

Statistical Methods for Data Analysis in Particle Physics

Luca Lista

Università Federico II, Napoli

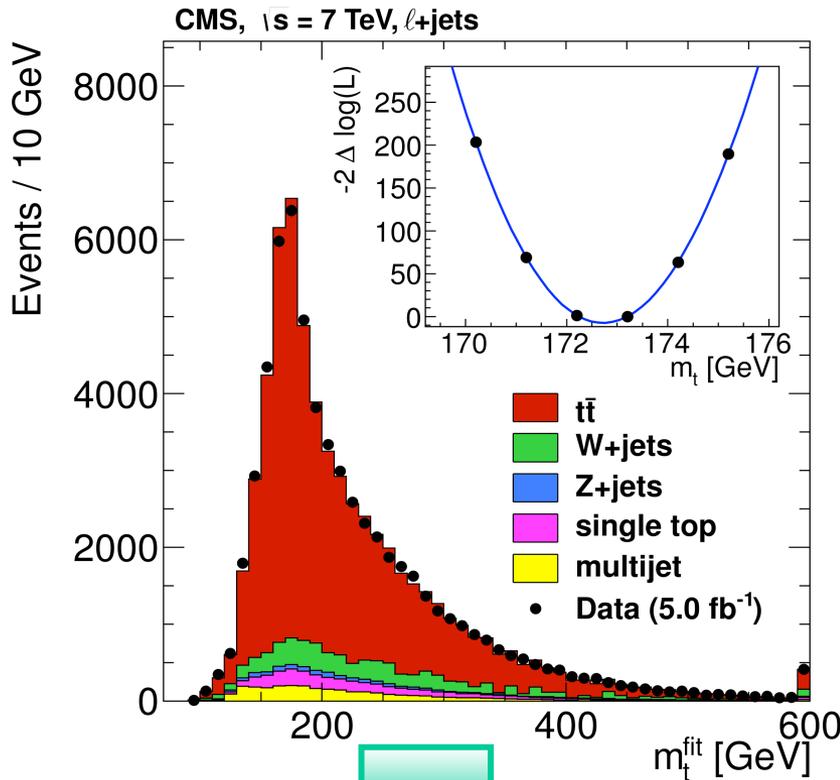
INFN Sezione di Napoli



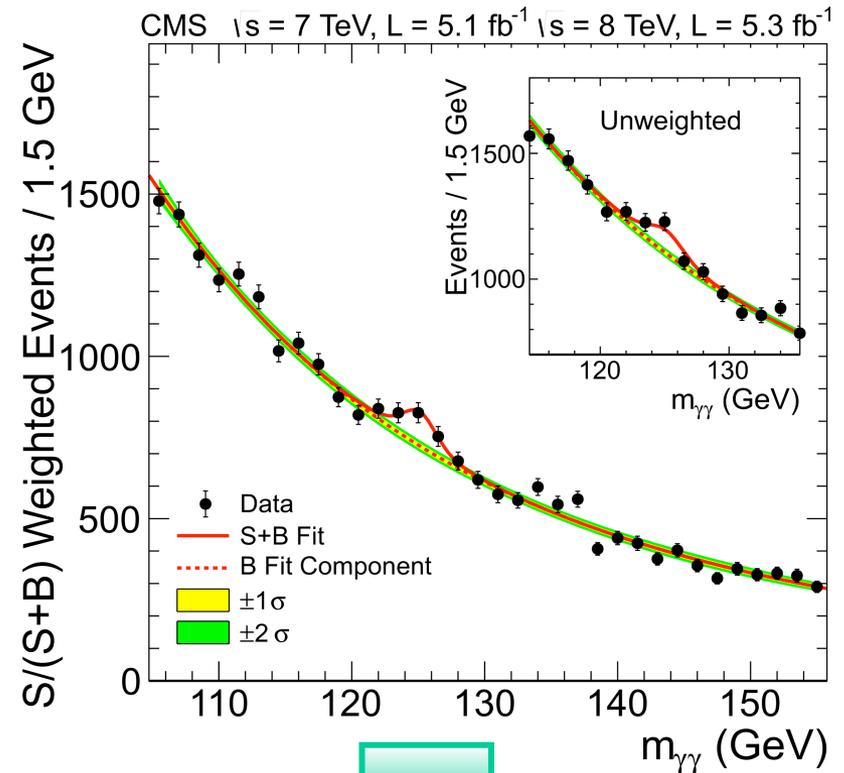
- Introduction to probability
 - Subjective/Bayesian probability vs Frequentist probability
 - Kolmogorov axiomatic approach to probability
 - Probability distributions: discrete, continuous, in more dimensions
 - Conditional probability, independent events and variables
 - The Bayes theorem
 - Examples of application of Bayes theorem
 - Bayes rule and likelihood function
 - Bayesian approach to probability as learning process
 - Inference with the Bayesian approaches
 - Choice of credible and confidence intervals: upper and lower limits
 - Error propagation with Bayesian estimates
 - Issues with prior PDF with the Bayesian approach
-

- Measurements

- Discoveries

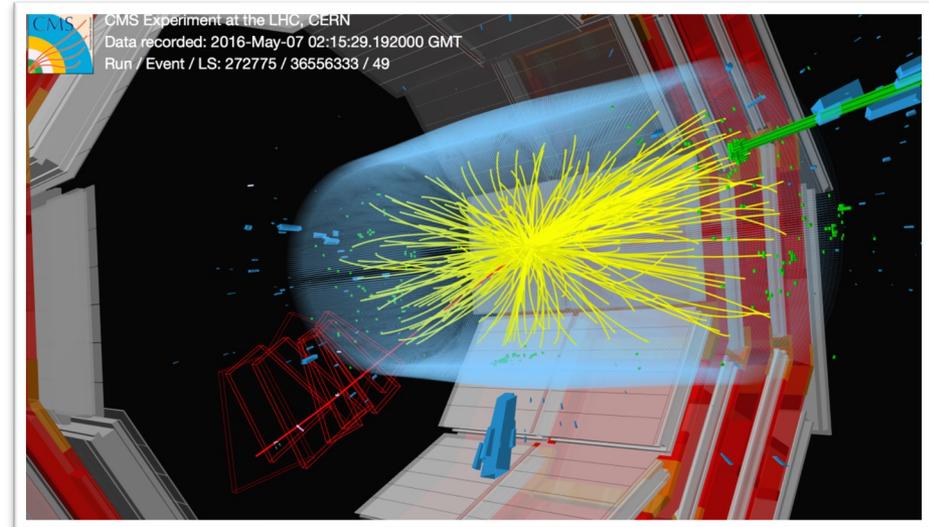


$m_t = 173.49 \pm 1.07 \text{ GeV}$



Higgs boson!

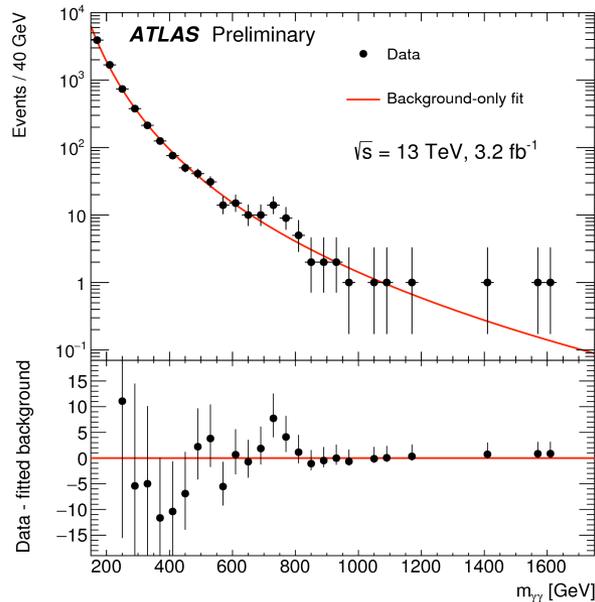
- Measurements are closely related to random processes and probability
- Repeated measurements always give the same result only for trivial cases
- Many physical processes have intrinsic randomness
 - Quantum Mechanics: $\mathcal{P} \propto |\mathcal{A}|^2$
- Detector response is somewhat random
 - Fluctuations, resolution, efficiency, ...



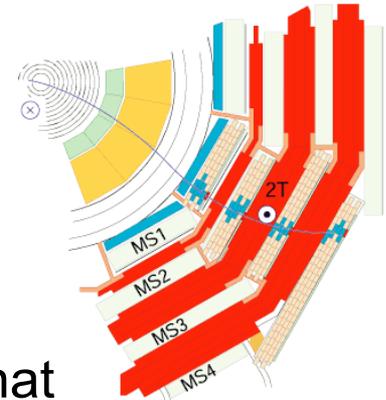
Repeatable experiments



- What's the probability to extract one ace in a deck of cards?
- What is the probability to win a lottery (or bingo, ...)?
- What is the probability that a pion is incorrectly identified as a muon in CMS?



- What is the probability that a fluctuation in the background can produce a peak in the $\gamma\gamma$ spectrum with a magnitude at least equal to what has been observed by ATLAS?
- **Note:** different question w.r.t.: what is the probability that the peak is due to a background fluctuation? (non repeatable!)



Non repeatable claims

DAY	HIGH/LOW	DESCRIPTION	PRECIP	WIND	HUMIDITY
TODAY May 20	69°/47°	☁️ Partly Cloudy	/ 0%	E 5 mph	54%
SAT May 21	78°/51°	☀️ Sunny	/ 0%	SSE 8 mph	54%
SUN May 22	73°/45°	☁️ Partly Cloudy	/ 20%	SW 9 mph	58%
MON May 23	57°/43°	☁️ Rain	/ 90%	SW 9 mph	67%
TUE May 24	66°/46°	☀️ Mostly Sunny	/ 20%	NE 7 mph	57%
WED May 25	71°/51°	☁️ Partly Cloudy	/ 20%	ESE 6 mph	60%

- Could be about **future events**:
 - what is the probability that **tomorrow it will rain in Geneva**?
 - what's the probability that **your favorite team will win next championship**?
- But also **past events**:
 - What's the probability that **dinosaurs went extinct because of an asteroid**?



- More in general, it's about **unknown events**:
 - What is the probability that **dark matter is made of particles heavier than 1 TeV**?
 - What is the probability that **climate changes are mainly due to human intervention**?



- Probability can be defined in different ways
- The applicability of each definition depends on the kind of claim we are considering to applying the concept of probability
- One subjective approach expresses the degree of belief/credibility of the claim, which may vary from subject to subject
- For repeatable experiments, probability may be a measure of how frequently the claim is true

- Probability determined by **symmetry** properties of a random device
- “**Equally undecided**” about event outcome, according to Laplace definition

$$\text{Probability: } P = \frac{\text{Number of favorable cases}}{\text{Number of total cases}}$$



Pierre Simon Laplace
(1749-1827)

$P = 1/2$



$P = 1/6$
(each dice)

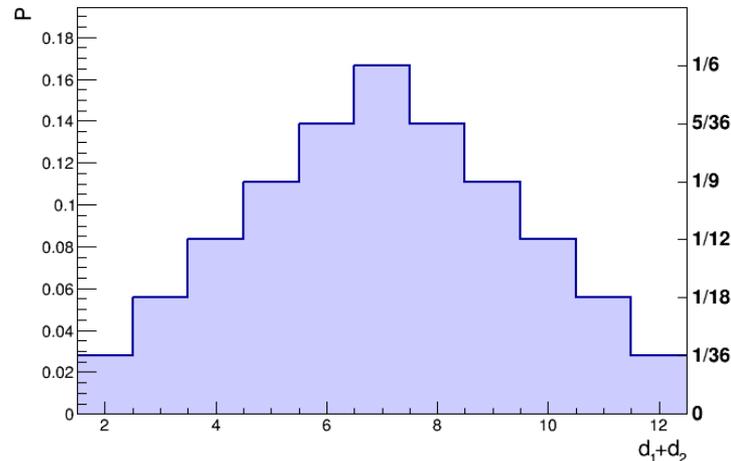


$P = 1/10$



$P = 1/4$

- Composite cases are managed via **combinatorial analysis**
- Reduce the (composite) event of interest into elementary equiprobable events (**sample space**)



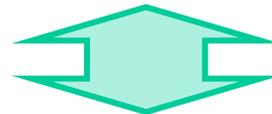
E.g:

- $2 = \{(1,1)\}$
- $3 = \{(1,2), (2,1)\}$
- $4 = \{(1,3), (2,2), (3,1)\}$
- $5 = \{(1,4), (2,3), (3,2), (4,1)\}$
- etc. ...

- Statements about an **event** can be defined via set algebra
 - **and/or/not** \Rightarrow intersection/union/complement

- E.g.: {

*“sum of two dices is even **and** greater than four”*

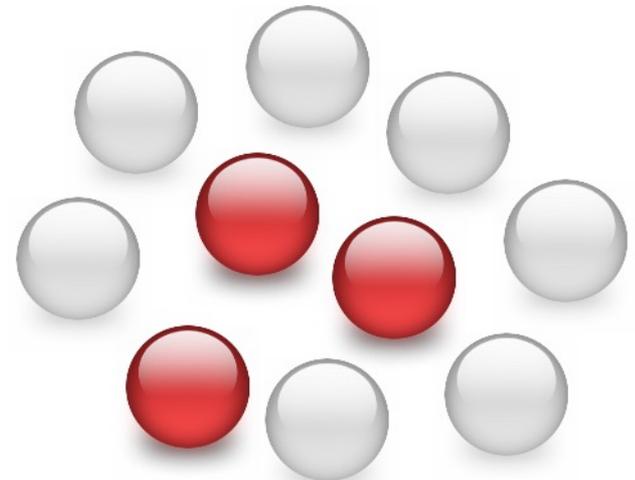


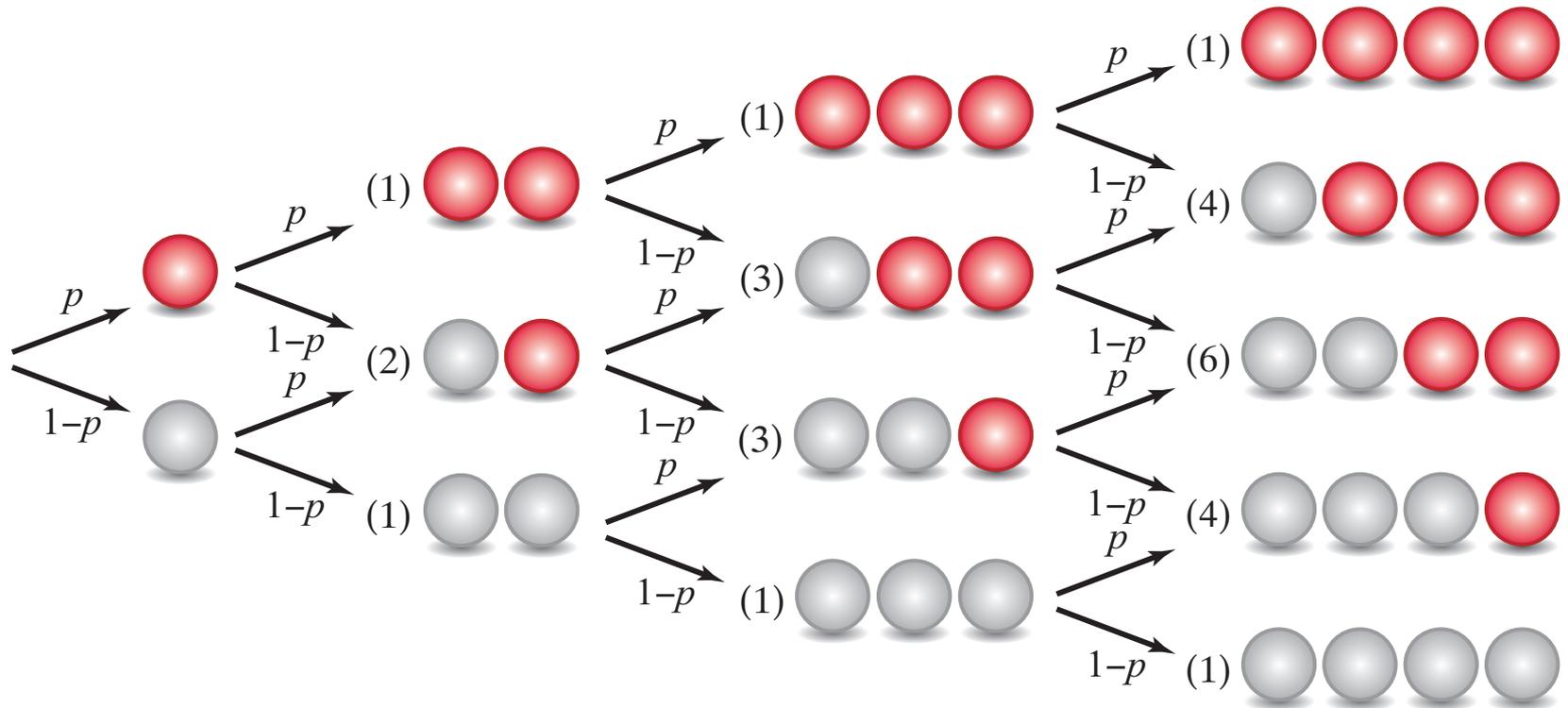
$$\{(d_1, d_2): \text{mod}(d_1 + d_2, 2) = 0\} \cap \{(d_1, d_2): d_1 + d_2 > 4\}$$

- Probability to extract n red balls over N trials, given the fraction p of red balls in a basket
- Each trial is called Bernoulli process

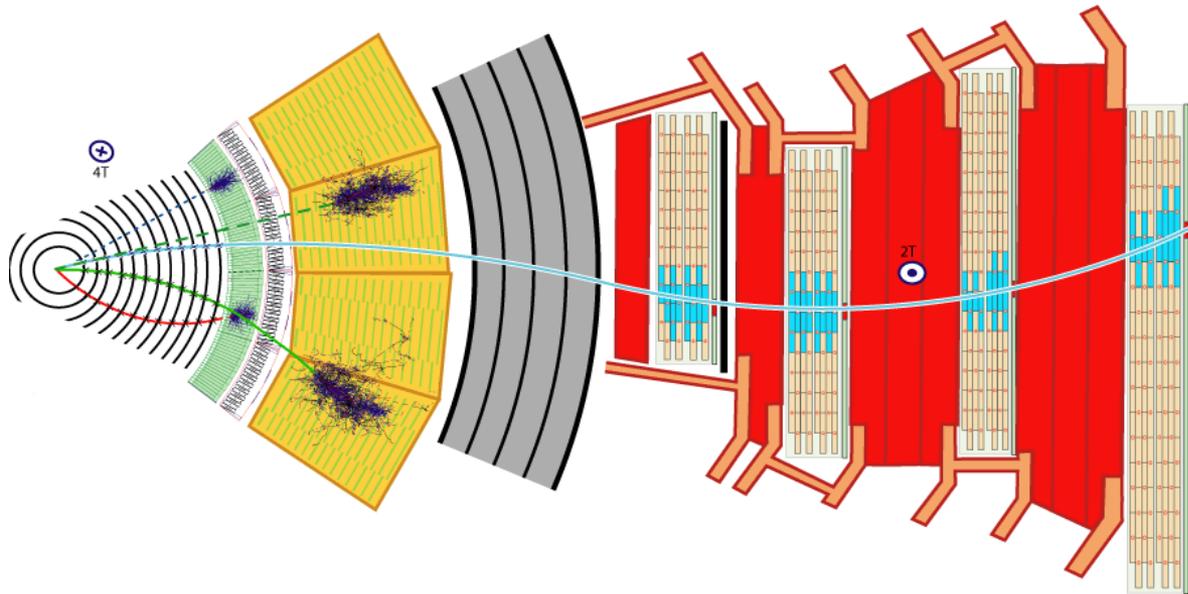
- Red:  $p = 3/10$
- White:  $1 - p = 7/10$

$$p = 3/10$$

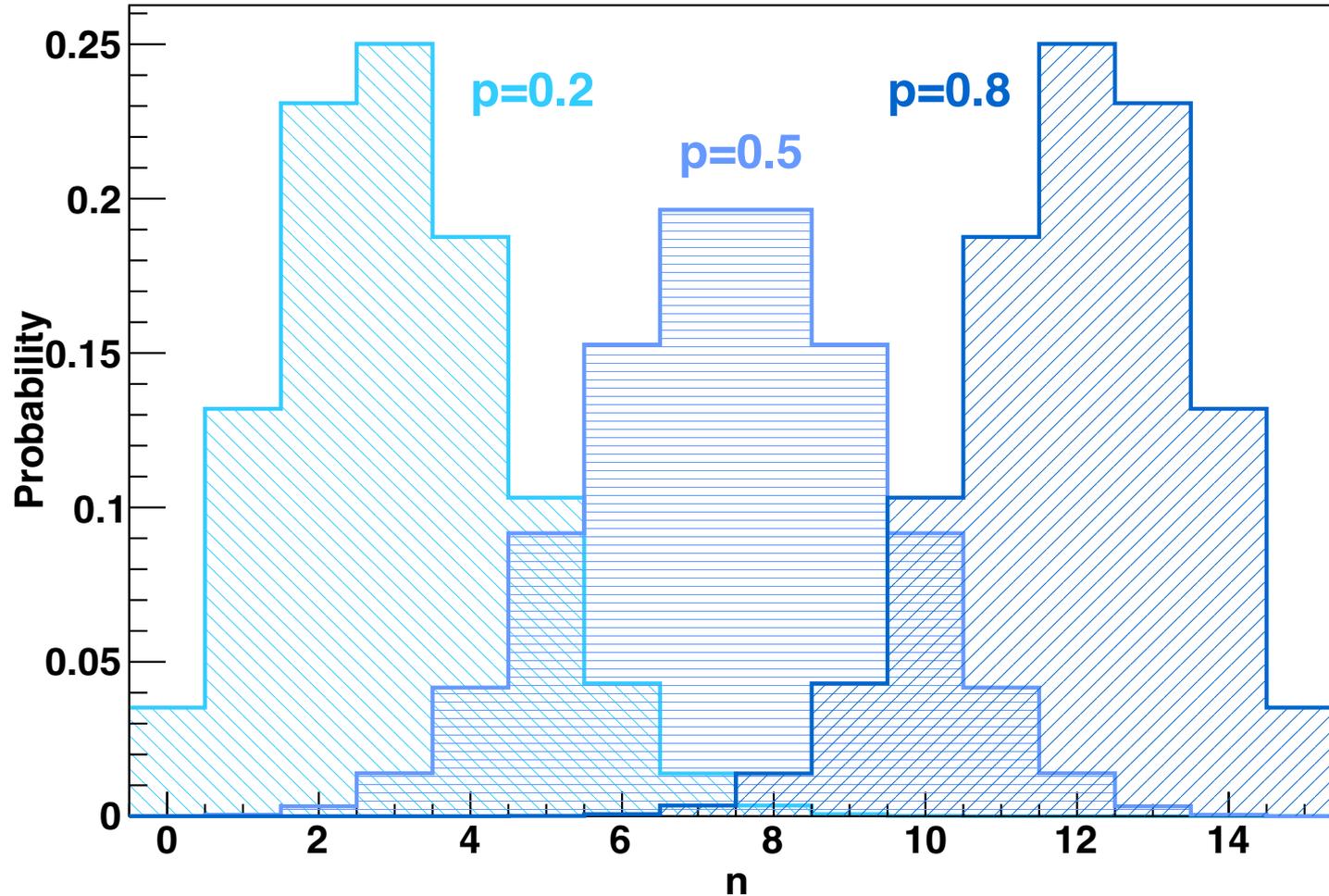




- Typical application in physics:
detector **efficiency** ($\varepsilon = p$)

