

STATISTICAL METHODS AND TOOLS

Giovanni Benato
giovanni.benato@gssi.it

Gran Sasso Science Institute

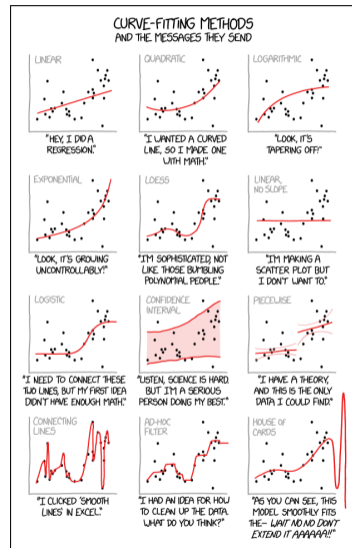
13th IDPASC School
September 19th, 2024

TABLE OF CONTENTS

- 1 INTRODUCTION
- 2 PROBABILITY THEORY
- 3 INFORMATION AND MEASUREMENT THEORY
- 4 POINT ESTIMATION
- 5 INTERVAL ESTIMATION
- 6 POINT AND INTERVAL ESTIMATION: BAYESIAN APPROACH

LITERATURE

- F. James, [Statistical methods in experimental physics](#)
→ Rigorous math, sometimes lacking “physical interpretation”
→ ~20 years old, so new methods are missing
- L. Lista, [Statistical methods for data analysis in particle physics](#)
→ Oriented to experimentalists; new methods also present
→ Concise theoretical construction;
- G. Cowan, [Statistical data analysis](#)
→ ~25 years old, so new methods are missing
- G. Cowan's [website](#)
→ Mostly slides, plus a lot of other references
- O. Behnke et al., [Data analysis in high-energy physics](#)
→ Very concise introduction; concentrates on modern methods



EXAMPLES

- A list of numerical examples can be found at this [link](#)
- All examples are written in C++/ROOT/BAT, but the code is commented so it should be possible to understand what it does

TABLE OF CONTENTS

- 1 INTRODUCTION
- 2 PROBABILITY THEORY
- 3 INFORMATION AND MEASUREMENT THEORY
- 4 POINT ESTIMATION
- 5 INTERVAL ESTIMATION
- 6 POINT AND INTERVAL ESTIMATION: BAYESIAN APPROACH

MATHEMATICAL PROBABILITY

- Let Σ be the set of all elementary and mutually exclusive events x_i
- We define the probability of occurrence of the event x_i to obey the *Kolmogorov axioms*:

$$P(x_i) \geq 0 \quad \forall i \quad (1)$$

$$P(x_i \vee x_j) = P(x_i) + P(x_j) \quad (2)$$

$$\sum_i P(x_i) = 1 \quad (3)$$

- Abstract definition
- Holds for any quantity that satisfies the 3 axioms

FREQUENTIST PROBABILITY

- Consider an experiment observing a series of N events
- Assume k events are of type X
- The frequentist probability for any **single** event to be of type X is:

$$P(X) = \lim_{N \rightarrow \infty} \frac{k}{N} = \lim_{N \rightarrow \infty} \frac{\# \text{ of favorable cases}}{\# \text{ of possible cases}} \quad (4)$$

- In principle, $P(X)$ can only be known for $N = \infty$, but often it can be computed (analytically or numerically) with high precision
- Can only be applied to **repeatable** experiments
 - Cannot predict if Italy will win the next World Cup
 - Cannot compute the probability that dinosaurs died by starvation

BAYESIAN PROBABILITY

- Probability defined as the **degree of belief** (DoB) that something will occur
- **Coherent bet**: which amount $F(x)$ are you willing to bet that x will occur, knowing that if you win you get a fixed amount K ?

$$P(x) = \frac{F(x)}{K} \Rightarrow \begin{cases} P(x) = 1 & \text{if we are sure } x \text{ will happen} \\ P(x) = 0 & \text{if we are sure } x \text{ won't happen} \\ 0 < P(x) < 1 & \text{otherwise} \\ \sum_i P(x_i) = 1 & \end{cases} \quad (5)$$

- The DoB is a property both of the observed system, as well as of the observer
- The DoB depends on the observer's knowledge, and will change according to it
- Can be applied to non-repeatable phenomena
→ Can be applied to the true value of a Physics theory!

BAYES THEOREM FOR DISCRETE EVENTS

- Let A and B be two sets of elementary events a_i and b_j
- Law of conditional probability:

$$P(A \cap B) = P(B|A) \cdot P(A) = P(A|B) \cdot P(B) \quad (6)$$

- Bayes theorem:

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)} \quad (7)$$

- More generally, if $A_i = A_0, \dots, A_n$ are **exclusive and exhaustive** sets, and if B is any event:

$$P(A_i|B) = \frac{P(B|A_i) \cdot P(A_i)}{\sum_j P(B|A_j) \cdot P(A_j)} \quad (8)$$

- $P(A_i) = \pi(A_i)$ is the **prior probability** of A_i
- $P(A_i|B)$ is the **posterior probability** of A_i after knowing that B has occurred

EXAMPLE: MEASURING PROTONS WITH A PARTICLE DETECTOR

- $P(B)$ = probability that any particle gives a triggered event
- $P(A)$ = probability of a proton hitting the detector
- $P(B|A)$ = probability of a proton giving a triggered event
- $P(A|B)$ = probability of a triggered event to be induced by a proton

EXAMPLE: DOGMAS

- Take a set of possible mutually exclusive events A_i
- Assume the following priors for each A_i :

$$\pi(A_i) = \begin{cases} 1 & \text{if } i = 0 \\ 0 & \text{if } i \neq 0 \end{cases}$$

- Applying Bayes theorem we get:

$$\text{If } i = 0 : P(A_0|B) = \frac{P(B|A_0) \cdot \pi(A_0)}{\sum_i P(B|A_i) \cdot \pi(A_i)} = \frac{P(B|A_0) \cdot 1}{P(B|A_0) \cdot 1} = 1 = \pi(A_0)$$

$$\text{If } i \neq 0 : 0 = P(A_i)$$

→ If I have a dogmatic belief about something, i.e. a belief that I cannot or am not willing to change for any reason, no experimental evidence will ever change my mind.

