



Combine Tutorial

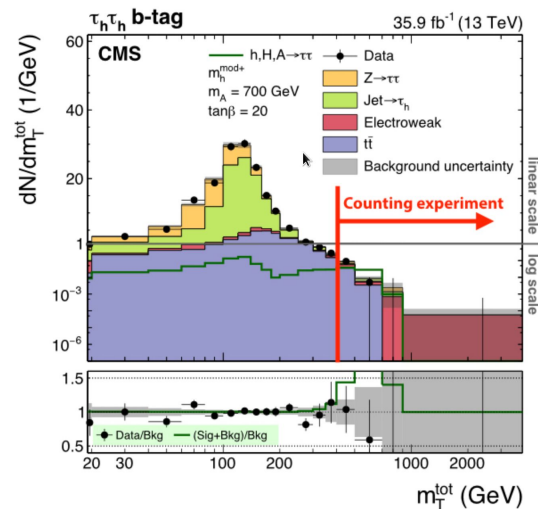
Terascale Statistics School
Kyle Cormier, Aliya Namigova, Nick Wardle

Introduction

Overview

Search for a heavy neutral higgs,
Which decays to $\tau\tau$

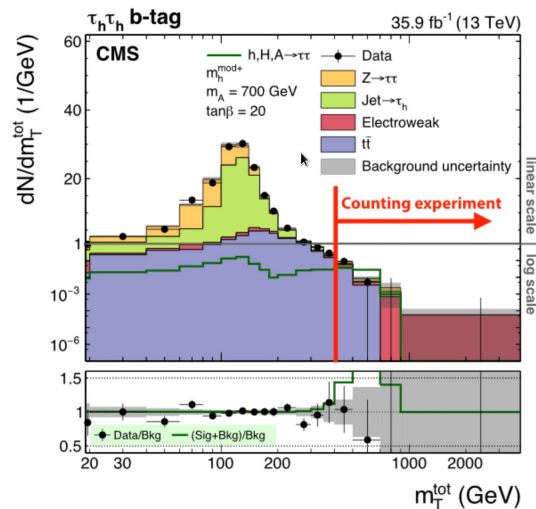
1. Simple Counting Experiment
2. Shape Based Analysis
3. Adding Control Regions
4. Physics models – beyond a single
Signal strength parameter



Overview

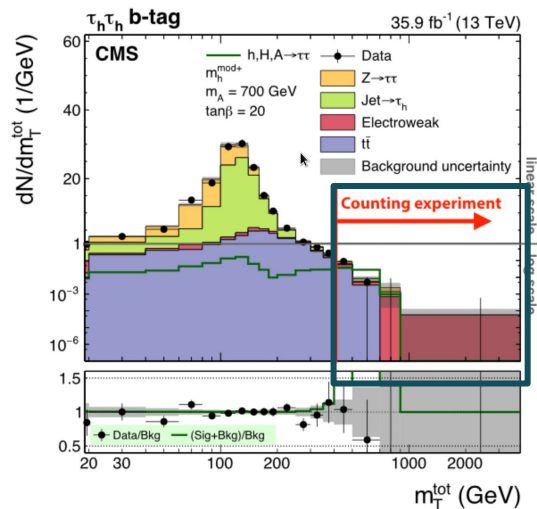
Search for a heavy neutral higgs,
Which decays to $\tau\tau$

- Limit Setting
- Significance Testing
- Asymptotic Calculations
- Toy-based Calculations
- Parameter Extraction
- Fit Debugging
- Fit Plotting
- Nuisance Parameter Checks
- Multi-Dimensional Likelihood Scans



Part 1 – Counting Experiment

Counting Experiment



Count number of events in high m_T region

imax

1

number of bins

jmax

4

number of processes minus 1

kmax

*

number of nuisance parameters

observation

bin

signal_region

backgrounds

signal

observation

10.0

bin

signal_region

signal_region

signal_region

signal_region

signal_reg

process

ttbar

diboson

Ztautau

jetFakes

bbHtautau

process

1

2

3

4

0

rate

4.43803

3.18309

3.7804

1.63396

0.711064

CMS_eff_b

lnN

1.02

1.02

1.02

-

1.02

CMS_eff_t

lnN

1.12

1.12

1.12

-

1.12

CMS_eff_t_highpt

lnN

1.1

1.1

1.1

-

1.1

acceptance_Ztautau

lnN

-

-

1.08

-

-

acceptance_bbH

lnN

-

-

-

-

1.05

acceptance_ttbar

lnN

1.005

-

-

-

-

norm_jetFakes

lnN

-

-

-

1.2

-

xsec_diboson

lnN

-

1.05

-

-

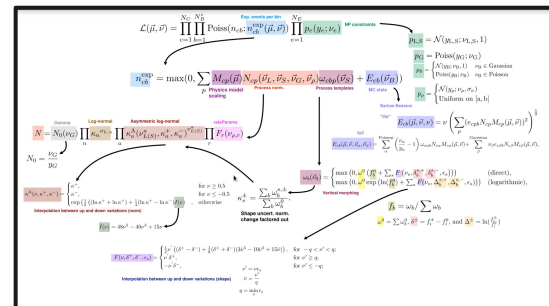
-

Systematic uncertainties

Neyman-Pearson doesn't usually help

We usually don't have explicit formulae for the pdfs $f(x|s)$, $f(x|b)$, so for a given x we can't evaluate the likelihood ratio

$$t(x) = \frac{f(x|s)}{f(x|b)}$$



Use a modified version of the likelihood ratio test statistic

Test statistic based on likelihood ratio

How can we choose a test's critical region in an 'optimal way', in particular if the data space is multidimensional?

Neyman-Pearson lemma states:

For a test of H_0 of size α , to get the highest power with respect to the alternative H_1 we need for all x in the critical region W

"likelihood ratio (LR)" $\rightarrow \frac{P(x|H_1)}{P(x|H_0)} \geq c_\alpha$

inside W and $\leq c_\alpha$ outside, where c_α is a constant chosen to give a test of the desired size.

Equivalently, optimal scalar test statistic is $t(x) = \frac{P(x|H_1)}{P(x|H_0)}$

N.B. any monotonic function of this is leads to the same test.

