frac{S}{S+B} G(m_i; m_H, \sigma) + \frac{B}{S+B} \alpha \, e^{-\alpha m_i}$$

60

# Poisson Example

Assume **Poisson distribution** with B = 0 :
Say we **observe n=5**, want to infer information on the parameter S

$$P(n\,;S) = e^{-S}\frac{S^n}{n!}$$

→ Try different values of S for a fixed data value n=5

→ Varying parameter, fixed data: **likelihood**

$$L(S\,;n=5) = e^{-S}\frac{S^5}{5!}$$



Observed
Value n=5

# Poisson Example

Assume **Poisson distribution** with B = 0 :
Say we **observe n=5**, want to infer information on the parameter **S**

$$P(n;S) = e^{-S} \frac{S^n}{n!}$$

→ Try different values of S for a fixed data value n=5

→ Varying parameter, fixed data: **likelihood**

$$L(S;n=5) = e^{-S} \frac{S^5}{5!}$$



**P(S = 0.5)**
**Low**
**likelihood**

**Read L(S; n=5) here**

**Observed**
**Value n=5**

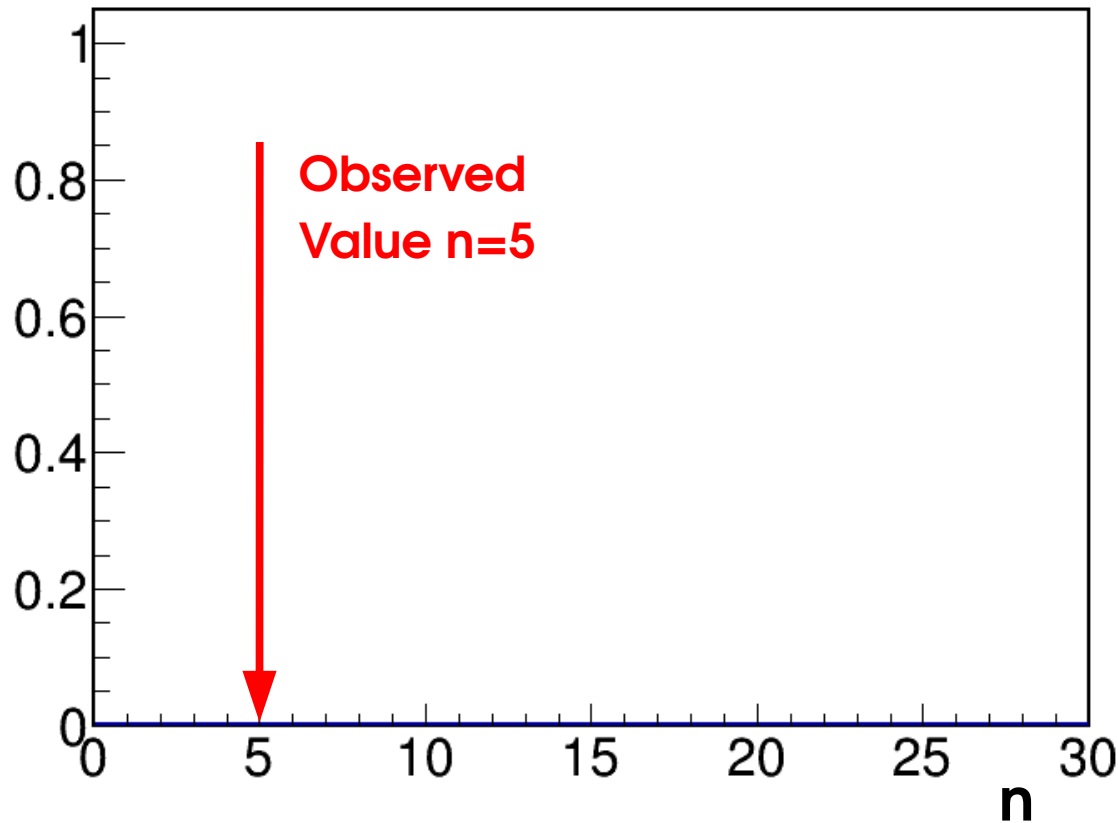# Poisson Example

Assume **Poisson distribution** with B = 0 :
Say we **observe n=5**, want to infer information on the parameter S