USE OF BAYES THEOREM IN PHYSICS

- Assume H_0 and H_1 are a complete set of hypotheses, i.e. a complete set of physics models describing a given physical phenomenon
- By convention
 - $H_0 =$ **background-only** hypothesis = the known physics is enough to explain the data
 - $H_1 =$ **signal + background** hypothesis = there is an additional component due to new physics and we know how to model it
- H_0 and H_1 will depend on some parameters, which might differ between the two hypotheses, and that we will generally indicate as $\vec{\theta} \in \Omega$
- Assume we perform n measurements of some physical observable x , which we will indicate as \vec{x}

BAYES THEOREM FOR PARAMETER ESTIMATION

$$P(\vec{\theta} | \vec{x}) = \frac{P(\vec{x} | \vec{\theta}) \pi(\vec{\theta})}{\int_{\Omega} P(\vec{x} | \vec{\theta}) \pi(\vec{\theta})} = \frac{P(\vec{x} | \vec{\theta}) \pi(\vec{\theta})}{P(\vec{x})} \quad (9)$$

- $P(\vec{\theta} | \vec{x})$ = Posterior probability for parameters $\vec{\theta}$ given the data \vec{x} and the model H_i
- $P(\vec{x} | \vec{\theta})$ = Probability of obtaining **exactly** the data \vec{x} given the parameters $\vec{\theta}$
- $\pi(\vec{\theta})$ = Prior probability of parameters $\vec{\theta}$ under the assumption of model H_i
- $P(\vec{x})$ = Probability of getting data \vec{x} given any possible value of $\vec{\theta}$
→ In the end, it's just a normalization factor.

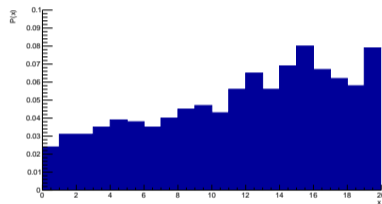
BAYES THEOREM FOR HYPOTHESIS TESTING

$$P(H_i | \vec{x}) = \frac{P(\vec{x} | H_i) \pi(H_i)}{\sum_i P(\vec{x} | H_i) \pi(H_i)} \quad (10)$$

- $P(H_i | \vec{x})$ = Posterior probability for hypothesis H_i after measuring data \vec{x}
- $\pi(H_i)$ = Prior probability for hypothesis H_i
→ This is the **subjective** part of the method
- $P(\vec{x} | H_i) = \int_{\Omega} P(\vec{x} | \vec{\theta}, H_i) \pi(\vec{\theta}) d\vec{\theta}$ = Probability of obtaining the data \vec{x} given all possible values of the parameters $\vec{\theta}$, assuming the model H_i

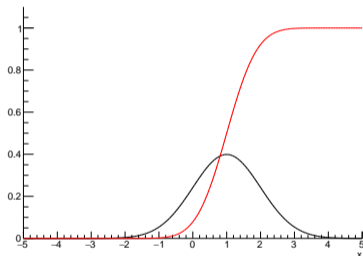
RANDOM VARIABLES AND PROBABILITY DISTRIBUTIONS

- Random event = event with more than one possible outcome, to which a probability may be associated
 - The outcome of a random event is unknown, only the probabilities of the possible outcomes are known
 - We can associate a random variable x to a random event X
- The possible outcomes $P(x_1), P(x_2), \dots$ of numerical values x_1, x_2, \dots form the **probability distribution**, obeying to the normalization condition $\sum_i P(x_i) = 1$
- What if a random variable covers a continuous interval?



PROBABILITY DENSITY FUNCTION

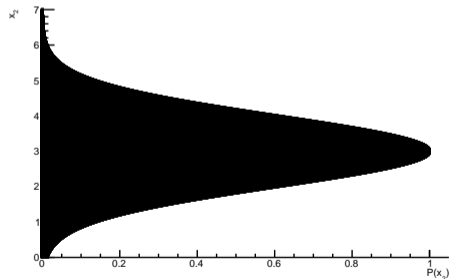
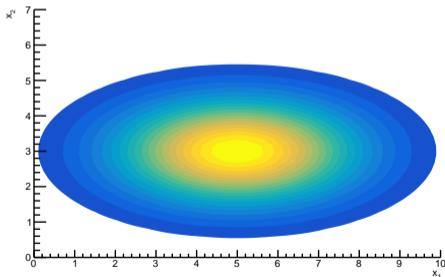
- Consider a sample space $\Omega \in \mathbb{R}^n$
- A random extraction (experiment) will lead to an outcome (measurement) corresponding to one point $\vec{x} \in \Omega$
- We can associate a **probability density** $f(\vec{x})$ to any point $\vec{x} \in \Omega$, with $f(\vec{x}) > 0$
- The probability of an event A is: $P(A) = \int_A f(\vec{x}) d\vec{x}$
 - $f(\vec{x})$ is the differential probability: $f(\vec{x}) = \frac{dP(\vec{x})}{d\vec{x}}$
 - $f(\vec{x})$ is normalized to 1: $\int_{\Omega} f(\vec{x}) d\vec{x} = 1$
- Cumulative distribution: $F(x) = \int_{-\infty}^x f(y) dy$
 - $F(\tilde{x}) = P(x \leq \tilde{x}) \forall \tilde{x}$



MARGINAL DISTRIBUTIONS

- Take an n-dimensional random variable $\vec{x} = (x_1, \dots, x_n)$ with PDF $f(\vec{x})$
- The marginal distribution for x_i is the 1-dim PDF:

$$f(x_i) = \int f(\vec{x}) dx_1 \cdots dx_{i-1} dx_{i+1} \cdots dx_n \quad (11)$$



EXAMPLE OF DISCRETE DISTRIBUTIONS

Binomial distribution: distribution of k successes out of n attempts, given a success probability p

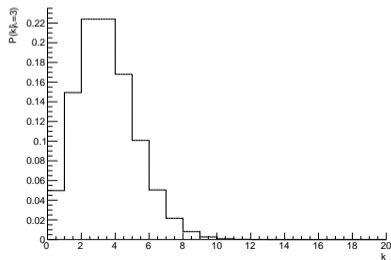
$$P(k | n, p) = \frac{n!}{k!(n-k)!} p^k (1-p)^{n-k}$$

Poisson distribution: limit of Binomial for $n \rightarrow \infty$, $p \rightarrow 0$ and $\lambda = np$ finite

$$P(k | \lambda) = \frac{e^{-\lambda} \lambda^k}{k!}$$

Properties of Poisson distribution

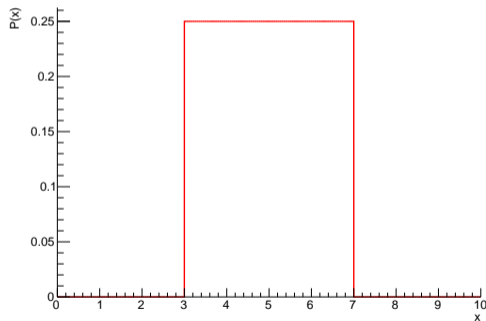
- $\lim_{\lambda \rightarrow \infty} P(k | \lambda) = \text{Gaus}(\mu = \lambda, \sigma = \sqrt{\lambda})$
- If k_1 and k_2 are Poisson-distributed with mean λ_1 and λ_2 , the sum $k_1 + k_2$ is Poisson-distributed with mean $\lambda_1 + \lambda_2$
- Randomly picking with probability ε from a Poisson process with mean λ_0 , gives a Poisson process with mean $\lambda = \varepsilon \lambda_0$