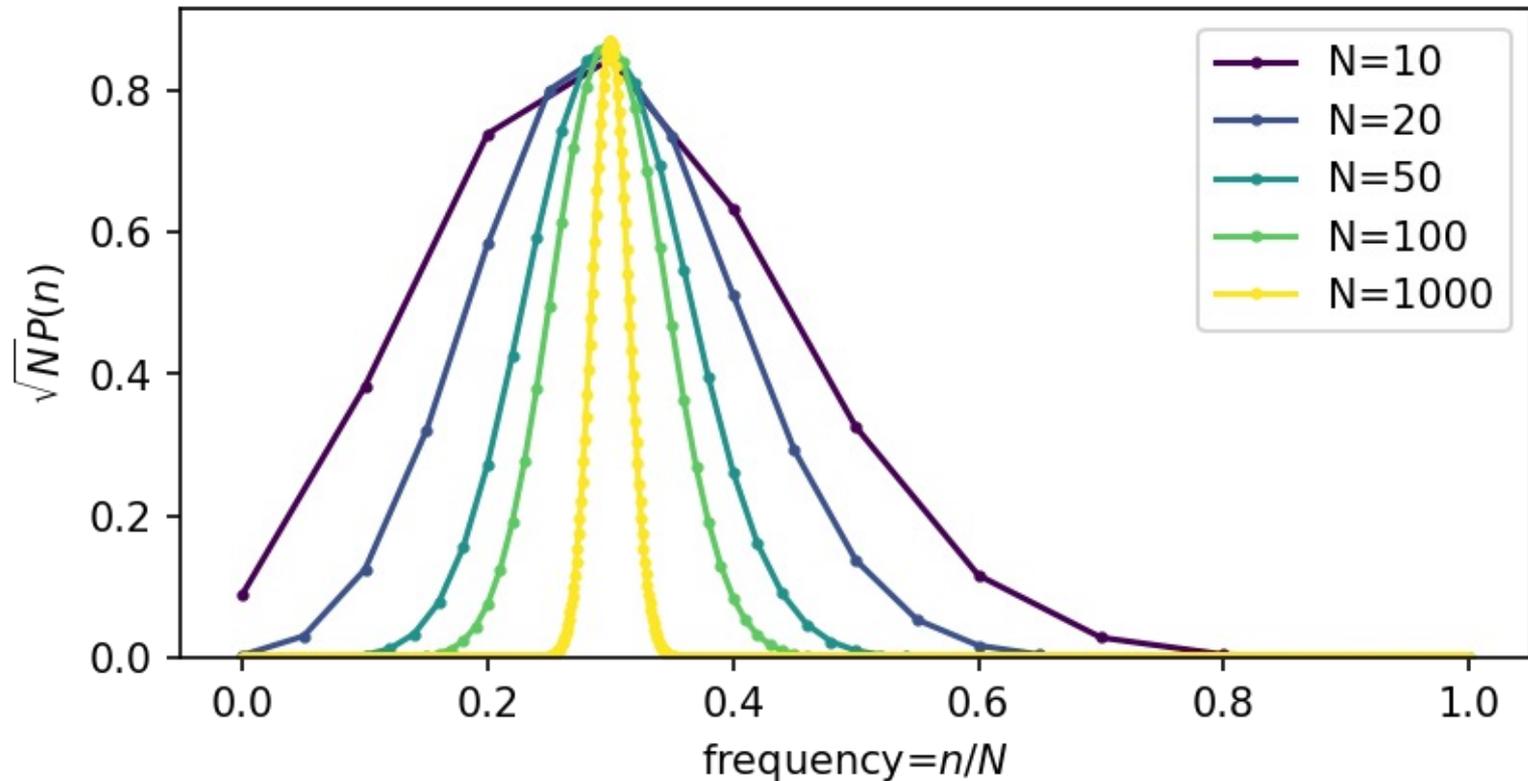


$$P(n; N, p) = \frac{N!}{n!(N - n)!} p^n (1 - p)^{N - n}$$



- $\forall \varepsilon \quad \lim_{N \rightarrow \infty} P \left(\left| \frac{n}{N} - p \right| < \varepsilon \right) = 1$

Binomial distribution

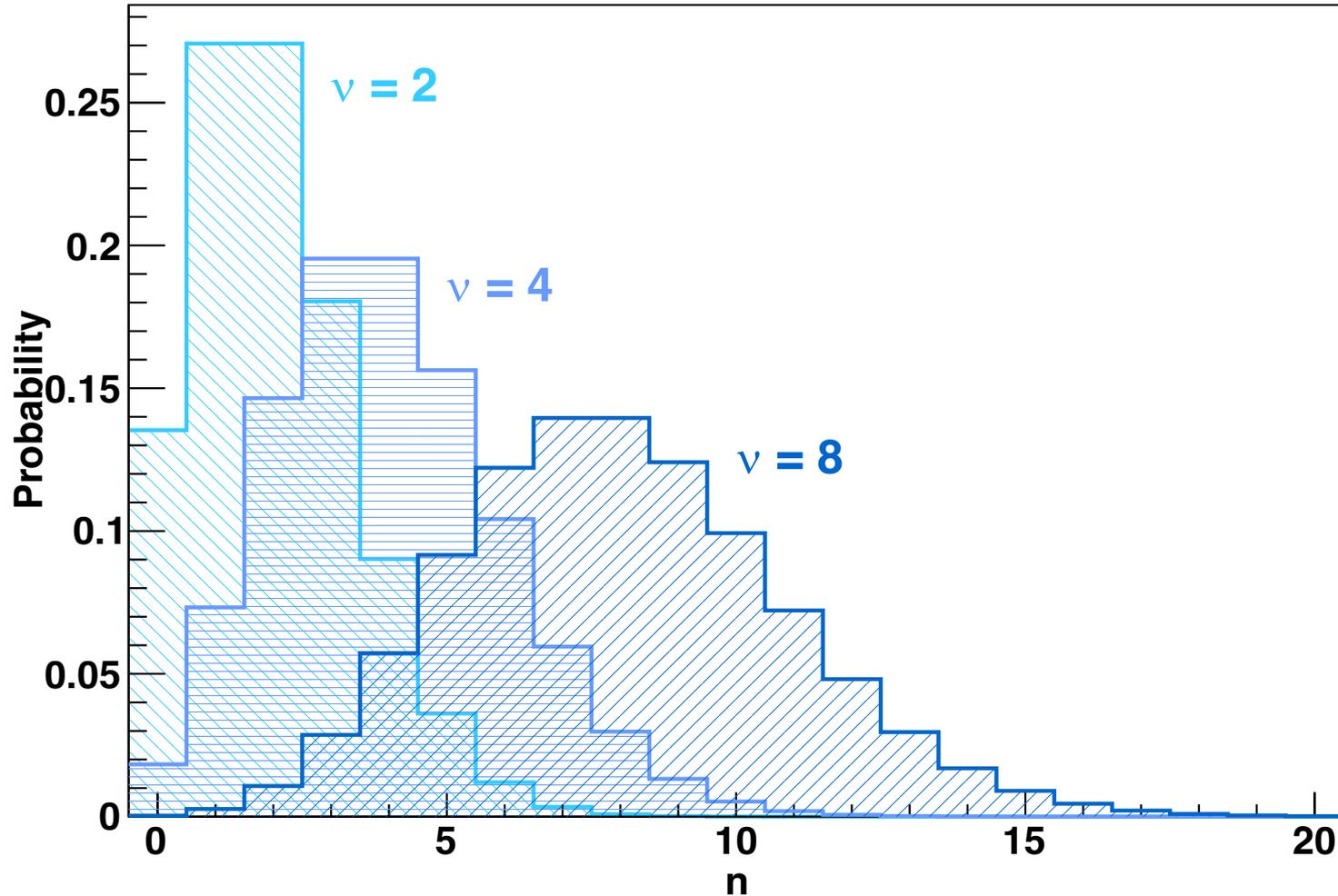


- Limit of Binomial for $N \rightarrow \infty$, $\nu = Np = \text{const.}$
- Distribution of the number of occurrences of random event **uniformly distributed** in a measurement range whose **rate** is known
 - E.g.: number of **rain drops** in a given area and in a given time interval, number of **cosmic rays** crossing a detector in a given time interval



*Siméon-Denis Poisson
(1781-1840)*

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



- Axiomatic probability definition
 - Terminology: Ω = sample space, F = event space, P = probability measure
 - Let $(\Omega, F \subseteq 2^\Omega, P)$ be a measure space that satisfy:

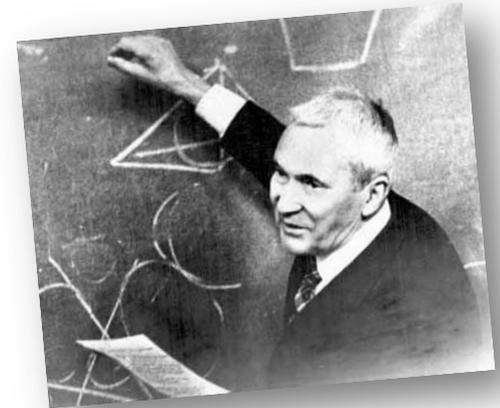
$$- 1 \quad P(E) \geq 0 \quad \forall E \in F$$

$$- 2 \quad P(\Omega) = 1 \quad (\text{normalization})$$

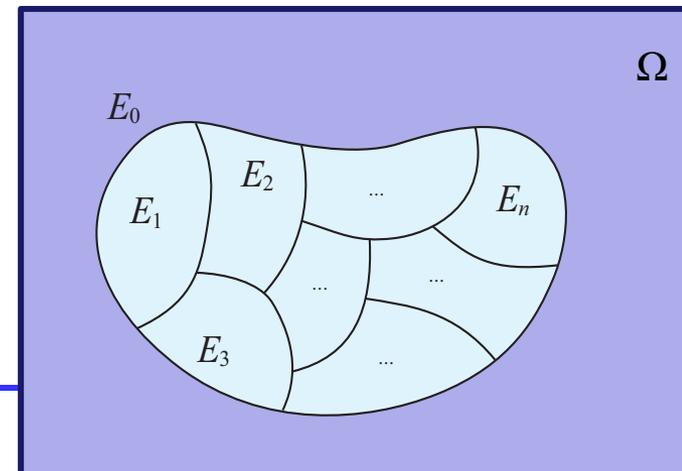
$$- 3 \quad \forall (E_1, \dots, E_n) \in F^n : E_i \cap E_j = \emptyset$$

$$P\left(\bigcup_{i=1, \dots, n} E_i\right) = \sum_{i=1, \dots, n} P(E_i)$$

- The same formalism applies to either frequentist and Bayesian probability



Andrej Nikolaevič Kolmogorov
(1903-1987)



- Expresses **one's degree of belief** that a claim is true
 - How strong? How much would you bet?
 - Applicable to all unknown events/claims, not only repeatable experiments
 - Each individual may have a different opinion/prejudice
- Quantitative rules exist about how subjective probability should be **modified after learning** about some observation/evidence
 - Consistent with **Bayes theorem** (→ will be introduced in next slides)
 - **Prior probability** → **Posterior probability** (following observation)
 - The more information we receive, the more Bayesian probability is insensitive on prior subjective prejudice (unless in pathological cases...)



- Given a **discrete random variable**, we can assign a probability to each individual value:

$$P(x) = P(\{x\})$$

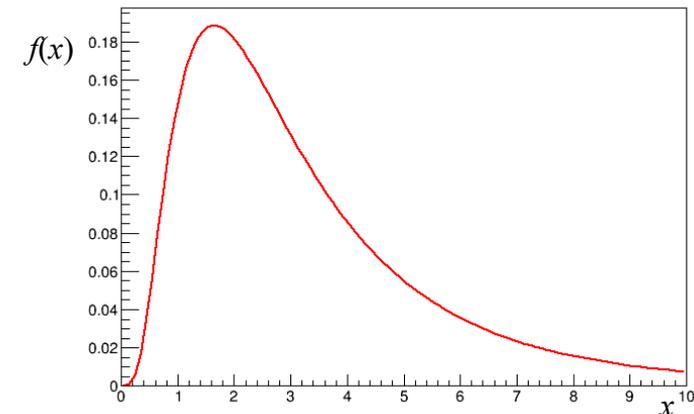
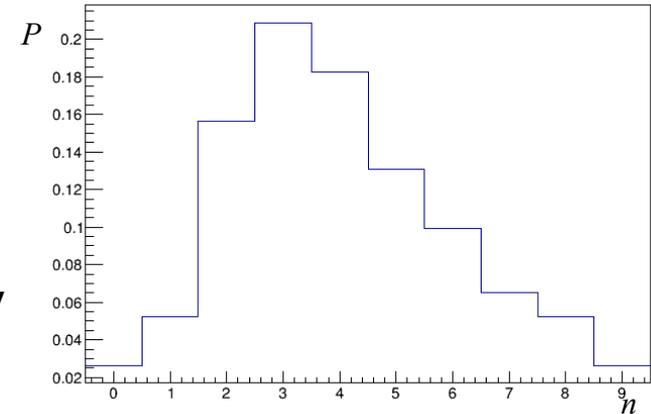
- In case of a **continuous variable**, the probability assigned to an individual value **may be zero**
- A **probability density** better quantifies the probability content (unlike $P(\{x\}) = 0$!):

$$\frac{dP(x)}{dx} = f(x)$$

- Discrete and continuous distributions can be combined using Dirac's delta functions.
- E.g.:

$$\frac{dP}{dx} = \frac{1}{2}\delta(x) + \frac{1}{2}f(x)$$

50% prob. to have zero ($P(\{0\}) = 0.5$), 50% distributed according to $f(x)$

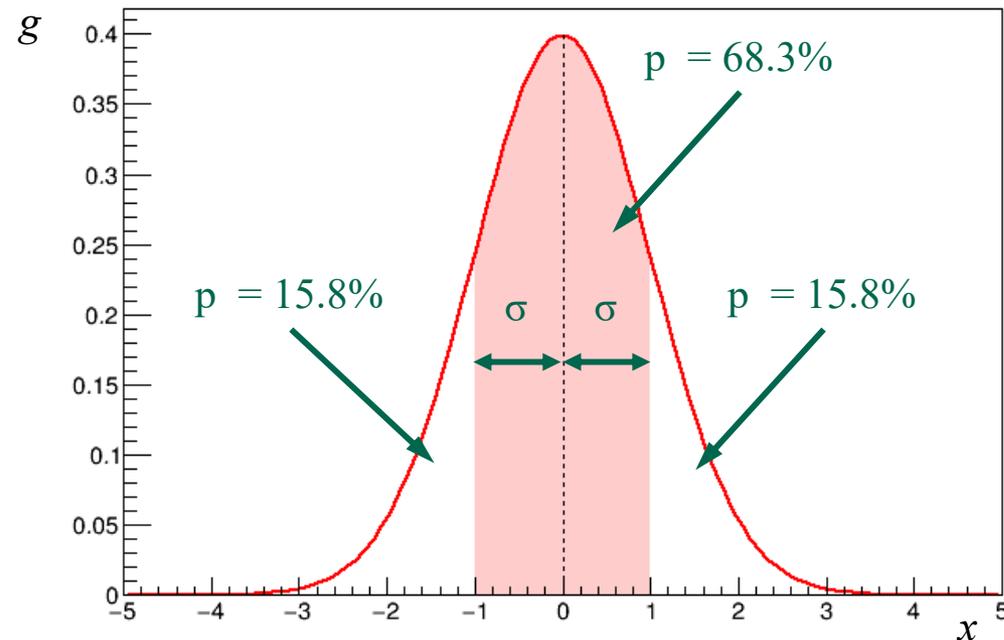


Gaussian distribution

- Many random variables in real experiments follow a **Gaussian distribution**
- **Central limit theorem**: approximate sum of multiple random contributions, regardless of the individual distributions
- Frequently used to model **detector resolution**

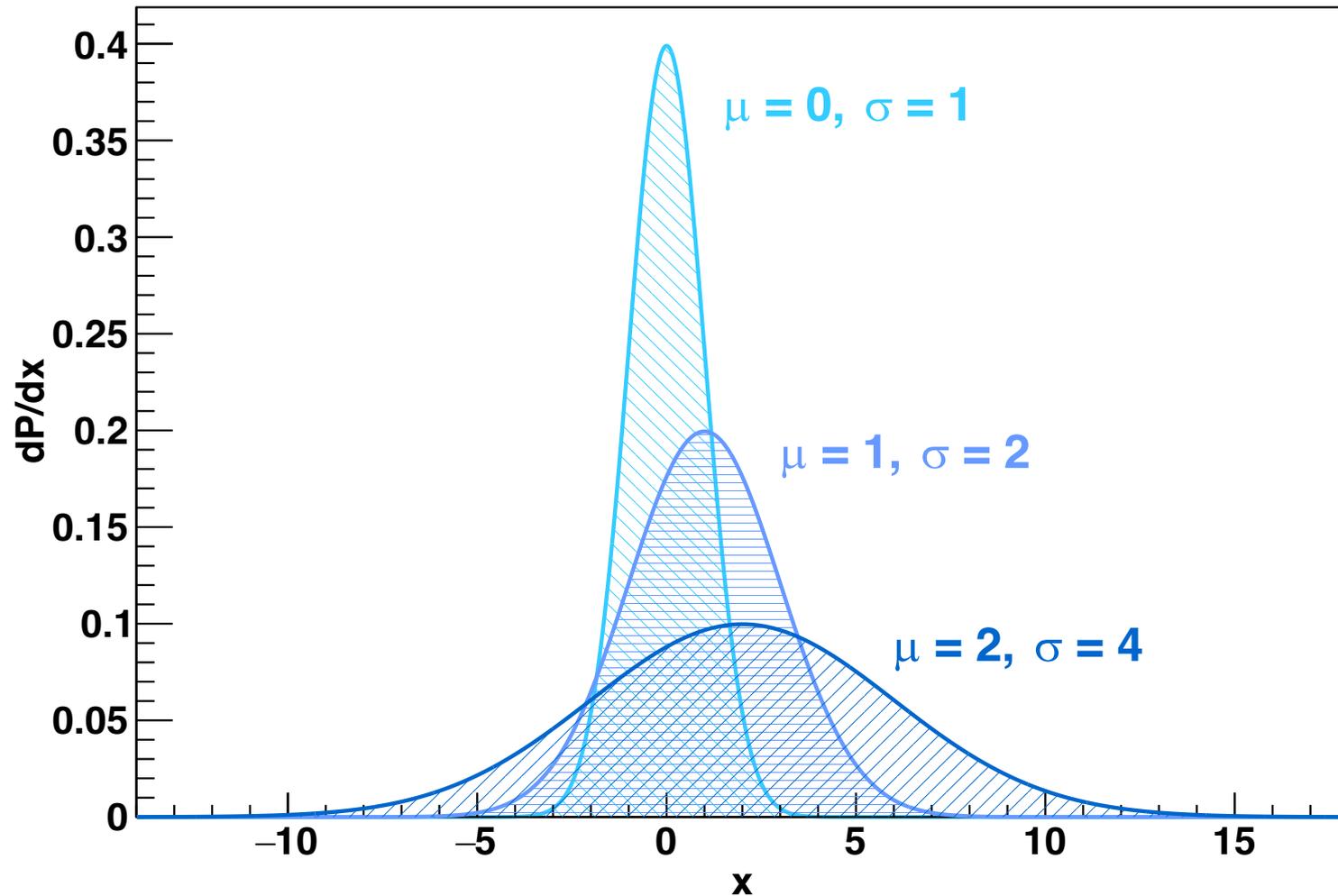


Carl Friedrich Gauss
(1777-1855)



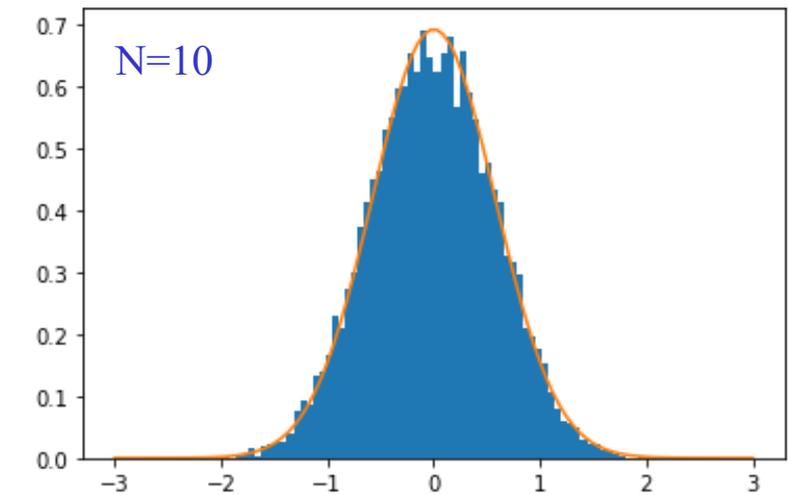
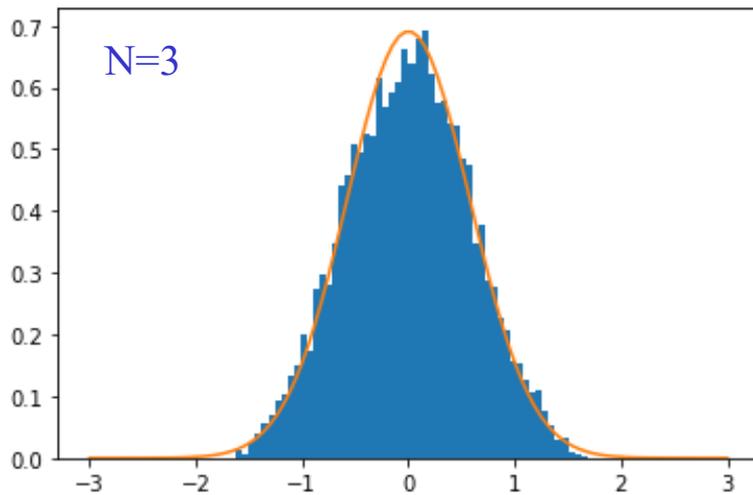
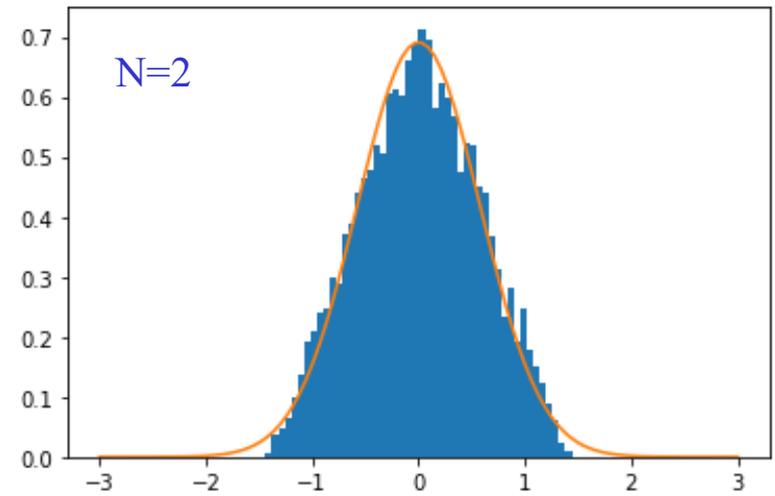
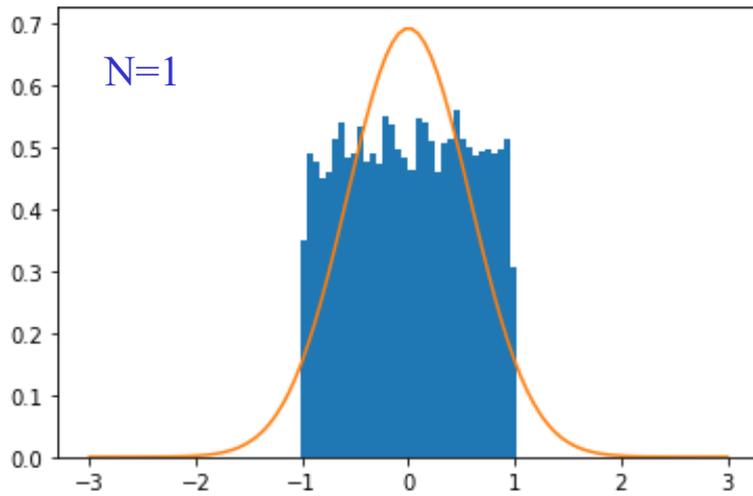
$$g(x; \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