$$P(n;S) = e^{-S}\frac{S^n}{n!}$$

→ Try different values of S for a fixed data value n=5

→ Varying parameter, fixed data: **likelihood**

$$L(S;n=5) = e^{-S}\frac{S^5}{5!}$$



**Read L(S; n=5) here**

P(S = 0.5) Low likelihood

Observed Value n=5

P(S = 5) High likelihood

# Poisson Example

Assume **Poisson distribution** with B = 0 :
Say we **observe n=5**, want to infer information on the parameter S

$$P(n;S) = e^{-S}\frac{S^n}{n!}$$

→ Try different values of S for a fixed data value n=5

→ Varying parameter, fixed data: **likelihood**

$$L(S;n=5) = e^{-S}\frac{S^5}{5!}$$



Read L(S; n=5) here

Observed
Value n=5

P(S = 0.5)
Low
likelihood

P(S = 5)
High
likelihood

P(S = 20)
Low
likelihood

# Poisson Example

Assume **Poisson distribution** with B = 0 :
Say we **observe n=5**, want to infer information on the parameter S

$$P(n;S) = e^{-S} \frac{S^n}{n!}$$

→ Try different values of S for a fixed data value n=5

→ Varying parameter, fixed data: **likelihood**

$$L(S;n=5) = e^{-S} \frac{S^5}{5!}$$



P(S = 0.5)
Low
likelihood

**Read L(S; n=5) here**

Observed
Value n=5

P(S = 5)
High
likelihood

L(S; n=5):
Likelihood
of S for n=5

S

n

# MLEs in Shape Analyses

**Binned shape analysis:**

$$L(S; n_i) = P(n_i; S) = \prod_{i=1}^{N} \text{Pois}(n_i; S f_i + B_i)$$

Maximize global L(S) (each bin may prefer a different **S**)
In practice easier to minimize

$$\lambda_{\text{Pois}}(S) = -2\log L(S) = -2\sum_{i=1}^{N}\log \text{Pois}(n_i; S f_i + B_i) \qquad \text{**Needs a computer...**}$$
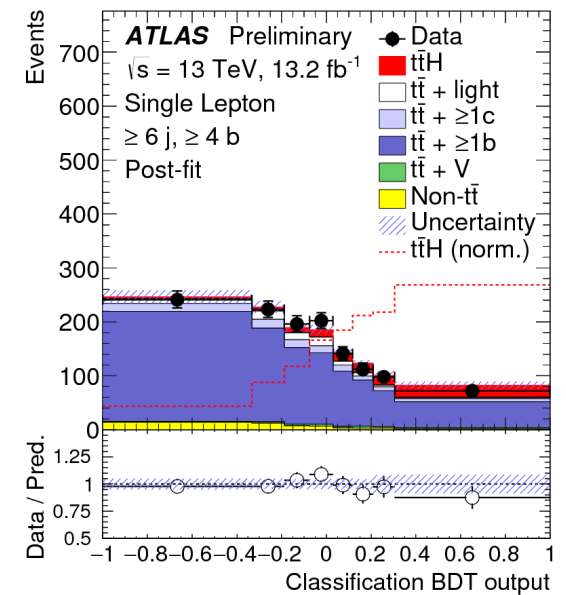
In the Gaussian limit

$$\lambda_{\text{Gaus}}(S) = \sum_{i=1}^{N} -2\log G(n_i; S f_i + B_i, \sigma_i) = \sum_{i=1}^{N}\left(\frac{n_i - (S f_i + B_i)}{\sigma_i}\right)^2 \qquad \chi^2 \text{ formula!}$$

➙ **Gaussian MLE** (min $\chi^2$ or min $\lambda_{\text{Gaus}}$) : *Best fit value* in a $\chi^2$ (Least-squares) fit

➙ **Poisson MLE** (min $\lambda_{\text{Pois}}$) : *Best fit value* in a *likelihood* fit (in `ROOT`, fit option "L")

In RooFit, $\boldsymbol{\lambda_{\text{Pois}}}$ ➙ `RooAbsPdf::fitTo()`, $\boldsymbol{\lambda_{\text{Gaus}}}$ ➙ `RooAbsPdf::chi2FitTo()`.

