



Statistical analysis tutorial

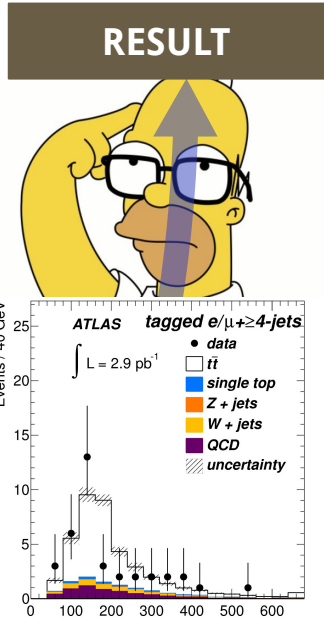
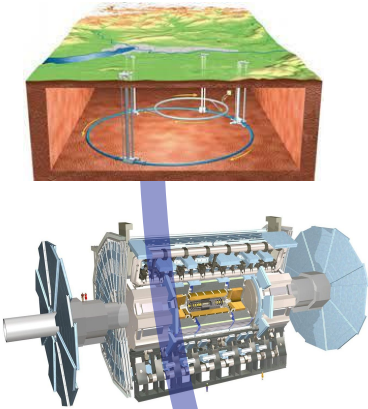
Michele Pinamonti (*INFN Trieste, ATLAS Udine Group*)

michele.pinamonti@ts.infn.it

CODATA-RDA Workshop - ICTP, Trieste 28/07/2022

Statistical analysis

- Why? What? ...
 - **why** do we need statistical analysis?
 - **what do we mean** by statistical analysis?



- **Statistical Analysis in particle-collider physics:**
 - the way to extract **quantitative information** from **collision data**
- ... and of course, what goes into the result section of your paper is the **quantitative information:**
 - we want to claim things like:
 - “ $X = Y \pm Z$ ”
 - “ $X > Y$ excluded at 95% confidence level”
 - “ X observed with a significance of 5σ ”

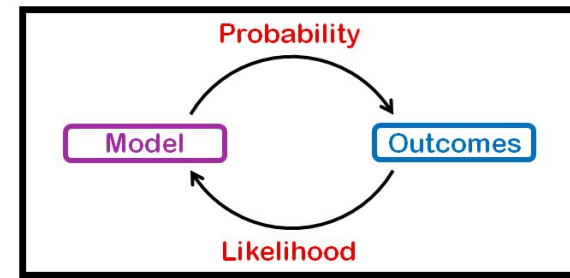


Statistical Analysis Basics (for HEP)



Maximum likelihood and Fits

- **Likelihood:**
 - defined as **probability** of observing a certain set of **data** given a model / hypothesis (with certain **parameter** values)

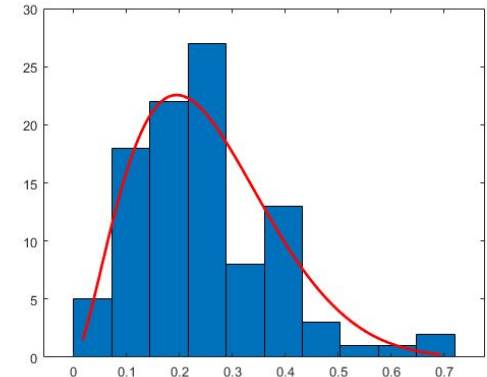


$$L(\vec{\theta}) = \text{Prob}(\vec{x}|\vec{\theta}) = \prod_i \text{Prob}(x_i|\vec{\theta})$$

if data points / measurements / observation are independent (i.e. uncorrelated)

Annotations: A red arrow points from the word 'probability' to the 'Prob' part of the equation. A purple arrow points from the word 'data' to the \vec{x} part. A blue arrow points from the word 'parameters' to the $\vec{\theta}$ part.

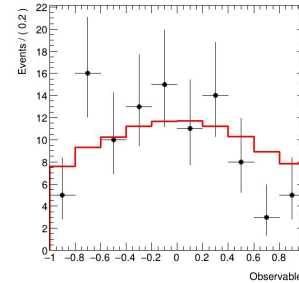
- **Maximum Likelihood principle:**
 - **estimated** value(s) of parameter(s) = value(s) **maximizing** the Likelihood
- **"Fit":**
 - **parameter estimation** procedure via Likelihood maximization



Types of likelihood

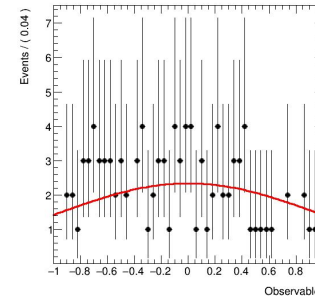
- **Binned Poisson Likelihood:**

$$\mathcal{L}(\vec{n}|\vec{\theta}) = \prod_{i \in bins} P(n_i|Y_i(\vec{\theta})) = \prod_{i \in bins} P(n_i|S_i(\vec{\theta}) + B(\vec{\theta}))$$



- **Unbinned Likelihood:**

$$\mathcal{L}(\vec{m}|\vec{\theta}) = P(n_{obs}|S + B) \times \prod_{i=1}^{n_{obs}} \frac{S \cdot \mathcal{P}_S(m_i, \vec{\theta}) + B \cdot \mathcal{P}_B(m_i, \vec{\theta})}{S + B}$$

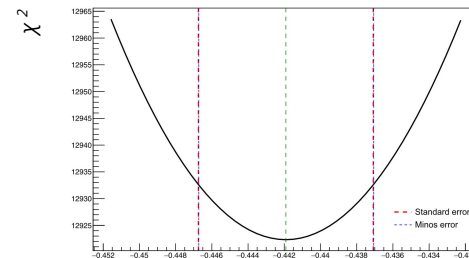


- **The χ^2 case:**

- binned Poisson L , for large $n \Rightarrow \mathcal{L}(\vec{n}|\vec{\theta}) \simeq \prod_{i \in bins} \mathcal{G}(n_i|\mu_i = Y_i(\vec{\theta}), \sigma_i = \sqrt{n_i}) = \prod_{i \in bins} \frac{1}{\sqrt{2\pi}\sigma_i} e^{-\frac{(n_i - \mu_i)^2}{2\sigma_i^2}}$

$$\Rightarrow -2 \log \mathcal{L}(\vec{n}|\vec{\theta}) = \sum_{i \in bins} \frac{(n_i - \mu_i)^2}{\sigma_i^2} + Const. = \chi^2$$