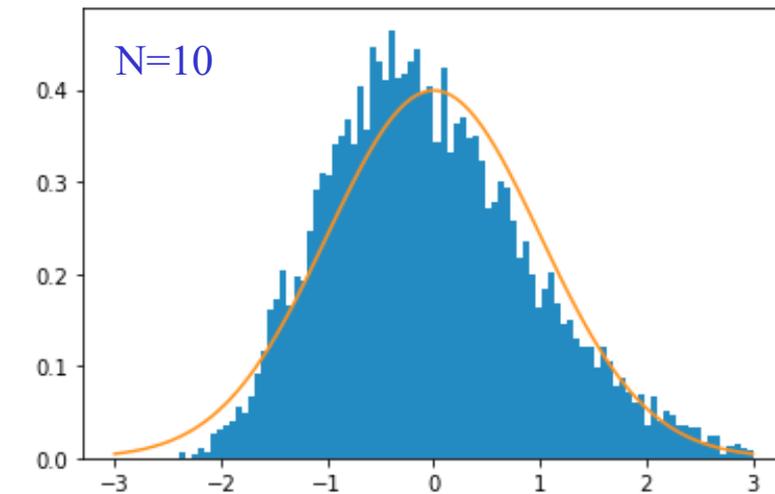
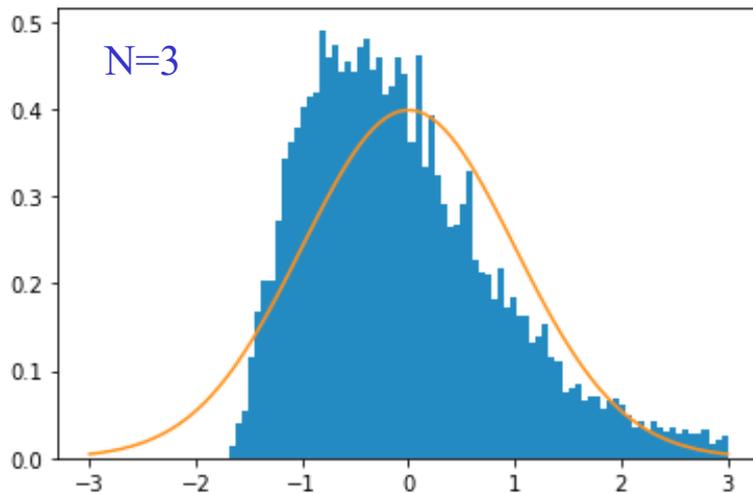
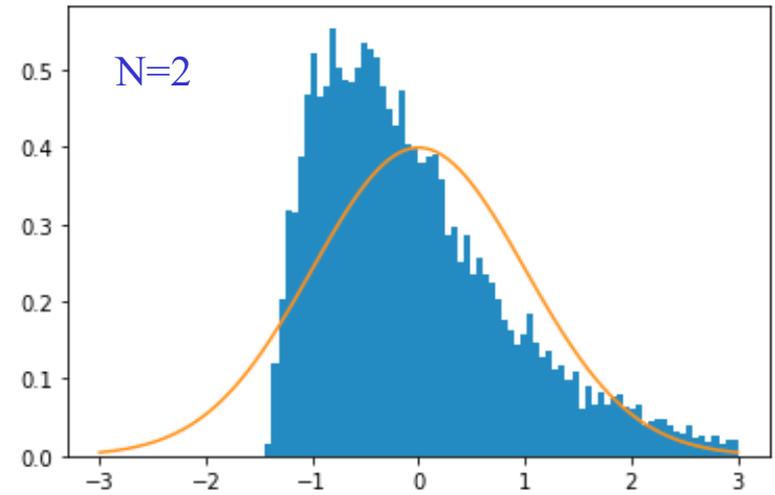
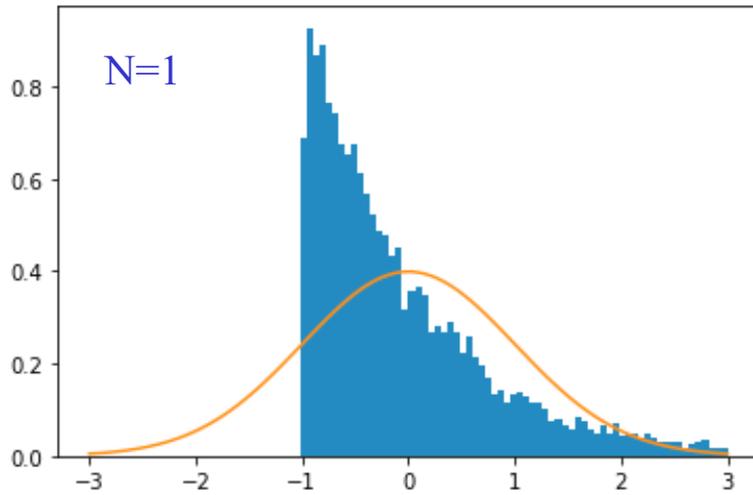
$n\sigma$	Prob.
1	0.683
2	0.954
3	0.997
4	$1 - 6.3 \times 10^{-5}$
5	$1 - 5.7 \times 10^{-7}$





- The distribution of the sum or average of N random variables all having the same distribution with finite variance tends to a Gaussian for large N





- Distributions of measured quantities in data:
 - are predicted by a theory model,
 - depend on some theory parameters,
 - e.g.: particle mass, cross section, etc.
- Given our data sample, we want to:
 - measure theory parameters,
 - e.g.: $m_t = 173.49 \pm 1.07 \text{ GeV}$
 - answer questions about the nature of data
 - Is there a Higgs boson? → Yes! (strong evidence? Quantify!)
 - Is there a Dark Matter? → No evidence, so far...
 - If not, what is the range of theory parameters compatible with the observed data? What parameter range can we exclude?

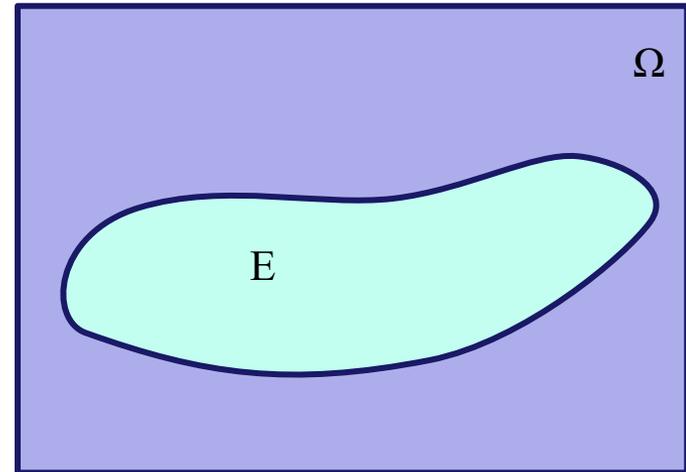
$$\begin{aligned}
 \mathcal{L} = & -\frac{1}{4} F_{\mu\nu} F^{\mu\nu} \\
 & + i\bar{\psi} \not{D} \psi + \text{h.c.} \\
 & + \chi_i y_{ij} \chi_j \phi + \text{h.c.} \\
 & + |D_\mu \phi|^2 - V(\phi)
 \end{aligned}$$

- In more dimensions (n random variables), PDF can be defined as:

$$\frac{d^n P}{dx_1 \cdots dx_n} = f(x_1, \cdots, x_n)$$

- The probability associated to an event E is obtained by integrating the PDF over the corresponding set in the sample space

$$P(E) = \int_E f(x_1, \cdots, x_n) d^n x$$



- Given a random variable x with distribution $f(x)$ we can define:

- Mean or expected value:

$$\left\{ \begin{array}{l} \mathbb{E}[x] = \langle x \rangle = \int x f(x) dx \\ \mathbb{E}[g(x)] = \langle g(x) \rangle = \int g(x) f(x) dx \end{array} \right.$$

- Variance:

$$\text{Var}[x] = \langle (x - \langle x \rangle)^2 \rangle = \overbrace{\langle x^2 \rangle}^{\text{r.m.s.: root mean square}} - \langle x \rangle^2$$

- Standard deviation:

$$\sigma_x = \sqrt{\text{Var}[x]} = \sqrt{\langle x^2 \rangle - \langle x \rangle^2}$$

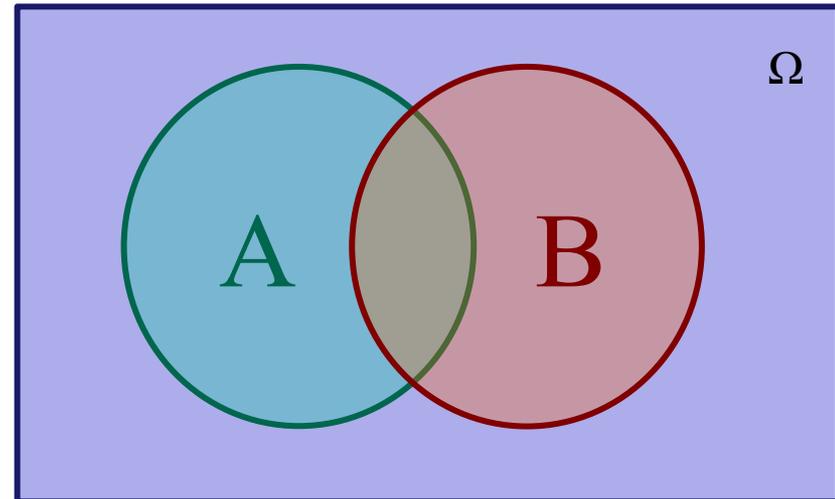
- Covariance and correlation coefficient of two variables x and y :

$$\text{cov}(x, y) = \langle (x - \langle x \rangle)(y - \langle y \rangle) \rangle$$

$$\rho_{xy} = \frac{\text{cov}(x, y)}{\sigma_x \sigma_y}$$

Note: integration extends over x and y in two dimensions when computing an average value!

- Probability of A , given B : $P(A | B)$, i.e.: probability that an event known to belong to set B also belongs to set A :
 - $P(A | B) = P(A \cap B) / P(B)$
 - Notice that:
 $P(A | \Omega) = P(A \cap \Omega) / P(\Omega)$
- Event A is said to be **independent** of B if the probability of A given B is equal to the probability of A :
 - $P(A | B) = P(A)$
- If A is independent of B then $P(A \cap B) = P(A) P(B)$
- \rightarrow If A is independent on B , B is independent on A



Independent variables



$$\frac{d^2 P}{dx dy} = f(x, y)$$

- 1D projections:
(marginal distributions)

$$\left\{ \begin{array}{l} f_x(x) = \int f(x, y) dy \\ f_y(y) = \int f(x, y) dx \end{array} \right.$$

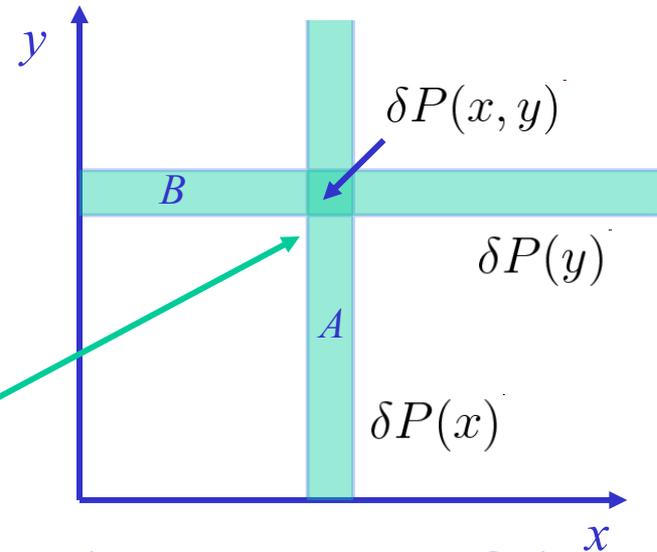
- x and y are independent if:

$$f(x, y) = f_x(x) f_y(y)$$

- We saw that A and B are independent events if:

$$P(A \cap B) = P(A)P(B)$$

- Where $A = \{x' : x < x' < x + \delta x\}$, $B = \{y' : y < y' < y + \delta y\}$



The Bayes theorem



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

$$P(A|B)P(B) = P(B|A)P(A)$$

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



Thomas Bayes (1702-1761)

- $P(A)$ = prior probability
- $P(A|B)$ = posterior probability

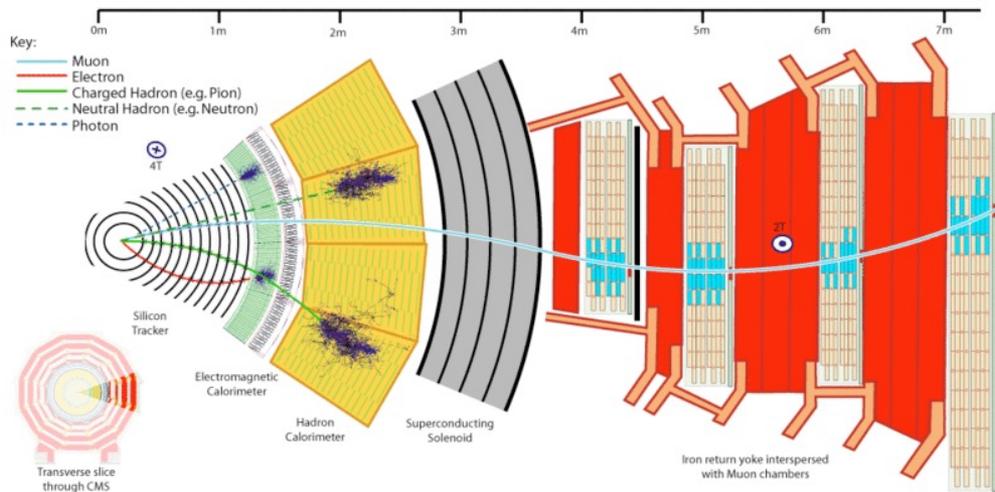


- Bayes theorem allows to determine **probability about hypotheses or claims H** that not related random variables, given an **observation or evidence E** :

$$P(H|E) = \frac{P(E|H)P(H)}{P(E)}$$

- $P(H)$ = **prior probability**
- $P(H|E)$ = **posterior probability**, given E
- The Bayes rule allows to define a **rational way** to modify one's prior belief once some observation is known

- A detector identifies **muons** with high efficiency, $\varepsilon = 95\%$
- A small fraction $\delta = 5\%$ of **pions** are incorrectly identified as muons (“fakes”)
- If a particle is identified as a **muon**, what is the probability it is really a **muon**?
 - The answer also depends on the composition of the sample!
 - i.e.: the fraction of **muons** and **pions** in the overall sample



This example is usually presented as an epidemiology case.

Naïve answers about fake positive probability are often wrong!

- Using Bayes theorem:

$$- P(\mu|+) = P(+|\mu) P(\mu) / P(+)$$

Law of total probability

- Where our inputs are:

$$- P(+|\mu) = \varepsilon = 0.95, P(+|\pi) = \delta = 0.05$$

+ denotes a positive id

- We can decompose $P(+)$ as:

$$- P(+)= P(+|\mu) P(\mu) + P(+|\pi) P(\pi)$$

normalization term

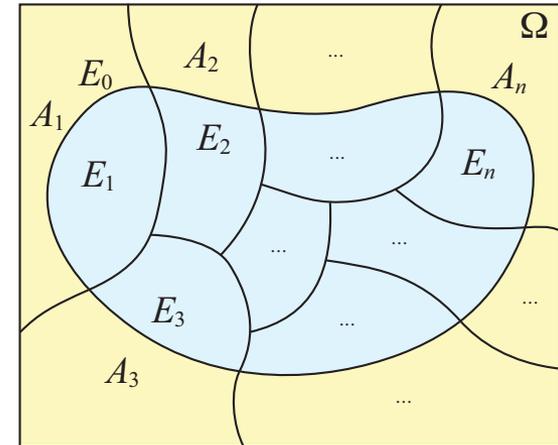
- Putting all together:

$$- P(\mu|+) = \varepsilon P(\mu) / (\varepsilon P(\mu) + \delta P(\pi))$$

- Assume we have a sample made of $P(\mu)=4\%$ muons and $P(\pi)=96\%$ pions, we have:

$$- P(\mu|+) = 0.95 \times 0.04 / (0.95 \times 0.04 + 0.05 \times 0.96) \cong 0.44$$

- Even if the selection efficiency is very high, the low sample purity makes $P(\mu|+)$ lower than 50%.



