# Statistics (for "daily" data analysis in physics)

Sara Bolognesi (IRFU, CEA)

# Significance



Number of events

S+B

В

# Significance



### Resolution

Back to histograms: **observed distribution = true distribution smeared, typically can be described by convolving with a Gaussian.** 



 $N_{j}^{reco} = N_{i}^{true} \cdot M_{ij}^{simu}$  where  $M_{ji}$  is a matrix which gives the probability for an event in true bin i to be reconstructed in bin j

Such matrix can be evaluated from MC (typically with cross-check, tuning from control samples)

- How to correct back from observed to "true"  $\rightarrow$  deconvolving detector effects With histograms is basically an algebric problem

# **Forward fitting**

Yesterday we considered the unfolding (inverting the matrix) The other possibility is **forward folding**, i.e. describe the true distribution as a function of unknown (to be measured) parameters and performing a fit to find the best values of the parameters which describe the observed data

 $N_i^{true} = N_j^{reco} \cdot M_{ij}^{simu} r_j$ 

r<sub>j</sub> = renormalize each bin with a semi-free term with prior value and uncertainty from MC (typically with Gaussian distribution) but to be tuned to data

# **Forward fitting**

The other possibility is **forward folding**, i.e. describe the true distribution as a function of unknown (to be measured) parameters and performing a fit to find the best values of the parameters which describe the observed data

$$N_i^{true} = N_j^{reco} \cdot M_{ij}^{simu} r_j$$

r<sub>j</sub> = renormalize each bin with a semi-free term **with prior value** and uncertainty from MC (typically with Gaussian distribution) but to be tuned to data

- The fit is an algorithm that change the MC expectations varying the parameters  $r_j$  until it find the 'best match' of MC expectation to data ('best match' = minimum of the likelihood)

# **Forward fitting**

The other possibility is **forward folding**, i.e. describe the true distribution as a function of unknown (to be measured) parameters and performing a fit to find the best values of the parameters which describe the observed data

 $N_i^{true} = N_j^{reco} \cdot M_{ij}^{simu} r_j$ 

r<sub>j</sub> = renormalize each bin with a semi-free term **with prior value** and uncertainty from MC (typically with Gaussian distribution) but to be tuned to data

- **The fit** is an algorithm that change the MC expectations varying the parameters  $r_j$  until it find the 'best match' of MC expectation to data ('best match' = minimum of the likelihood)

- Likelihood ~ function which described how well the data match with my model/expectations. Actually in frequentist terms: how probable is to observe my data, given the model

$I(N^{data} \cdot N^{simu} \cdot f(\alpha))$		$u \cdot f(\alpha)$	built in such a way to be
$L(1)_j$	, ⊥ <b>v</b> j	$(\alpha_k))$	minimal when

 $N_{i}^{data} \sim N_{i}^{simu} \cdot f(\alpha_{k})$ 

- $\alpha_{k}$  are parameters describing 'freedom' in the expectation:
  - parameters you want to measure (aka parameters of interest)

- systematic uncertainties on the model, both the physics model and the detector model (aka nuisances parameters)

 likelihood function written in a statistically correct way to consider statistical uncertainty in data and prior knowledge/uncertainty on nuisances

# Systematic uncertainties (aka nuisance parameters)

• The expectations and their dependence on nuisances  $N_i^{simu} f(\alpha_k)$  typically can be

- in form of a **full analytical description** (typically unpractical since it is difficult to encode in a single analytical function all the detector effects and their possible variations)

- in form of a **simulated histogram** which is reproduced with full simulation at each variation of all the parameters (typically unpractical since it is computationally expensive to perform a full simulation for each fit iteration)

- in form of a simulated histogram with **parametrization of uncertainties in** form of reweigthing of the histogram

# Systematic uncertainties (aka nuisance parameters)

• The expectations and their dependence on nuisances  $N_i^{simu} f(\alpha_k)$  typically are

- in form of a **full analytical description** (typically unpractical since it is difficult to encode in a single analytical function all the detector effects and their possible variations)