## In both cases, MLE ⟺ *Best Fit*

ATLAS Preliminary
√s = 13 TeV, 13.2 fb⁻¹
Single Lepton
≥ 6 j, ≥ 4 b
Post-fit

Data, tt̄H, tt̄ + light, tt̄ + ≥1c, tt̄ + ≥1b, tt̄ + V, Non-tt̄, Uncertainty, tt̄H (norm.)

Classification BDT output

# H→γγ

$$L(S, B; m_i) = e^{-(S+B)} \prod_{i=1}^{n_{evts}} S\, P_{sig}(m_i) + B\, P_{bkg}(m_i)$$



*Estimate the MLE $\hat{S}$ of $S$ ?*

→ Perform (likelihood) best-fit of model to data

⇒ fit result for S is the desired $\hat{S}$.

In particle physics, often use the *MINUIT* minimizer within ROOT.

# MLE Properties

- **Asymptotically Gaussian** 
  $$P(\hat{\mu}) \propto \exp\left(-\frac{(\hat{\mu}-\mu^*)^2}{2\sigma_{\hat{\mu}}^2}\right) \quad \text{for } n \to \infty$$

  **and unbiased** $\langle\hat{\mu}\rangle = \mu^*$ for $n \to \infty$

  Standard deviation of the distribution of $\hat{\mu}$

  for large enough datasets

- **Asymptotically Efficient** : $\sigma_{\hat{\mu}}$ is the **lowest possible value** (in the limit n→∞)

  among consistent estimators.

  → MLE captures all the available information in the data

- Also **consistent**: $\hat{\mu}$ converges to the true value for large n, $\quad \hat{\mu} \xrightarrow{n\to\infty} \mu^*$

- **Log-likelihood :** Can also **minimize** $\lambda$ = -2 log L

  → Usually more efficient numerically

  → For Gaussian L, $\lambda$ is parabolic:

- Can **drop multiplicative constants in L** (additive constants in $\lambda$)

# Extra: Fisher Information

**Fisher Information:**

$$I(\mu) = \left\langle \left( \frac{\partial}{\partial \mu} \log L(\mu) \right)^2 \right\rangle = -\left\langle \frac{\partial^2}{\partial \mu^2} \log L(\mu) \right\rangle$$

Measures the **amount of information** available in the measurement of $\mu$.

**Gaussian likelihood:** $\qquad I(\mu) = \dfrac{1}{\sigma_{Gauss}^2}$

$\rightarrow$ smaller $\sigma_{Gauss}$ $\Rightarrow$ more information.

**Cramer-Rao bound:** $\qquad \mathrm{Var}(\widetilde{\mu}) \geq \dfrac{1}{I(\mu)}$

For any estimator $\widetilde{\mu}$.

$\rightarrow$ cannot be more precise than allowed by information in the measurement.