$$P(E_0) = \sum_{i=1}^n P(E_0|A_i)P(A_i)$$

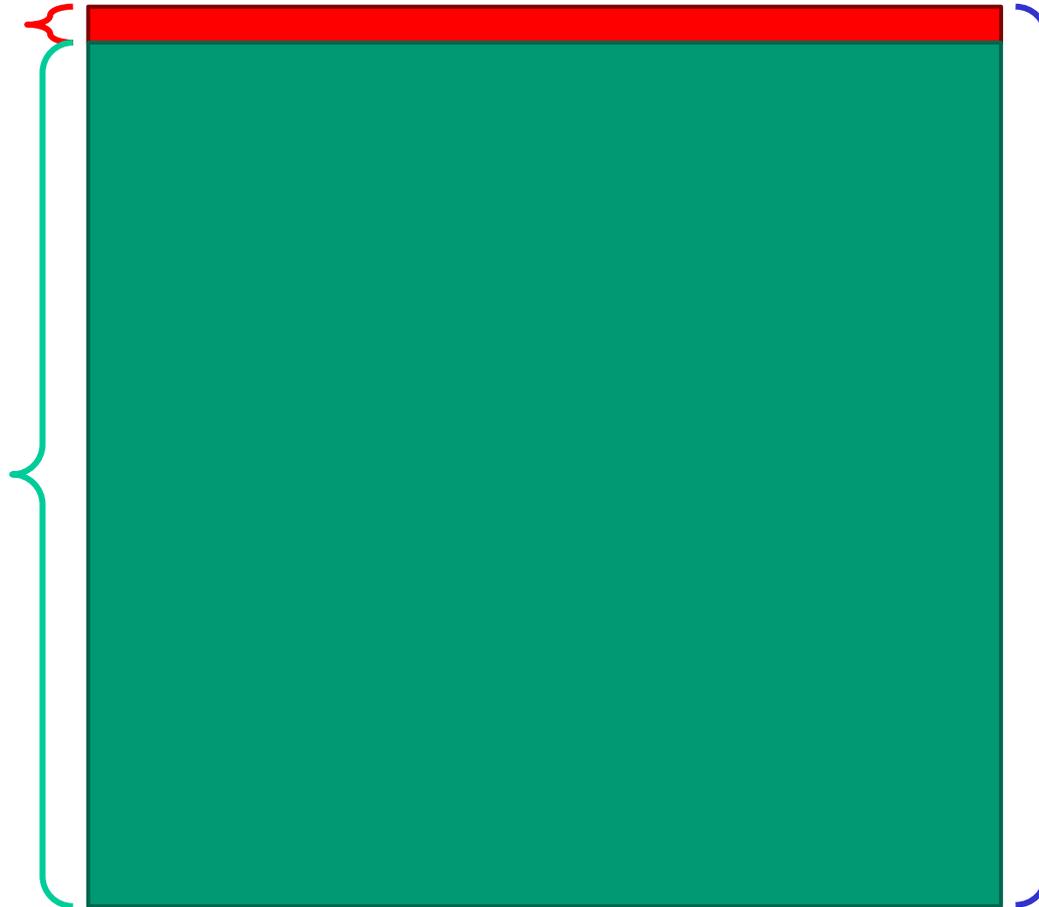
$E_0 = '+'$, $A_i = \mu, \pi$

Before any muon id. information

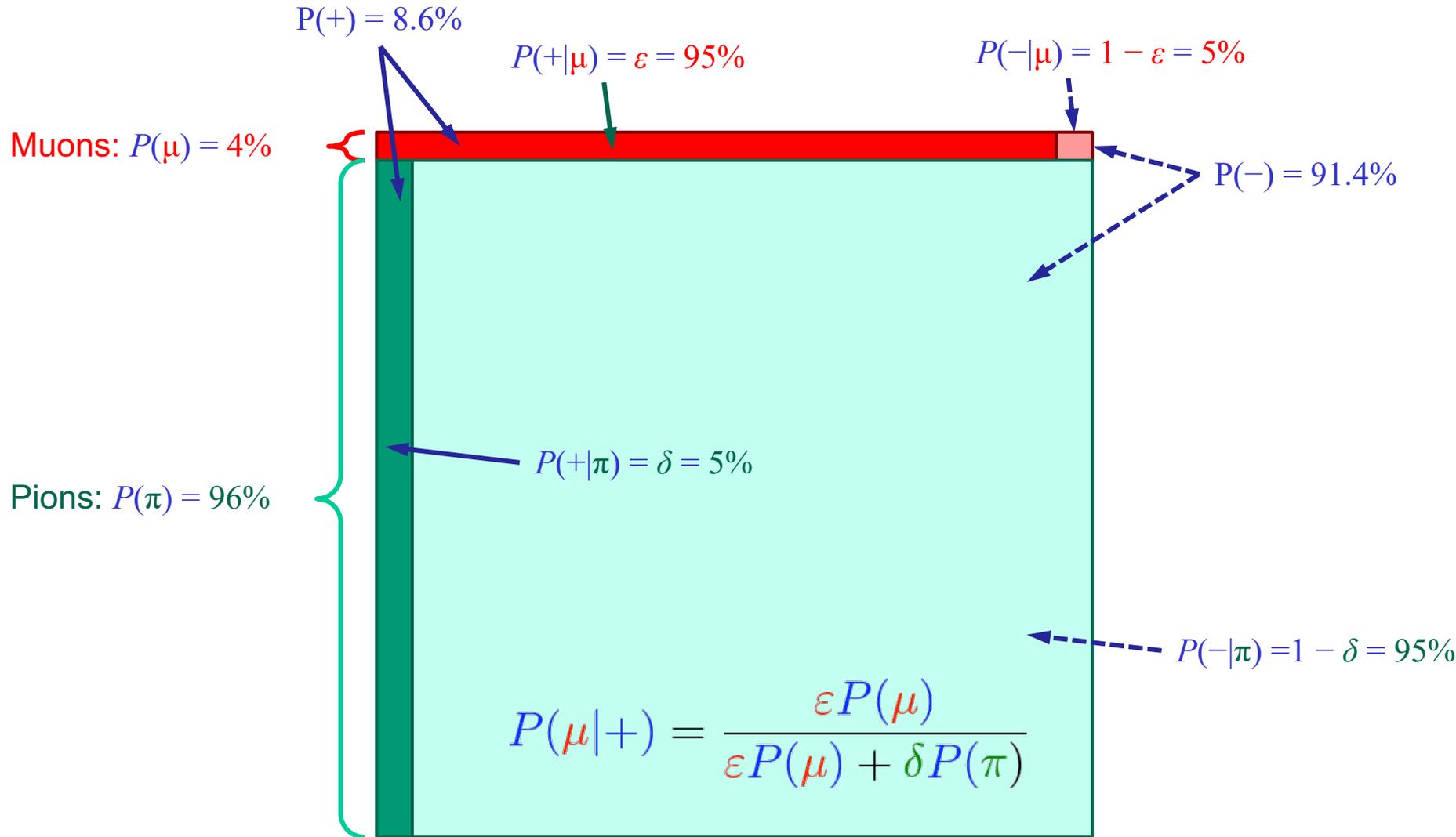


Muons: $P(\mu) = 4\%$

Pions: $P(\pi) = 96\%$



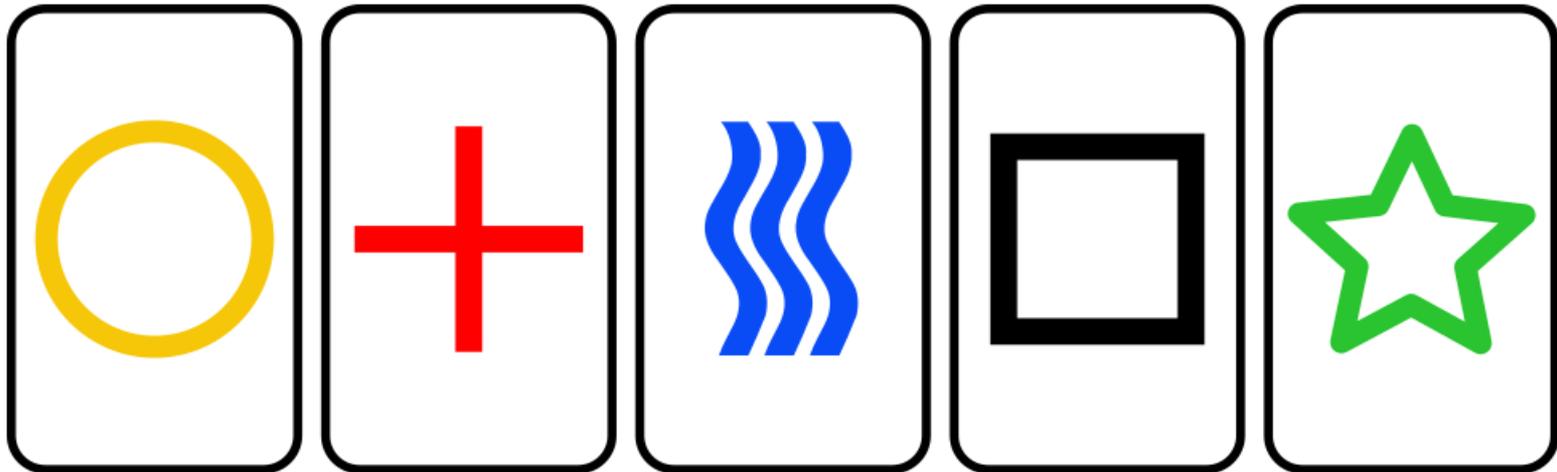
All particles:
 $P(\Omega) = 100\%$

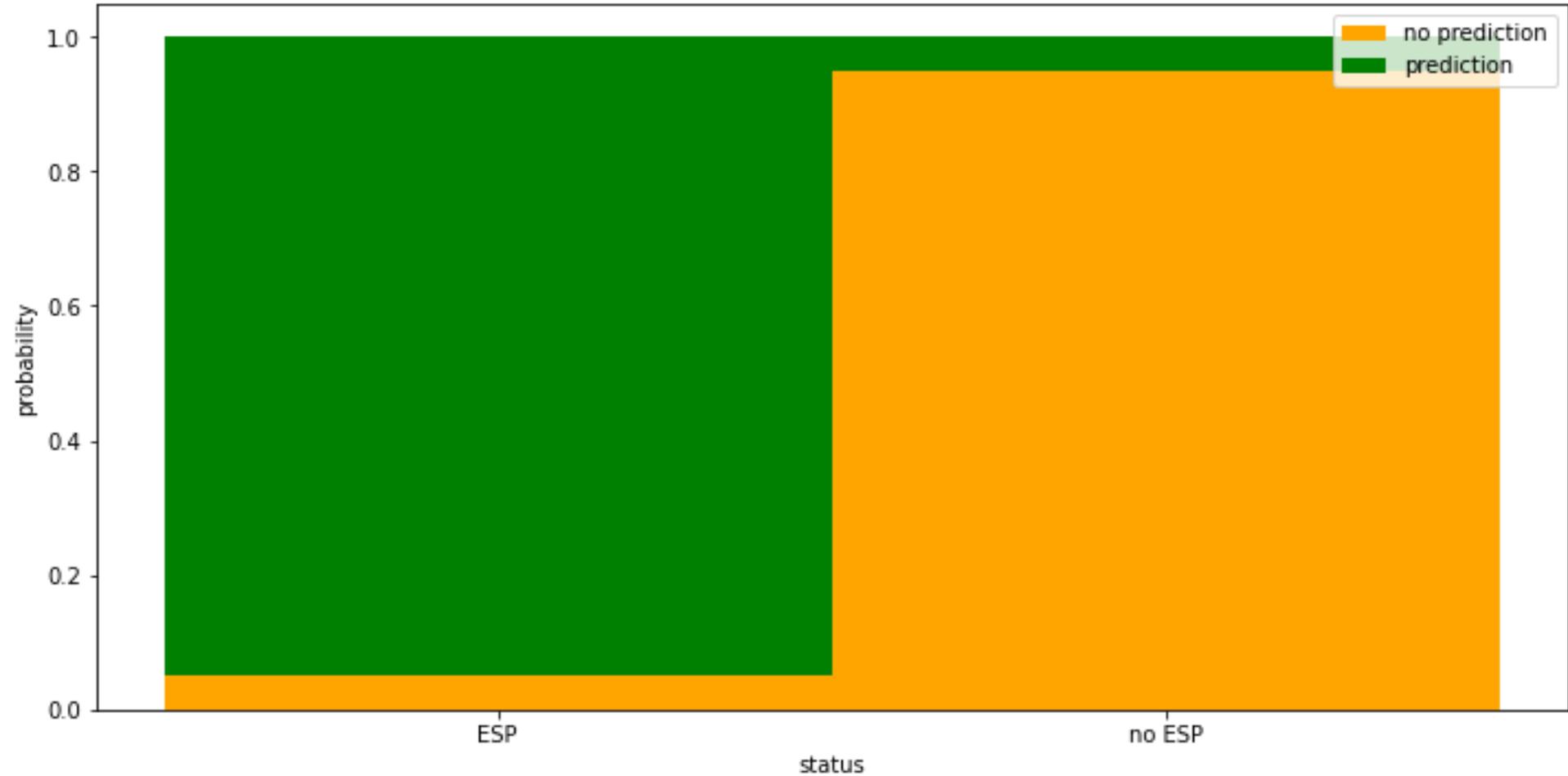


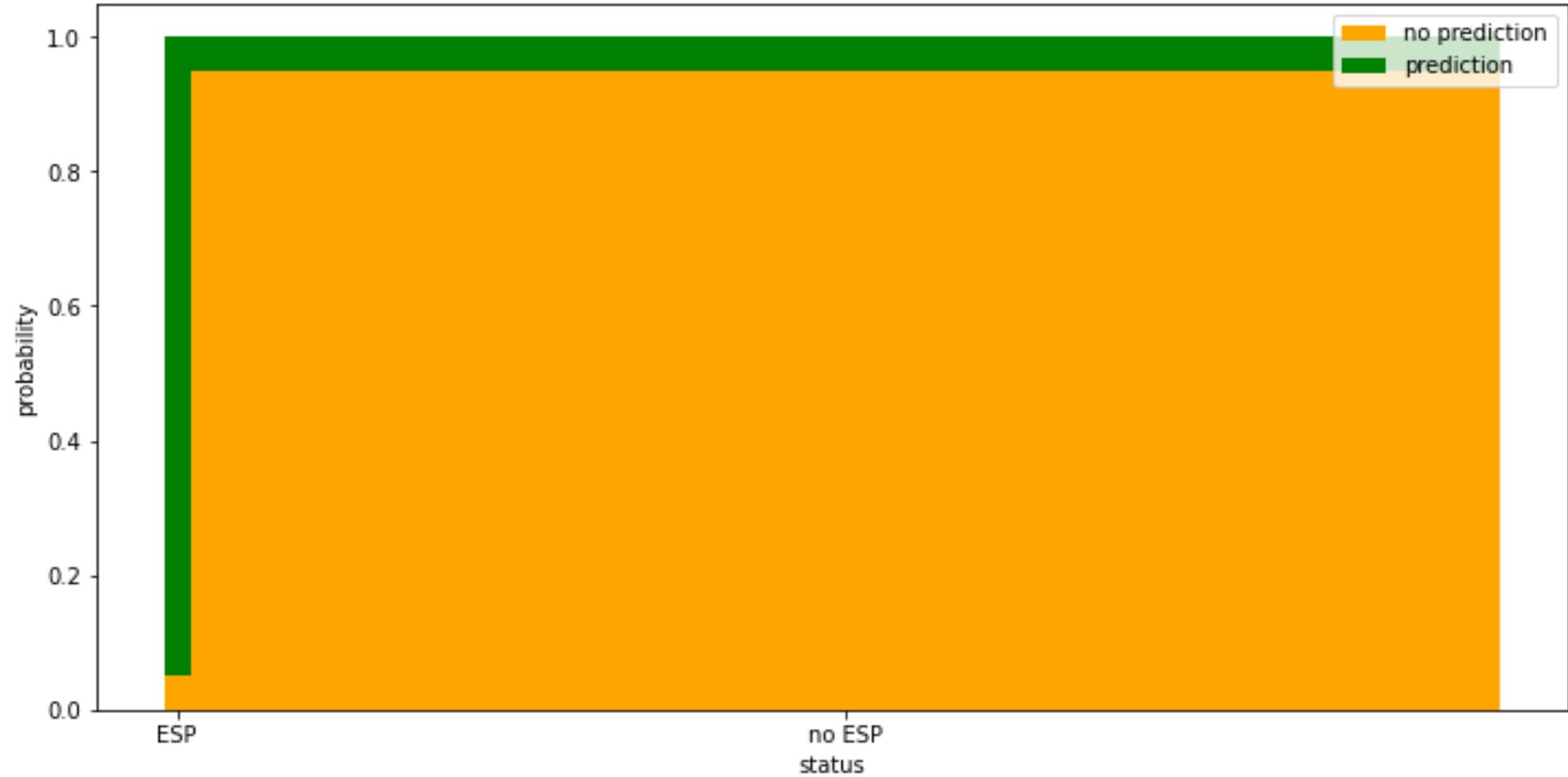
The same approach to unknowns

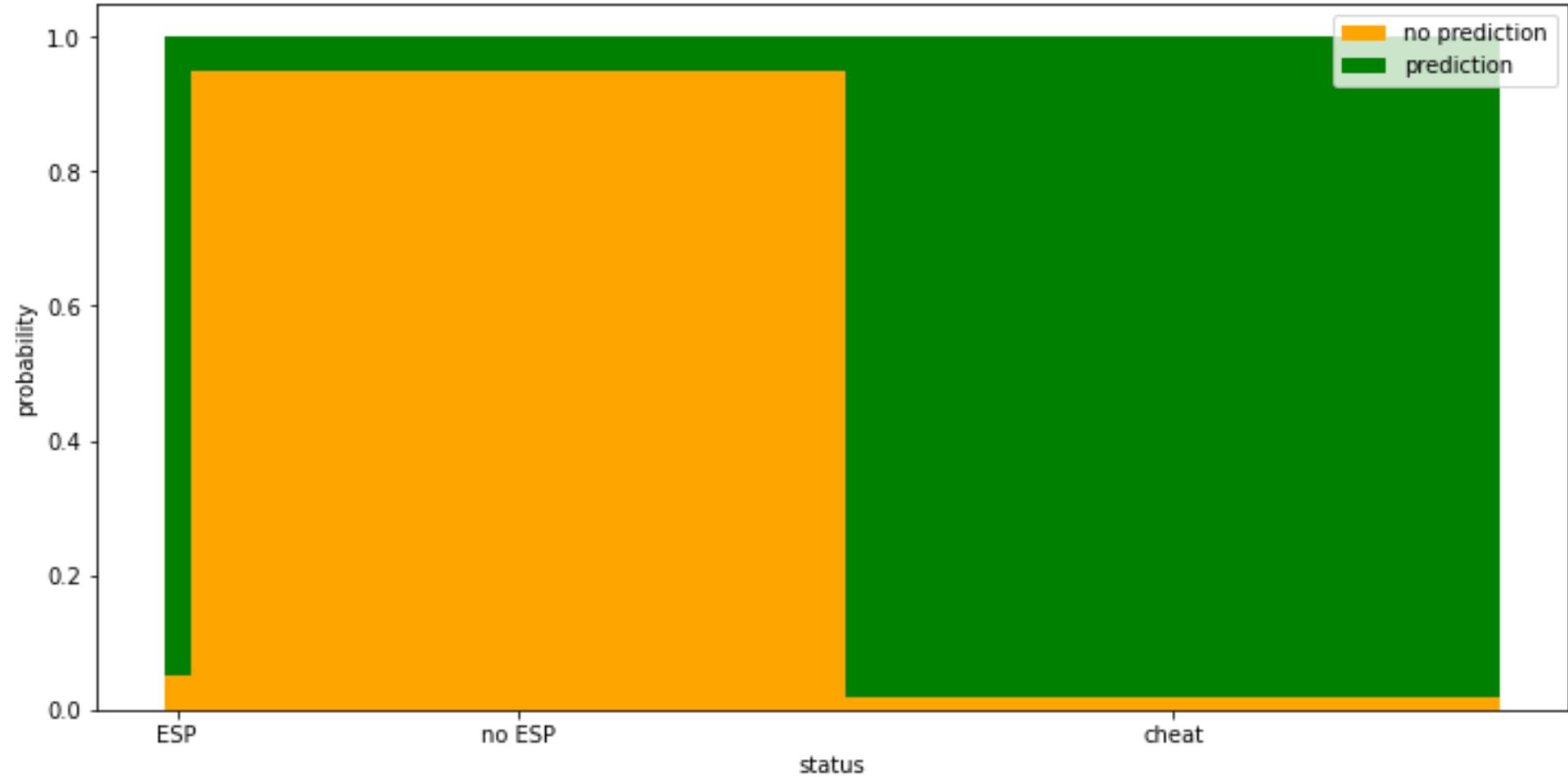


- ESP: extra-sensory perception:
Supernatural (?) prediction of a set of cards











- Note that claiming that all scientific evidences are fakes is the approach of **conspiracy theorists**
- The main difference is that the prior chosen by conspiracy theorists ignore most of the evidence in favour of fantasy inventions
- **Scientific reasoning** is closely related to the Bayesian approach, provided that evidences are not discarded on purpose

- In many cases, the outcome of our experiment can be modeled as a set of random variables x_1, \dots, x_n whose distribution takes into account:
 - **intrinsic sample randomness** (quantum physics is intrinsically random),
 - **detector effects** (resolution, efficiency, ...).
- Theory and detector effects can be described according to some parameters $\theta_1, \dots, \theta_m$, whose values are, in most of the cases, unknown
- The overall PDF, evaluated at our observation x_1, \dots, x_n , is called likelihood function:

$$L = f(x_1, \dots, x_n; \theta_1, \dots, \theta_m)$$

- In case our sample consists of N **independent measurements** (collision events) the likelihood function can be written as:

$$L = \prod_{i=1}^N f(x_1^i, \dots, x_n^i; \theta_1, \dots, \theta_m)$$

- Given a set of measurements x_1, \dots, x_n , Bayesian posterior PDF of the unknown parameters $\theta_1, \dots, \theta_m$ can be determined as:

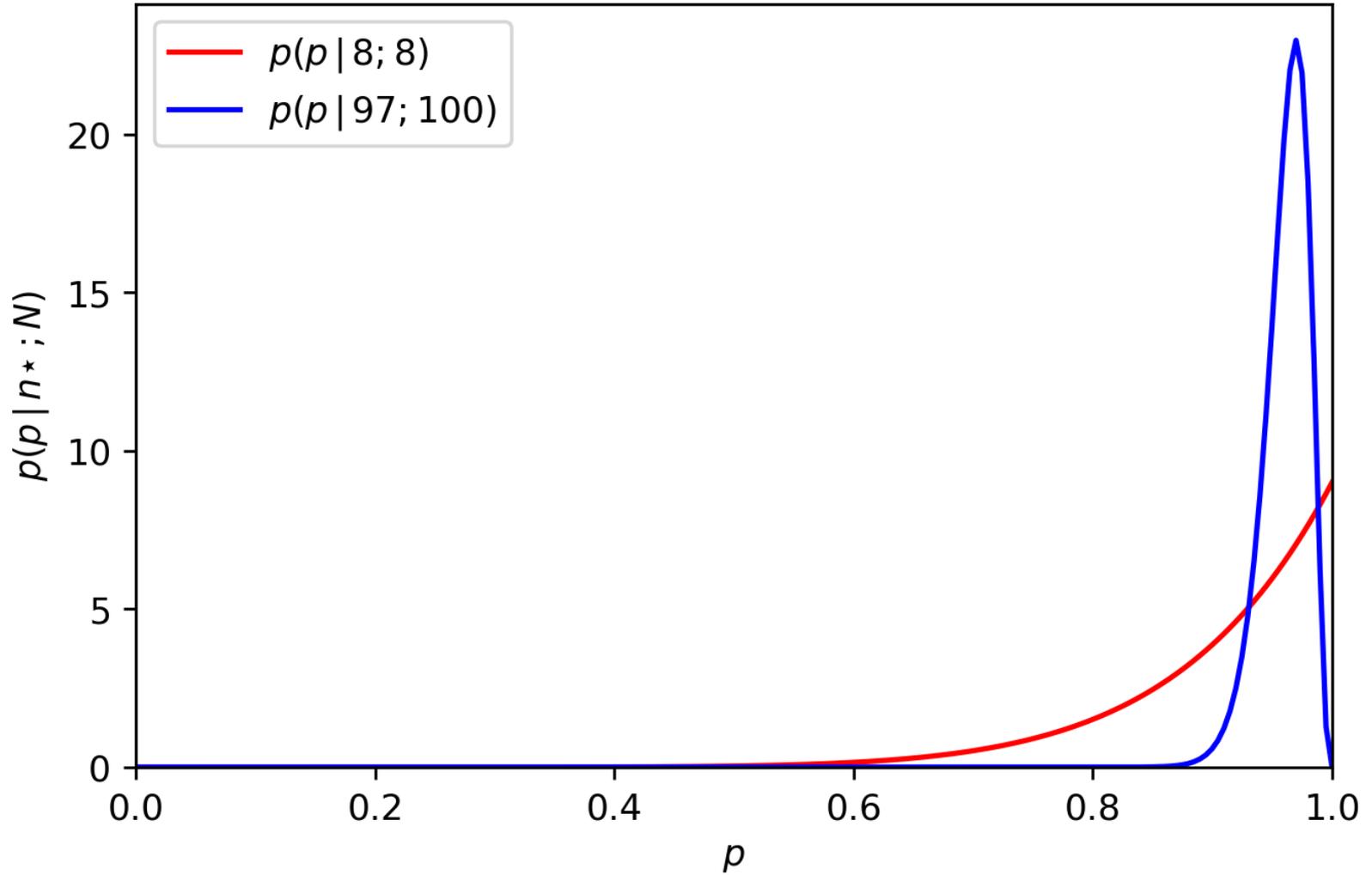
$$P(\theta_1, \dots, \theta_m | x_1, \dots, x_n) = \frac{L(x_1, \dots, x_n; \theta_1, \dots, \theta_m) \pi(\theta_1, \dots, \theta_m)}{\int L(x_1, \dots, x_n; \theta_1, \dots, \theta_m) \pi(\theta_1, \dots, \theta_m) d^m \theta}$$

- Where $\pi(\theta_1, \dots, \theta_m)$ is the subjective prior probability
- The denominator $\int L(x, \theta) \pi(\theta) d^m \theta$ is a normalization factor
- The observation of x_1, \dots, x_n modifies the prior knowledge of the unknown parameters $\theta_1, \dots, \theta_m$
- If $\pi(\theta_1, \dots, \theta_m)$ is sufficiently smooth and L is sharply peaked around the true values $\theta_1, \dots, \theta_m$, the resulting posterior will not be strongly dependent on the prior's choice



- The number of positive/negative feedbacks of an online seller can be assumed to follow a binomial distribution
- If a seller has 100% of positive feedbacks (8/8) and another one has 97% (97/100), which one is more reliable?
- The posterior is a special case of the so-called Beta distribution, prop to: $p^n(1-p)^{n-N}$:

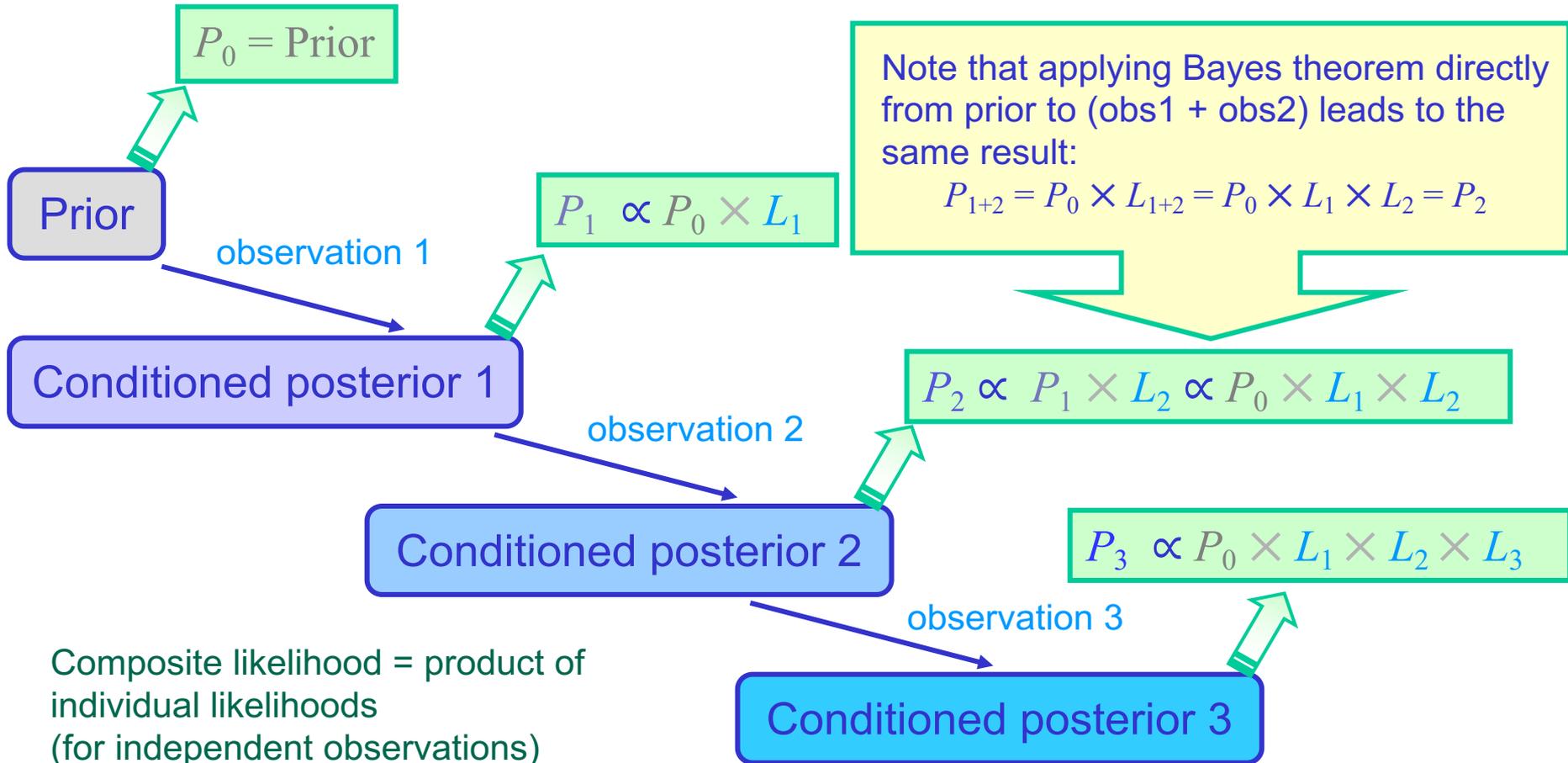
$$p(x; \alpha, \beta) = \frac{x^{\alpha-1}(1-x)^{\beta-1}}{\Gamma(\alpha)\Gamma(\beta)/\Gamma(\alpha+\beta)}$$