- in form of a **simulated histogram** which is reproduced with full simulation at each variation of all the parameters (typically unpractical since it is computationally expensive to perform a full simulation for each fit iteration)

- in form of a simulated histogram with **parametrization of uncertainties in** form of reweigthing of the histogram

• Typically  $\alpha_k$  are not completely free: they are known with a certain precision from control samples or from simulation

 $\rightarrow$  included in the likelihood with a 'penalty term' which makes the likelihood large (i.e. makes bad data-MC match) if the nuisance parameters value move away from the 'prior' estimated value

 $\rightarrow$  you need to decide how well you know this prior value and what is the distribution of its uncertainty

# Systematic uncertainties (aka nuisance parameters)

• The expectations and their dependence on nuisances  $N_i^{simu} f(\alpha_k)$  typically are

- in form of a **full analytical description** (typically unpractical since it is difficult to encode in a single analytical function all the detector effects and their possible variations)

- in form of a **simulated histogram** which is reproduced with full simulation at each variation of all the parameters (typically unpractical since it is computationally expensive to perform a full simulation for each fit iteration)

- in form of a simulated histogram with **parametrization of uncertainties in** form of reweigthing of the histogram

• Typically  $\alpha_k$  are not completely free: they are known with a certain precision from control samples or from simulation

 $\rightarrow$  included in the likelihood with a 'penalty term' which makes the likelihood large (i.e. makes bad data-MC match) if the nuisance parameters value move away from the 'prior' estimated value

 $\rightarrow$  you need to decide how well you know this prior value and what is the distribution of its uncertainty

• Typically  $\alpha_k$  prior knoweldge/uncertainty is assumed Gaussian but not always obvious. For instance theoretical uncertainty  $\rightarrow$  you can use other distributions (eg flat)

# Likelihood

Likelihood ~ function which described how well the data match with my model/expectations. Actually in frequentist terms: how probable is to observe my data, given the model

$$\begin{split} L\left(N_{j}^{data}; N_{j}^{simu}\right) &= \\ \sum_{j}^{recobins} 2\left(N_{j}^{simu} - N_{j}^{data} + N_{j}^{data} \ln\left(\frac{N_{j}^{data}}{N_{j}^{simu}}\right)\right) &+ \dots \end{split}$$

**Statistical term:** minimum when data ~ simu and written in a statistical correct way for Gaussian (Poisson) uncertainties

All derivations here (\*)

# Likelihood

Likelihood ~ function which described how well the data match with my model/expectations. Actually in frequentist terms: how probable is to observe my data, given the model

$$L(N_{j}^{\textit{data}};N_{j}^{\textit{simu}}\cdot f(\alpha_{k})) =$$

$$\sum_{j}^{recobins} 2\left(N_{j}^{simu} \cdot f(\alpha) - N_{j}^{data} + N_{j}^{data} \ln\left(\frac{N_{j}^{data}}{N_{j}^{simu} \cdot f(\alpha)}\right)\right) + \sum_{k,i} \left(\alpha_{k} - \alpha_{k}^{prior}\right) M_{ki} \left(\alpha_{i} - \alpha_{i}^{prior}\right)$$

"Chi square" multidimensional term considering possible correlations in prior knowledge of nuisances. Large if value of  $\alpha$  away from prior (some freedom with  $\sigma_{\alpha}$ )

<sup>1D:</sup> 
$$\chi^2 = \frac{(\alpha - \alpha^{prior})}{\sigma_{\alpha}}$$

The minimization algorithm will change  $\alpha$  value until finding the value which make data ~ MC at smallest possible expense of deviation from  $\alpha$  prior

### **Toys and Asimov**

You exercise/tune/develop your fit on Monte Carlo samples:

- you produce a sample of simulated events which is your reference MC sample to evaluate  $N_i^{simu} f(\alpha_k)$ 

- you produce many other samples of simulated events with small variations: eg, statistical fluctuations as expected in data, small change of systematic value (eg slightly larger detector efficiency, resolution...)