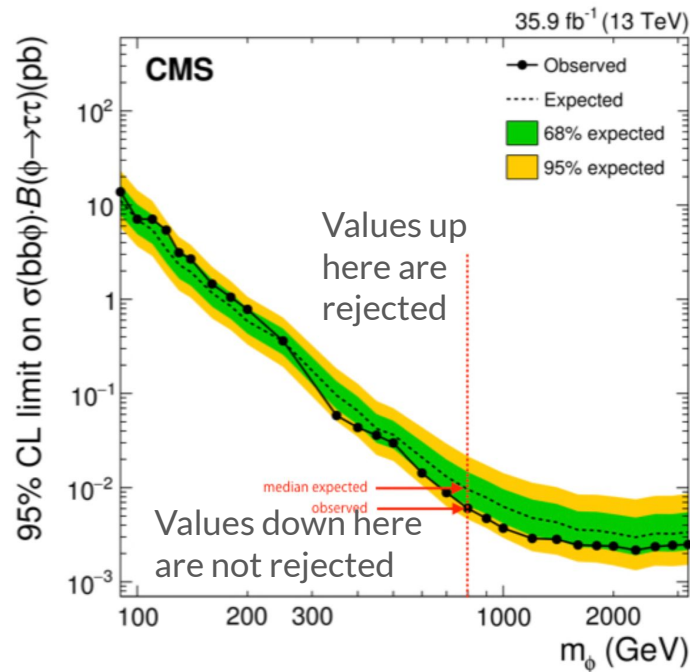
G. Cowan / RHUL Physics

Terascale Statistics 2024 / Lecture 1

41

We use:

$$\tilde{q}_\mu = \begin{cases} -2 \log \left(\frac{\mathcal{L}(\mu)}{\mathcal{L}(\hat{\mu}=0)} \right) & \hat{\mu} < 0 \\ -2 \log \left(\frac{\mathcal{L}(\mu)}{\mathcal{L}(\hat{\mu})} \right) & 0 < \hat{\mu} < \mu \\ 0 & \mu < \hat{\mu} \end{cases} \quad (\text{And use CL}_s \text{ criterion})$$



Asymptotic vs Using Toys

Wilks' Theorem

Wilks' Theorem: if the hypothesized $\mu_i(\theta)$, $i = 1, \dots, N$, are true for some choice of the parameters $\theta = (\theta_1, \dots, \theta_M)$, then in the large sample limit (and provided regularity conditions are satisfied)

$$t_\mu = -2 \ln \frac{L(\boldsymbol{\mu}(\hat{\boldsymbol{\theta}}))}{L(\hat{\boldsymbol{\mu}})}$$

MLE of $(\theta_1, \dots, \theta_M)$ follows a chi-square distribution for $N - M$ degrees of freedom.

MLE of (μ_1, \dots, μ_N)

The regularity conditions include: the model in the numerator of the likelihood ratio is “nested” within the one in the denominator, i.e., $\boldsymbol{\mu}(\theta)$ is a special case of $\boldsymbol{\mu} = (\mu_1, \dots, \mu_N)$.

Proof boils down to having all estimators \sim Gaussian.

S.S. Wilks, *The large-sample distribution of the likelihood ratio for testing composite hypotheses*, Ann. Math. Statist. **9** (1938) 60-2.

Under some conditions the distribution of the test statistic is known analytically:

→ Use asymptotic approximation

Otherwise:

Generate many sets of pseudodata to get an empirical distribution of the test statistic

Asymptotic Approximation

```
<<< Combine >>>
<<< v9.1.0 >>>
>>> Random number generator seed is 123456
>>> Method used is AsymptoticLimits

-- AsymptoticLimits ( CLs ) --
Observed Limit: r < 10.8183
Expected 2.5%: r < 7.0537
Expected 16.0%: r < 9.8108
Expected 50.0%: r < 14.5625
Expected 84.0%: r < 22.3988
Expected 97.5%: r < 33.5971
```

Empirical Distribution

```
-- Hybrid New --
Limit: r < 11.1291 +/- 0.163054 @ 95% CL
```

Expected

2.5%	Limit: r < 5.46875 +/- 0.15625 @ 95% CL
16.0%	Limit: r < 10.4676 +/- 0.123997 @ 95% CL
50.0%	Limit: r < 14.5396 +/- 0.136762 @ 95% CL
84.0%	Limit: r < 21.7222 +/- 0.271188 @ 95% CL
97.5%	Limit: r < 33.2392 +/- 1.62741 @ 95% CL

Empirical test statistic distributions

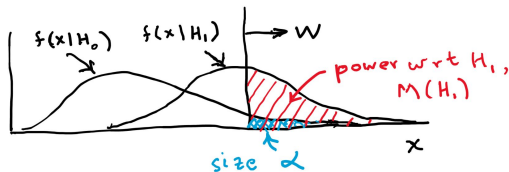
Can directly look at the distributions of the test statistics under the background-only and signal_background hypothesis

Definition of a test (2)

But in general there are an infinite number of possible critical regions that give the same size α .

Use the alternative hypothesis H_1 to motivate where to place the critical region.

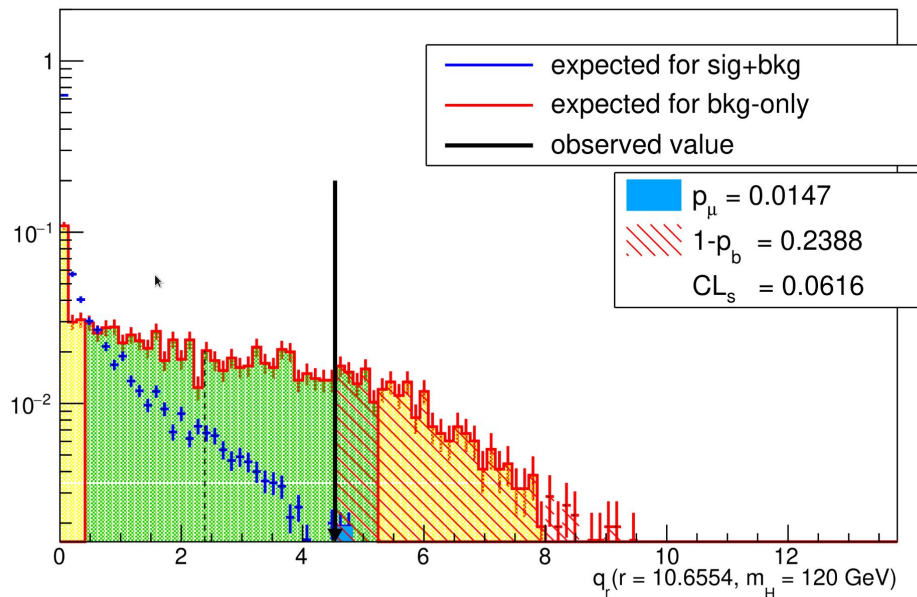
Roughly speaking, place the critical region where there is a low probability (α) to be found if H_0 is true, but high if H_1 is true:



G. Cowan / RHUL Physics

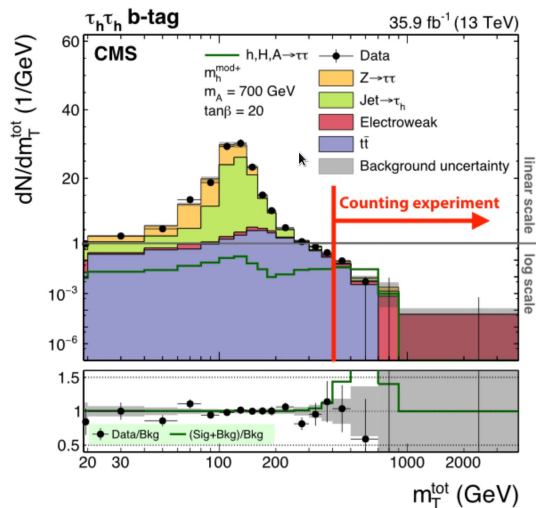
Terascale Statistics 2024 / Lecture 1

10



This point ($r=10.6554$) is (just barely) not rejected

Part 2 - Shape Experiment



Use the full histogram!

imax 1
 jmax 1
 kmax *

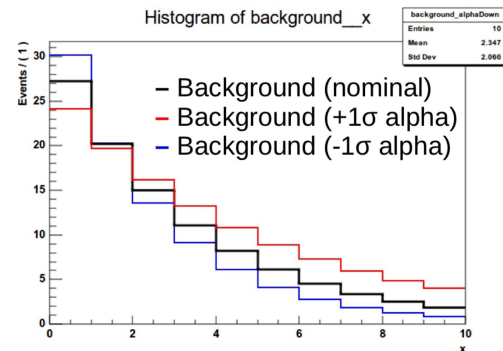
Input histograms

```
shapes * * simple-shapes-TH1_input.root $PROCESS $PROCESS_$SYSTEMATIC
shapes signal * simple-shapes-TH1_input.root $PROCESS$MASS $PROCESS$MASS_$SYSTEMATIC
```

bin bin1
 observation 85

bin	bin1	bin1
process	signal	background
process	0	1
rate	10	100

lumi	lnN	1.10	1.0
bgnorm	lnN	1.00	1.3
alpha	shape	-	1



```

imax      1 number of bins
jmax      4 number of processes minus 1
kmax      * number of nuisance parameters
-----
shapes bbHtautau * datacard_part2.shapes.root signal_region/$PROCESS$MASS signal_region/$PROCESS$MASS_$SYSTEMATIC
shapes * * datacard_part2.shapes.root signal_region/$PROCESS signal_region/$PROCESS_$SYSTEMATIC
-----
bin          signal_region
observation  3416.0
-----
bin          signal_region  signal_region  signal_region  signal_region  signal_region
process      ttbar          diboson       Ztautau       jetFakes       bbHtautau
process      1              2              3              4              0
rate         683.017        96.5185       742.649       2048.94       0.913183
-----
CMS_eff_b    lnN          1.02          1.02          1.02          -              1.02
CMS_eff_t    lnN          1.12          1.12          1.12          -              1.12
acceptance_Ztautau lnN          -              -              1.08          -              -
acceptance_bbH   lnN          -              -              -              -              1.05
acceptance_ttbar lnN          1.005         -              -              -              -
lumi_13TeV     lnN          1.025         1.025         1.025         -              1.025
norm_jetFakes  lnN          -              -              -              1.2           -
xsec_Ztautau   lnN          -              -              1.04          -              -
xsec_diboson   lnN          -              1.05          -              -              -
xsec_ttbar     lnN          1.06          -              -              -              -
# These ones are new
top_pt_ttbar_shape shape      1              -              -              -              -
CMS_scale_t_1prong0pi0_13TeV shape      1              1              1              -              1
CMS_scale_t_1prong1pi0_13TeV shape      1              1              1              -              1
CMS_scale_t_3prong0pi0_13TeV shape      1              1              1              -              1
CMS_eff_t_highpt shape      1              1              1              -              1

```

Shape-based analysis improved limits over simple counting experiment

```
-- AsymptoticLimits ( CLs ) --  
Observed Limit: r < 7.9771  
Expected 2.5%: r < 4.7720  
Expected 16.0%: r < 6.8417  
Expected 50.0%: r < 10.5312  
Expected 84.0%: r < 16.9959  
Expected 97.5%: r < 26.5059
```

Shape-based

```
<<< Combine >>>  
<<< v9.1.0 >>>  
>>> Random number generator seed is 123456  
>>> Method used is AsymptoticLimits  
  
-- AsymptoticLimits ( CLs ) --  
Observed Limit: r < 10.8183  
Expected 2.5%: r < 7.0537  
Expected 16.0%: r < 9.8108  
Expected 50.0%: r < 14.5625  
Expected 84.0%: r < 22.3988  
Expected 97.5%: r < 33.5971
```

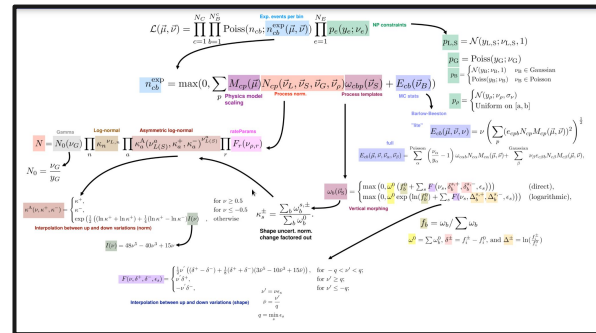
Counting

Fit diagnostics

Neyman-Pearson doesn't usually help

We usually don't have explicit formulae for the pdfs $f(\mathbf{x}|s)$, $f(\mathbf{x}|b)$, so for a given \mathbf{x} we can't evaluate the likelihood ratio

$$t(\mathbf{x}) = \frac{f(\mathbf{x}|s)}{f(\mathbf{x}|b)}$$



In this case we do have explicit formulas for the pdf, we constructed it with combine!