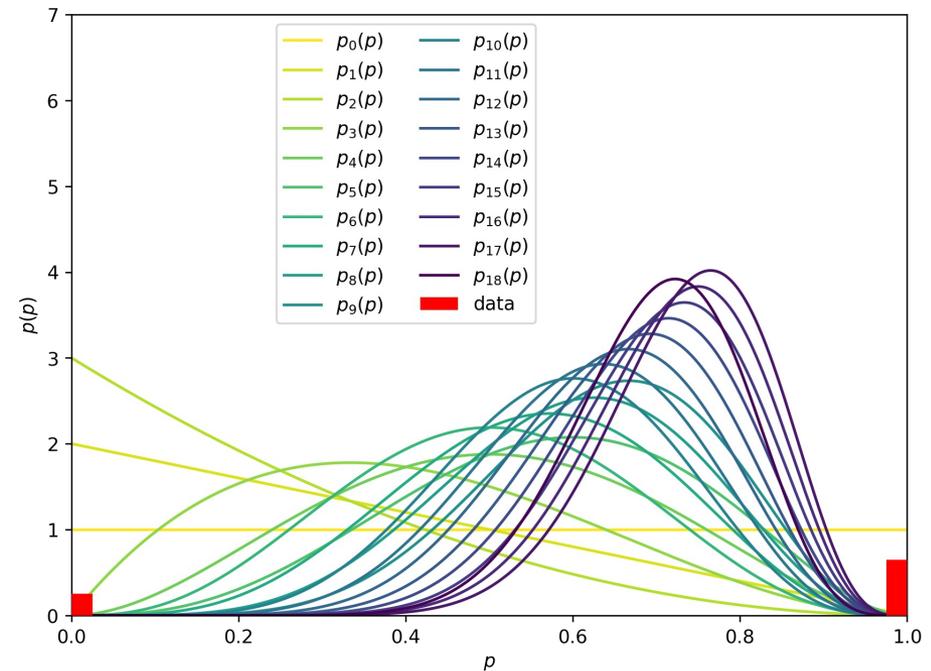
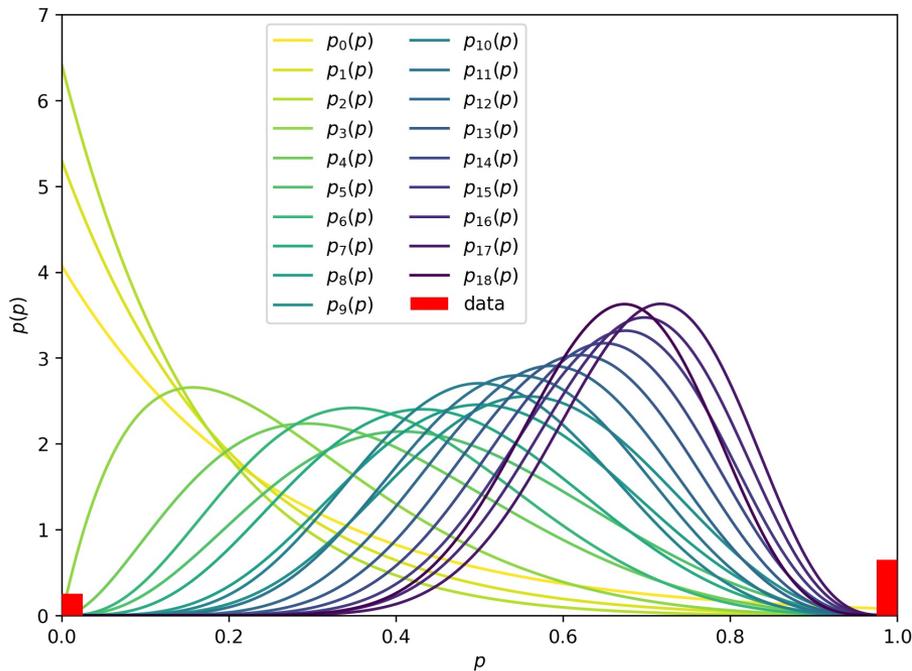
Repeated use of Bayes theorem



- Bayes theorem can be applied sequentially for repeated independent observations (posterior PDF = learning from experiments)

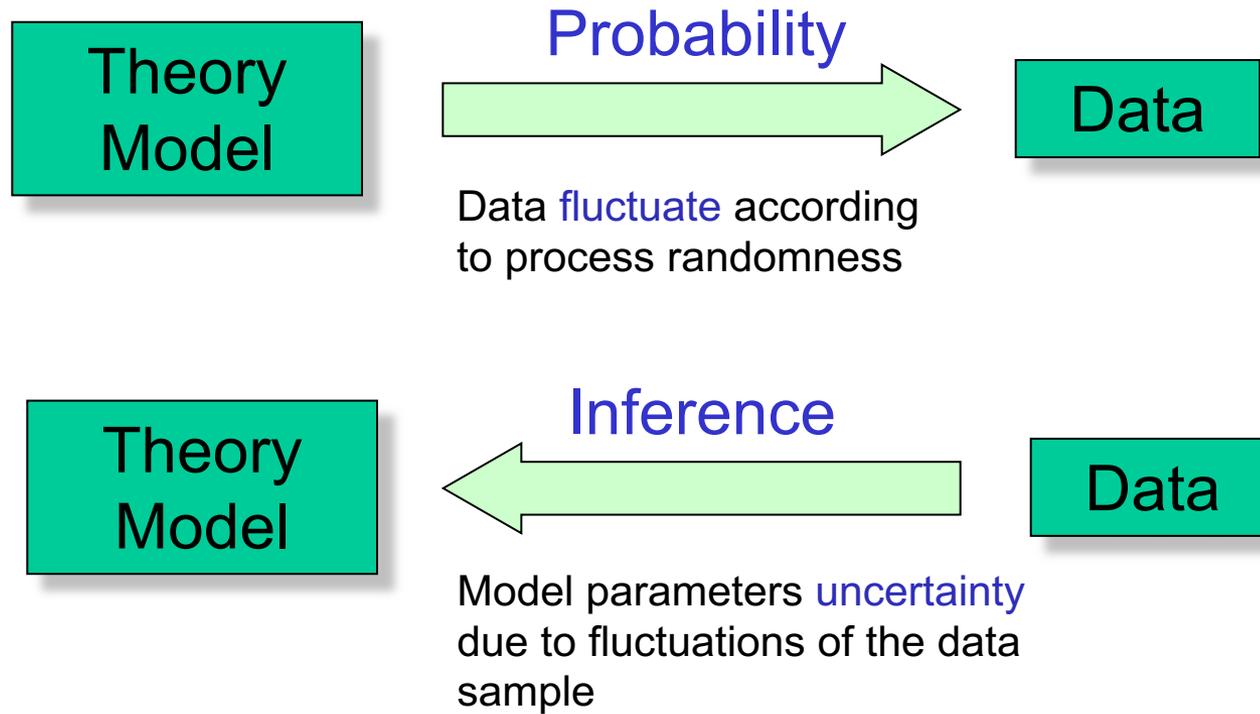


- Inference of a Binomial parameter as repeated application of Bayes rule for many Bernoulli extractions:
- $1 \rightarrow p_{i+1} = p_i \times p$
- $0 \rightarrow p_{i+1} = p_i \times (1 - p)$





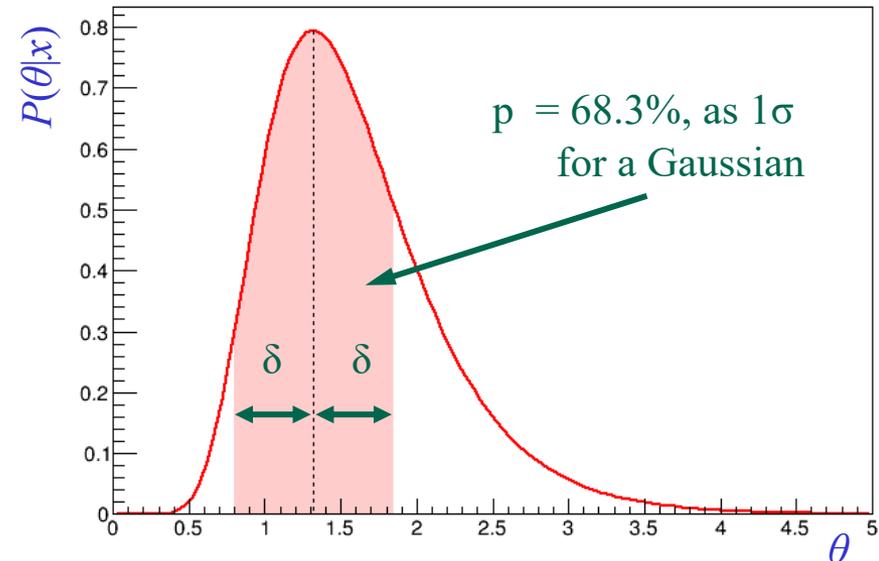
- Determining information about unknown parameters using probability theory



- The posterior PDF provides all the information about the unknown parameters (let's assume here it's just a single parameter θ for simplicity)

$$P(\theta|x) = \frac{L(x; \theta)\pi(\theta)}{\int L(x; \theta)\pi(\theta)d\theta}$$

- Given $P(\theta|x)$, we can determine:
 - The **most probable value** (best estimate)
 - **Intervals** corresponding to a specified probability
- Notice that if $\pi(\theta)$ is a constant, the most probable value of θ correspond to the **maximum of the likelihood function**



- Posterior PDF, assuming the prior to be $\pi(s)$:

$$P(s|n) = \frac{\frac{s^n e^{-s}}{n!} \pi(s)}{\int_0^\infty \frac{s'^n e^{-s'}}{n!} \pi(s') ds'}$$

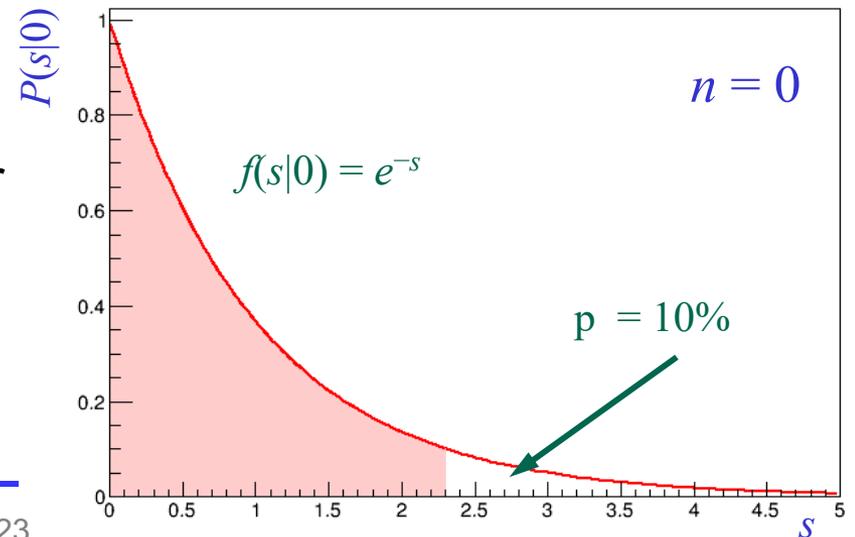
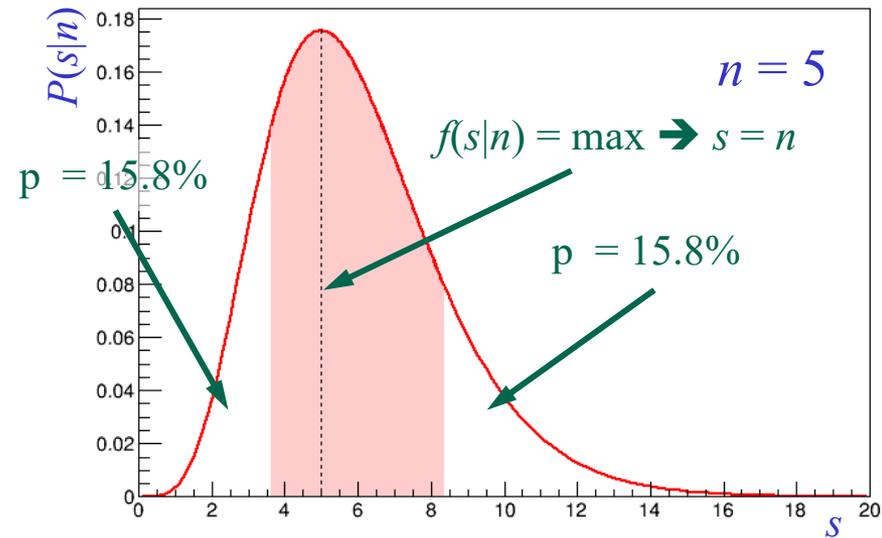
- If $\pi(s)$ is uniform:

$$P(s|n) = \frac{s^n e^{-s}}{n!}$$

- Note: $\langle s \rangle = n + 1$, $\text{Var}[s] = n + 1$
- For $n = 0$, one may quote an upper limit at 90% or 95% CL:
 - $s < 2.303$ (90% CL)
 - $s < 2.996$ (95% CL)

}

zero observed events

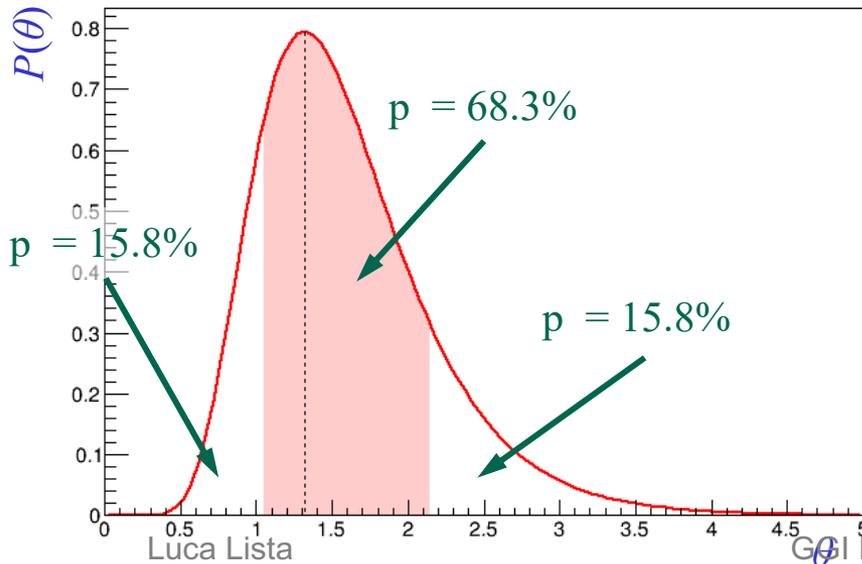


Choice of 68% prob. intervals

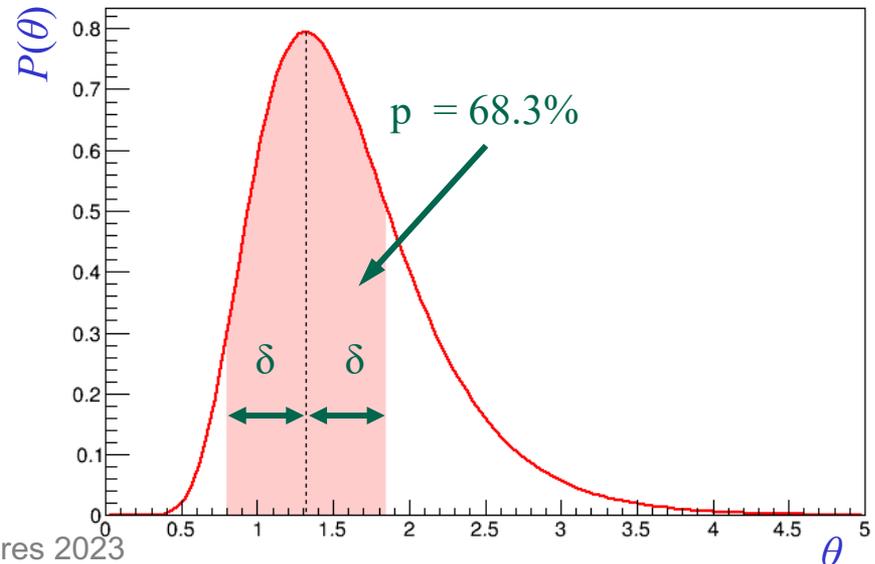


- Different **interval choices** are possible, corresponding to the same probability level (usually 68%, as 1σ for a Gaussian)
 - Equal areas in the right and left tails
 - Symmetric interval
 - Shortest interval
 - ...
- } All equivalent for a symmetric distribution (e.g. Gaussian)
- Reported as $\theta = \hat{\theta} \pm \delta$ (sym.) or $\theta = \hat{\theta}_{-\delta_2}^{+\delta_1}$ (asym.)

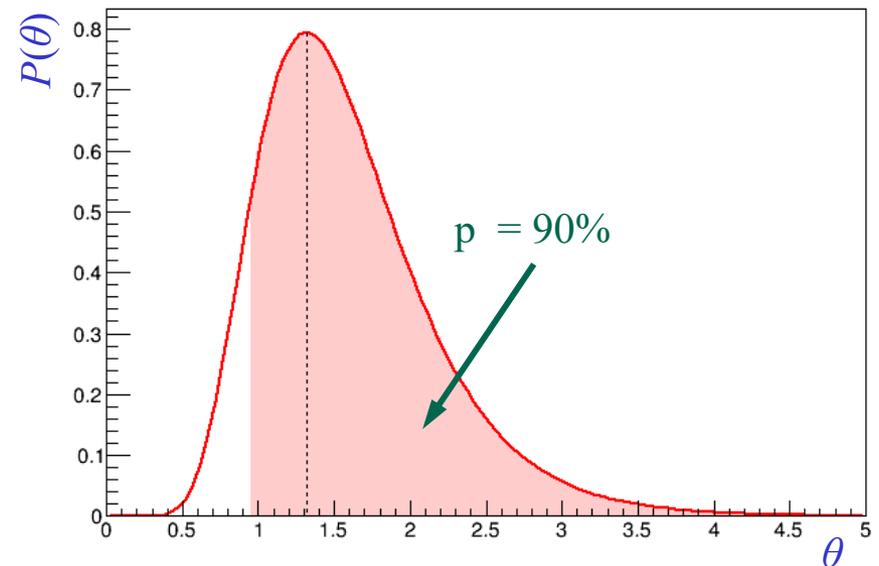
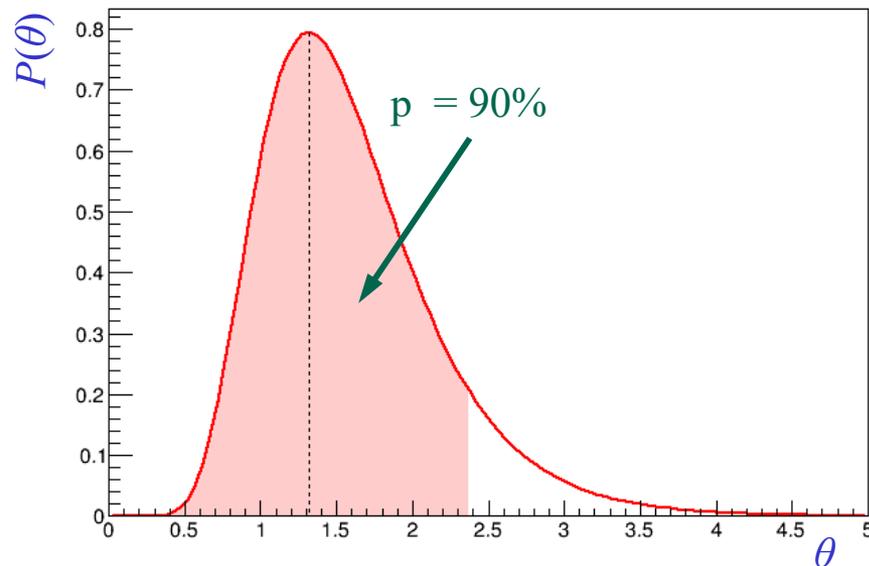
Equal tails interval



Symmetric interval



- A **fully asymmetric interval** choice is obtained setting one extreme of the interval to the lowest or highest allowed range
- The other extreme indicates an **upper or lower limits** to the “allowed” range
- For upper or lower limits, usually a probability of **90%** or **95%** is preferred to the usual 68% adopted for central intervals
- Reported as: $\theta < \theta^{\text{up}}$ (90% CL) or $\theta > \theta^{\text{lo}}$ (90% CL)



- Applying a parameter transformation, say $\eta = H(\theta)$, results in a transformed central value and transformed uncertainty interval
- The error propagation can be done transforming the posterior PDF, then computing the interval on the transformed PDF:

$$f'(\eta) = \int \delta(\eta - H(\theta)) f(\theta) d\theta$$

- Transformations for cases with more than one variable proceed in a similar way:

$$\eta = H(\theta_1, \theta_2) : \longrightarrow \downarrow$$

$$f'(\eta) = \int \delta(\eta - H(\theta_1, \theta_2)) f(\theta_1, \theta_2) d\theta_1 d\theta_2$$

$$\eta_1 = H_1(\theta_1, \theta_2), \eta_2 = H_2(\theta_1, \theta_2) : \longrightarrow \downarrow$$

$$f'(\eta_1, \eta_2) = \int \delta(\eta_1 - H_1(\theta_1, \theta_2)) \delta(\eta_2 - H_2(\theta_1, \theta_2)) f(\theta_1, \theta_2) d\theta_1 d\theta_2$$

- If the prior PDF is uniform in a choice of variable, it won't be uniform when applying coordinate transformation
- Given a prior PDF in a random variable, there is always a transformation that makes the PDF uniform
- The problem is: chose one metric where the PDF is uniform
- **Harold Jeffreys'** prior: chose the prior form that is **invariant under parameter transformation**

$$\pi(\theta) \propto \sqrt{\det I(\vec{\theta})} \quad I_{ij}(\vec{\theta}) = \left\langle \frac{\partial \ln L(\vec{x}; \vec{\theta})}{\partial \theta_i} \frac{\partial \ln L(\vec{x}; \vec{\theta})}{\partial \theta_j} \right\rangle$$

- Some commonly used cases:

- Poissonian mean:

$$\pi(\mu) \propto 1/\sqrt{\mu}$$

- Poissonian mean with background b :

$$\pi(\mu) \propto 1/\sqrt{\mu + b}$$

- Gaussian mean:

$$\pi(\mu) \propto 1$$

- Gaussian standard deviation:

$$\pi(\mu) \propto 1/\sigma$$

- Binomial parameter:

$$\pi(\mu) \propto 1/\sqrt{\varepsilon(1 - \varepsilon)}$$

Note: the previous simple Poissonian example was obtained with $\pi(\mu) = \text{const.}$!