**Efficient** estimators reach the bound : **e.g. MLE in the large dataset limit.**

---

**Gaussian case**:
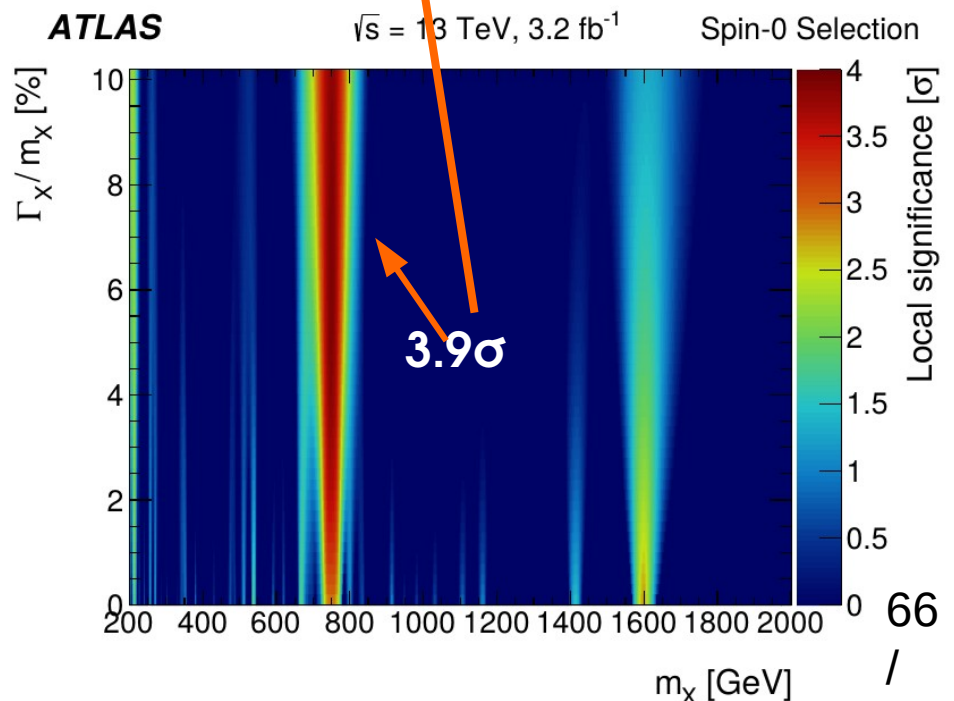