BUT we want to make sure the model we've constructed is sensible and the fits are running well!

Fit Parameter Values

RooFitResult: minimized FCN value: -2.55338e-05, estimated distance to minimum: 7.54243e-06
covariance matrix quality: Full, accurate covariance matrix
Status : MINIMIZE=0 HESSE=0

Floating Parameter	FinalValue +/-	Error
CMS_eff_b	-4.5380e-02 +/-	9.93e-01
CMS_eff_t	-2.6311e-01 +/-	7.33e-01
CMS_eff_t_highpt	-4.7146e-01 +/-	9.62e-01
CMS_scale_t_1prong0pi0_13TeV	-1.5989e-01 +/-	5.93e-01
CMS_scale_t_1prong1pi0_13TeV	-1.6426e-01 +/-	4.94e-01
CMS_scale_t_3prong0pi0_13TeV	-3.0698e-01 +/-	6.06e-01
acceptance_Ztautau	-3.1262e-01 +/-	8.62e-01
acceptance_bbH	-2.8676e-05 +/-	1.00e+00
acceptance_ttbar	4.9981e-03 +/-	1.00e+00
lumi_13TeV	-5.6366e-02 +/-	9.89e-01
norm_jetFakes	-9.3327e-02 +/-	2.56e-01
r	-2.7220e+00 +/-	2.59e+00
top_pt_ttbar_shape	1.7586e-01 +/-	7.00e-01
xsec_Ztautau	-1.6007e-01 +/-	9.66e-01
xsec_diboson	3.9758e-02 +/-	1.00e+00
xsec_ttbar	5.7794e-02 +/-	9.46e-01

name	b-only fit	s+b fit	rho
CMS_eff_b	-0.04, 0.99	-0.05, 0.99	+0.01
CMS_eff_t	* -0.24, 0.73*	* -0.26, 0.73*	+0.06
CMS_eff_t_highpt	* -0.56, 0.94*	* -0.47, 0.96*	+0.02
CMS_scale_t_1prong0pi0_13TeV	* -0.17, 0.58*	* -0.16, 0.59*	-0.04
CMS_scale_t_1prong1pi0_13TeV	! -0.12, 0.45!	! -0.16, 0.49!	+0.20
CMS_scale_t_3prong0pi0_13TeV	* -0.31, 0.61*	* -0.31, 0.61*	+0.02
acceptance_Ztautau	* -0.31, 0.86*	* -0.31, 0.86*	-0.05
acceptance_bbH	+0.00, 1.00	-0.00, 1.00	+0.05
acceptance_ttbar	+0.01, 1.00	+0.00, 1.00	+0.00
lumi_13TeV	-0.05, 0.99	-0.06, 0.99	+0.01
norm_jetFakes	! -0.09, 0.26!	! -0.09, 0.26!	-0.05
top_pt_ttbar_shape	* +0.24, 0.69*	* +0.18, 0.70*	+0.22
xsec_Ztautau	-0.16, 0.97	-0.16, 0.97	-0.02
xsec_diboson	+0.03, 1.00	+0.04, 1.00	-0.02
xsec_ttbar	+0.08, 0.95	+0.06, 0.95	+0.02



Part 3

Rate Parameters

Can add rate parameters which scale certain processes

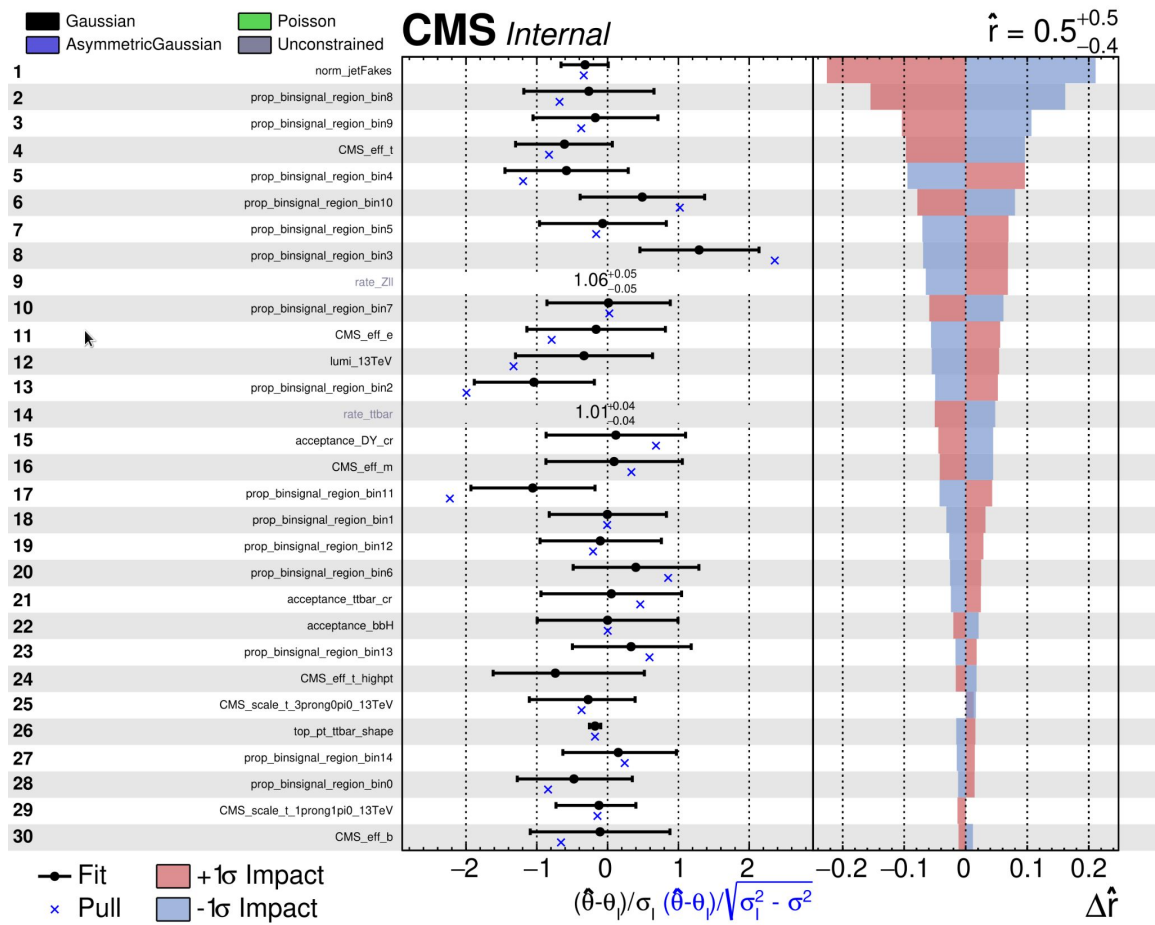
```
[name] rateParam [channel] [process] [init] [min,max]
```