- Maximizing Likelihood = Minimizing χ^2 or $-2 \log L$



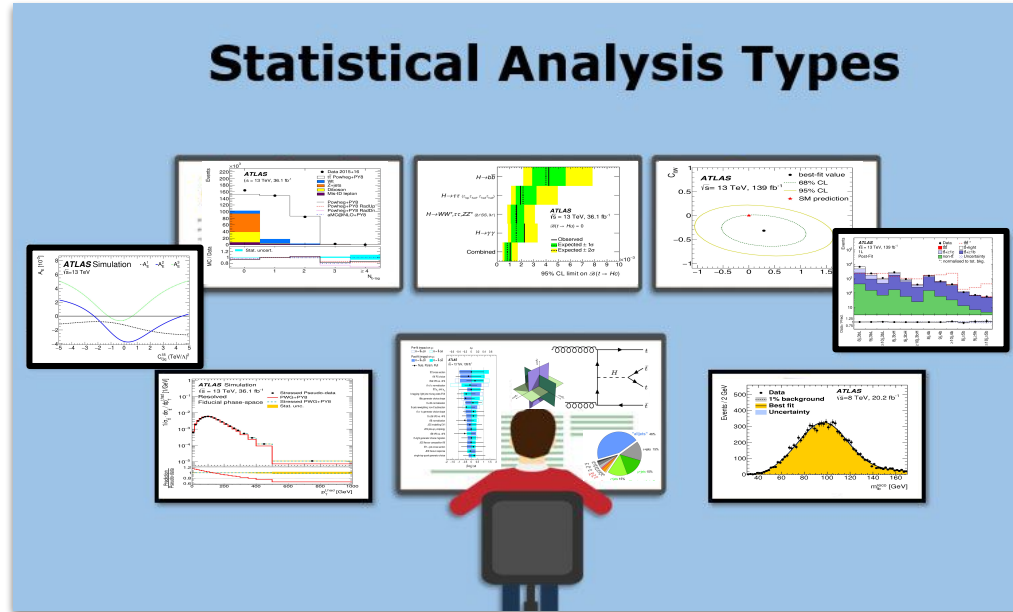
Types of measurements

- In HEP data analysis, **different types of measurement**:

- **searches** for a new process
- **cross-section** measurements:
 - *total cross-section*
(full / fiducial phase-space...)
 - *differential cross-section*
 - *ratios of cross-sections...*
- other **parameter estimation**:
 - *usually "shape analyses"*
(e.g. *top mass, top width...*)
- **EFT** fits / limits ...

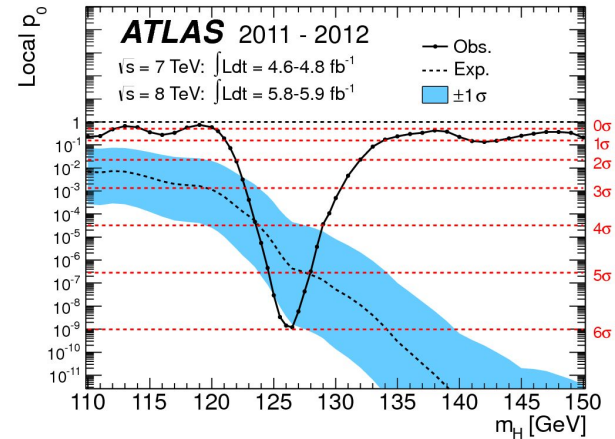
- Also, measurements can take as **inputs**:

- **binned** data \Rightarrow histogram counts are the inputs
- **unbinned** data \Rightarrow individual events as "input measurements"
- other existing **measurements** / differential cross-section bins \Rightarrow "2 step" analysis

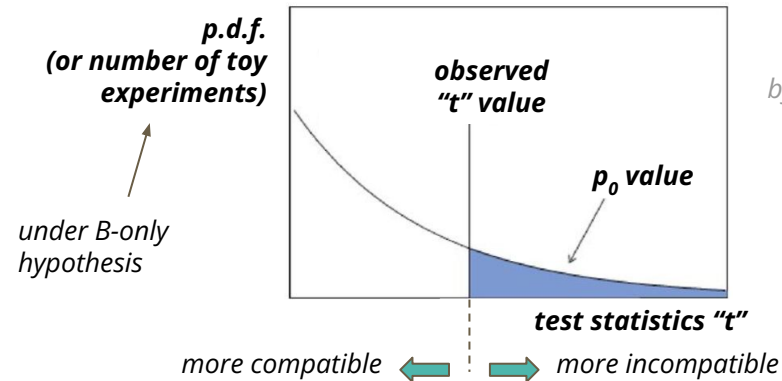


Searches: discovery significance

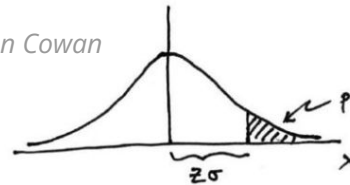
- **Observing** a new process (*in the Frequentist language*)
= seeing data **incompatible** with **background-only** hypothesis
(or “null hypothesis”)
- **How to quantify it?**
 - define “**test statistics**”, quantifying agreement btw. data and a prediction (e.g. likelihood, χ^2 , LH-ratio...)
 - define **p_0 -value** = probability of seeing **worse agreement** than observed one, in the background-only hypothesis
 - *i.e. “probability that what we see is a fake signal”*



- turn p_0 into number of Gaussian std.dev, define **significance “Z”** in terms of *number of sigmas*:
 - $5\sigma = \sim 3 \cdot 10^{-7}$