EXAMPLE OF CONTINUOUS DISTRIBUTIONS

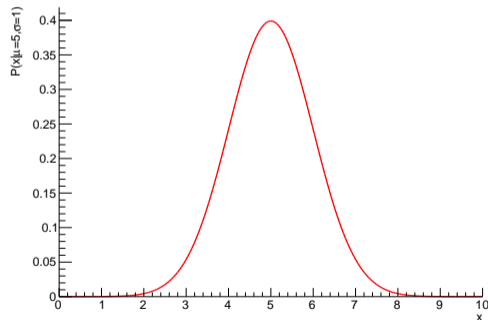
Uniform distribution

$$f(x) = \begin{cases} \frac{1}{b-a} & \text{if } a \leq x \leq b \\ 0 & \text{otherwise} \end{cases}$$



Normal (Gaussian) distribution

$$G(x | \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$



PSEUDORANDOM NUMBERS AND MC METHODS

- A pseudorandom number generator is an algorithm that generates a sequence of numbers distributed according to some PDF and that resemble very closely an actual distribution of random numbers with the same PDF
→ The true issue lies on the specification and implementation of that **very closely**
- Properties of pseudorandom number generators:
 - Each extraction must be statistically independent from the previous ones

$$f(x_i | x_{i-m}) = f(x_i)$$

- All extractions must follow the same PDF

$$f(x_i) = f(x_j) \quad \forall i, j$$

- After a given period p , the sequence will repeat itself with $x_{i+p} = x_i$, so the sequence is truly random up to $n \leq p$
- We should be able to replicate the exact same sequence for debugging purpose (seeding)

MARKOV-CHAIN MC

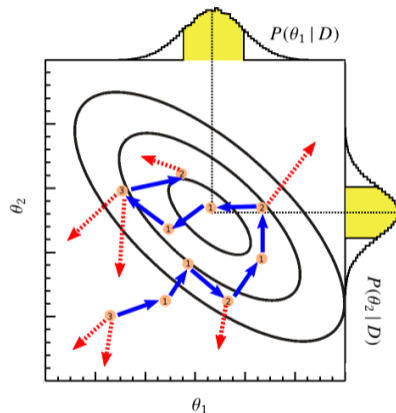
- A Markov-Chain Monte Carlo (MCMC) is a sequence of random variables x_1, \dots, x_n whose PDF obeys to

$$f(x_{i+1} | x_0, \dots, x_i) = f(x_{i+1} | x_i)$$

→ Can be more efficient in sequency highly peaked PDFs

- Example of MCMC: Metropolis-Hastings

- 1 Pick a point \vec{x}_1 uniformly distributed in space Ω and evaluate $f(\vec{x}_0)$
- 2 Generate a second point \vec{x} according to a predefined PDF $q(\vec{x}, \vec{x}_0)$ called “proposal function” or “step function”, and evaluate $f(\vec{x})$
- 3 Generate uniform number $u \in [0, 1[$
- 4 If $\frac{f(\vec{x})q(\vec{x}_0, \vec{x})}{f(\vec{x}_0)q(\vec{x}, \vec{x}_0)} > u$, accept the point and set $\vec{x}_2 = \vec{x}$, otherwise reject \vec{x}
- 5 Iterate back to (2)



PSEUDORANDOM VS MCMC-GENERATED NUMBERS

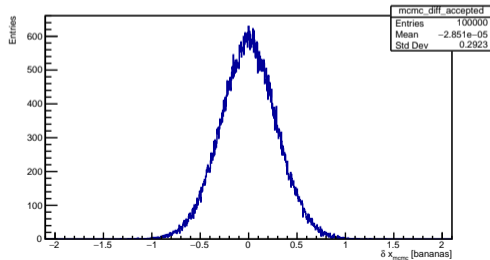
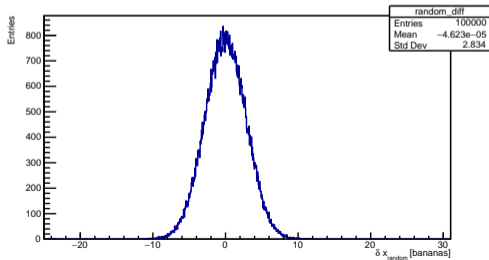
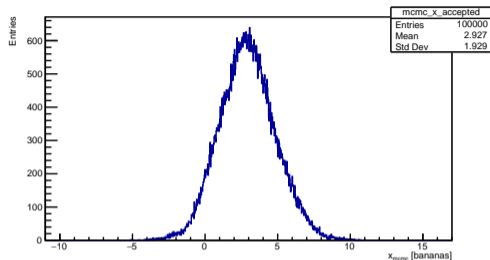
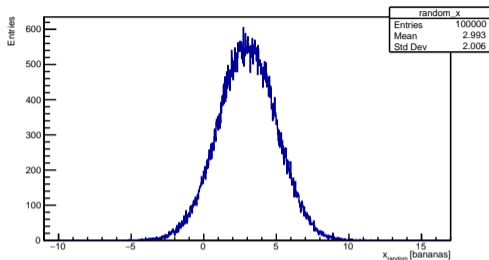


TABLE OF CONTENTS

- 1 INTRODUCTION
- 2 PROBABILITY THEORY
- 3 INFORMATION AND MEASUREMENT THEORY**
- 4 POINT ESTIMATION
- 5 INTERVAL ESTIMATION
- 6 POINT AND INTERVAL ESTIMATION: BAYESIAN APPROACH

LIKELIHOOD

- Take a real random variable \vec{x} with PDF $f(\vec{x} | \vec{\theta})$, where $\vec{\theta}$ is a set of real parameters
- Suppose we make a set of n measurements of \vec{x} : $\vec{x} = \vec{x}_1, \vec{x}_2, \dots, \vec{x}_n$
- The joint PDF of all n measurements becomes:

$$P(\vec{x} | \vec{\theta}) = P(\vec{x}_1, \dots, \vec{x}_n | \vec{\theta}) = \prod_{i=1}^n f(\vec{x}_i | \vec{\theta})$$

- Since the values \vec{x}_i are fixed (they are measured), **P is no longer a PDF**, but just a function of $\vec{\theta}$, and we denote it as \mathcal{L} :

$$\mathcal{L}(\vec{\theta}) = \mathcal{L}(\vec{x} | \vec{\theta}) = \prod_{i=1}^n f(\vec{x}_i | \vec{\theta}) \quad (12)$$

STATISTIC (SINGULAR, NOT PLURAL)

- A statistic is any new random variable $t = t(\vec{x}_1, \dots, \vec{x}_n)$
- For example, the average $\langle \vec{x} \rangle$ is a statistic
- In practice, we use a statistic whenever we need to map a highly-dimensional random variable into a one-dimensional one
- A statistic $t = t(\vec{x})$ is sufficient for θ if the conditional PDF $f(\vec{x} | t)$ is independent of θ
 - If t is a sufficient statistic, any monotonic function of t is also a sufficient statistic
 - t contains the same information about θ as the original data \vec{x}