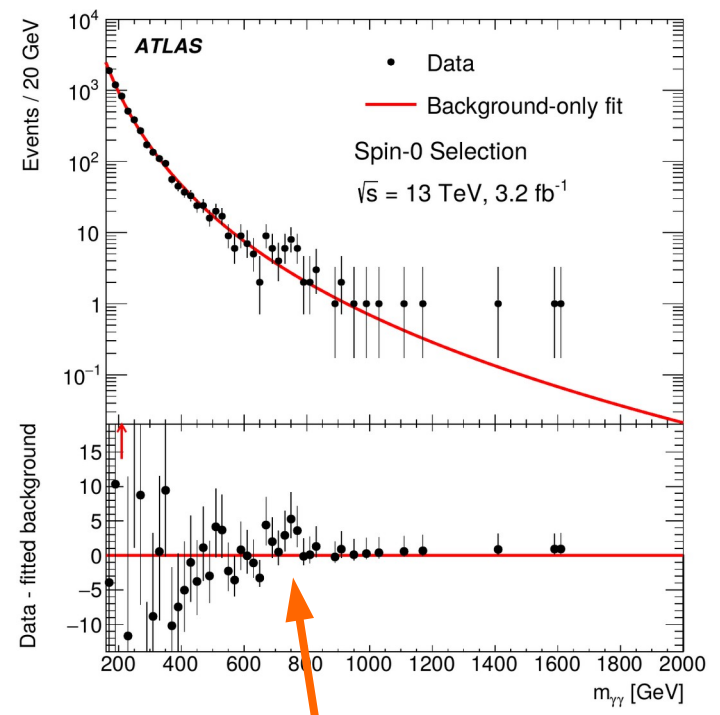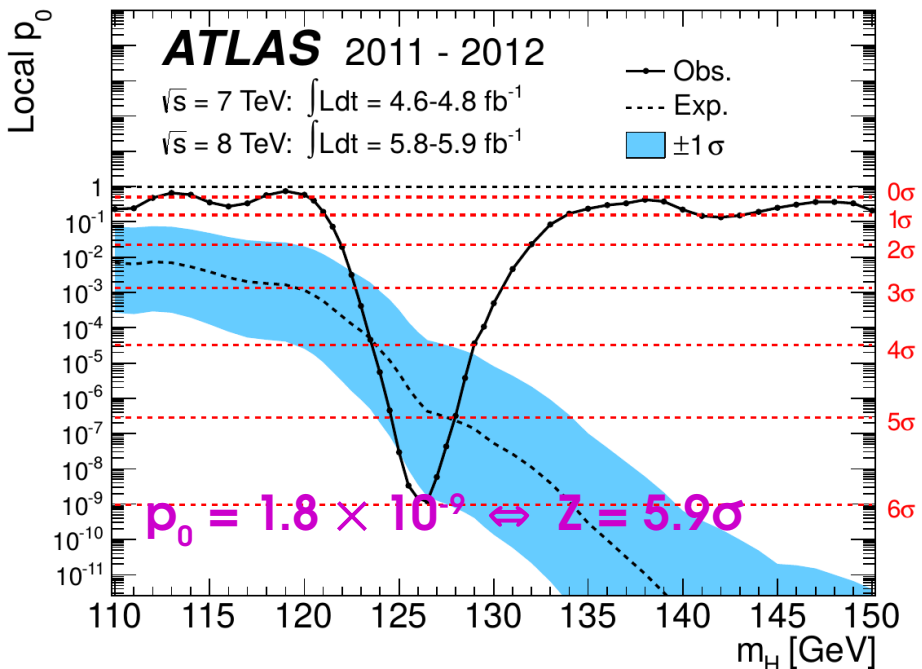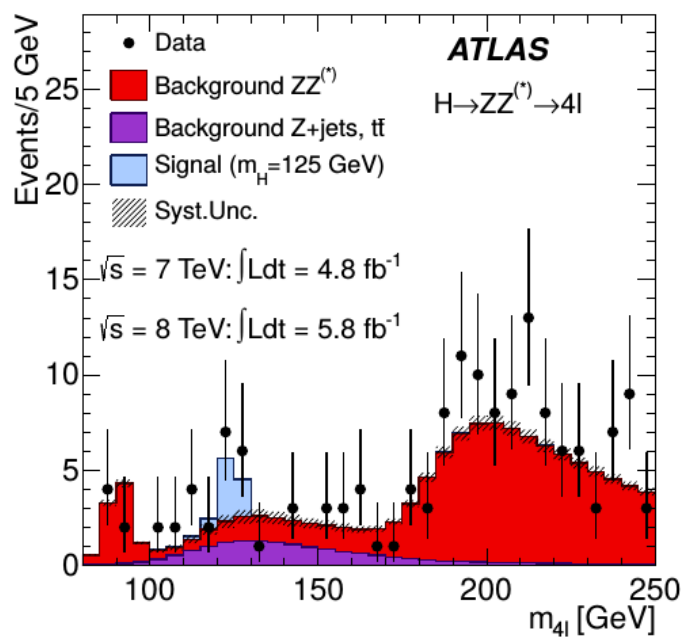- For a Gaussian estimator $\widetilde{\mu}$