To allow the rates of the ttbar and Z->ll process in the control regions to influence those in the signal region,
Connect them to each other via a rate parameter

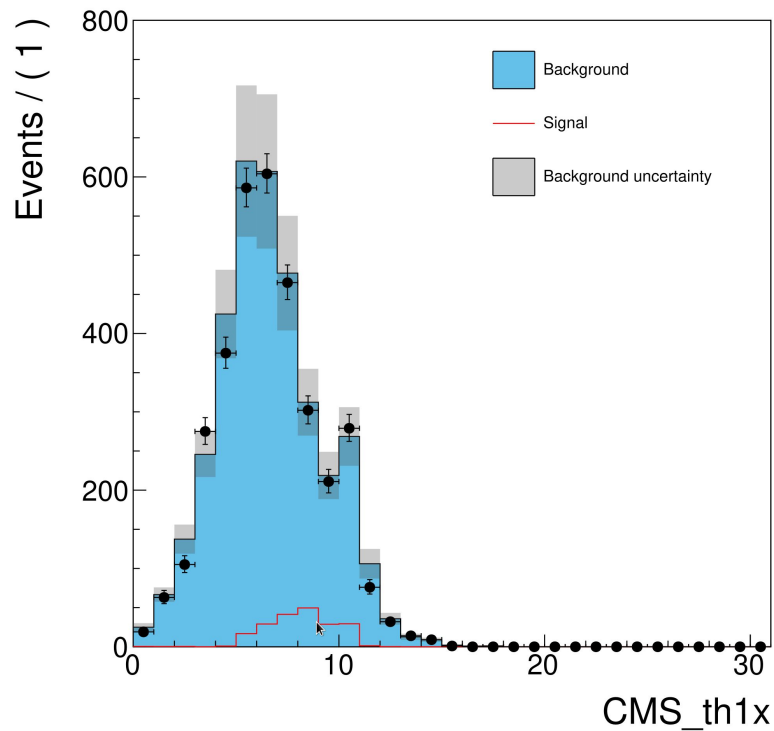
```
rate_ttbar rateParam * ttbar 1  
rate_Zll rateParam * Ztautau 1  
rate_Zll rateParam * Zmumu 1
```

Impacts

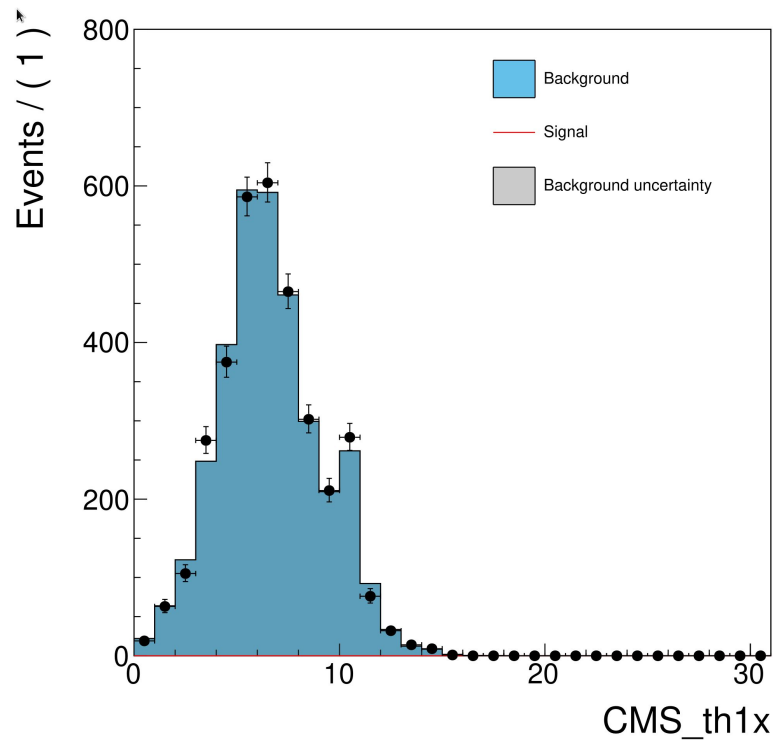


Visualizing fits

Prefit



background -only fit



Significance

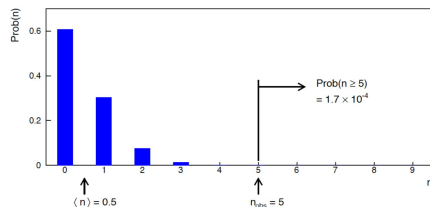
Poisson counting experiment: discovery p -value

Suppose $b = 0.5$ (known), and we observe $n_{\text{obs}} = 5$.

Should we claim evidence for a new discovery?