 $\rightarrow$  You analyse these samples as they were many examples of actual data

#### **Asimov fit** = fit of the reference MC sample to itself: both $N_i^{simu} f(\alpha_{\nu})$ and $N_i^{simu} f(\alpha_{\nu})$ from

the same MC reference sample

→ the fit must converge to your expectation by definition (basic closure test) → you can use it to estimate the **expected sensitivity** (i.e. postfit precision on parameter of interest)

#### **Toys** = fit of the reference MC sample to the 'varied' samples of MC

 $\rightarrow$  all the fits **should** converge

 $\rightarrow$  you can use it to look how your data may look like: in principle data should look like one of those sample  $\rightarrow$  how data fit is similar to them? (**P-value, Confidence level...**)

MINUIT (or any other algorithm) will find the minimum for you



MINUIT (or any other algorithm) will find the minimum for you



#### How to define "1 sigma" error on $\alpha$ ?

If the likelihood is a  $\chi^2$ , ie all your uncertainties have a Gaussian distribution then you have the simple  $\chi^2$  rules

$$L_{min}$$
 + 1  $\rightarrow \alpha_{min}$  +/-  $\delta \alpha$ 

MINUIT (or any other algorithm) will find the minimum for you



Typically real world is never perfectly Gaussian

→ toys: run many fits on MC by changing the prior values of your parameters around true values  $\rightarrow$  look at distribution of L<sub>min</sub>-L<sub>true</sub>

#### How to define "1 sigma" error on $\alpha$ ?

If the likelihood is a  $\chi^2$ , ie all your uncertainties have a Gaussian distribution then you have the simple  $\chi^2$  rules

$$L_{min}$$
 + 1  $\rightarrow \alpha_{min}$  +/-  $\delta \alpha$ 



MINUIT (or any other algorithm) will find the minimum for you



Typically real world is never perfectly Gaussian

→ toys: run many fits on MC by changing the prior values of your parameters around true values → look at distribution of  $L_{min}$ - $L_{true}$ e.g. integrate over 68% of your results to know the  $\Delta L$ ~'1 $\sigma$ ' error

#### How to define "1 sigma" error on $\alpha$ ?

If the likelihood is a  $\chi^2$ , ie all your uncertainties have a Gaussian distribution then you have the simple  $\chi^2$  rules

$$L_{min}$$
 + 1  $\rightarrow \alpha_{min}$  +/-  $\delta \alpha$ 



# Many dimensions

Typically the likelihood is multidimensional (since you have many unknown parameters  $\alpha_{i}$ )



In general if **correlations** are present between parmeters  $\rightarrow$  non-circular projection (eg ellipses)



18

## Linear correlations

Imagine repeating the measurement of two variables (x,y) many times

- if the two measurement are independent  $\rightarrow$  uncorrelated



- if the two measurement are positively correlated:  $y \sim \rho x$ 



(e.g. two xsec measurement at same experiment share same uncertainty on L: luminosity)

- if the two measurement are negatively correlated  $y \sim -\rho x$ 



 $\langle cov(x,y) \rangle < 0$  (e.g. rate of  $v_{e} \rightarrow e$  and  $v_{\mu} \rightarrow \mu$  are linked by  $\mu$ -e mis-identification)

0.4

-ż

<u>ک</u> 0.2

p(X)

### Relations

Typically the likelihood is multidimensional and you (almost certainly) have relations between various free parameters



The projection of "1 sigma" contour on a set of measured variables may not be circular Eg here relation such that  $\alpha_2 \sim \alpha_1^2$ 

(By the way: if  $x \sim y^2$  then mathematically their correlation = 0 !)