$$P(\widetilde{\mu}) \propto \exp\left( -\frac{(\widetilde{\mu} - \mu^*)^2}{2\sigma_{\widetilde{\mu}}^2} \right)$$

- **MLE**: Var($\hat{\mu}$) = $\sigma_{\hat{\mu}}^2$

**Cramer-Rao: Var($\widetilde{\mu}$) $\geq$ $\sigma_{Gauss}^2$ = $\sigma_{\widetilde{\mu}}^2$**

# Some Examples

## Higgs Discovery: Phys. Lett. B 716 (2012) 1-29



$$p_0 = 1.8 \times 10^{-9} \Leftrightarrow Z = 5.9\sigma$$



3.9σ

# Upper Limit Pathologies

**Upper limit**: $S_{up} \sim \hat{S} + 1.64\, \sigma_s$.

**Problem**: for negative $\hat{S}$, get **very** good observed limit.
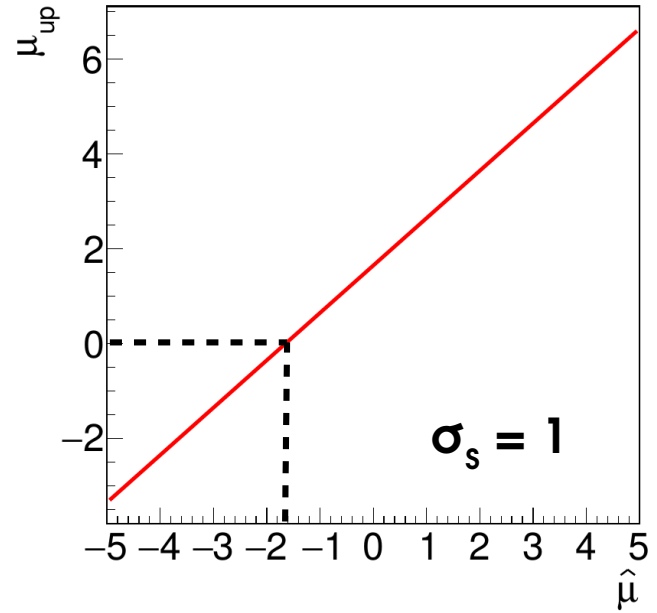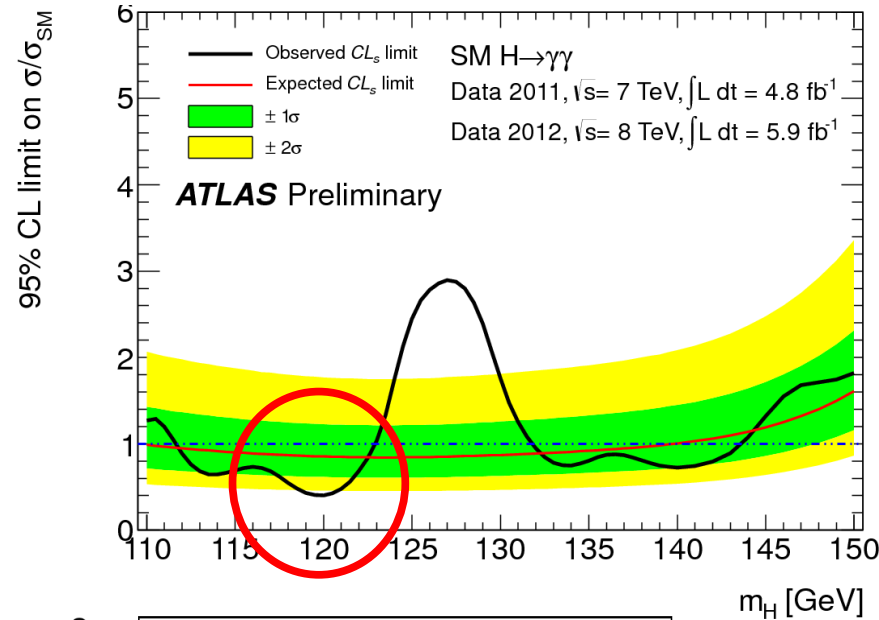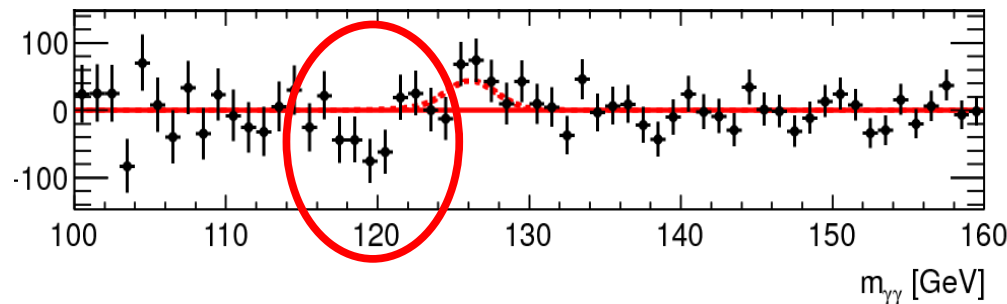
→ For $\hat{S}$ sufficiently negative, even $S_{up} < 0$ !

How can this be ?

→ **Background modeling issue ?...** Or:

→ This is a **95%** limit ⇒ **5% of the time, the limit wrongly excludes the true value**, *e.g.* S*=0.

**Options**

→ **live with it**: sometimes report limit $< 0$

→ **Special procedure to avoid these cases**, since if we assume S must be >0, we know a priori this is just a fluctuation.





$\sigma_s = 1$

# CL$_s$

Usual solution in HEP : **CL$_s$**.