Give p -value for hypothesis $s = 0$, suppose relevant alt. is $s > 0$.

$$\begin{aligned} p\text{-value} &= P(n \geq 5; b = 0.5, s = 0) \\ &= 1.7 \times 10^{-4} \neq P(s = 0)! \end{aligned}$$



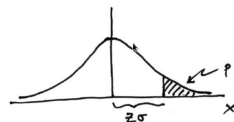
G. Cowan / RHUL Physics

Terascale Statistics 2024 / Lecture 1

19

Significance from p -value

Often define significance Z as the number of standard deviations that a Gaussian variable would fluctuate in one direction to give the same p -value.



$$p = \int_Z^{\infty} \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx = 1 - \Phi(Z)$$

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

in ROOT:

```
p = 1 - TMath::Freq(Z)
Z = TMath::NormQuantile(1-p)
```

in python (scipy.stats):

```
p = 1 - norm.cdf(Z) = norm.sf(Z)
Z = norm.ppf(1-p)
```

Result Z is a “number of sigmas”. Note this does not mean that the original data was Gaussian distributed.

G. Cowan / RHUL Physics

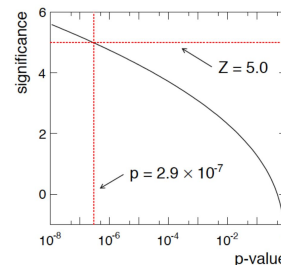
Terascale Statistics 2024 / Lecture 1

20

Poisson counting experiment: discovery significance

Equivalent significance for $p = 1.7 \times 10^{-4}$: $Z = \Phi^{-1}(1 - p) = 3.6$

Often claim discovery if $Z > 5$ ($p < 2.9 \times 10^{-7}$, i.e., a “5-sigma effect”)



In fact this tradition should be revisited: p -value intended to quantify probability of a signal-like fluctuation assuming background only; not intended to cover, e.g., hidden systematics, plausibility signal model, compatibility of data with signal, “look-elsewhere effect” (~multiple testing), etc.

G. Cowan / RHUL Physics

Terascale Statistics 2024 / Lecture 1

21

We calculate the p -value of the modified Likelihood Ratio test-statistic q_0 and quote a significance

$$q_0 = \begin{cases} 0 & \hat{\mu} < 0 \\ -2 \log \left(\frac{\mathcal{L}(\mu_{\text{NP}}=0)}{\mathcal{L}(\hat{\mu}_{\text{NP}})} \right) & \hat{\mu} \geq 0 \end{cases}$$

Significance

Simple Asymptotic Calculation – assume known distribution of test-statistic

```
<<< Combine >>>
<<< v9.1.0 >>>
>>> Random number generator seed is 123456
>>> Method used is Significance

-- Significance --
Significance: 1.11273
```

Can also instruct combine to give the p-value directly with the `--pvalue` flag

```
-- Significance --
p-value of background: 0.132912
Done in 0.00 min (cpu), 0.00 min (real)
```

Significance

Can also check expected significance for various signal strengths, e.g. with `-t -1 --expectSignal 1.5`

```
-- Significance --  
Significance: 3.52007
```

Can also use fit model after a first fit to the data to get model parameters with `--toysFrequentist`

```
-- Significance --  
Significance: 3.13954
```


Signal Strength measurement

Maximum Likelihood Estimators (MLEs)

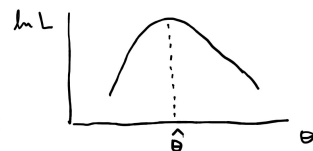
We *define* the maximum likelihood estimators or MLEs to be the parameter values for which the likelihood is maximum.

Maximizing L
equivalent to
maximizing $\log L$

$$\hat{\theta} = \operatorname{argmax}_{\theta} L(\theta)$$

Could have multiple maxima (take highest).

MLEs not guaranteed to have any 'optimal' properties, (but in practice they're very good).



G. Cowan / RHUL Physics

Terascale Statistics 2024 / Lecture 1

25

Example of variance by graphical method

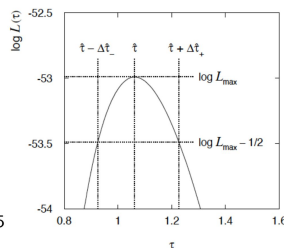
ML example with exponential:

$$\hat{\tau} = 1.062$$

$$\Delta \hat{\tau}_{-} = 0.137$$

$$\Delta \hat{\tau}_{+} = 0.165$$

$$\hat{\sigma}_{\hat{\tau}} \approx \Delta \hat{\tau}_{-} \approx \Delta \hat{\tau}_{+} \approx 0.15$$

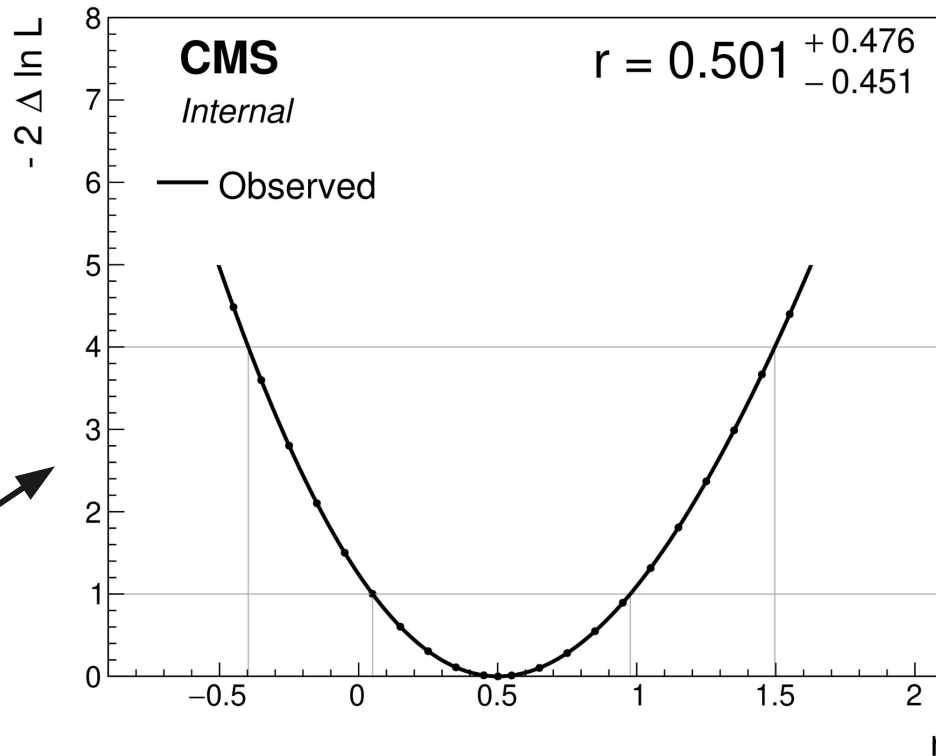


Not quite parabolic $\ln L$ since finite sample size ($n = 50$).