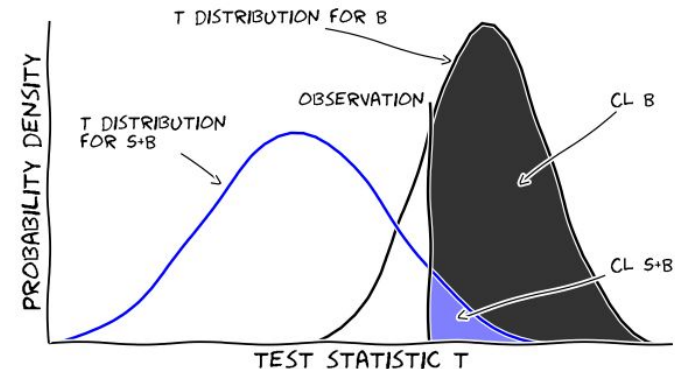
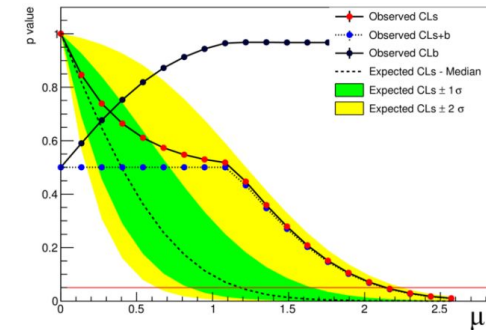
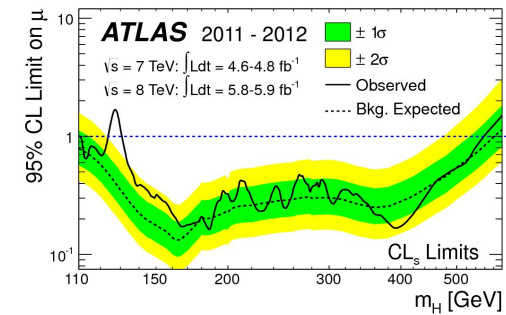


by Glen Cowan

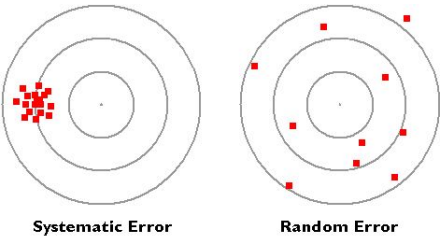


Searches: exclusion limits

- If no evidence for signal, setting **exclusion limits**
- Usually limits set on **signal strength** $\mu = \sigma^{\text{obs}} / \sigma^{\text{theory}}$:
 - values of $\mu >$ quoted value **excluded at 95% confidence-level**
- **Operationally:**
 - define **test-statistics** (as before), t (data, μ)
 - **scan** values of μ , get t^{obs} for each μ
 - assign prob. of seeing worse t than t^{obs} , assuming that value of μ
 - **find μ** for which **prob. = 5%** (i.e. 1 - 95%, corresponding to 2σ)
- What does **CL_s** mean?
 - description above defines "**CL_{s+b}**"
 - can then define "**CL_b**" as follows:
 - get t^{obs} for each μ (as before)
 - define CL_b as prob. of seeing worse t , in the **B-only hypothesis** ($\mu=0$)
 - then define **CL_s = CL_{s+b} / CL_b**



Systematic Uncertainties



Inclusion of systematic uncertainties

- In particle collision physics we distinguish:
 - **statistical uncertainty:**
 - result of **stochastic fluctuations** in data
 - consequence of limited size of analysed dataset
 - **systematic uncertainties:**
 - everything that is not a statistical uncertainty
 - uncertainties associated with measurement apparatus, assumptions made, or model used
 - *Statistical* uncertainty usually **intrinsically included** in inference method (e.g. in χ^2 fit)
 - *Systematic* uncertainties: **non-obvious inclusion** in and **propagation** through statistical analysis
 - **Side considerations:**
 - in our world, systematic uncertainties are uncertainties on $\text{Prob}(x, \theta)$, i.e. uncertainties **on expected values** (e.g. exp. S+B), **not on data** (!)
 - systematics divided into multiple independent / uncorrelated “**sources**”
- ➡ Fully **uncorrelated** between subsequent measurements
- ➡ Fully **correlated** between subsequent measurements

The Profile Likelihood formalism

- More and more common approach for including systematics in HEP statistical analysis:
 - include systematic uncertainties as unknown parameters in the model
 - nuisance parameters modifying expectations in a parametric way
 - prior probabilities on values of nuisance parameters to reflect limited knowledge

- **The binned profile-likelihood:**

$$L(\vec{n} \mid \vec{\theta}, \vec{k}) = \prod_i P(n_i \mid S_i(\vec{\theta}, \vec{k}) + B_i(\vec{\theta}, \vec{k})) \times \prod_j G(\theta_j)$$

data → \vec{n} Poisson ↓ $P(n_i \mid \dots)$ Gaussian (or other pdf...) ↓ $G(\theta_j)$

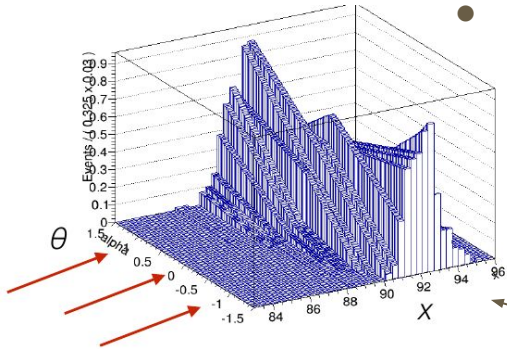
→ → → $\vec{\theta}$ → → \vec{k}

constrained parameters: nuisance parameters (**NPs**) associated to systematic uncertainties

data events in bin i prediction in bin i (signal+background) constraint term for nuisance parameter j

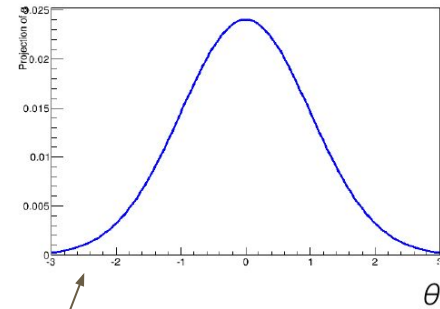
unconstrained parameters: parameter of interest (**POI** or “ μ ”) + unconstrained nuisance parameters (e.g. background normalization parameters)