- **Problematic with PDF in more than one dimension!**

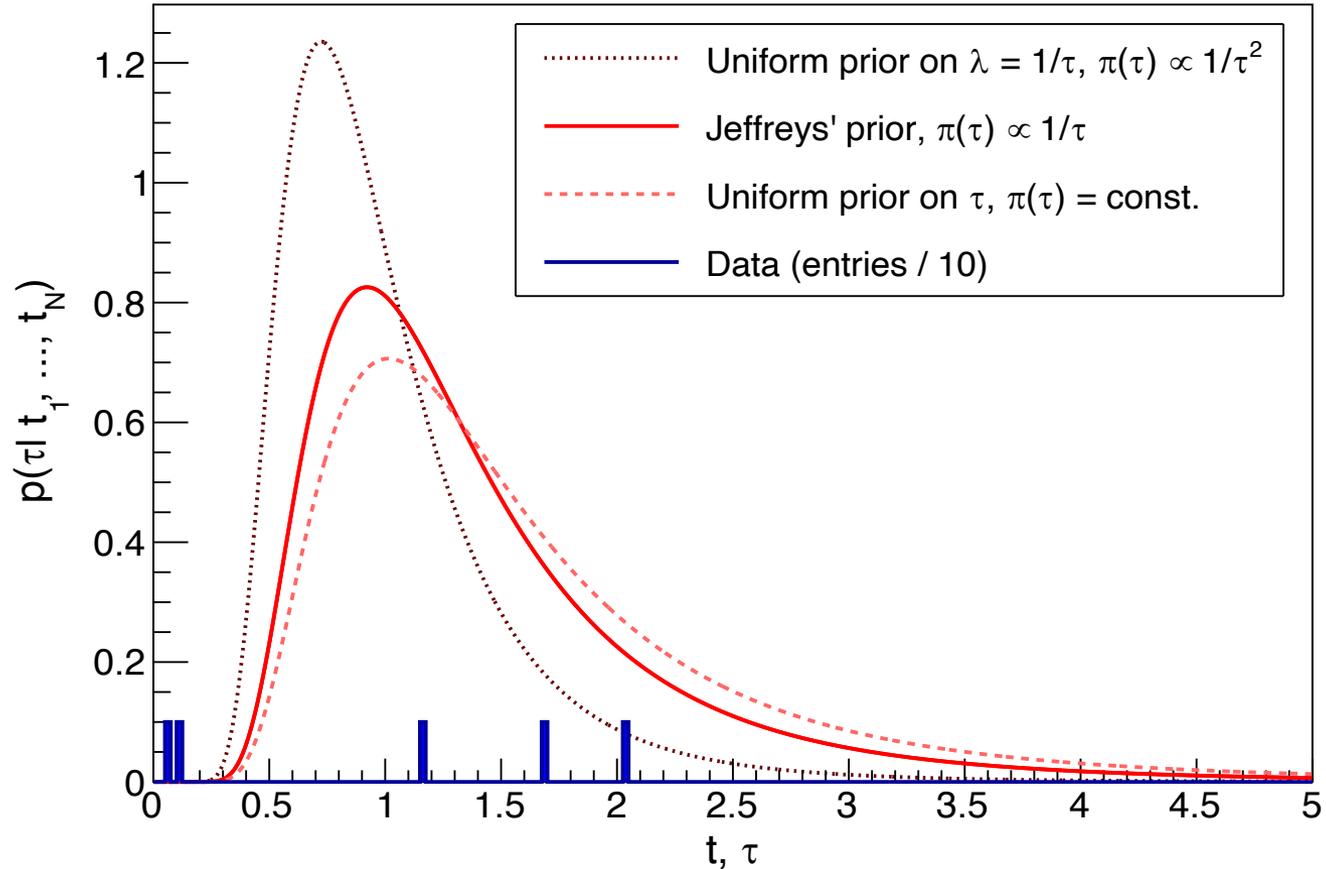
$$L(t_1, \dots, t_N | \tau) = \frac{1}{\tau^N} \prod_{i=1}^N e^{-t_i/\tau}$$

- The posterior for τ is:

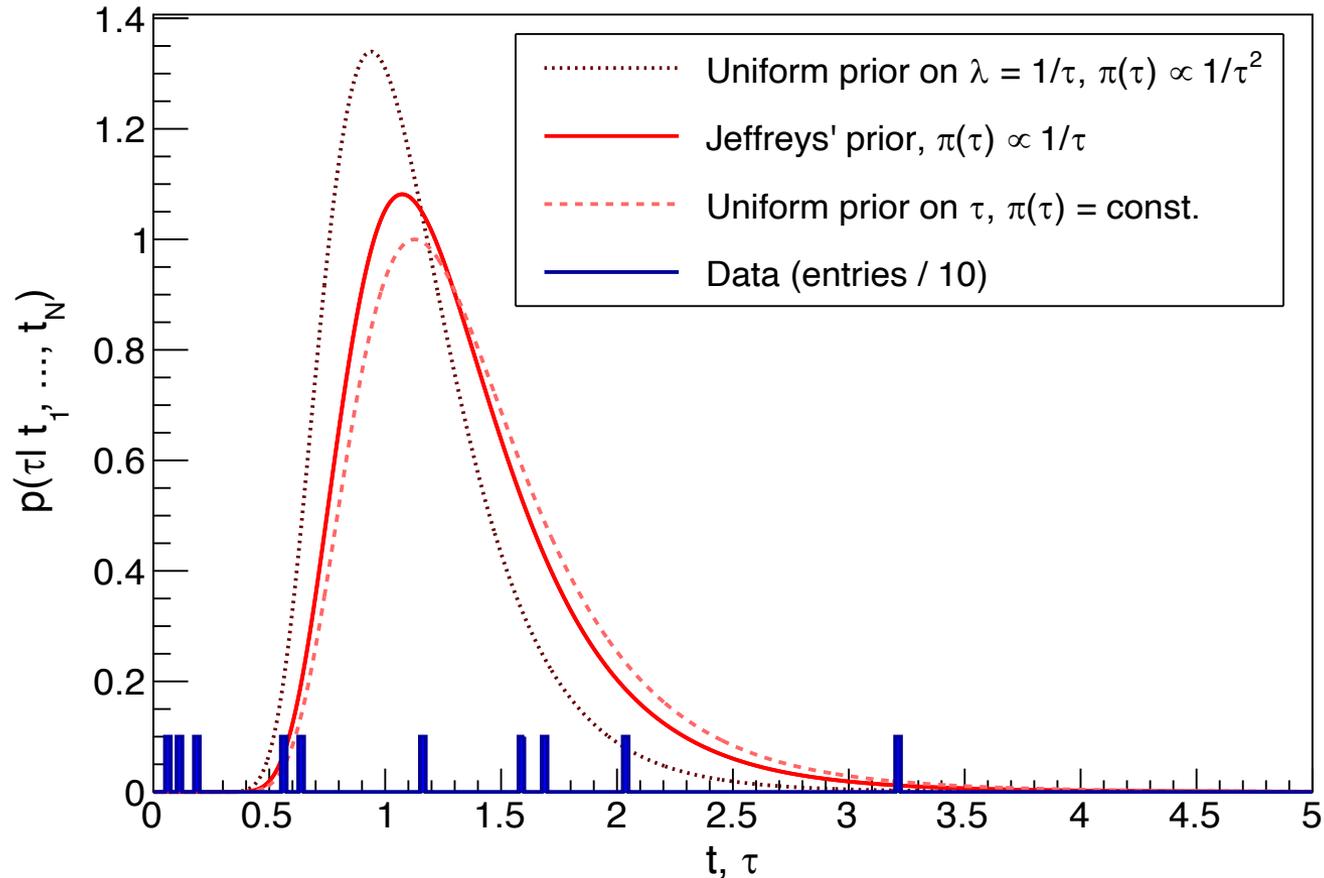
$$p(\tau | t_1, \dots, t_N) = \frac{\pi(\tau) e^{-\sum_{i=1}^N t_i/\tau} / \tau^N}{\int \pi(\tau') e^{-\sum_{i=1}^N t_i/\tau'} / \tau'^N d\tau'}$$

- The prior can be chosen in a different way:
 - Uniform in τ , $\pi(\tau) = \text{const.}$
 - Uniform in $\lambda = 1/\tau$, $\pi(\tau) = 1/\tau^2$
 - Jeffrey's prior, $\pi(\tau) = 1/\tau$
- All choices give as posterior a gamma distribution with different parameters ($k = N, N + 2, N + 1$)
 - $p(\tau | t_1, \dots, t_N) = C \tau^k e^{-\sum_{i=1}^N t_i/\tau}$
- The maximum of the PDF is at $\tau = \sum_{i=1}^N t_i / k$ which is equal to $N\bar{t}/k$, that, for large N , tends to \bar{t} , regardless of the prior choice.

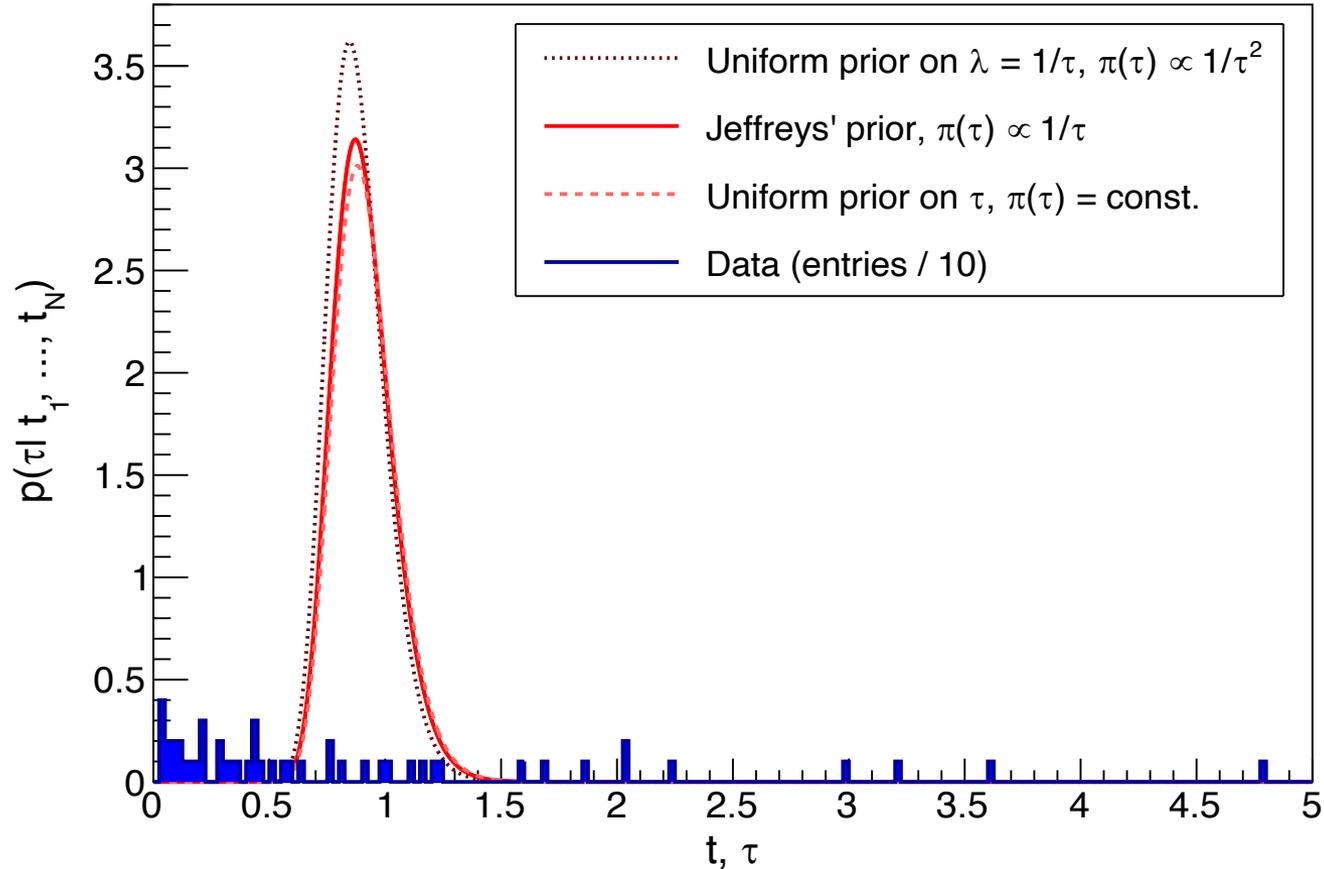
$$n = 5$$



$$n = 10$$



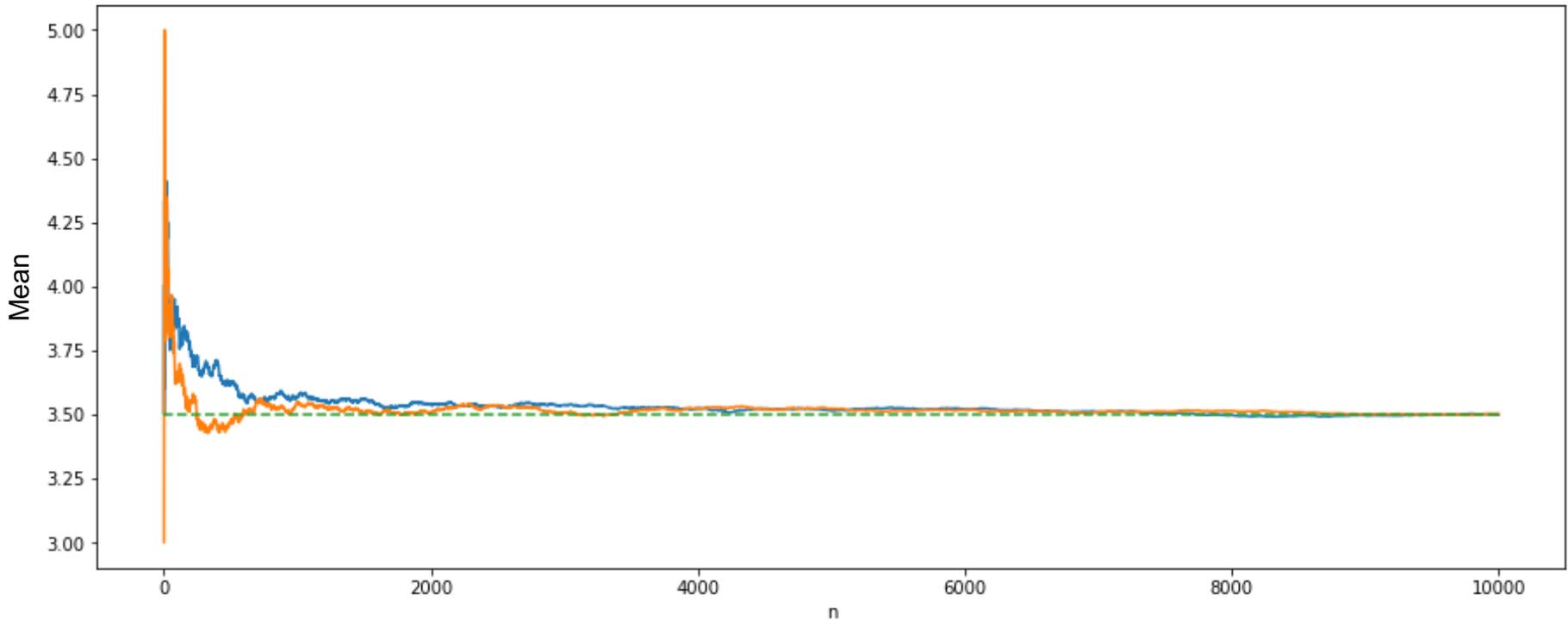
$n = 50$





- Law of large numbers and frequentist probability
 - Inference with the frequentist approach: coverage
 - The maximum likelihood method
 - Extended likelihood functions
 - Binned and unbinned fits
 - Properties of frequentist estimators
 - Neyman's confidence interval
 - Binomial confidence intervals according to Clopper and Pearson
 - Error estimates with the maximum likelihood method
 - Issues with asymmetric uncertainties
 - Two-dimensional uncertainty contours
 - Binned fits: the minimum chi-squared method
 - Goodness of the fit with chi-squared test
 - Baker-Cousins binned likelihood ratio fits
 - Combination of measurements and the BLUE method
-

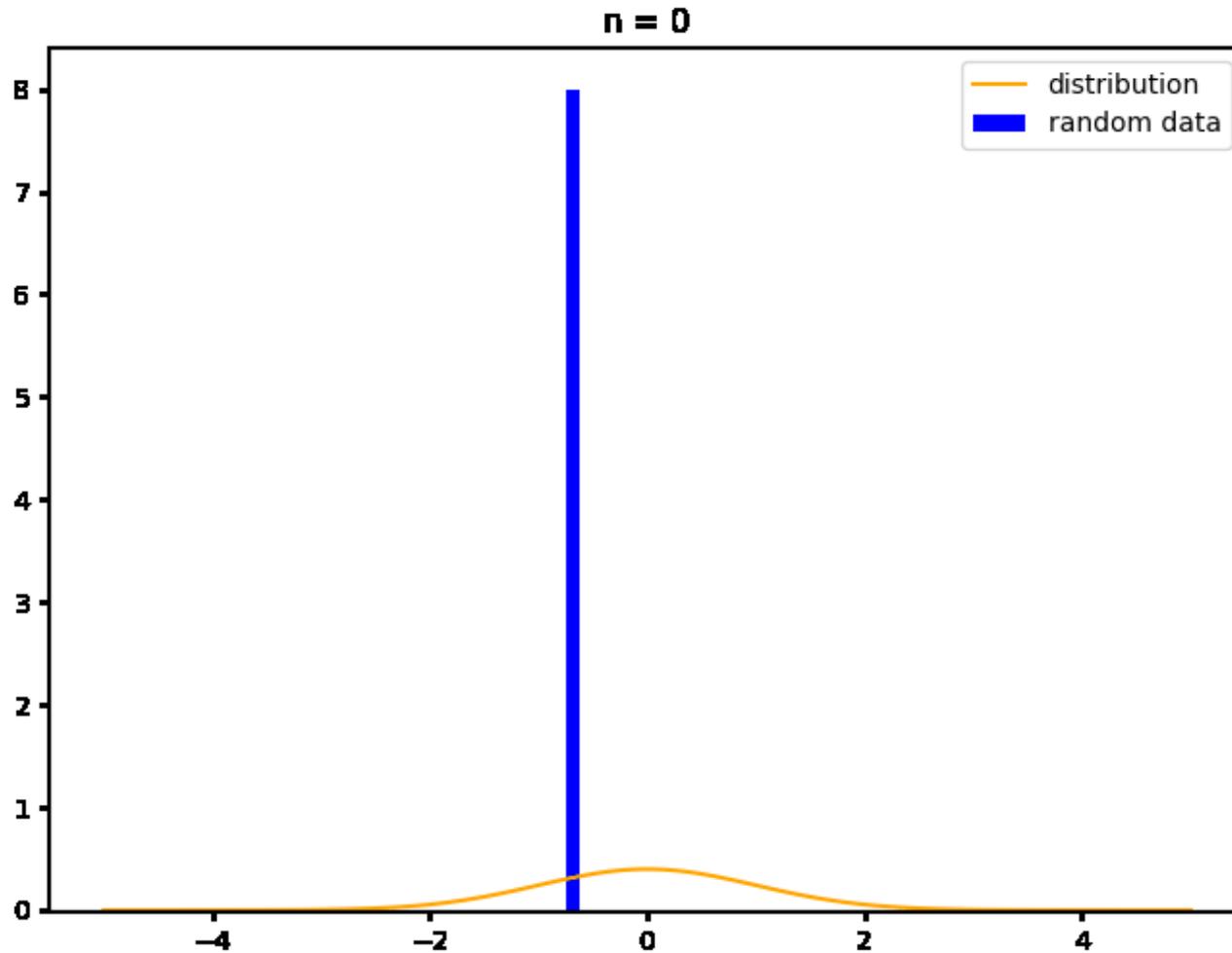
- The mean of a large number of random extractions following the same distribution tends to the expected value
- This law assumes an underlying probability model



- Probability P = frequency of occurrence of an event in the limit of very large number ($N \rightarrow \infty$) of repeated trials

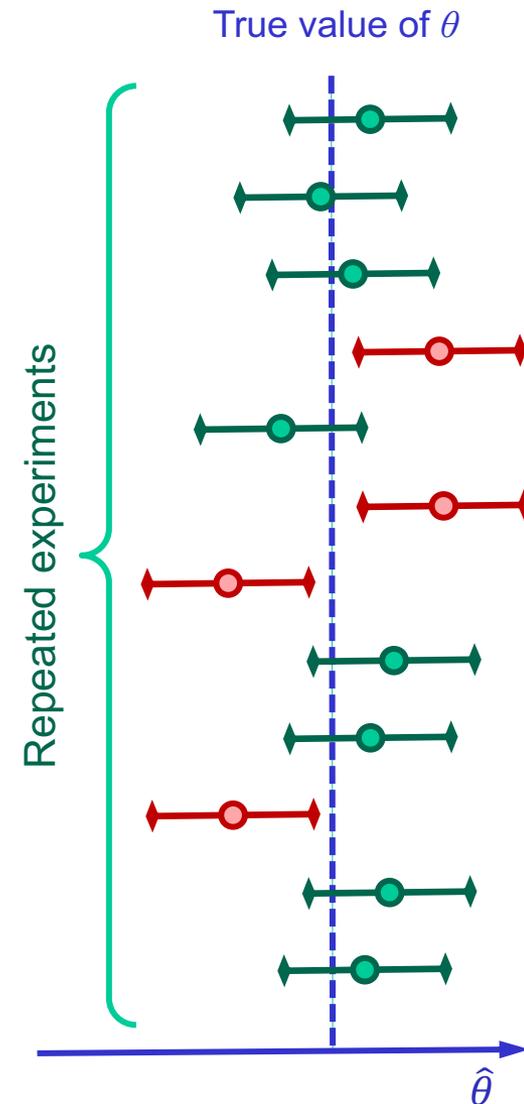
$$\text{Probability: } P = \lim_{N \rightarrow \infty} \frac{\text{Number of favorable cases}}{N = \text{Number of trials}}$$

- Exactly realizable only with an **infinite number of trials**
 - Conceptually may be unpleasant
 - Pragmatically acceptable by physicists
- Only applicable to repeatable experiments



- Assigning a probability level of an unknown parameter makes no sense in the frequentist approach
 - Parameters are not random variables!
- A frequentist inference procedure determines a **central value** and an **uncertainty interval** that depend on the observed measurements
- The **central value and interval extremes are random variables**
- **No subjective element** is introduced in the determination
- The function that returns the central value given an observed measurement is called **estimator**
- Different estimator choices are possible, the most frequently adopted is the **maximum likelihood estimator** because of its statistical properties discussed in the following