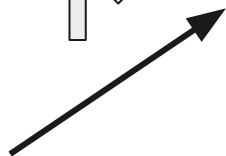
G. Cowan / RHUL Physics

Terascale Statistics 2024 / Lecture 1

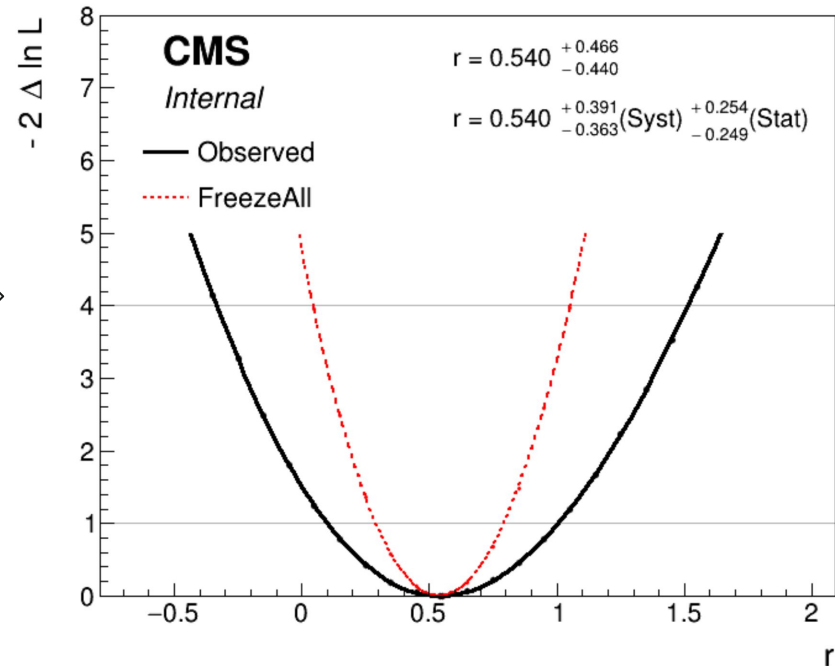
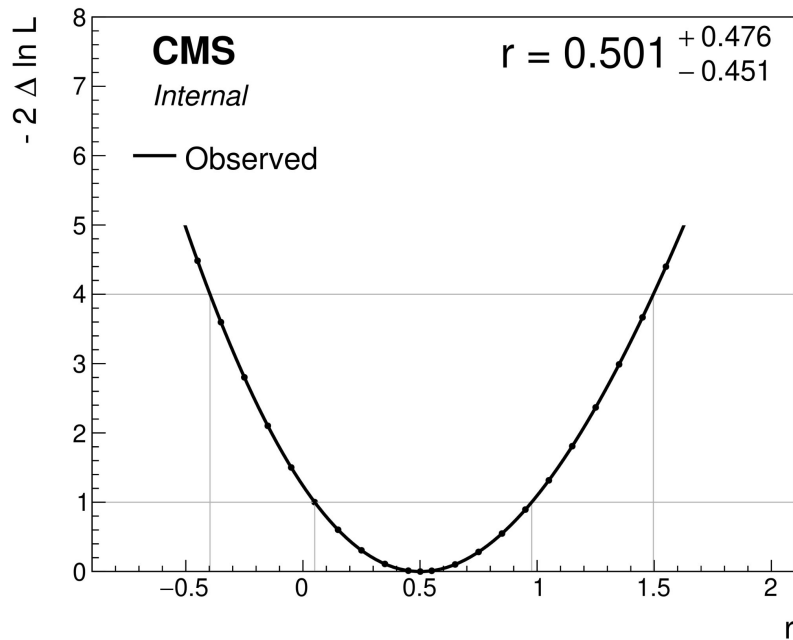
33