Nuisance parameters and systematic uncertainties



- Each (*independent*) source of **systematic uncertainty** included in the likelihood as constrained **nuisance parameters** (NPs):
 - affecting $S+B$ prediction in a **coherent way**
 - effect **interpolated** and **extrapolated** from **3 discrete values** (0 = nominal, 1 = “up” var., -1 = “down” var.) to range of **continuous values**

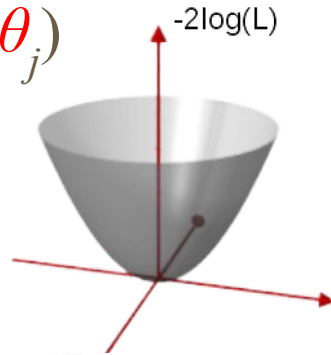
normal distribution



$$L(\vec{n} | \vec{\theta}, \vec{k}) = \prod_i P(n_i | S_i(\vec{\theta}, \vec{k}) + B_i(\vec{\theta}, \vec{k})) \times \prod_j G(\theta_j)$$

- The fit procedure becomes a **multi-dimensional Likelihood maximisation** problem

- the fit **result** is not just the value (and uncertainty) on parameter(s) of interest (POI), but a **set of values** for all the parameters, including nuisance parameters: $(\hat{\mu}, \hat{\theta}_0, \dots, \hat{\theta}_{N-1}) : \mathcal{L}(\hat{\mu}, \hat{\theta}) = \max$



Profile likelihood ratio and asymptotic regime

- **Neyman-Pearson lemma:**

- the **likelihood ratio** $\lambda = L(H_1)/L(H_0)$ is the **optimal discriminator** when testing hypothesis H_1 vs. H_0 (e.g. $H_1 =$ presence of signal ($\mu > 0$), H_0 no signal ($\mu = 0$))
- in the case of our profile likelihood, can build profile likelihood ratio, as a function of POI:

Profile likelihood ratio
only dependent on μ

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

Maximize L for a given μ
'conditional' likelihood

Maximize L
'unconditional' likelihood

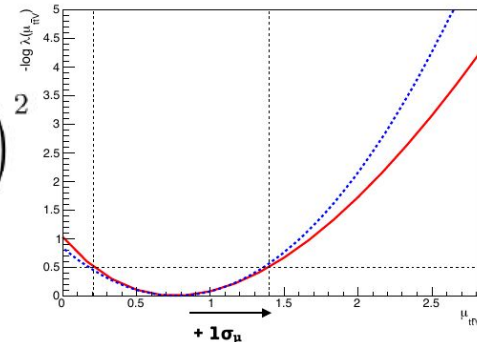
- maximizing λ vs. $\mu =$ maximizing L vs. (μ, θ)

- **Wilks' Theorem:** in large statistics data samples, λ distribution follows **χ^2 distribution:**

$$-2 \log \lambda(\mu) = -2(\log L(\mu, \hat{\theta}) - \log L(\hat{\mu}, \hat{\theta})) = \left(\frac{\mu - \hat{\mu}}{\sigma_\mu} \right)^2$$

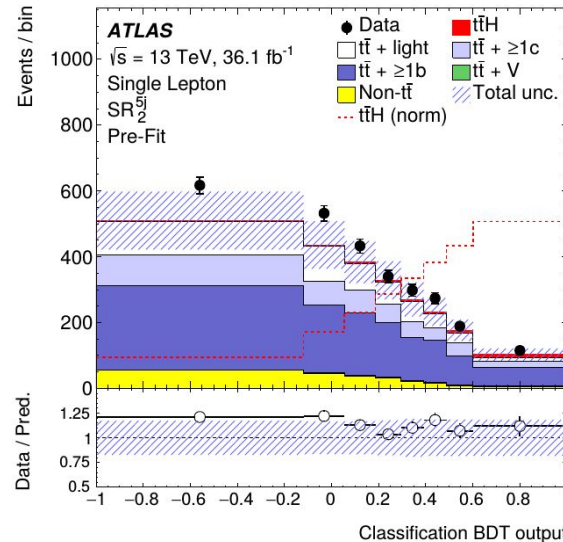
⇒ can get the **uncertainty on μ (including effect of all systematics!!)**

- large-statistics means $> \sim O(10)$ events
- saves from running very time consuming pseudo-experiments

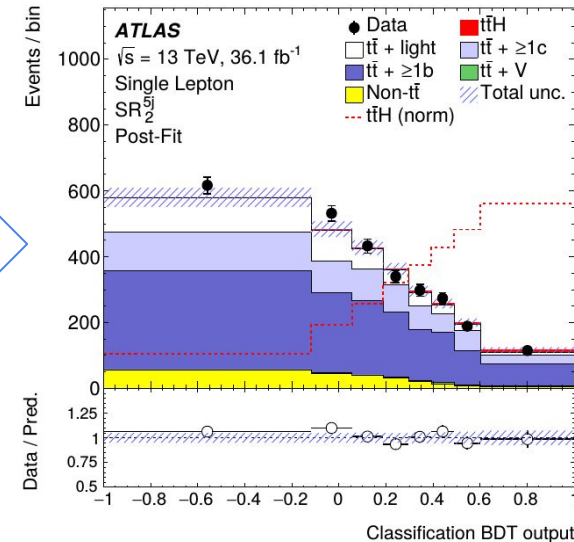


Profiling, pre-fit and post-fit

- Profile likelihood fit can:
 - **change background prediction**, if best-fit θ values different from θ_0
 - **reduce uncertainty** on background, through:
 - **constraint** of NPs
("improved knowledge" of parameters that are affected by systematic uncertainties, i.e. data have enough statistical power to further constraint the NP)
 - **correlations** between NPs