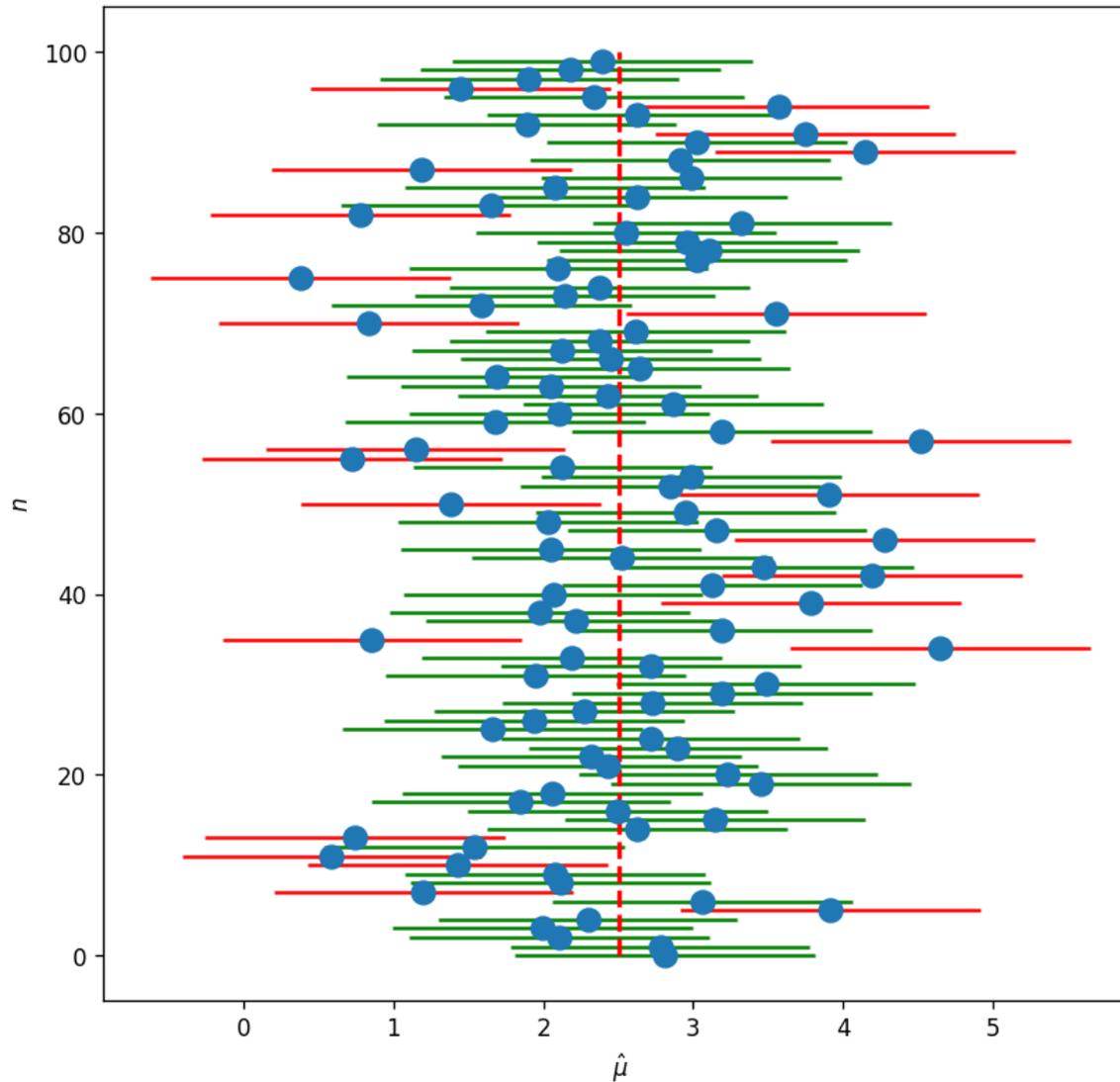
- Repeating the experiment will result each time in a different **data sample**
- For each data sample, the **estimator** returns a different **central value $\hat{\theta}$**
- An **uncertainty interval $[\hat{\theta} - \delta, \hat{\theta} + \delta]$** can be associated to the estimator's value $\hat{\theta}$
- Some of the confidence intervals contain the fixed and unknown true value of θ , corresponding to a fraction equal to 68% of the times, in the limit of very large number of experiments (**coverage**)





- An **estimator** is a function of a given set of measurements that provides an approximate value of a **parameter of interest** which appears in our PDF model (“**best fit**”)
- Simplest example:
 - Assume a Gaussian PDF with a *known* σ and an *unknown* μ
 - A single experiment provides a measurement x
 - We estimate μ as $\hat{\mu} = x$
 - The distribution of $\hat{\mu}$ (repeating the experiment many times) is the original Gaussian
 - 68.3% of the experiments (in the limit of large number of repetitions) will provide an estimate within: $\mu - \sigma < \hat{\mu} < \mu + \sigma$
- We can quote:

$$\mu = x \pm \sigma$$



- The **maximum-likelihood estimator** is the most adopted parameter estimator
- The “**best fit**” **parameters** correspond to the set of values that maximizes the likelihood function
 - Good statistical properties (→ next slides)
- The maximization can be performed analytically only in the simplest cases, and numerically for most of realistic cases

- **Minuit** is historically the most widely used minimization engine in High Energy Physics
 - F. James, 1970's; rewritten in C++ and released under CERN's ROOT framework





- If we have n independent measurements all modeled with (or approximated to) the same Gaussian PDF, we have:

$$-2 \ln L = \underbrace{\sum_{i=1}^n \frac{(x_i - \mu)^2}{\sigma^2}}_{\text{(example of a } \chi^2 \text{ variable)}} + n(\ln 2\pi + 2 \ln \sigma)$$

- An analytical minimization of $-2 \ln L$ w.r.t μ (assuming σ^2 is known) gives the **arithmetic mean** as ML estimate of μ :

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i$$

- If σ^2 is also **unknown**, the ML estimate of σ^2 is:

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})^2$$

- The above estimate can be demonstrated to have an unpleasant feature, called *bias* (\rightarrow next slide)

- **Consistency**: for large number of measurements the estimator $\hat{\theta}$ should converge, in probability, to the true value θ .
 - ML estimators are consistent
- **Bias**: the bias of a parameter is the average value of its deviation from the true value

$$b(\theta) = \langle \hat{\theta} - \theta \rangle = \langle \hat{\theta} \rangle - \theta$$

- ML estimators may have a bias, but the bias decreases with large number of measurements (if the fit model is correct...!)
- E.g.: in the case of the estimate of a Gaussian's σ^2 , the unbiased estimate is the well known:

$$\hat{\sigma}_{\text{unbias.}}^2 = \frac{n}{n-1} \hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \hat{\mu})^2$$

← ML method underestimates the variance σ^2

- The **variance** of any consistent estimator is subject a **lower bound** (Cramér-Rao bound):

$$V[\hat{\theta}] \geq \frac{\left(1 + \frac{\partial b(\theta)}{\partial \theta}\right)^2}{\mathbb{E}\left[\left(\frac{\partial^2 \ln L(\theta)}{\partial \theta^2}\right)^2\right]} = \frac{\left(1 + \frac{\partial b(\theta)}{\partial \theta}\right)^2}{\mathbb{E}\left[-\frac{\partial^2 \ln L(\theta)}{\partial \theta^2}\right]} = V_{CR}$$

bias of θ

- Efficiency** can be defined as the ratio of Cramér-Rao bound and the estimator's variance:

$$\varepsilon(\hat{\theta}) = \frac{V_{CR}}{\text{Var}[\hat{\theta}]}$$

- Efficiency for ML estimators tends to 1 for large number of measurements

$$\lim_{n \rightarrow \infty} V[\hat{\theta}] = -\frac{1}{\mathbb{E}\left[\left(\frac{\partial^2 \ln L(\theta)}{\partial \theta^2}\right)^2\right]} \cong -\frac{1}{\frac{\partial^2 \ln L(\theta)}{\partial \theta^2} \Big|_{\theta=\hat{\theta}}}$$

- I.e.: ML estimates have, asymptotically, the smallest possible **variance**



- A **parabolic approximation** of $-2\ln L$ around the minimum is equivalent to a **Gaussian approximation**
 - Sufficiently accurate in many but not all cases

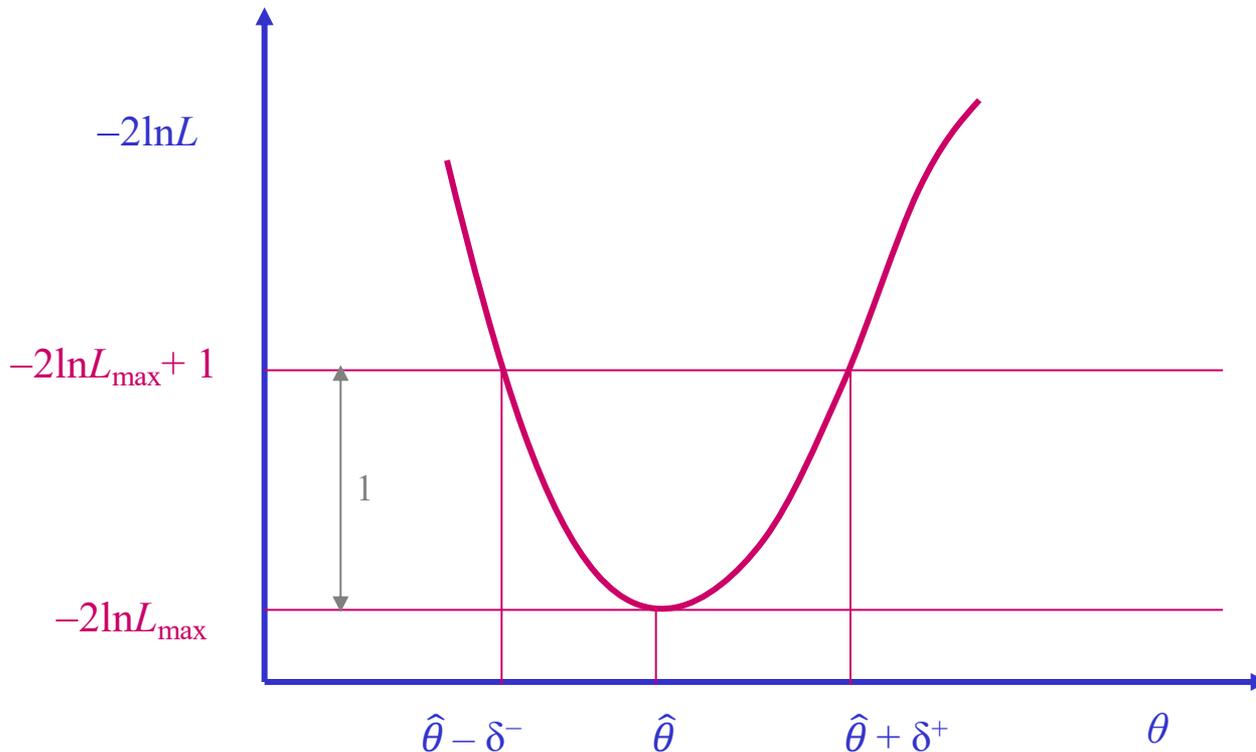
$$-2 \ln L = \sum_{i=1}^n \frac{(x_i - \mu)^2}{\sigma^2} + \text{const.}$$

- Estimate of the covariance matrix from 2nd order partial derivatives w.r.t. fit parameters at the minimum:

$$V_{ij}^{-1} = - \left. \frac{\partial^2 \ln L}{\partial \theta_i \partial \theta_j} \right|_{\theta_k = \hat{\theta}_k}$$

- Implemented in Minuit as MIGRAD/HESSE function

- Another approximation alternative to the parabolic one may be to evaluate the excursion range of $-2\ln L$.
- Error ($n\sigma$) determined by the range around the maximum for which $-2\ln L$ increases by $+1$ ($+n^2$ for $n\sigma$ intervals)



- Errors can be asymmetric
- For a Gaussian PDF the result is identical to the 2nd order derivative matrix
- Implemented in Minuit as MINOS function

- The probability distribution for a single measurement is:

$$p(t; \tau) = \frac{1}{\tau} e^{-t/\tau}$$

- And the likelihood function is

$$L(t_1, \dots, t_n; \tau) = \frac{1}{\tau^n} \left(\prod_{i=1}^n e^{-t_i/\tau} \right)$$

- Minimization gives:

$$\hat{\tau} = \frac{1}{n} \sum_{i=1}^n t_i$$

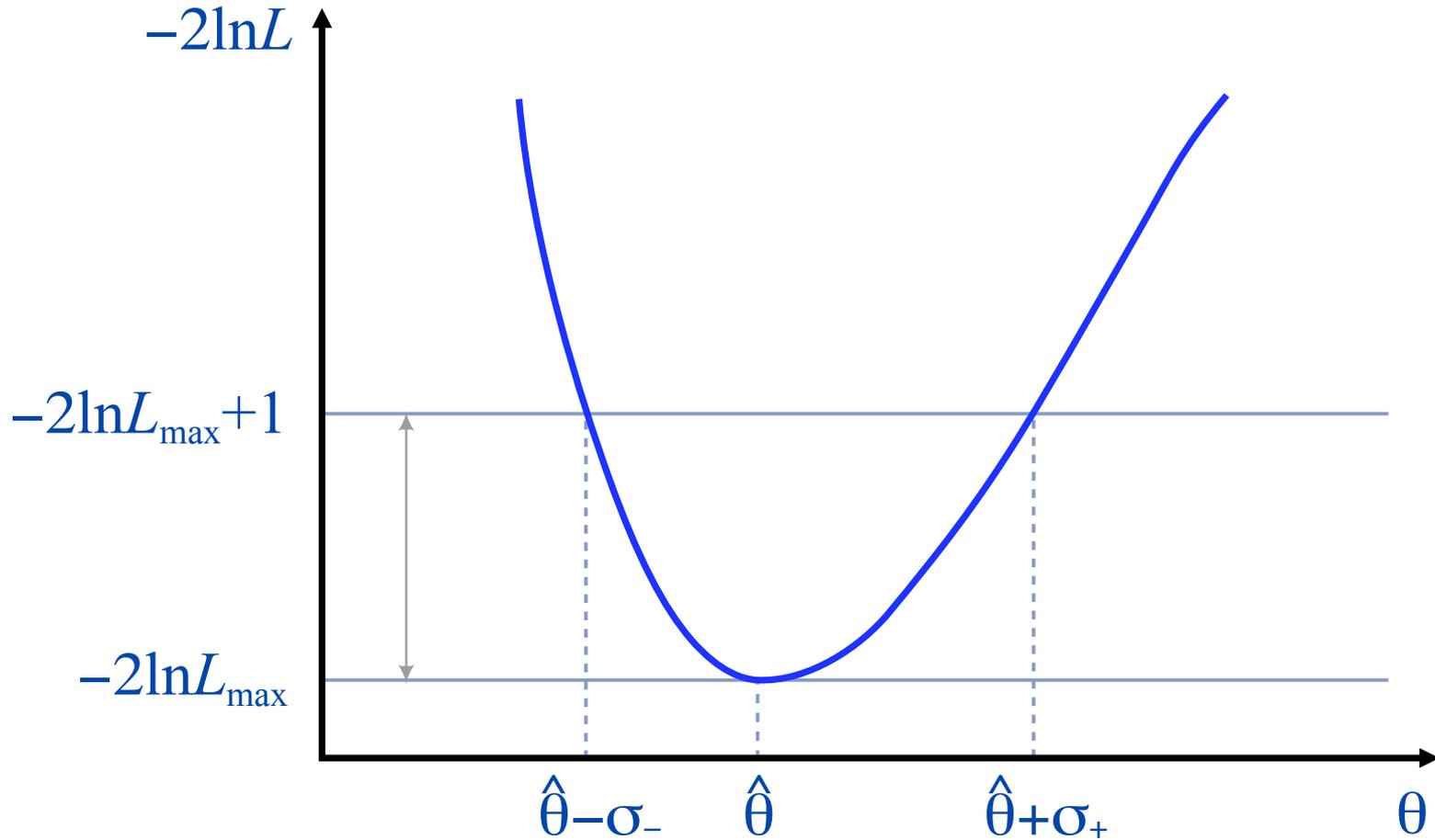
- Which is clearly unbiased
- The distribution of $\hat{\tau}$ is a gamma distribution scale parameter τ and shape parameter 1:

$$p(\hat{\tau}) = \hat{\tau}^{2n} e^{-\hat{\tau}/\tau} / (\tau^{2n} (2n - 1)!)$$

- With uncertainty given by the square root of the variance equal to:

$$\sigma_{\hat{\tau}} = \tau / \sqrt{n}$$

- This is also equal to the Cramer-Rao bound



INFN 2D intervals

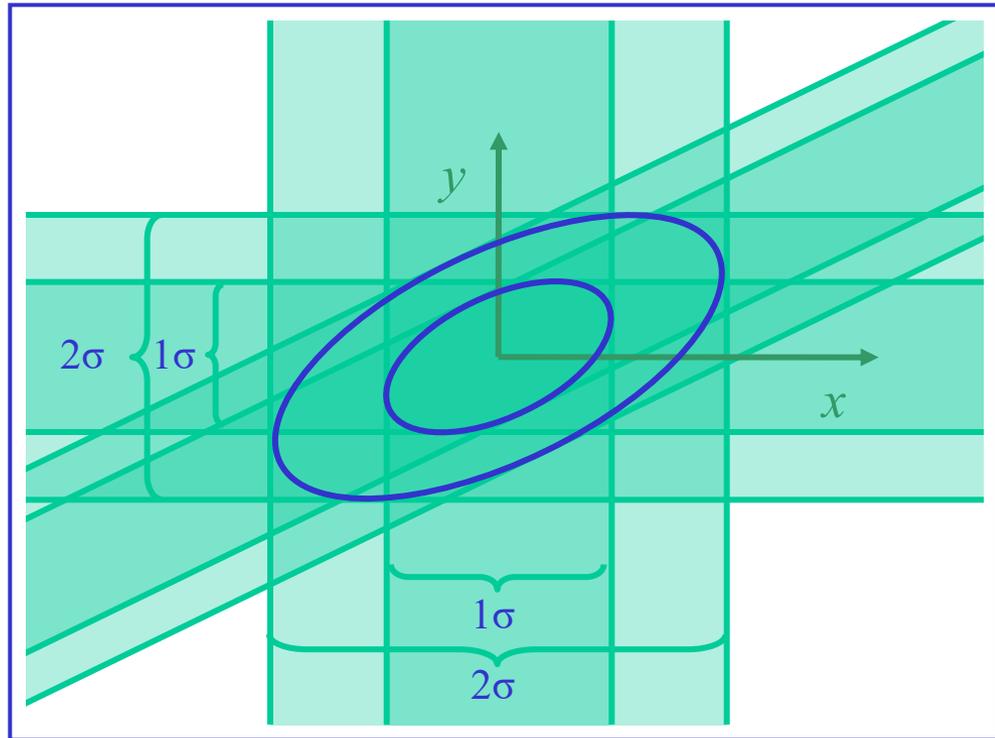


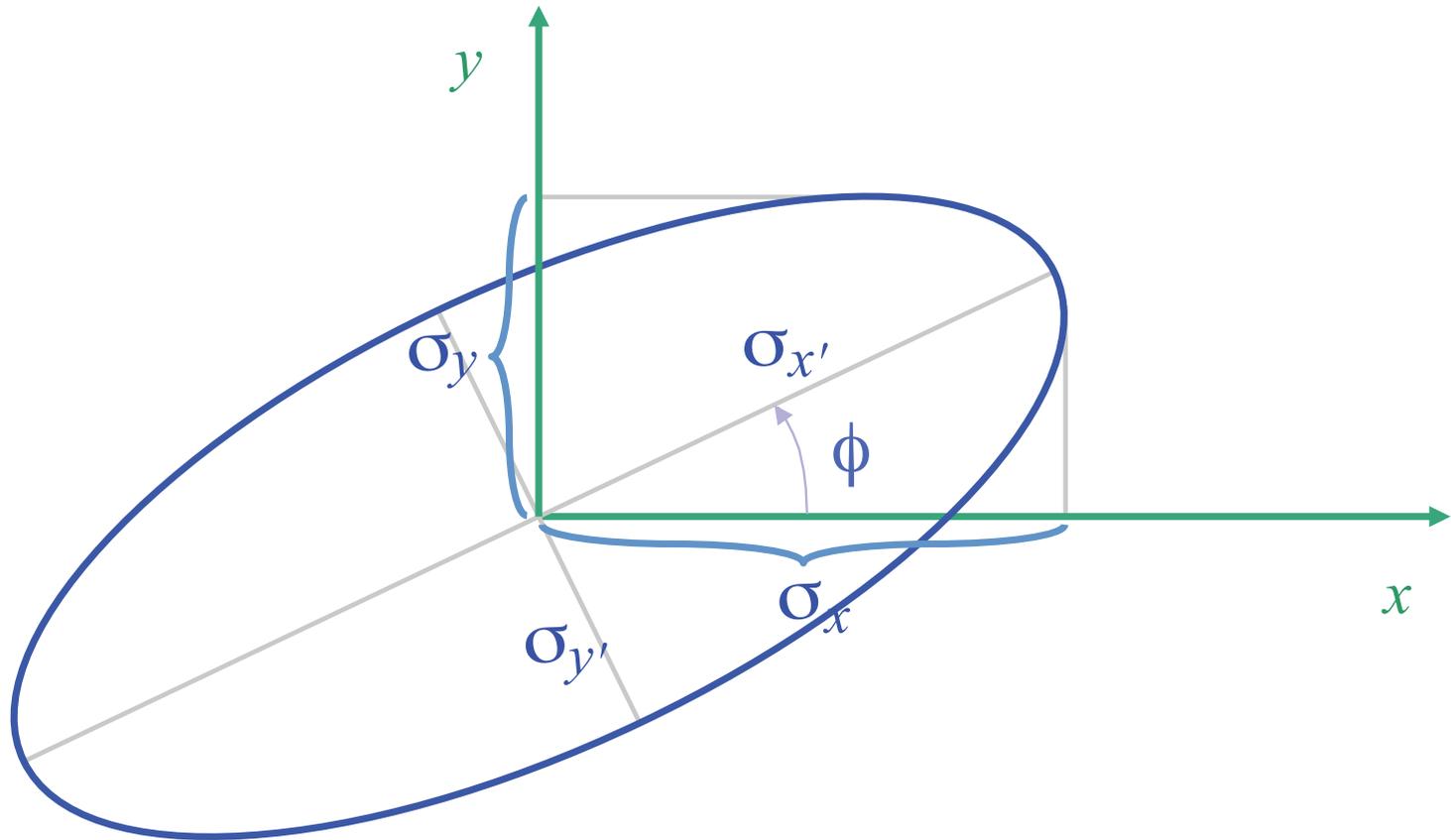
- In more dimensions one can determine 1σ and 2σ contours
- Note: different probability content in 2D compared to one dimension
- 68% and 95% contours are usually preferable

$$P_{1D}(n\sigma) = \sqrt{\frac{2}{\pi}} \int_0^n e^{-\frac{x^2}{2}} dx = \text{erf}\left(\frac{n}{\sqrt{2}}\right)$$

$$P_{2D}(n\sigma) = \int_0^n e^{-\frac{r^2}{2}} r dr = 1 - e^{-\frac{n^2}{2}}$$

Width	P_{1D}	P_{2D}
1σ	0.6827	0.3934
2σ	0.9545	0.8647
3σ	0.9973	0.9889
1.515σ		0.6827
2.486σ		0.9545
3.439σ		0.9973





- Given a sample of N measurements of the variables (x_1, \dots, x_n) , the likelihood function expresses the probability density of the sample, as a function of the unknown parameters:

$$L = \prod_{i=1}^N f(x_1^i, \dots, x_n^i; \theta_1, \dots, \theta_m)$$

- If the size N of the sample is also a random variable, the **extended likelihood** function is usually also used:

$$L = P(N; \theta_1, \dots, \theta_m) \prod_{i=1}^N f(x_1^i, \dots, x_n^i; \theta_1, \dots, \theta_m)$$

- Where $P(N; \theta_1, \dots, \theta_m)$ is in practice always a **Poisson** distribution whose expected rate is a function of the unknown parameters
- In many cases it is convenient to use $-\ln L$ or $-2\ln L$: $\prod_i \rightarrow \sum_i$

- For Poissonian signal and background processes:

$$L(x_i; s, b, \theta) = \frac{(s + b)^n e^{-(s+b)}}{n!} \prod_{i=1}^n (f_s P_s(x_i; \theta) + f_b P_b(x_i; \theta))$$

$$\left. \begin{aligned} f_s &= \frac{s}{s+b} \\ f_b &= \frac{b}{s+b} \end{aligned} \right\} = \frac{e^{-(s+b)}}{n!} \prod_{i=1}^n (s P_s(x_i; \theta) + b P_b(x_i; \theta))$$

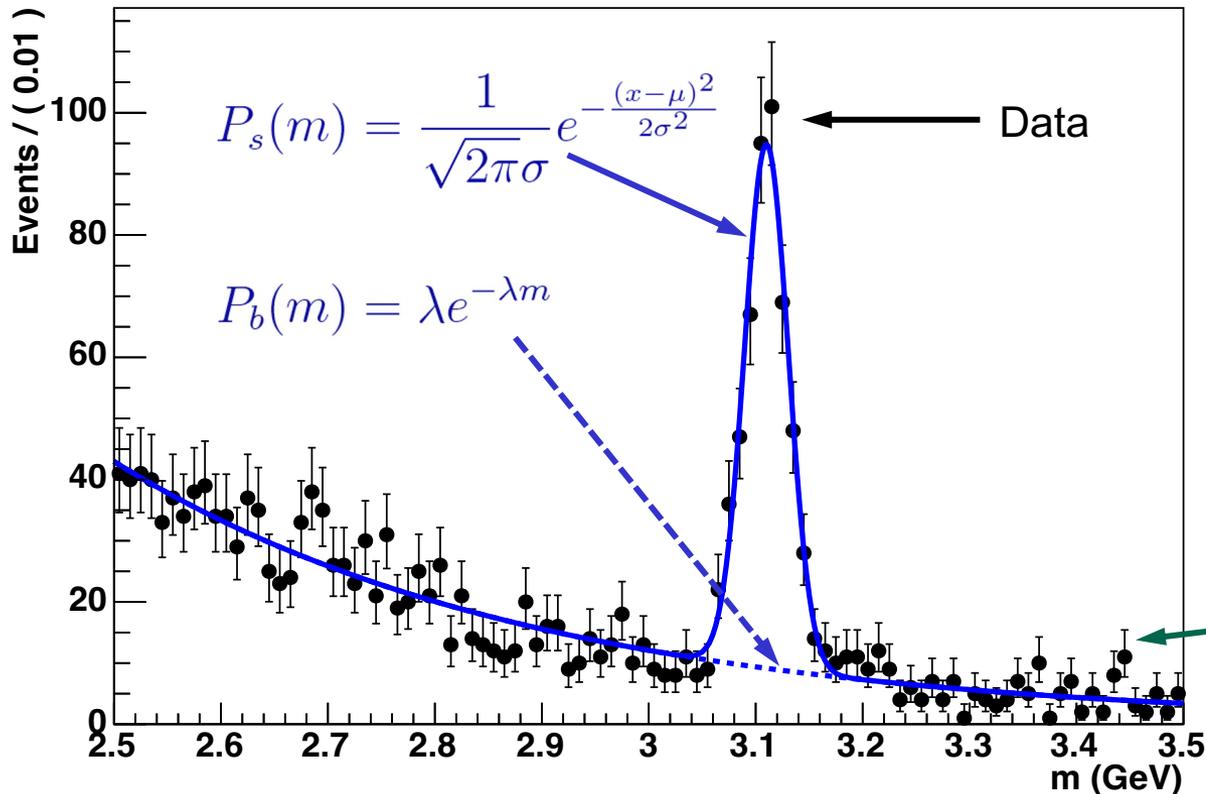
- We can fit simultaneously s , b and θ minimizing: constant!

$$-\ln L = s + b - \sum_{i=1}^n \ln(s P_s(x_i; \theta) + b P_b(x_i; \theta)) + \ln n!$$

- Sometimes s is replaced by μs_0 , where s_0 is the theory estimate and μ is called **signal strength**

- $P_s(m)$: Gaussian peak
- $P_b(m)$: exponential shape

Exponential decay parameter λ , Gaussian mean μ and standard deviation σ can be fit together with sig. and bkg. yields s and b .