ADDRESSED QUESTION VS REQUIRED METHODS

Question	Method
Based on the measured data \vec{x} , what is the single value $\hat{\theta}$ that is closest to the true unknown of θ ? Or, what is the most probable value for θ ?	Point estimation
Based on the measured data \vec{x} , what is the range of values that is most likely to enclose the true unknown value of θ ? Or what interval encloses the true value with a given amount of probability?	Interval estimation
Is our model $f(\vec{x} \vec{\theta})$ good enough to describe the measured data?	Goodness of fit
Assuming we have two alternative models H_0 and H_1 to describe the observed process, which of the two agrees better with the measured data?	Hypothesis testing

ADDRESSED QUESTION VS REQUIRED METHODS

Method	Solution	Algorithm
Point estimation	Find the parameter values $\hat{\theta}$ that maximize \mathcal{L} or $P(\vec{\theta} \vec{x})$	Minimizer algorithm
Interval estimation	Study the tails of \mathcal{L} or $P(\vec{\theta} \vec{x})$	Study all possible combinations of $\vec{\theta}$ giving \vec{x} with Toy-MC, or map $P(\vec{\theta} \vec{x})$ with Markov Chain MC
Goodness of fit	Quantify the probability of a random fluctuation to give a worse fit	Analytical method (e.g. χ^2) or toy-MC
Hypothesis testing	Compare probability of two models to give a better or worse fit	Toy-MC plus some method to compare the validity of the alternative hypotheses

TABLE OF CONTENTS

- 1 INTRODUCTION
- 2 PROBABILITY THEORY
- 3 INFORMATION AND MEASUREMENT THEORY
- 4 POINT ESTIMATION
- 5 INTERVAL ESTIMATION
- 6 POINT AND INTERVAL ESTIMATION: BAYESIAN APPROACH

ESTIMATOR AND THEIR PROPERTIES

- **Estimate of an unknown parameter:** mathematical procedure to determine the central value of the parameter as a function of the observed data sample
- **Estimator:** function of the data sample returned by the estimate
- Properties of estimators:
 - Consistency: does the estimator converge in probability to the true value $\bar{\theta}$ of the unknown parameter θ ?
 - Unbiasedness: expected value of the deviation of the parameter estimate from the true value $\bar{\theta}$
 - Efficiency: variance of the estimator with respect to the ideal case (Cramer-Rao inequality)
 - Robustness: (in)sensitivity to small deviations from the assumed PDF model

MAXIMUM-LIKELIHOOD ESTIMATOR

- The **maximum-likelihood estimate** of the parameter θ is that value $\hat{\theta}$ for which $\mathcal{L}(\vec{x} | \theta)$ is maximal, given the observed data \vec{x}
- Since the logarithm is monothonic, finding $\max \mathcal{L}$ is equivalent to finding $\max(\ln(\mathcal{L}))$ or $\min(-\ln(\mathcal{L}))$
→ Instead of maximizing $\prod f(x | \theta)$, we can maximize $\sum \ln f(x | \theta)$, which often is numerically simpler!

EXTENDED LIKELIHOOD

- Suppose we perform n measurements of a random variable x with PDF $f(x | \theta)$
- Suppose the number of observations n is itself a random variable with PDF $P(n | \theta)$
- We can extend the likelihood to include $P(n | \theta)$:

$$\mathcal{L} = P(n | \theta) \prod_{i=1}^n f(x_i | \theta) \quad (13)$$

- In many cases, $P(n | \theta)$ is a Poisson distribution with mean $\lambda = \lambda(\theta)$:

$$\mathcal{L} = \frac{e^{-\lambda} \lambda^n}{n!} \prod_{i=1}^n f(x_i | \theta) \quad (14)$$

EXAMPLE: MAXIMUM-LIKELIHOOD FOR SIGNAL + BACKGROUND

- Assume we want to describe the data with a signal (new physics) plus background (known physics). The overall PDF is:

$$f(x | s, b, \theta) = \frac{s}{s+b} f_s(x | \theta) + \frac{b}{s+b} f_b(x | \theta) \quad (15)$$

where s and b are the expected number for signal and background counts

→ The expectation value for n is $\lambda = s + b$

- The likelihood becomes:

$$\mathcal{L} = \frac{e^{-(s+b)} (s+b)^n}{n!} \prod_{i=1}^n \frac{sf_s(x | \theta) + bf_b(x | \theta)}{s+b} = \frac{e^{-(s+b)}}{n!} \prod_{i=1}^n [sf_s(x | \theta) + bf_b(x | \theta)] \quad (16)$$

- Taking the logarithm:

$$\ln \mathcal{L} = -s - b - \ln n! + \sum_{i=1}^n \ln [sf_s(x | \theta) + bf_b(x | \theta)] \quad (17)$$

BINNED LIKELIHOOD

- Suppose the number of measurement n is so large that computing $\sum \ln f(x_i | \theta)$ takes too long
- We can simplify the problem by binning x in $m \ll n$
→ The likelihood becomes a multinomial times the PDF of n :

$$\mathcal{L} = P(n | \theta) \left[\prod_{i=1}^m \frac{p_i^{k_i}}{k_i!} \right] n! \quad (18)$$

where:

- i is now the bin index (not the event index!)
- k_i is the number of counts in bin i
- p_i is the probability associated to bin i
- $\sum_i k_i = n$

BINNED LIKELIHOOD FOR POISSON-DISTRIBUTED n

- Take $\lambda_i = \lambda \int_{\delta x_i} f(x | \theta) dx =$ expectation value for bin i
→ $\sum_i \lambda_i = \lambda$
- Take $p_i = \frac{\lambda_i}{\lambda}$
- The likelihood becomes:

$$\mathcal{L} = \frac{e^{-\lambda} \lambda^n}{n!} \prod_{i=1}^m \left(\frac{\lambda_i}{\lambda} \right)^{k_i} \frac{1}{k_i!} = e^{-\lambda} \prod_{i=1}^m \frac{\lambda_i^{k_i}}{k_i!} = \prod_{i=1}^m \frac{e^{-\lambda_i} \lambda_i^{k_i}}{k_i!} \quad (19)$$

→ The binned version of an extended \mathcal{L} with Poisson-distributed n is the product of a Poisson term for each bin!