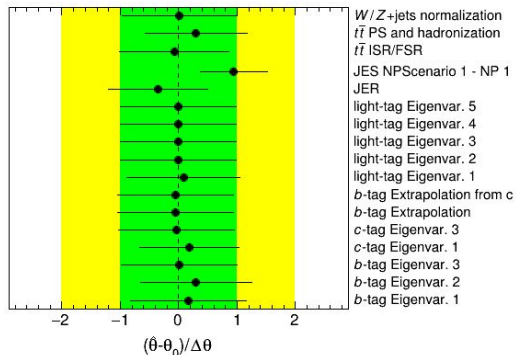


FIT

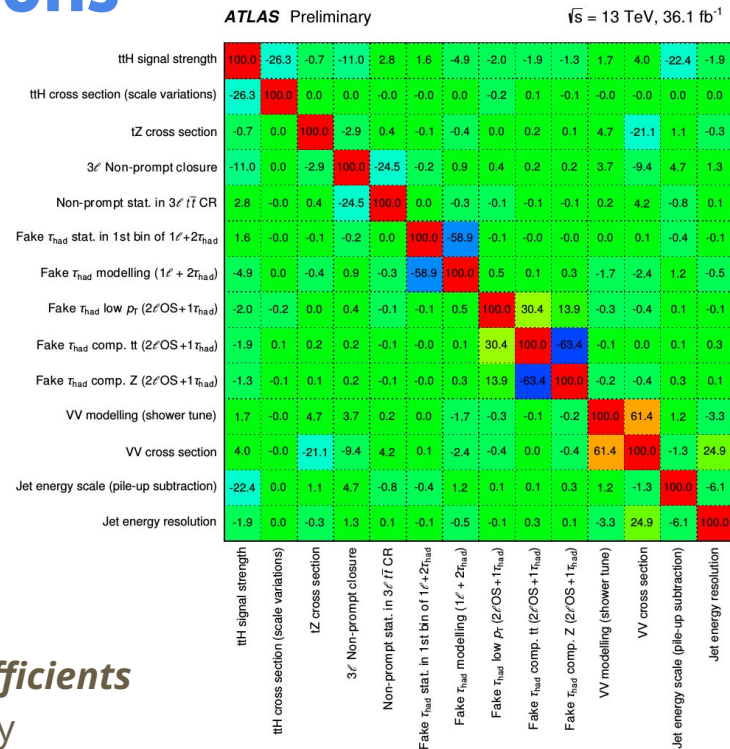


NP pulls, constraints and correlations

- Useful to **monitor** NP **pulls** and **constraints**:
 - they are "*nuisance*", but they can be important!

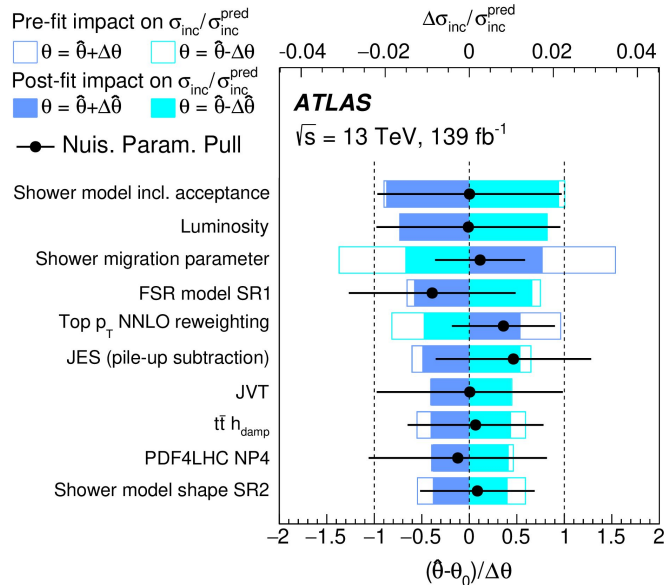


- Important to consider also NP **correlations**:
 - uncertainties on NPs (*and POI*) extracted from **covariance matrix**, which includes **correlation coefficients**
 - correlation **built by the fit**, even if completely independent / uncorrelated sources of uncertainty before the fit (*correlation in the improved knowledge of the parameters*)
 - (anti-)correlations can **reduce** total post-fit uncertainty!



Impact of systematics

1. "Ranking plot" shows *pre-fit* and *post-fit* **impact** of **individual NP** on the determination of μ :
 - **each NP fixed** to ± 1 pre-fit and post-fit error
 - fit re-done with $N-1$ parameters
 - impact = difference in **central value** of μ
2. "Grouped impact table" reports *contributions* to *total uncertainty* from **groups** of syst.:



"which systematics are more important?"

- fix a group of NPs to post-fit values
- repeat the fit, get reduced error on μ
- impact = difference in quadrature btw. original and reduced error on μ
- get stat. uncertainty by fixing all NPs

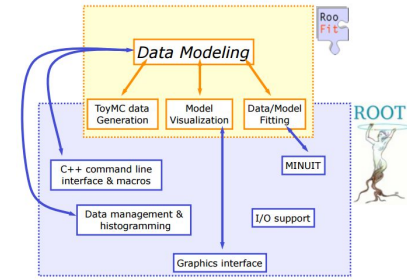
Category	$\frac{\Delta\sigma_{\text{fid}}}{\sigma_{\text{fid}}} [\%]$	$\frac{\Delta\sigma_{\text{inc}}}{\sigma_{\text{inc}}} [\%]$
Signal modelling		
$t\bar{t}$ shower/hadronisation	± 2.8	± 2.9
$t\bar{t}$ scale variations	± 1.4	± 2.0
...		
Total systematic uncertainty	± 4.3	± 4.6
Data statistical uncertainty	± 0.05	± 0.05
Total uncertainty	± 4.3	± 4.6

Tools for statistical analysis (with Profile Likelihood)



Profile likelihood - Implementation in ROOT

- **Roofit:** toolkit to extend **ROOT** providing language to describe data models
 - model distribution of observable x in terms of parameters θ using probability density function PDF
- **Roostats:** project to provide advanced stat. techniques for LHC collaborations
 - built on top of **Roofit**
- **Rooworkspace:** generic container class for all **Roofit** objects, containing:
 - full model configuration
(i.e. all information to run statistical calculations)
 - PDF and parameter/observables descriptions uncertainty/shape of nuisance parameters
 - (multiple) data sets
- **HistFactory:** tool for creating **Roofit** workspaces formatted for use with **Roostats** tools
 - meant for analyses based on template histograms