With a
sign flip

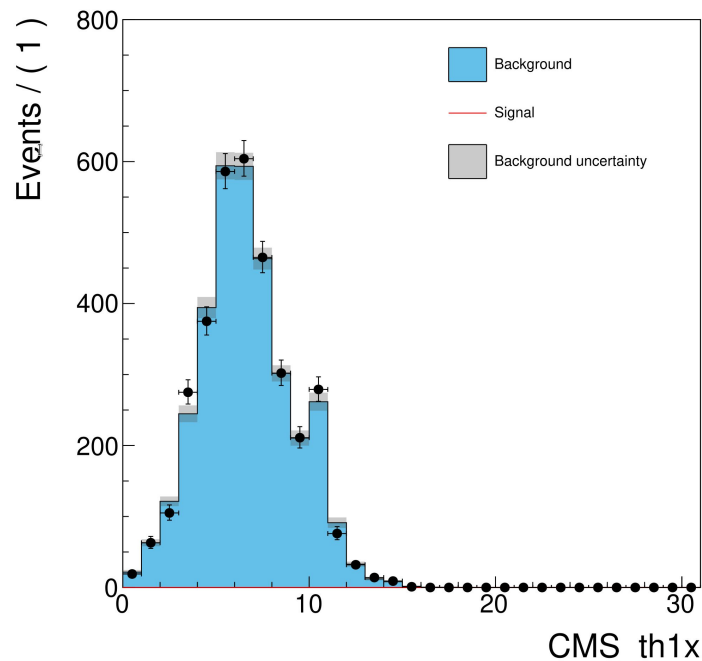
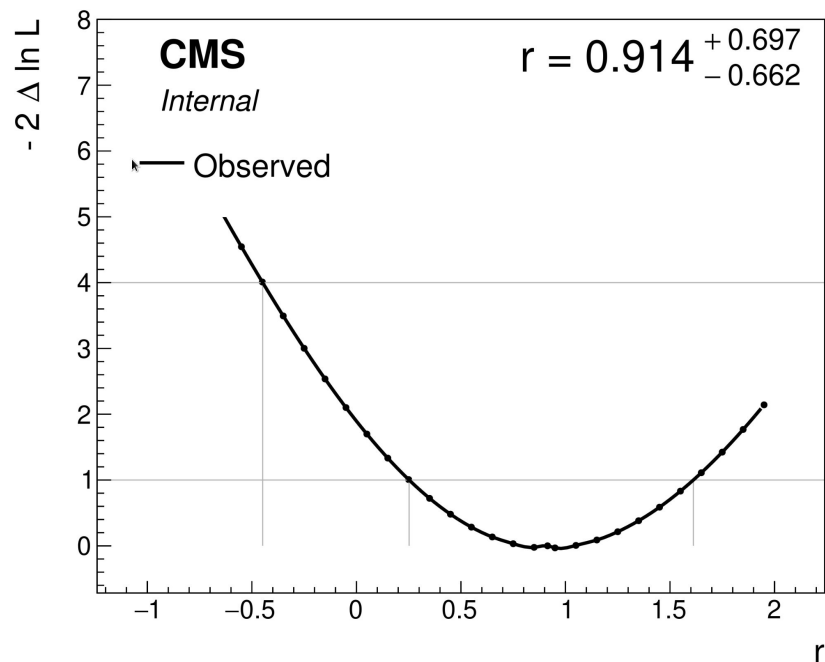


Uncertainty Breakdowns



Channel Masks

Can also mask particular channels to investigate the fit, e.g. masking the control regions:



Part 4

Two Parameters of Interest

```
from HiggsAnalysis.CombinedLimit.PhysicsModel import PhysicsModel
```

```
class DASModel(PhysicsModel):
```

```
    def doParametersOfInterest(self):
```

```
        """Create POI and other parameters, and define the POI set."""
```

```
        self.modelBuilder.doVar("r[0,0,10]")
```

```
        self.modelBuilder.doSet("POI", "", ".join(["r"])
```

```
    def getYieldScale(self, bin, process):
```

```
        "Return the name of a RooAbsReal to scale this yield by or the two special values 1 and 0
```

```
        if self.DC.isSignal[process]:
```

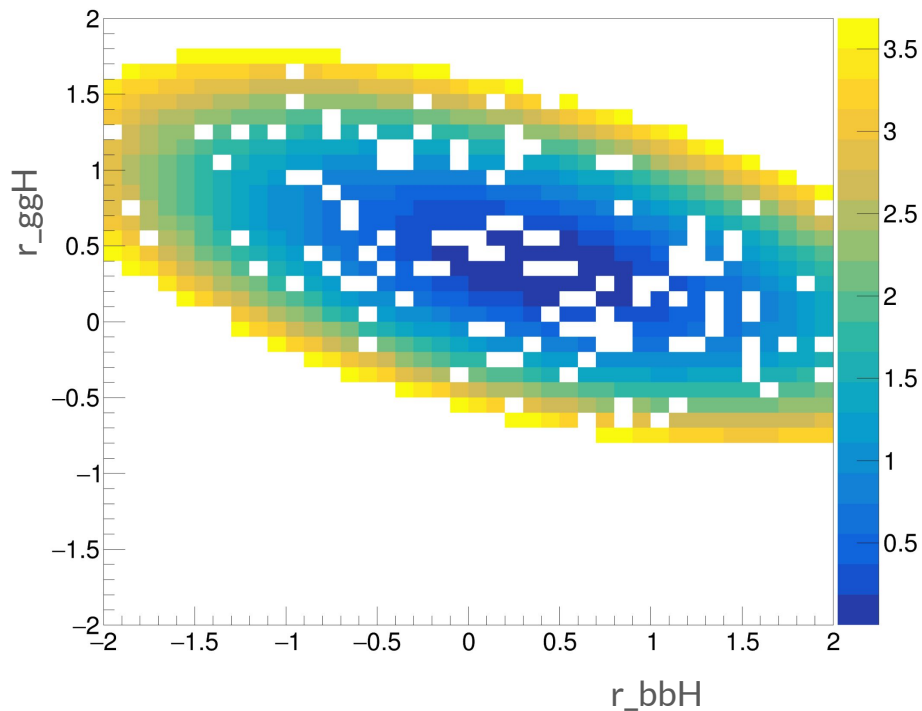
```
            print("Scaling %s/%s by r" % (bin, process))
```

```
            return "r"
```

```
        return 1
```

```
dasModel = DASModel()
```

2-Dimensional NLL map



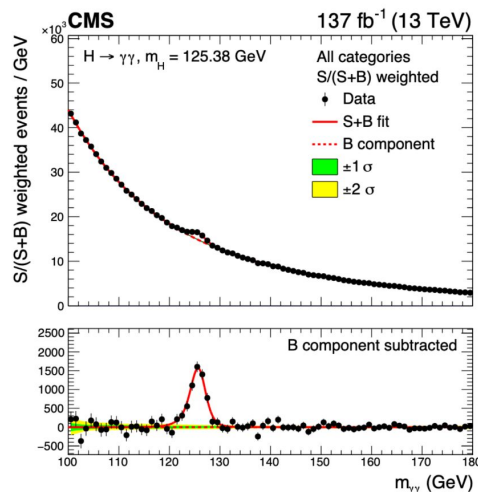
Prep for tomorrow

Tomorrow – Parametric Models

If you're not very familiar with ROOT and RooFit You might want to check this pre-tutorial:

<https://cms-analysis.github.io/HiggsAnalysis-CombinedLimit/latest/part5/roofit/>

Analysis types: input



Parametric model:

- Suitable for the analysis with analytically described bkg and signal: e.g. Gaussian signal on smoothly falling polynomial background
- In most cases used in the analyses with data-driven bkg description
- The systematic uncertainties assigned on the parameters of the model