# Multiple minima

The likelihood may have multiple minima which appears as 'islands' in the projected contours

 $\rightarrow$  careful to explore the likelihood in all the domain where minima could be present (otherwise you may overestimate the power of your experimental data)









#### **Practical example 1**

#### **Neutrino oscillation:**

fit to likelihood of near and far detector data to extract best value and uncertainty of parameters of interest which dictate the oscillation probability

# Practical example: v oscillation



Oscillation probability estimated by comparing  $v_{\mu}$  and  $v_{e}$  rate between near and far detectors:

$$P(\nu_{\alpha} \rightarrow \nu_{\beta}) = \frac{\sin^{2}(2\theta)}{\exp^{2}\left(1.27 \frac{\Delta m_{ji}^{-}[ev^{-}]L[Km]}{E_{\nu}[GeV]}\right)}$$
(simplified 2-flavors approximation amplitude frequency)

Parameters of interest: mixing angle  $\theta$ , mass difference  $\Delta m^2$  between neutrino mass-eigenstates

### Practical example: v oscillation



The **oscillation probability**  $\nu_{\alpha} \rightarrow \nu_{\alpha'}$  which you want to estimate: it depends on the parameters you want to measure (mixing angle  $\theta$ , mass difference  $\Delta m^2$ )

#### $\nu$ oscillation: near detector

Predictions depend on number of produced neutrinos and their probability to interact.

 $N_{\nu_{\alpha}}^{ND}(E_{\nu}) = \phi(E_{\nu}) \times \sigma(E_{\nu}) dE_{\nu}$ 

flux= number of neutrinos produced by the accelerator

**cross-section** = probability of interaction of the neutrinos in the material of the detector

### $\nu$ oscillation: near detector

Predictions depend on number of produced neutrinos and their probability to interact.

$$N_{\nu_{\alpha}}^{ND}(E_{\nu}) = \phi(E_{\nu}) \times \sigma(E_{\nu}) dE_{\nu}$$

flux= number of neutrinos produced by the accelerator

**cross-section** = probability of interaction of the neutrinos in the material of the detector

Detector effects:



#### $\nu$ oscillation: near detector fit

$$N_{\nu_{\alpha}}^{ND}(E_{\nu}) \approx \phi_{\nu_{\alpha}}^{ND}(E_{\nu}) \times \sigma_{\nu_{\alpha}}^{ND}(E_{\nu}) \times \frac{1}{\epsilon^{ND}} \times p^{ND}$$

Model implemented in **MC simulation predicts expectations for flux and cross-section and detector effects:** uncertainties described by nuisance parameters constrained in a fit to near detector data

Likelihood fit to find nuisance values which best  $N_i^{data} \sim N_i^{simu} \cdot f(b_k, x_i, d_j)$  Likelinood iit to match the data:  $-2\log L_{ND} = 2\sum_{i=0}^{N_{bins}} \left( N_i^{pred}(\vec{b}, \vec{x}, \vec{d}) - N_i^{obs} + N_i^{obs} \log \frac{N_i^{obs}}{N_i^{pred}(\vec{b}, \vec{x}, \vec{d})} \right)$  $\sum \Delta \vec{b}_i \left( V_b^{-1} \right)_{ij} \Delta \vec{b}_j^T$ flux (b) Nuisances (b, x, d)  $+\sum \sum \Delta \vec{x}_i \left( V_x^{-1} \right)_{ii} \Delta \vec{x}_j^T$ are not completely cross-section free: prior knowledge (X) from simulation and  $+\sum_{a}^{N_{d}}\sum_{a}^{N_{a}}\Delta\vec{d_{i}}\left(V_{d}^{-1}\right)_{ij}\Delta\vec{d_{j}}^{T}$ control-samples in detector (d) "penalty terms"

#### $\nu$ oscillation: near detector fit

$$N_{\nu_{\alpha}}^{ND}(E_{\nu}) \approx \phi_{\nu_{\alpha}}^{ND}(E_{\nu}) \times \sigma_{\nu_{\alpha}}^{ND}(E_{\nu}) \times \frac{1}{\epsilon^{ND}} \times p^{ND}$$

Model implemented in **MC simulation predicts expectations for flux and cross-section and detector effects:** uncertainties described by nuisance parameters constrained in a fit to near detector data

 $N_{j}^{data} \sim N_{i}^{simu} \cdot f(\alpha_{k}) \qquad \Longrightarrow \qquad \alpha_{k} \pm \delta \alpha_{k} \twoheadrightarrow \alpha_{k} M_{ik} \alpha_{i}$ 

#### Actually the flux and xsec uncertainties are strongly anticorrelated



Number of events in data ~ flux times xsec  $\rightarrow$  if you increase xsec then you need to decrease xsec and viceversa...