Mathematical concept	Roofit class
variable x	<code>RoorealVar</code>
function $f(x)$	<code>RoorealAbs</code>
PDF $f(x)$	<code>RoorealPdf</code>
space point \vec{x}	<code>RoorealArgSet</code>
integral $\int_{x_{\min}}^{x_{\max}} f(x) dx$	<code>RoorealIntegral</code>
list of space points	<code>RoorealAbsData</code>

Practical part



Repository and environment



- GitHub repository: <https://github.com/pinamont/statistics-tutorial>

- The whole tutorial will be run through Jupyter notebooks (python and ROOT/C++ based)



- 2 available options:

- **Binder** 

- **SWAN+cern-box**  

- **Goal:** guide you through **what's actually done** to publish your results
 - with some **exercises** to get acquainted with the machinery
 - we'll choose **dynamically** what to cover (*raise your hands!*)
 - you may use the rests as a **reference** (*& feel free to contact us!*)

Setting up the environment

- Go to the [GitHub repository](#)



- Choose one of the 2 options:

- Binder:**



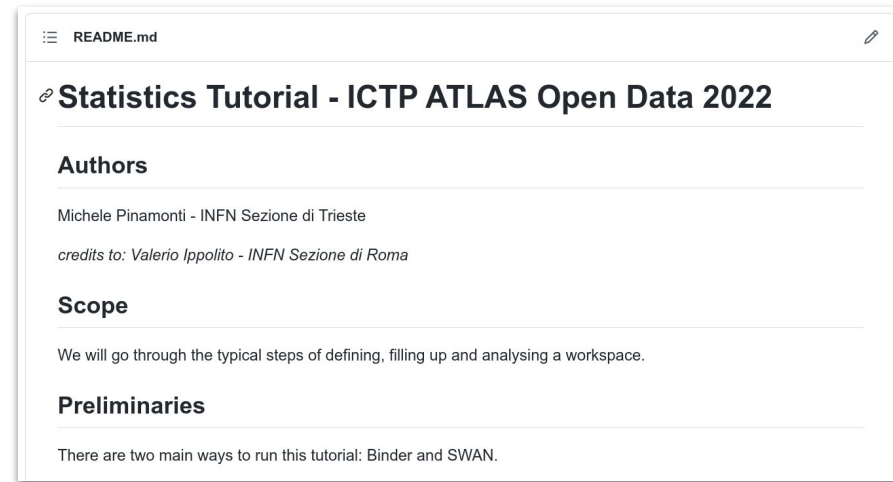
- no CERN account needed
 - could take more time to load...

- SWAN:**



- CERN account needed (and cern-box / eos space set up)
 - should be faster to start

- Follow instructions on the **README file** for setting up environment, according to chosen option
- Once ready, try running the `hello_world.ipynb` notebook



- Caveat:**

- exercise doesn't seem to work with ROOT version 26.04 (set by default in SWAN)
 - setting-up ROOT version 24.06 in Binder
 - following instructions on README for SWAN should work as well (setting-up 24.06)

Binder and SWAN interfaces

The Binder interface is divided into two main sections. On the left is a file browser with a search bar and a list of files and folders. On the right is the 'Launcher' section with four categories of tools: Notebook, Console, Other, and Terminal/Text File.

Name	Last Modified
create_data	27 minutes ago
data	27 minutes ago
fit	27 minutes ago
limit	27 minutes ago
p_values	27 minutes ago
systematics	27 minutes ago
environment.yml	27 minutes ago
hello_world.ipynb	27 minutes ago
README.md	27 minutes ago

Launcher

- Notebook**
 - Python 3 (ipykernel)
 - ROOT C++
- Console**
 - Python 3 (ipykernel)
 - ROOT C++
- Other**
- Terminal**
- Text File**

The SWAN interface shows a project page for 'statistics-tutorial1'. It features a dark header with navigation options and a table listing the project's files and folders.

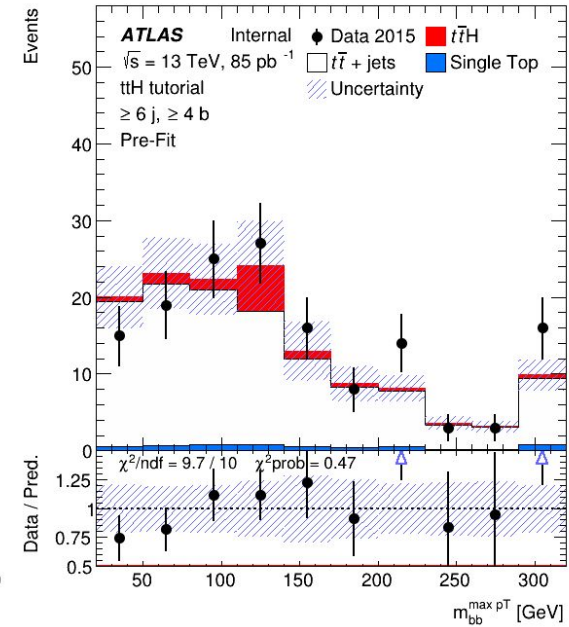
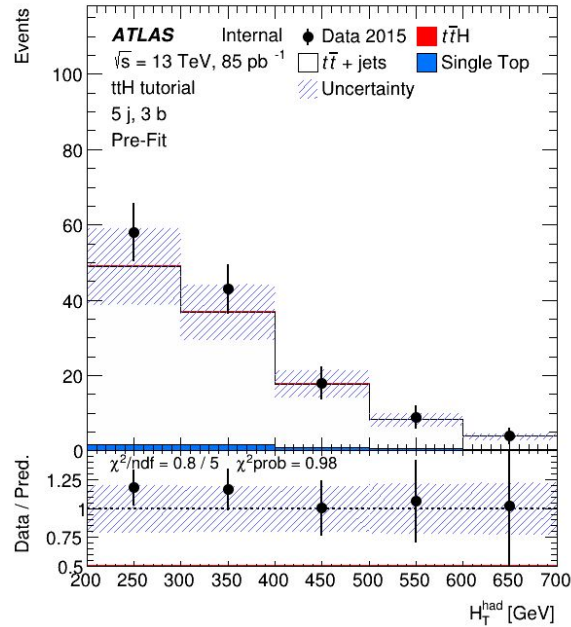
NAME	SIZE	STATUS	MODIFIED
create_data			2 giorni fa
data			2 giorni fa
fit			2 giorni fa
limit			2 giorni fa
p_values			2 giorni fa
systematics			2 giorni fa
hello_world.ipynb	2.16 kB		2 giorni fa
environment.yml	128 B		2 giorni fa
README.md	4.1 kB		2 giorni fa