LEAST-SQUARES (OR χ^2) METHOD

- Take a set of n measurements with Gaussian uncertainties $y_i \pm \sigma_i$
- Assume each y_i corresponds to a perfectly known x_i following the relation $y = y(x, \theta)$
→ $y(x, \theta)$ is just a function, not a PDF
- The likelihood will be a product of Gaussian PDFs around the curve $y(x, \theta)$:

$$\mathcal{L}(\vec{y} | \theta) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma_i} \exp\left(-\frac{(y_i - y(x_i, \theta))^2}{2\sigma_i^2}\right) \quad (20)$$

- Let's minimize $-2 \ln \mathcal{L}$ instead:

$$-2 \ln \mathcal{L} = \sum_{i=1}^n \frac{(y_i - y(x_i, \theta))^2}{2\sigma_i^2} + 2 \sum_{i=1}^n \ln(2\pi\sigma_i^2) \quad (21)$$

→ The first term is a sum of standard-normal variables, so it follows a χ^2 distribution

→ The second term is a constant and can be dropped

- If θ is not known, $\min(-2 \ln \mathcal{L})$ follows a χ^2 distribution only if all y_i are uncorrelated

LEAST-SQUARES METHOD FOR HISTOGRAMS

- Assume we have a set of measurements x_i following a PDF $f(x_i | \theta)$
- Assume we bin the data, and every bin has a large number of entries k_i so that the corresponding Poisson distribution can be approximated with a Gaussian
- We can approximate the likelihood with:

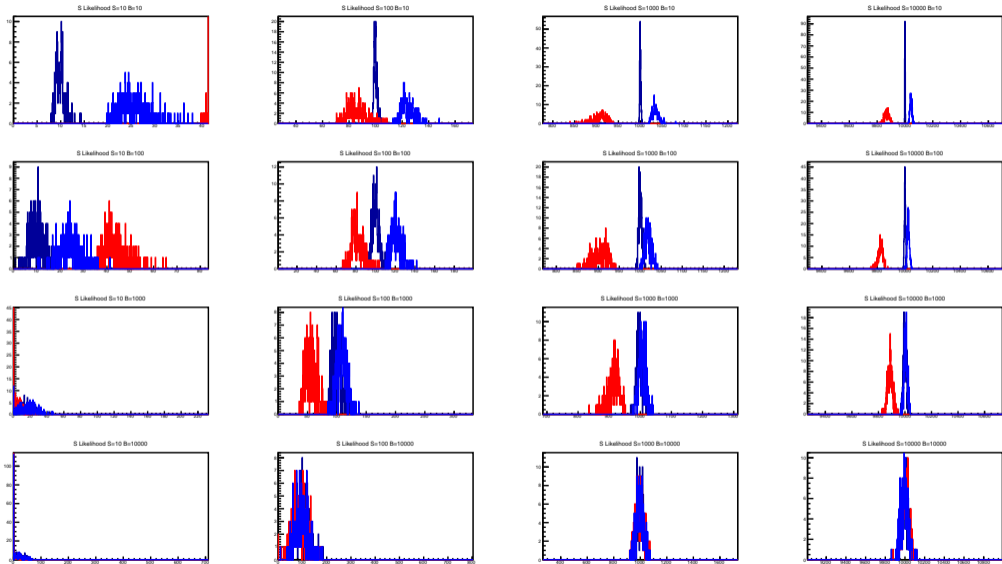
$$\mathcal{L} \simeq \prod_{\text{bin } i=1}^m \frac{1}{\sqrt{2\pi}\sigma_i} \exp\left(-\frac{(k_i - \lambda_i(\theta))^2}{2\sigma_i^2}\right) \Rightarrow -2 \ln \mathcal{L} = \sum_{i=1}^m \ln(2\pi\sigma_i^2) + \sum_{i=1}^m \frac{(k_i - \lambda_i(\theta))^2}{\sigma_i^2} \quad (22)$$

- The problem is that we don't know $\sigma_i \dots$
 - The first term varies slowly, so we can neglect it
 - **Neyman's χ^2** : set $\sigma_i^2 = k_i$
 - Standard for many predefined fit methods, e.g. in ROOT/Minuit
 - Fails miserably if just one bin is empty
 - **Pearson's χ^2** : set $\sigma_i^2 = \lambda_i(\theta)$
 - Works also for empty bins, but is unreliable for small λ_i
- Take-home message: **NEVER USE A χ^2 FIT ON A HISTOGRAM!**

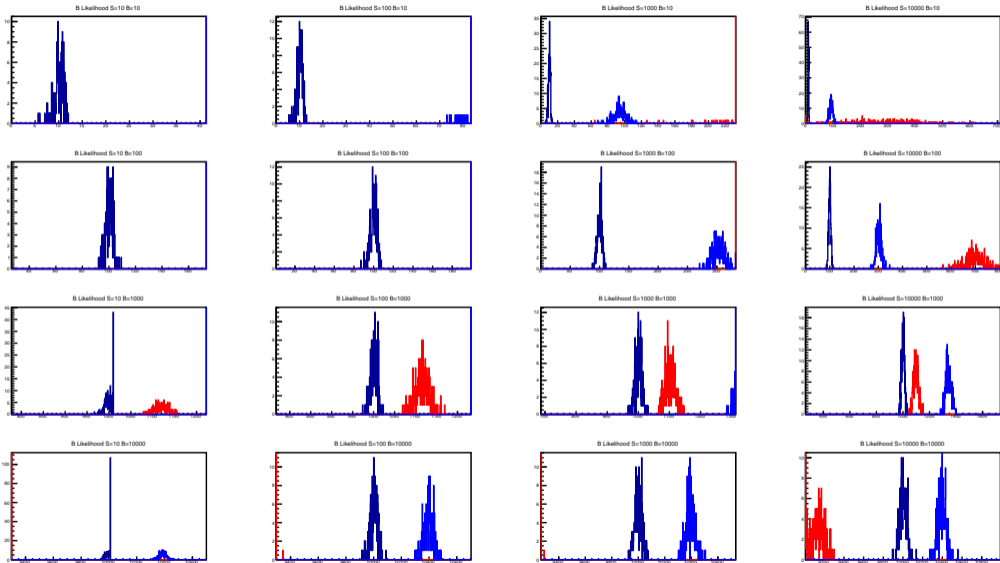
EXAMPLE: \mathcal{L} VS NEYMAN- χ^2 VS PEARSON- χ^2

- Generate toy-MC spectra with S Gaussian-distributed signal events and B uniformly distributed background events
- Fit each toy-MC with \mathcal{L} , Neyman- χ^2 and Pearson- χ^2
- Compare distribution of best-fit values for 10^4 toy-MC, and for various combination of $S, B = 10, 100, 1000, 10000$

EXAMPLE: \mathcal{L} VS NEYMAN- χ^2 VS PEARSON- χ^2



EXAMPLE: \mathcal{L} VS NEYMAN- χ^2 VS PEARSON- χ^2



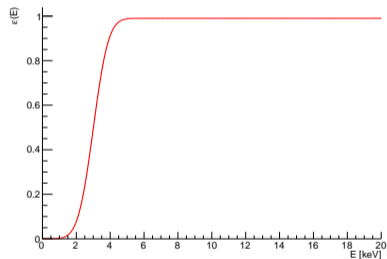
CHOOSING THE CORRECT PDF

- Goal: study the trigger efficiency curve for a detector
- Assume we know that the efficiency curve can be modeled with an error function (cumulative of a Gaussian):

$$\varepsilon = \frac{p}{2} \left[1 + \operatorname{erf} \left(\frac{E - \mu}{\sigma} \right) \right] \quad (23)$$

where $p = 0.99$, $\mu = 3 \text{ keV}$, $\sigma = 1 \text{ keV}$

- Suppose we inject a known number n of pulser events at energies $E_i = 1, 2, 3, \dots, 20 \text{ keV}$, and we want to run a Bayesian fit to reconstruct the PDF of p
- What is the correct likelihood? What PDF are the data following?