→ Compute modified p-value

$$p_{CL_s} = \frac{p_{S_0}}{(1 - p_B)}$$

**The usual p-value under** H(S=S$_0$) (=5%)

**The p-value computed under H(S=0)**

⇒ **Rescale** exclusion at S$_0$ by exclusion at S=0.
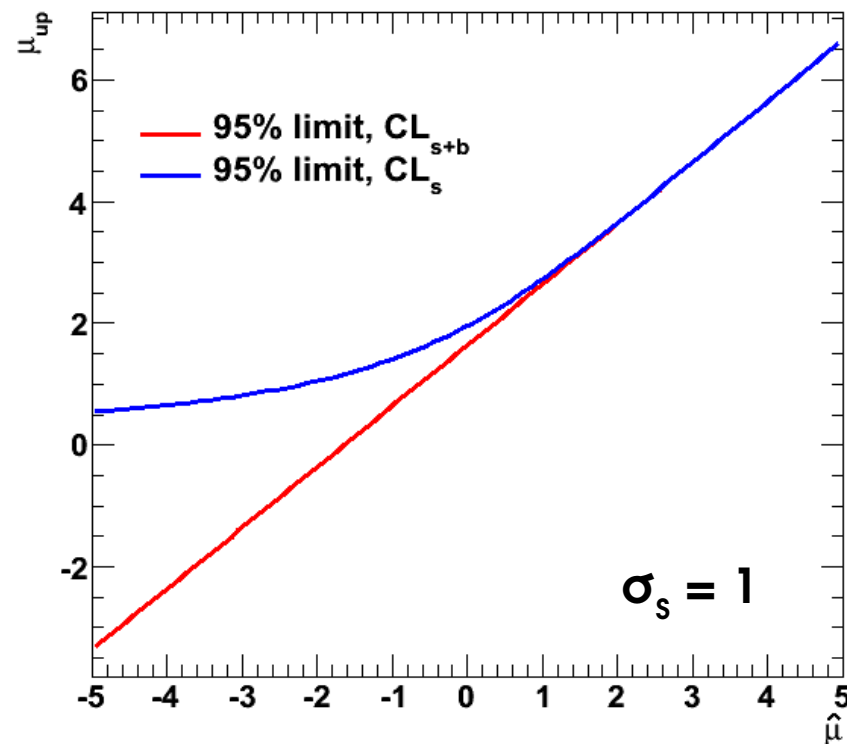
→ Somewhat ad-hoc, but good properties...

**Ŝ compatible with 0** : $p_B \sim O(1)$

$p_{CLs} \sim p_{S0} \sim$ **5%, no change**.

**Far-negative Ŝ** : $1 - p_B \ll 1$

$p_{CLs} \sim p_{S0}/(1-p_B) \gg$ **5%**

→ lower exclusion ⇒ higher limit, usually >0 as desired



95% limit, CL$_{s+b}$
95% limit, CL$_s$

$\sigma_s = 1$

**Drawback**: *overcoverage*

→ limit is claimed to be 95% CL, but actually >95% CL for small 1-$p_B$.

# CL$_s$ : Gaussian Bands

Usual Gaussian counting example with known B:
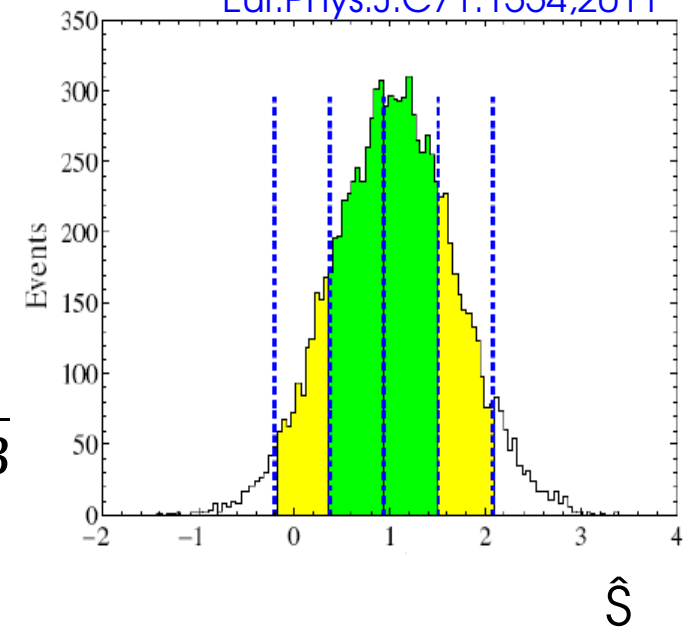
95% CL$_s$ upper limit on S:

$$S_{up} = \hat{S} + \left[\Phi^{-1}\left(1 - 0.05\ \Phi\left(\hat{S}/\sigma_s\right)\right)\right]\sigma_S \qquad \text{with} \quad \sigma_S = \sqrt{B}$$
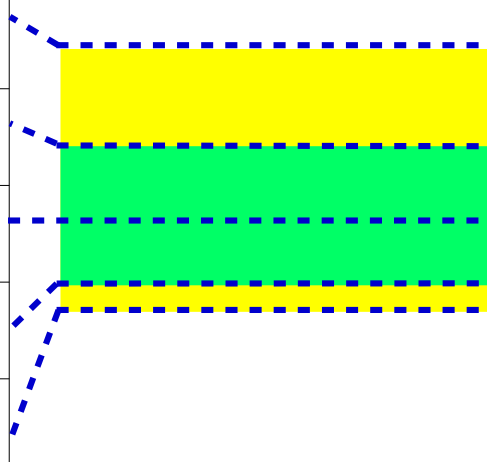
Compute expected bands for S=0:

→ **Asimov dataset ⇔ Ŝ = 0** : 
$$S^0_{up,exp} = 1.96\ \sigma_S$$

→ **± nσ bands**:
$$S^{\pm n}_{up,exp} = \left(\pm n + \left[1 - \Phi^{-1}\left(0.05\ \Phi(\mp n)\right)\right]\right)\sigma_S$$

Ŝ

| n | $S_{exp}^{\pm n}$ / $\sqrt{B}$ |
|---|---|
| +2 | 3.66 |
| +1 | 2.72 |
| 0 | 1.96 |
| -1 | 1.41 |
| -2 | 1.05 |

**CLs** :
- Positive bands somewhat reduced,
- Negative ones more so

Band width from $\sigma^2_{S,A} = \dfrac{S^2}{q_S(\textbf{Asimov})}$ depends on S, for non-Gaussian cases,different values for each band...
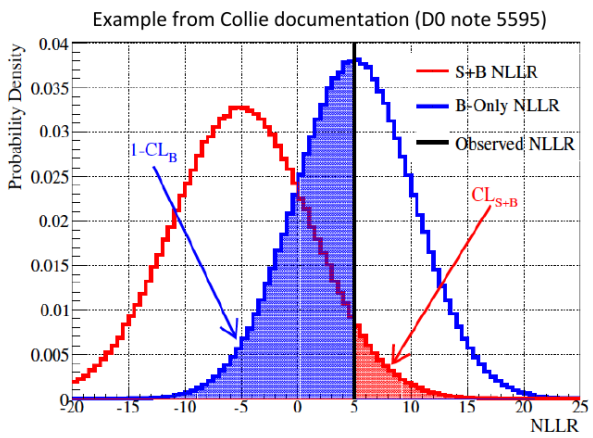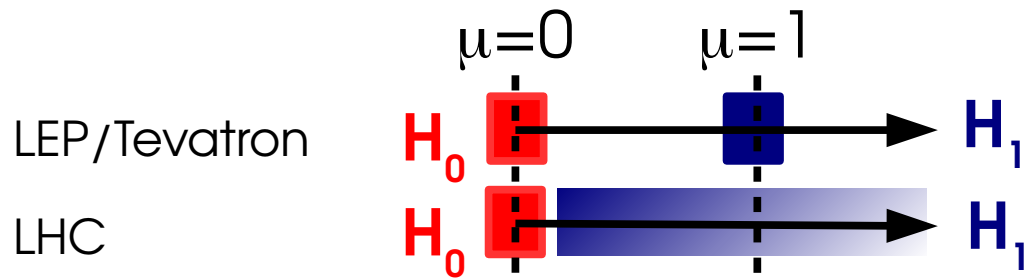
# Comparison with LEP/TeVatron definitions

Likelihood ratios are not a new idea:

- **LEP**: Simple LR with NPs from MC
  - Compare μ=0 and μ=1
- **Tevatron**: PLR with profiled NPs

$$q_{LEP} = -2\log\frac{L(\mu=0, \widetilde{\theta})}{L(\mu=1, \widetilde{\theta})}$$

$$q_{Tevatron} = -2\log\frac{L(\mu=0, \hat{\hat{\theta}}_0)}{L(\mu=1, \hat{\hat{\theta}}_1)}$$

Both compare to **μ=1** instead of best-fit **μ̂**



μ=0    μ=1

LEP/Tevatron    $H_0$ ──────► $H_1$

LHC    $H_0$ ──────► $H_1$

→ Asymptotically:

- **LEP/Tevaton**: q linear in μ ⇒ **~Gaussian**
- **LHC**: q quadratic in μ ⇒ **~χ2**

→ Still use TeVatron-style for discrete cases



Example from Collie documentation (D0 note 5595)