Tutorial

- Tutorial structured as a **set of notebooks**, each performing a single action:
 - `create_data/create_workspace.ipynb` → create a RooWorkspace from existing histograms
→ will use output of this notebook for all other operations
 - simplified version `create_workspace_minimal.ipynb` also available
 - `create_data/inspect_workspace.ipynb` → inspect what's inside the workspace we just created
 - `fit/simple_fit.ipynb` → perform a fit and print fit results
 - `fit/postfit_plots.ipynb` → visualize projection of fit results to expected distributions
 - `systematics/ranking.ipynb` → breakdown of impact of systematics - method 1
 - `systematics/impact_table.ipynb` → breakdown of impact of systematics - method 2
 - `limit/toys.ipynb` → perform exclusion limit extraction
 - `p_values/pvalues.ipynb` → p-value and significance calculation

Our example workspace

- We'll use as an exercise a set of inputs (histograms):
 - ATLAS ttH search ($H \rightarrow bb$), part of real fitting exercise with very first 2015 data
 - $t\bar{t}+(b)$ -jets selection (1-lepton channel)
- Two statistically independent datasets ("regions" or "channels", as you wish):
 - "5 j, 3 b"
 - Control Region, enriched in $t\bar{t} + (b)$ jets
 - " ≥ 6 j, ≥ 4 b"
 - Signal Region



Backup

p_0 -value and discovery significance

- **Observing** a new process
= seeing data **incompatible** with **background-only** hypothesis (“null hypothesis”)
- **How to quantify it?**
 - define “**test statistics**”, quantifying data-prediction agreement
 - define **p_0 -value** = probability of seeing **worse agreement** (in B-only hypothesis)
 - turn p_0 into number of Gaussian std.dev, define **significance “Z”** in terms of *number of sigmas*

One-sided tests-statistics

- in the case of profile-likelihood ratio:

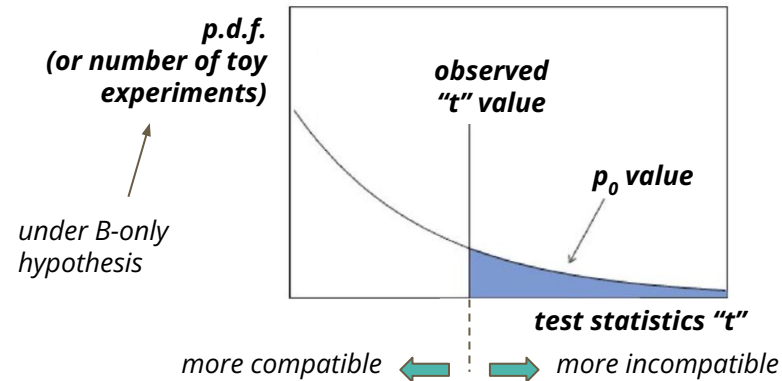
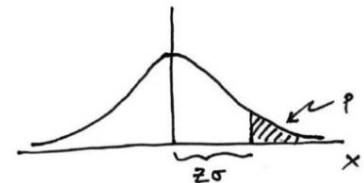
$$q_0 = \begin{cases} -2\ln\lambda(0) & \hat{\mu} \geq 0 \\ 0 & \hat{\mu} < 0 \end{cases}$$

$$p_0 = \int_{q_{0,\text{obs}}}^{\infty} f(q_0|0) dq_0$$

$$Z_0 = \Phi^{-1}(1 - p_0)$$

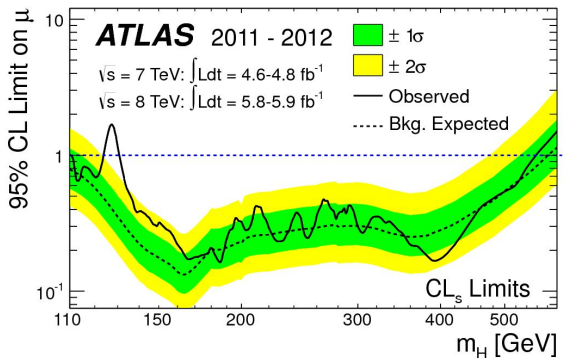
reminder:

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

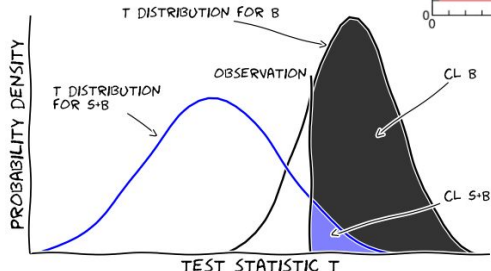
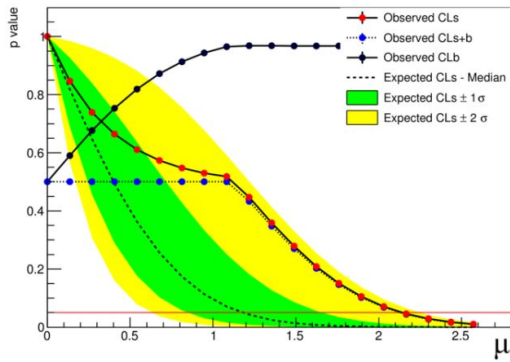


Exclusion limits

- No evidence \Rightarrow **exclusion limits**
 - usually on **signal strength** $\mu = \sigma^{\text{obs}} / \sigma^{\text{theory}}$
- Define **test-statistics** (as before), **t (data, μ)**
 - **scan** values of μ , get t^{obs} for each μ
 - assign prob. of seeing worse t than t^{obs} , assuming that value of μ
 - **find μ** for which **prob. = 5%** (i.e. 1 - 95%, corresponding to 2σ)