CHOOSING THE CORRECT PDF

- For each energy E_i , we inject a **known** number of pulser events n
 - The number of detected events k_i follows a **binomial** distribution
 - We must use a binomial likelihood!

$$\mathcal{L} = \prod_i \frac{n!}{k_i!(n - k_i)!} \varepsilon(E_i | p, \mu, \sigma)^{k_i} (1 - \varepsilon(E_i | p, \mu, \sigma))^{n - k_i}$$

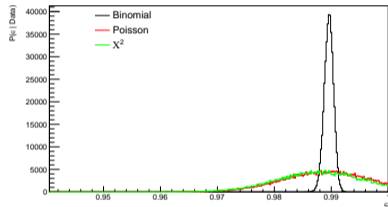
- What happens if we use a Poisson likelihood, or a χ^2 instead?

CHOOSING THE CORRECT PDF

- For each energy E_i , we inject a **known** number of pulser events n
 - The number of detected events k_i follows a **binomial** distribution
 - We must use a binomial likelihood!

$$\mathcal{L} = \prod_i \frac{n!}{k_i!(n - k_i)!} \varepsilon(E_i | p, \mu, \sigma)^{k_i} (1 - \varepsilon(E_i | p, \mu, \sigma))^{n - k_i}$$

- What happens if we use a Poisson likelihood, or a χ^2 instead?



SIMULTANEOUS FITS

- Suppose we have two sets of data \vec{x} and \vec{y} with PDFs that have some parameters in common:

$$f(\vec{x} | \vec{\theta}, \vec{\nu}) \quad \text{and} \quad f(\vec{y} | \vec{\theta}, \vec{\phi})$$

- The likelihood will simply be the product of the two likelihoods:

$$\mathcal{L}(\vec{x}, \vec{y} | \vec{\theta}, \vec{\nu}, \vec{\phi}) = \mathcal{L}(\vec{x} | \vec{\theta}, \vec{\nu}) \cdot \mathcal{L}(\vec{y} | \vec{\theta}, \vec{\phi}) \quad (24)$$

TABLE OF CONTENTS

- 1 INTRODUCTION
- 2 PROBABILITY THEORY
- 3 INFORMATION AND MEASUREMENT THEORY
- 4 POINT ESTIMATION
- 5 INTERVAL ESTIMATION**
- 6 POINT AND INTERVAL ESTIMATION: BAYESIAN APPROACH

CONFIDENCE INTERVAL

- Goal: we want to find the range $\theta_a \leq \theta \leq \theta_b$ that contains the true value θ_0 with some probability β (usually 68%)
- Given an observation x from a PDF $f(x | \theta)$, the probability content β of the region $[a, b]$ in **x-space** is:

$$\beta = P(a \leq x \leq b) = \int_a^b f(x | \theta) dx$$

- However, we want to find a corresponding interval $[\theta_a, \theta_b]$ in **θ -space** so that $P(\theta_a \leq \theta \leq \theta_b) = \beta$
 - Such interval is called **confidence interval**
 - A method that provides such an interval is said to have the property of **coverage**
- Notice that:
 - The true value θ_0 is and will always remain an unknown constant
 - θ_a and θ_b must be functions of x , not of θ

CONFIDENCE INTERVAL FOR NORMALLY-DISTRIBUTED DATA

- Let $f(x | \vec{\theta})$ be a normal distribution with $\vec{\theta} = (\mu, \sigma)$
- If μ and σ are known:

$$\beta = P(a \leq x \leq b) = \Phi\left(\frac{b - \mu}{\sigma}\right) - \Phi\left(\frac{a - \mu}{\sigma}\right) \quad \text{where } \Phi \text{ is the cumulative}$$

- If μ and σ are unknown, we can't compute the probability content of $[a, b]$, but we can compute the probability β that x lies in some interval $[\mu + c, \mu + d]$. By defining $y = \frac{x - \mu}{\sigma}$ we get:

$$\beta = P(\mu + c \leq x \leq \mu + d) = \Phi\left(\frac{d}{\sigma}\right) - \Phi\left(\frac{c}{\sigma}\right)$$

- Then we can invert to obtain:

$$\beta = P(x - d \leq \mu \leq x - c)$$

→ This is still a probability statement about x , while μ is still unknown!

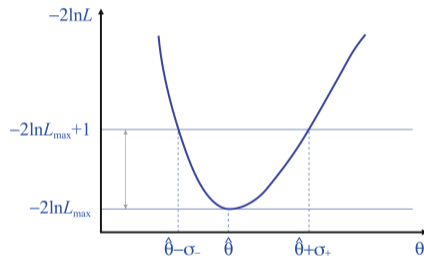
→ There is infinite possible choices for such interval!

LOG-LIKELIHOOD SCAN

- Assume we have a 1-dim Gaussian likelihood:

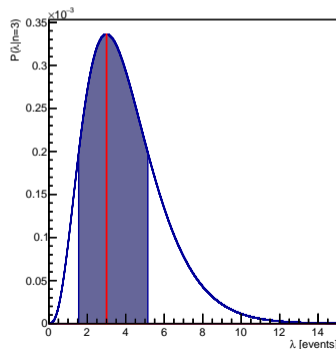
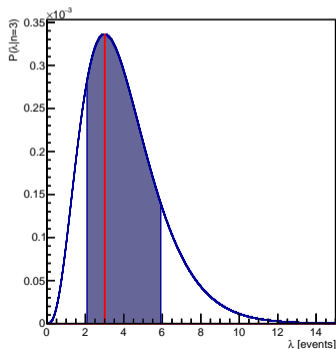
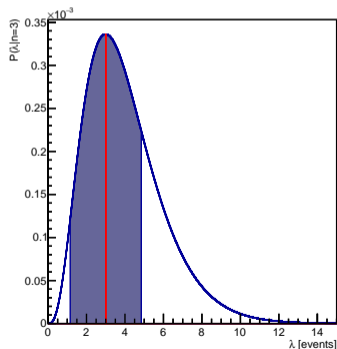
$$-2 \ln \mathcal{L} = -2 \ln \mathcal{L}_{\max} + \frac{(\mu - x)^2}{\sigma^2}$$

- The intercepts at $-2 \ln \mathcal{L}_{\max} + 1$ correspond to the $\pm 1\sigma$ interval
- The intercepts at $-2 \ln \mathcal{L}_{\max} + 4$ correspond to the $\pm 2\sigma$ interval
- For a non-Gaussian likelihood, $-2 \ln \mathcal{L}$ is not parabolic, but the max- \mathcal{L} estimate is invariant under reparameterization
 - The intercept at $-2 \ln \mathcal{L}_{\max} + k$ will maintain the same coverage as in the Gaussian case!