#### v oscillation: far detector fit

$$N_{\nu_{\alpha'}}^{FD}(E_{\nu}) \approx P_{\nu_{\alpha} \rightarrow \nu_{\alpha'}}(E_{\nu}) \times \phi_{\nu_{\alpha'}}^{FD}(E_{\nu}) \times \sigma_{\nu_{\alpha'}}^{FD}(E_{\nu}) \times \frac{1}{\epsilon^{FD}} \times p^{FD}$$

 $N_{j}^{data} \sim N_{i}^{simu} \cdot f(\alpha_{k}, \beta)$ 

where  $\alpha$  are nuisances of flux and xsec strongly constrained by ND + nuisances on detector systematics (efficiency and purity)

 $\beta$  = oscillation parameters. Described by standard oscillation formulas (PMNS)



2D likelihood over 2 parameters of interest



#### v oscillation: far detector fit

$$N_{\nu_{\alpha'}}^{FD}(E_{\nu}) \approx P_{\nu_{\alpha} \rightarrow \nu_{\alpha'}}(E_{\nu}) \times \phi_{\nu_{\alpha'}}^{FD}(E_{\nu}) \times \sigma_{\nu_{\alpha'}}^{FD}(E_{\nu}) \times \frac{1}{\epsilon^{FD}} \times p^{FD}$$

 $N_{j}^{data} \sim N_{i}^{simu} \cdot f(\alpha_{k},\beta)$ 

where  $\alpha$  are nuisances of flux and xsec strongly constrained by ND + nuisances on detector systematics (efficiency and purity)

 $\beta$  = oscillation parameters. Described by standard oscillation formulas (PMNS)





# Where did the nuisances $\alpha$ go?

The likelihood depends both on the parameters of interest you want to measure (PMNS parameters:  $\beta$ ) and the nuisances parameters describing just systematics effects (flux, xsec, detector:  $\alpha$ )

#### How to "project" the likelihood on $\beta$ ? Profiling or marginalizing on nuisance $\alpha$

- When we minimize a likelihood, we can just add our nuisance parameters to the list of things to minimize

- Find a global minimum across all parameters

- Look at the variation of the parameter of interest at the best estimate of the nuisance parameters

$$L(\beta) \cdot d \alpha \sim L(\alpha_{\min}, \beta)$$

Alternatively (Bayesian) we integrate (or marginalize) over the nuisance parameters

$$L(\beta) = \int L(\alpha, \beta) d\alpha$$

# **Profiling vs marginalization?**

Profiling ~ marginalization, if error on  $\beta$  ~ constant over  $\alpha$  nuisances



If error on POI  $\beta$  changes with  $\alpha$  values and/or non linear correlation then results can be widely different!





#### Practical example 2

Higgs spin-parity:

fit to likelihood of newly discovered "Higgs-like" resonance decay kinematics to do hypothesis testing on spin/parity of the resonance

Collapse all the information on the kinematics of the final state Higgs decay in a single discriminant (you do not need to know how is built, consider it as an observable)



Collapse all the information on the kinematics of the final state Higgs decay in a single discriminant (you do not need to know how is built, consider it as an observable)



P(x|H1, H0) = probability of observing data given alternative hypothesis 1 (alternative = 0-) or given hypothesis 0 (baseline = Standard Model)

$$\begin{array}{ll} = \text{these are the likelihood !} & L^{SM}(N_{j}^{data}\,;N_{j}^{simuSM}\!\cdot\!f(\,\alpha_{k})) \\ & L^{alt}(N_{j}^{data}\,;N_{j}^{simualt}\!\cdot\!f(\,\alpha_{k})) \end{array} \end{array}$$

Expected distribution of test statistics q over many toys in Monte Carlo varying the systematic uncertainties (nuisances)



Expected distribution of test statistics q over many toys in Monte Carlo varying the systematic uncertainties (nuisances)



$$q = -2\ln(L^{alt}/L^{SM})$$

Value of test statistics from a likelihood fit to data is SM-like. By how much?