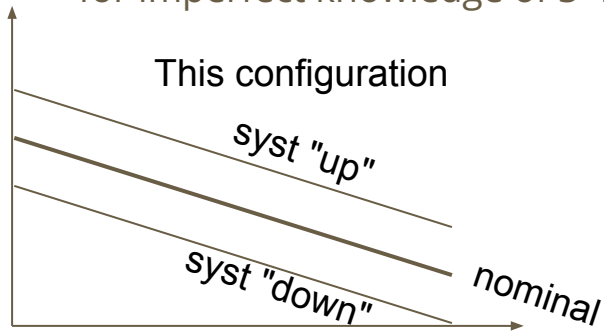
- What does **CL_s** mean?
 - description above defines "**CL_{s+b}**"
 - can then define "**CL_b**" as follows:
 - get t^{obs} for each μ (as before)
 - define CL_b as prob. of seeing worse t , in the **B-only hypothesis** ($\mu=0$)
 - then define **CL_s = CL_{s+b} / CL_b**



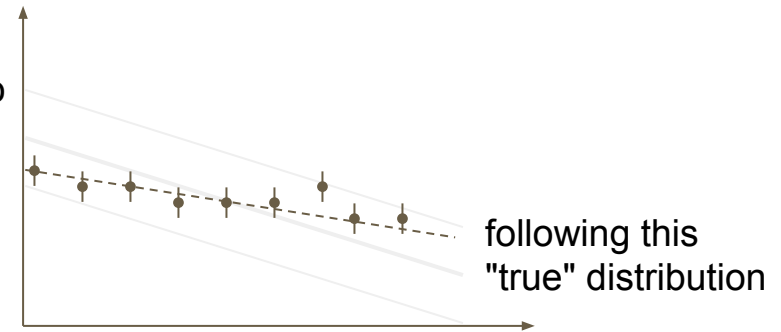
Profiling pitfalls



- The profile likelihood approach is **valid** with some **assumptions**
 - in particular, assumed that "*nature*" can be described by the model with **a single combination of values** for the parameters
- Cannot just take *large uncertainties* hoping that they are enough to cover for imperfect knowledge of S+B expectation!



will not be able to fit these points



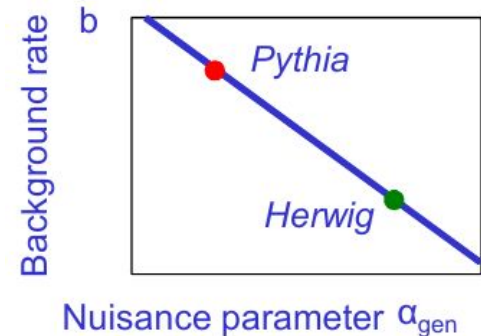
- "**Flexibility**" / "**granularity**" of the systematics model needs to be considered

Theory modeling systematics

- **Experimental systematics** nowadays often well suited for profile likelihood application:
 - come from calibrations \Rightarrow gaussian constraint appropriate
 - broken-down into several independent/uncorrelated components (JES, b -tagging...)
- Different situation for **theory systematics**:
 - **difficulty 1**: what is the **distribution** of the subsidiary measurement?
 - **difficulty 2**: what are the **parameters** of the systematic?
 - can a combination of the included parameters describe **any possible** configuration?
 - is **any allowed value** of the parameter physically meaningful?

See: https://indico.cern.ch/event/287744/contributions/1641261/attachments/535763/738679/Verkerke_Statistics_3.pdf

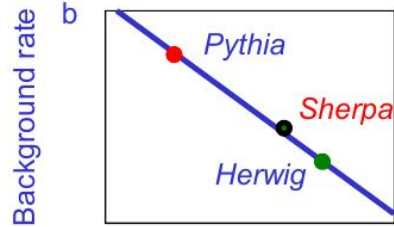
- **The obviously tricky case: "two point" systematics**
 - e.g. Herwig vs. Pythia as "parton shower and hadronization model uncertainty", as a single NP



Theory modeling systematics

One-bin case:

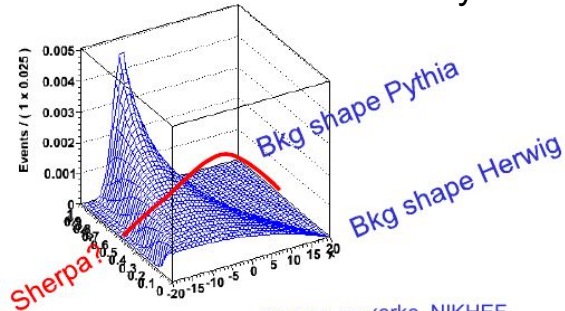
- reasonable to think that "Sherpa" can be between Herwig and Pythia



Nuisance parameter α_{gen}

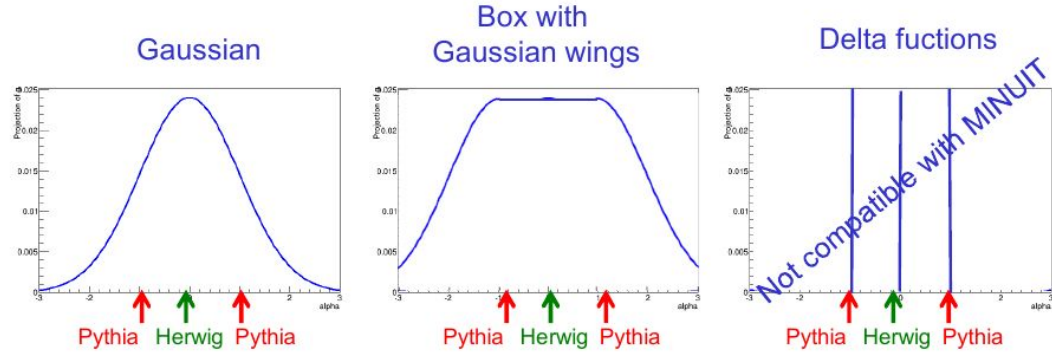
Shape case:

- Sherpa can be different from linear combination of Py and Her...



wouter.verkerke, NIKHEF

Which prior?



Pre-fit / non-constrained NP could be fine to cover for all possible models...