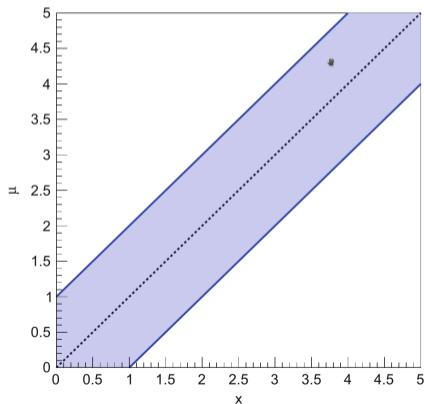
ORDERING RULES

- **Central interval:** $[x_L, x_U] = [\bar{x} - \delta, \bar{x} + \delta]$
- **Equal areas:** $\int_{-\infty}^{x_L} f(x | \theta) dx = \int_{x_U}^{+\infty} f(x | \theta) dx = \frac{1-\beta}{2}$
- **Shortest interval:** $f(x_L | \theta) = f(x_U | \theta) \quad \wedge \quad \int_{x_L}^{x_U} f(x | \theta) dx = \beta$
- **Lower limit:** $\int_{-\infty}^{x_L} f(x | \theta) dx = 1 - \beta$
- **Upper limit:** $\int_{x_U}^{+\infty} f(x | \theta) dx = 1 - \beta$



NEYMAN CONFIDENCE BELT

- Take a variable x with a PDF $f(x | \theta)$
- Suppose we can take a function of the data $t(x)$ so that $\beta = P(t_1(\theta) \leq t \leq t_2(\theta))$
- Suppose that we can compute t_1 and t_2 for each value of θ
 - The space between the curves $t_1(\theta)$ and $t_2(\theta)$ is the Neyman confidence belt
- For the case $t = x$, and Gaussian-distributed x :



FLIP-FLOPPING

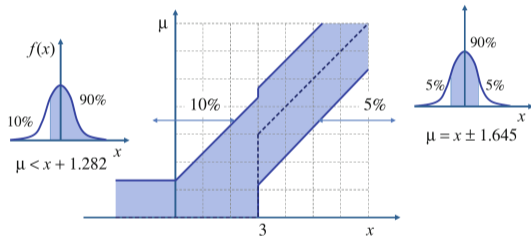
- Take a variable x with a Gaussian PDF with known σ
- Suppose Physics tells us that $\mu \geq 0$
- Suppose the observable x is subject to instrumental fluctuations and can be negative
- We can decide to quote:

$$\hat{\mu}(x) = \begin{cases} x & \text{if } \frac{x}{\sigma} \geq 3 \\ 0 & \text{otherwise} \end{cases}$$

- The corresponding interval will be:

$$[\mu_1, \mu_2] = \begin{cases} [\hat{\mu} - 1.65\sigma, \hat{\mu} + 1.65\sigma] & \text{if } \frac{x}{\sigma} \geq 3 \\ [0, \hat{\mu} + 1.282\sigma] & \text{otherwise} \end{cases}$$

- For some values we have a coverage of 85% only!



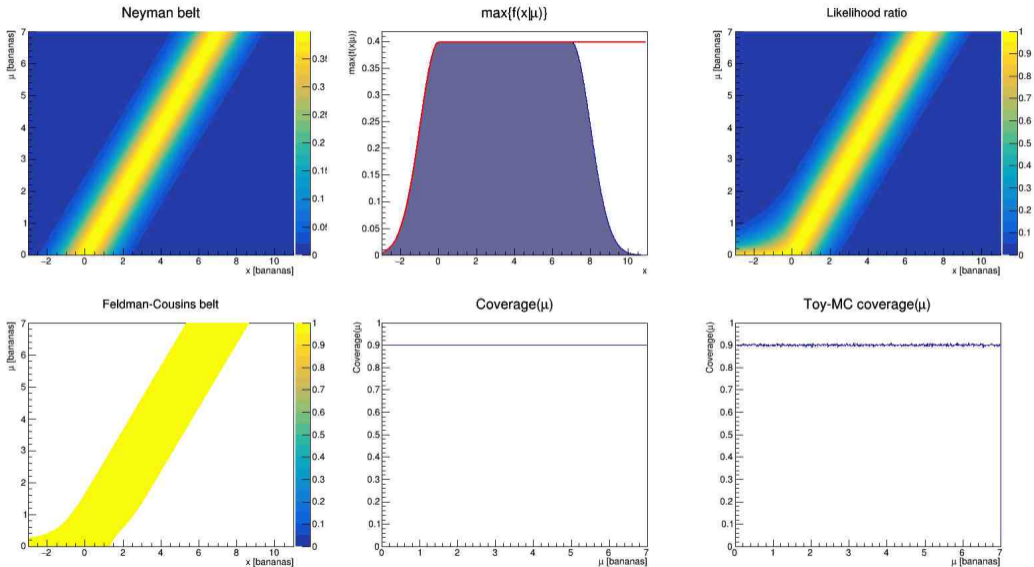
UNIFIED FELDMAN-COUSINS APPROACH

- Solution to flip-flopping issue: invent an ordering rule that allows us to smoothly switch from a two-sided interval to an upper limit with no dependence on the measured data
- Use the \mathcal{L} -ratio as an ordering rule

$$\lambda(x | \theta_0) = \frac{\mathcal{L}(x | \theta_0)}{\mathcal{L}(x | \hat{\theta})}$$

- Procedure:
 - ① Find the best-fit $\hat{\theta}$
 - ② Fix θ to some value θ_0 , then start from $\hat{x} = \max_x \lambda(x | \theta_0)$ and move left and right until we get a coverage β on $\mathcal{L}(x | \theta_0)$
 - ③ Repeat for all values of θ (in the physical range)

FELDMAN-COUSINS FOR GAUSSIAN CASE



FELDMAN-COUSINS FOR KATRIN DATA

- In Phys. Rev. Lett. 123 (2019) 221802, the KATRIN experiment reported the following result on the electron neutrino mass:
 $m^2 = -1.0^{+0.9}_{-1.1} \text{ eV}^2 \rightarrow$ This is their fit observable
 $m_\nu < 1.1 \text{ eV} \rightarrow$ This is their parameter of interest
- We can approximate the PDF of m_ν^2 with a Gaussian:

$$f(m^2|m_\nu) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(m^2 - m_\nu^2)^2}{2\sigma^2}\right)$$

with $\sigma = 1$, which is very close to the asymmetric uncertainties quoted above

- By doing this, we get $m_\nu < 0.91 \text{ eV}$, which is pretty close to their published value!