Expected distribution of test statistics q over many toys in Monte Carlo varying the systematic uncertainties (nuisances)



$$q = -2\ln\left(L^{alt}/L^{SM}\right)$$

Value of test statistics from a likelihood fit to data is SM-like. By how much?

$$\frac{P(q \ge q_{obs}; H_{alt})}{P(q \ge q_{obs}; H_{SM})} < \alpha$$

alternative signal hypotheses is excluded or not with a given confidence level  $(1 - \alpha)$ .  $\rightarrow$  alternative H: 0- excluded at 99.9%

# Hypothesis testing



**Type I error:** reject baseline (null) hypothesis when it is true ( $\alpha$ )

**Type II error:** fail to reject baseline (null) hypothesis when the alternative hypothesis is actually true ( $\beta$ )

# Hypothesis testing



**Type I error:** reject baseline (null) hypothesis when it is true ( $\alpha$ )

**Type II error:** fail to reject baseline (null) hypothesis when the alternative hypothesis is actually true ( $\beta$ )

**P-value:** probability, assuming H, to observe data with equal or lesser compatibility with H relative to the data we got.

# Hypothesis testing



 $\mu - \sigma$ 

 $\mu + \sigma$ 

**Type I error:** reject baseline (null) hypothesis when it is true ( $\alpha$ )

**Type II error:** fail to reject baseline (null) hypothesis when the alternative hypothesis is actually true ( $\beta$ )

**P-value:** probability, assuming H, to observe data with equal or lesser compatibility with H relative to the data we got.

#### We often use the Gaussian distribution as an intuitive "metric"

Even if our test statistics does not have a Gaussian distribution,

often we translate the "confidence level" or a "pvalue" as the number of standard deviation that a Gaussian variable would fluctuate in one direction to give the same p-value 42 (5 $\sigma$  discovery ~ 3x10<sup>-7</sup> probability)

### **Frequentist vs Bayesian**

The likelihood is the probability of observing data, given a certain hypothesis

$$L(N_{data}; N_{simu} \cdot f(\alpha)) = P(data|H(\alpha))$$

often used as probability of hypothesis given data but it is not correct. What we would like is the posterior PDF of  $H(\alpha)$  = probability of H (or a value) given the data

### Frequentist vs Bayesian

The likelihood is the probability of observing data, given a certain hypothesis

$$L(N_{data}; N_{simu} \cdot f(\alpha)) = P(data|H(\alpha))$$

often used as probability of hypothesis given data but it is not correct. What we would like is the posterior PDF of  $H(\alpha)$  = probability of H (or a value) given the data

$$P(H(\alpha)|data) = \frac{P(data|H(\alpha)) \cdot P(H(\alpha))}{P(data)}$$

This is our *f*requentist likelihood

prior probability of the data: since this doesn t depend on  $\alpha$  it is essentially a normalisation constant

prior probability of  $\alpha$  , i.e. encompassing our knowledge of  $\alpha$  before the measurement

### Frequentist vs Bayesian

The likelihood is the probability of observing data, given a certain hypothesis

$$L(N_{data}; N_{simu} \cdot f(\alpha)) = P(data|H(\alpha))$$

often used as probability of hypothesis given data but it is not correct. What we would like is the posterior PDF of  $H(\alpha)$  = probability of H (or a value) given the data

$$P(H(\alpha)|data) = \frac{P(data|H(\alpha)) \cdot P(H(\alpha))}{P(data)}$$
This is our frequentist prior probability of the data: since this doesn't depend on prior probability of  $\alpha$ , i.e. encompassing our knowledge of the data is the maximum of the data is the data is the maximum of the data is the

 $\alpha$  before the measurement

There is some arbitrariness on how to chose the functional form the prior  $P(H(\alpha))$ 

normalisation constant

 $\alpha$  it is essentially a

A good experiment (with large sensitivity to H( $\alpha$ )) does not depend on the choice of the prior ... but then it means that you can choose a flat prior P(H( $\alpha$ ))~constant so:  $P(H(\alpha)|data) \propto P(data|H(\alpha))$