FELDMAN-COUSINS FOR KATRIN DATA

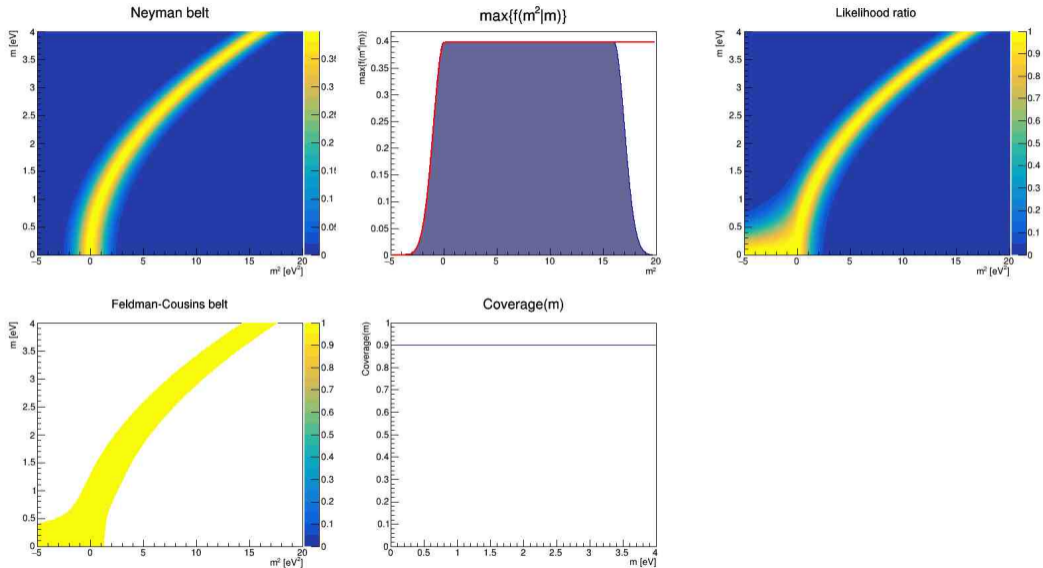


TABLE OF CONTENTS

- 1 INTRODUCTION
- 2 PROBABILITY THEORY
- 3 INFORMATION AND MEASUREMENT THEORY
- 4 POINT ESTIMATION
- 5 INTERVAL ESTIMATION
- 6 POINT AND INTERVAL ESTIMATION: BAYESIAN APPROACH

BAYESIAN PARAMETER ESTIMATION

- Take the usual n measurement of a variable \vec{x} with PDF $f(\vec{x} | \vec{\theta})$
- Before the measurement, our degree of belief on the parameters is $\pi(\vec{\theta})$
- The probability of obtaining exactly the data \vec{x} is

$$P(\vec{x} | \vec{\theta}) = \mathcal{L}(\vec{x} | \vec{\theta})$$

- We can use the Bayes theorem to find the posterior probability for the parameters $\vec{\theta}$ given the data \vec{x} :

$$P(\vec{\theta} | \vec{x}) = \frac{\mathcal{L}(\vec{x} | \vec{\theta}) \pi(\vec{\theta})}{\int_{\Omega} \mathcal{L}(\vec{x} | \vec{\theta}) \pi(\vec{\theta})}$$

- The global mode of $P(\vec{\theta} | \vec{x})$ represent the most probable combination of all parameters

BAYESIAN INTERVAL ESTIMATION

- Suppose we are just interested in one of the parameters, e.g. θ , out of a list of parameters $(\theta, \vec{\nu})$
 - What is the most probable value for θ given the data \vec{x} ?
 - What is the interval that contains the most probable value of θ with 68% probability?
- The answer is: **marginalization**

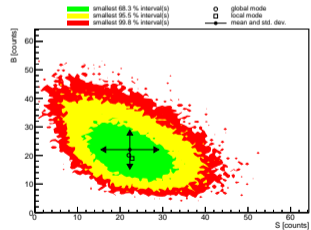
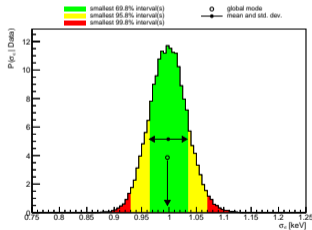
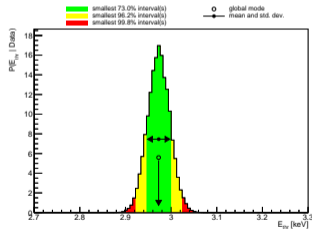
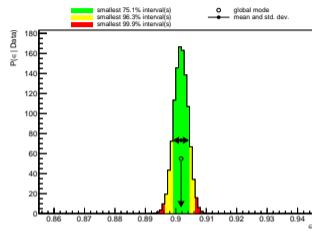
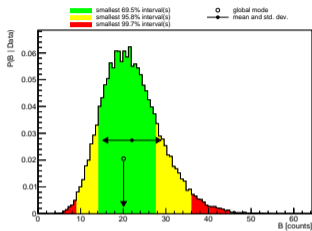
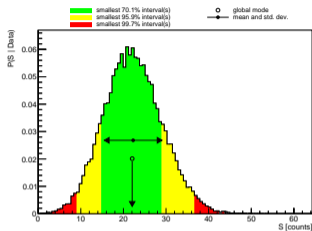
$$P(\theta | \vec{x}) = \int_{\Omega_{\nu}} P(\theta, \vec{\nu} | \vec{x}) d\vec{\nu}$$

- From $P(\theta | \vec{x})$ we can quote:
 - The marginalized mode $\hat{\theta}$, which can also be at the boundary of the physics region
 - The central/shortest interval $[\theta_1, \theta_2]$, or an upper/lower limit (typically at 90% or 95%)
 - These intervals are called **credible intervals** because they state that, based on our current knowledge, we believe the true value of θ is in that range with the specified probability.
 - The concept of coverage does not apply here!

EXAMPLE: BAYESIAN SIMULTANEOUS FIT

- Suppose we want to run a Bayesian fit of the pulser data described above and physics data, recorded between 10 and 20 keV, and featuring
 - a Gaussian signal with $\mu = 15$ keV and $\sigma = 1.2$ keV
 - a flat background
- How does the marginalized distribution of each parameter look like?
- How do the correlation plots of all pairs of parameters look like?

EXAMPLE: BAYESIAN SIMULTANEOUS FIT



BAYESIAN FITS: PRACTICAL TIPS

- We need to **map** the posterior $P(\vec{\theta} | \vec{x})$ and integrate it over nuisance parameters, however $P(\vec{\theta} | \vec{x})$ might be complicated to integrate over $d\vec{v}$ or $d\vec{\theta}$
- Solution:
 - ① Map $P(\vec{\theta} | \vec{x})$ using a MCMC with n tested samples
 - ② The denominator $\int \mathcal{L}\pi d\vec{\theta}$ is just a constant, so forget about it
 - ③ Build the marginalized PDF of θ by histogramming its tested values
 - ④ Renormalize the marginalized PDF by dividing over n
 - ⑤ Extract $\hat{\theta}$ and the credible interval from the histogram

Questions?