The additional parameters, beyond the **parameters of interest** (s in this case), used to model background, resolution, etc. are examples of **nuisance parameters**

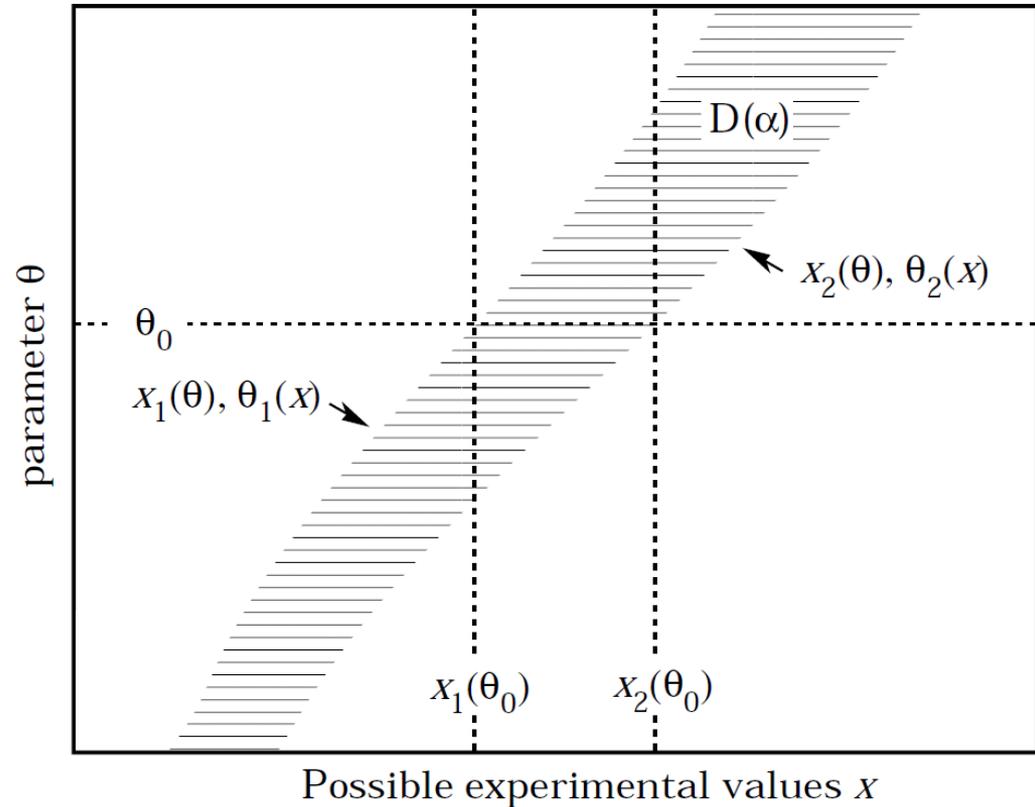
In the plot, data are accumulated into **bins** of a given width

Error bars usually represent uncertainty on each bin count (in this case: Poissonian)

Procedure to determine frequentist confidence intervals

- Scan the allowed range of an unknown parameter θ
- Given a value of θ compute the interval $[x_1, x_2]$ that contain x with a probability $1 - \alpha$ equal to 68% (or 90%, 95%)
- **Choice of interval needed!**
- Invert the **confidence belt**: for an observed value of x , find the interval $[\theta_1, \theta_2]$
- A fraction of the experiments equal to $1 - \alpha$ will measure x such that the corresponding $[\theta_1, \theta_2]$ contains (“covers”) the true value of θ (“coverage”)
- **Note:** the random variables are $[\theta_1, \theta_2]$, not θ !

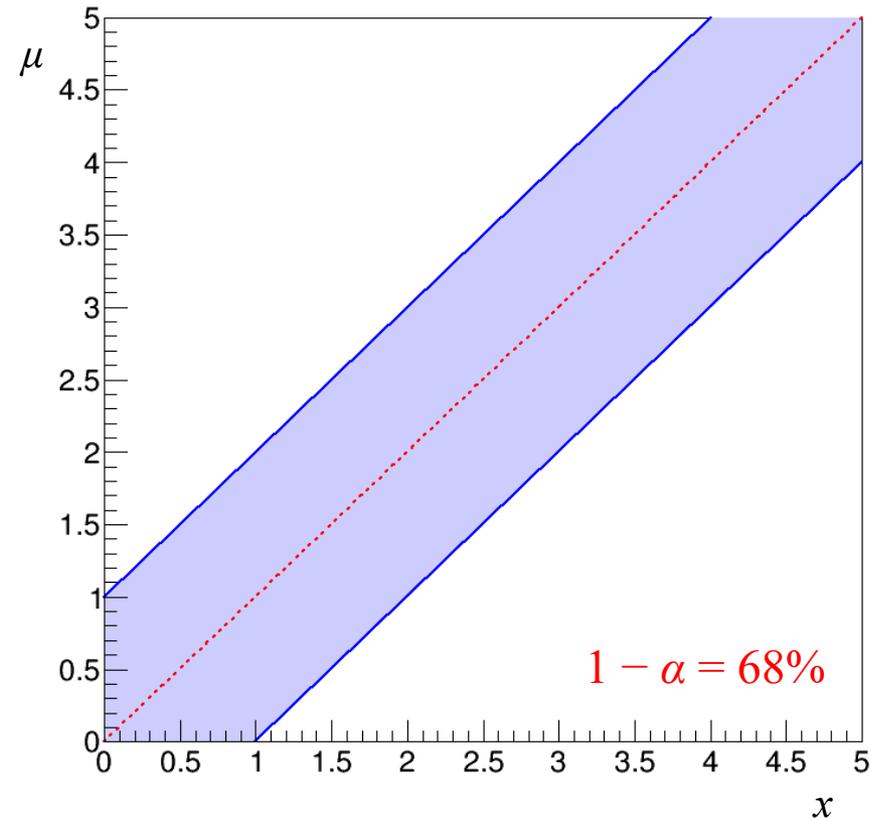
Plot from PDG statistics review



$\alpha =$ significance level

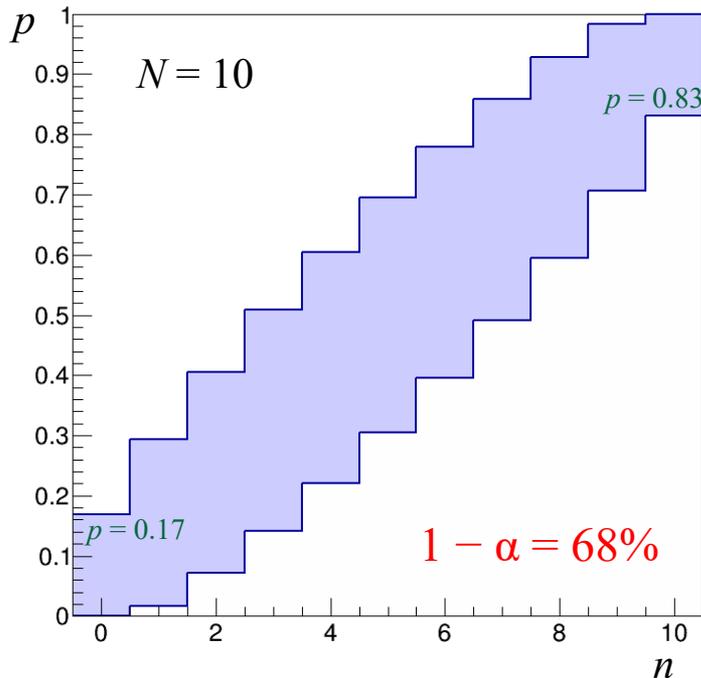
- Assume a Gaussian distribution with unknown average μ and known $\sigma = 1$
- The belt inversion is trivial and gives the expected result:
Central value $\hat{\mu} = x$,
 $[\mu_1, \mu_2] = [x - \sigma, x + \sigma]$
- So we can quote:

$$\mu = x \pm \sigma$$



- The Neyman's belt construction may only guarantee **approximate coverage** in case of **discrete variables**
- For a Binomial distribution: find the interval $\{n_{\min}, \dots, n_{\max}\}$ such that:

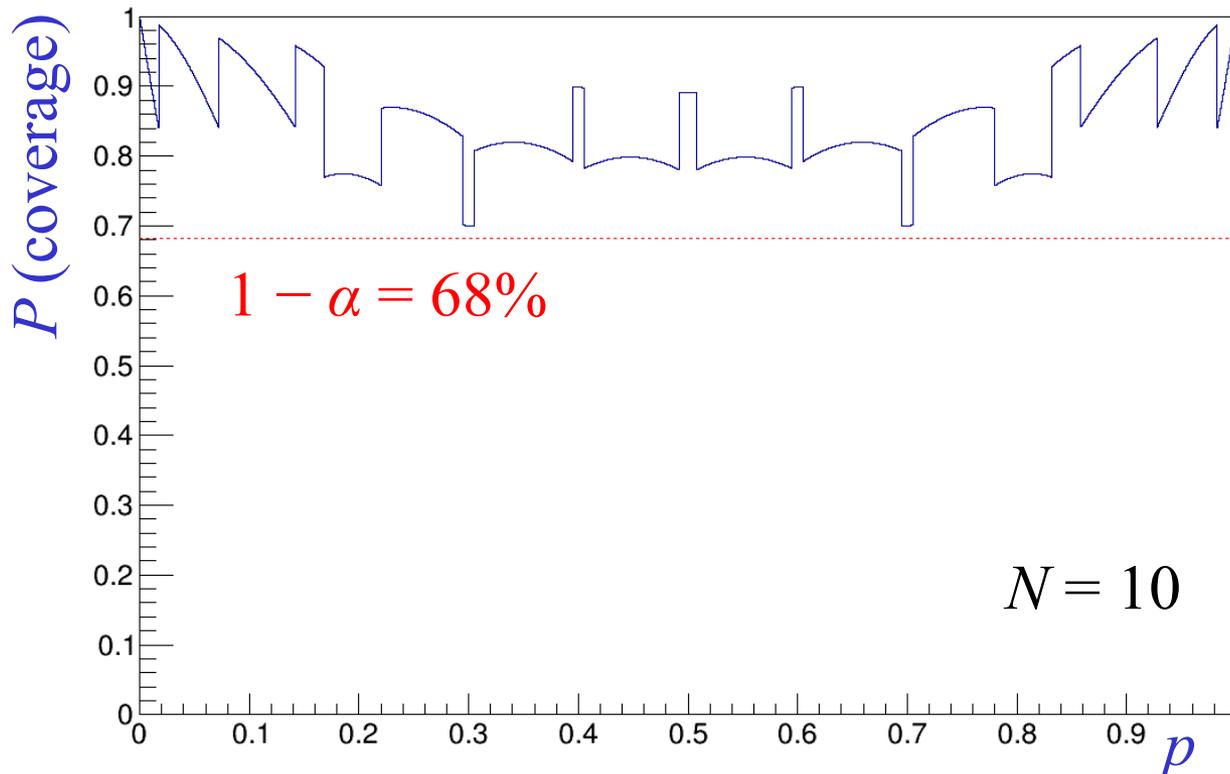
$$\sum_{n=n_{\min}}^{n=n_{\max}} \frac{N!}{n!(N-n)!} p^n (1-p)^{N-n} \geq 1 - \alpha$$



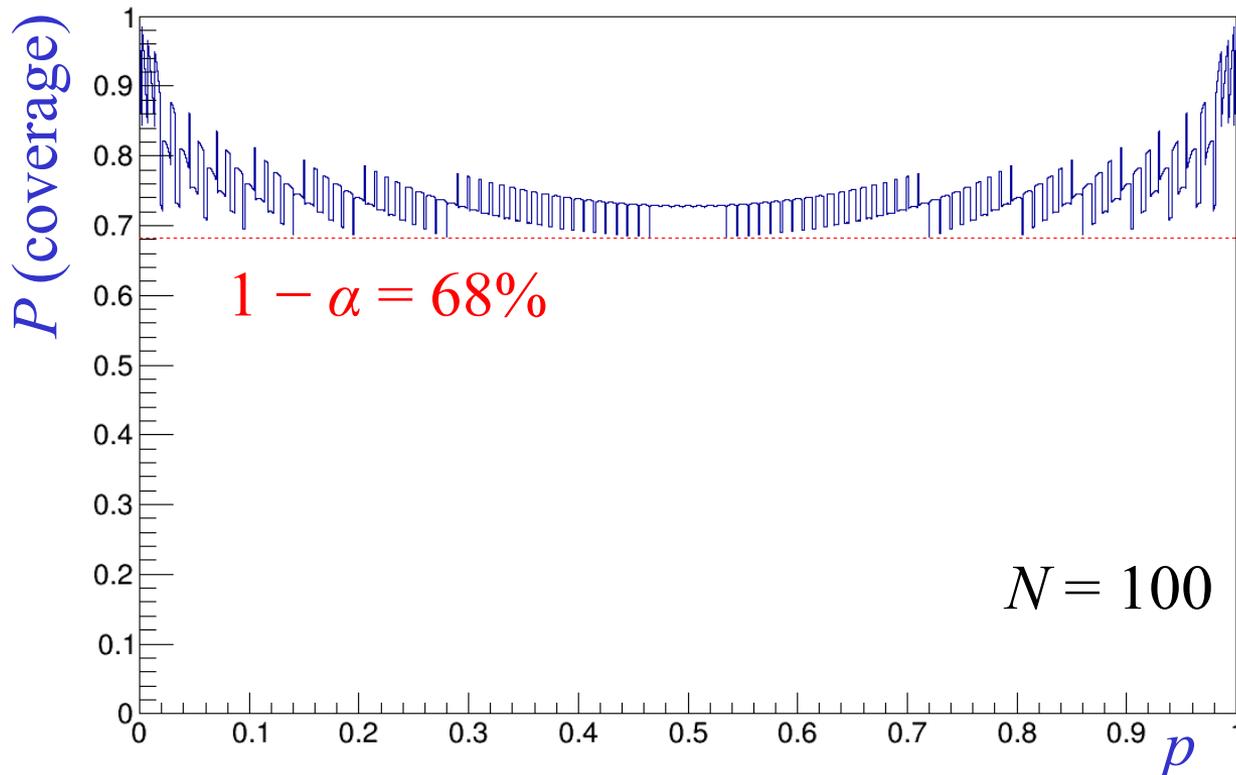
- Clopper and Pearson** (1934) solved the belt inversion problem for central intervals
- For an observed $n = k$, find lowest p^{lo} and highest p^{up} such that:
- $P(n \leq k | N, p^{\text{lo}}) = \alpha/2$, $P(n \geq k | N, p^{\text{up}}) = \alpha/2$
- E.g.: $n = N = 10$, $P(N|N) = p^N = \alpha/2$, hence:
 $p^{\text{lo}} = \sqrt[10]{\alpha/2} = 0.83$ (68% CL), 0.74 (90% CL)
- A frequently used approximation, which **fails** for $n = 0$, N is:

$$\hat{p} = \frac{n}{N}, \quad \sigma_{\hat{p}} \simeq \sqrt{\frac{\hat{p}(1-\hat{p})}{N}}$$

- CP intervals are often defined as “exact” in literature
- Exact coverage is often impossible to achieve for discrete variables



- For larger N the “ripple” gets closer to the nominal 68% coverage



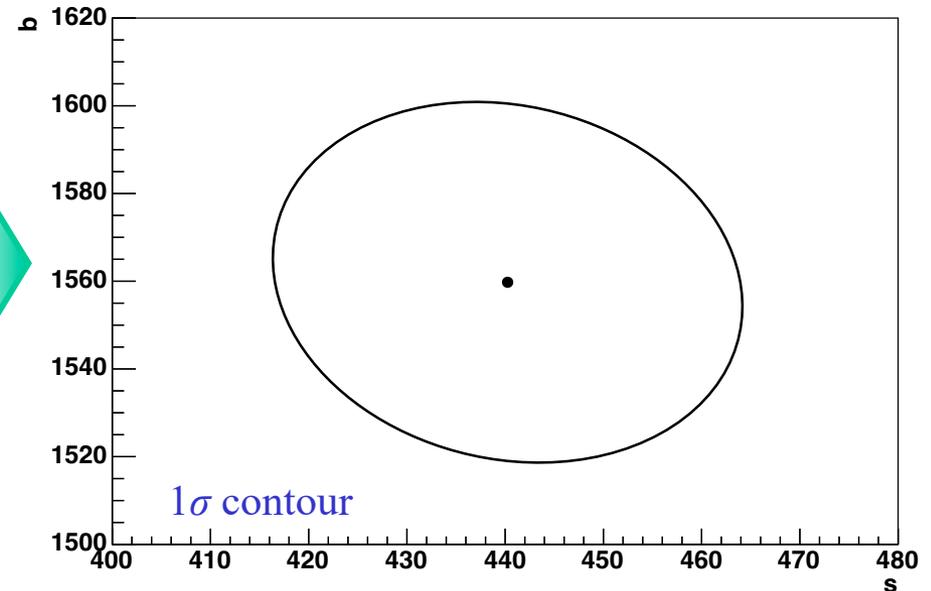
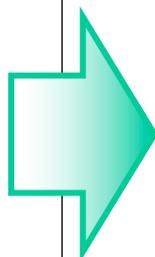
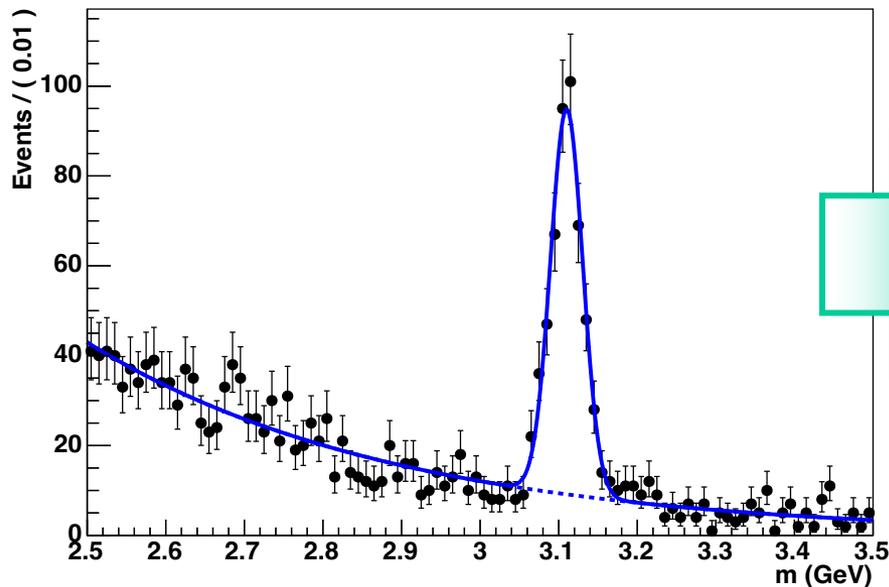


- Note: if the true value is $p = 0$ (similarly for $p = 1$), the observed value is always $n = 0$, therefore the confidence interval is $[0, pup[$
- The true value is therefore contained in the confidence interval with 100% probability instead of 68% (or 90%, or whatever)
- This is against the definition of frequentist coverage, but it is unavoidable for discrete variable
- This feature is typical of cases with low number of counts, for instance, for Poissonian counting experiments

- From previous fit example:
 - $P_s(m)$: Gaussian peak
 - $P_b(m)$: exponential shape

Exponential decay parameter, Gaussian mean and standard deviation are fit together with s and b yields.

The contour shows for this case a mild correlation between s and b



- Assume we estimate from a fit the parameter set:
 $\theta = (\theta_1, \dots, \theta_n)$ and we know their **covariance matrix** Θ_{ij}
- We want to determine a new set of parameters that are functions of θ :
 $\eta = (\eta_1, \dots, \eta_m)$.
- For small uncertainties, a linear approximation maybe sufficient
- A Taylor expansion around the central values of θ gives, using the error matrix Θ_{ij} :

$$H_{ij} = \sum_{k,l} \frac{\partial \eta_i}{\partial \theta_k} \frac{\partial \eta_j}{\partial \theta_l} \Theta_{kl}$$

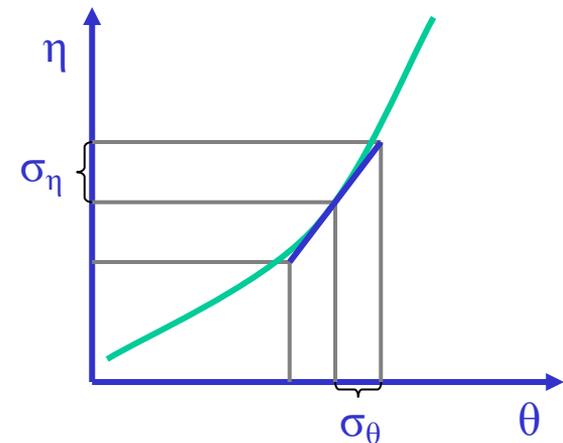
- Few examples in case of **no correlation**:

$$\sigma_{x+y} = \sigma_{x-y} = \sqrt{\sigma_x^2 + \sigma_y^2}$$

$$\frac{\sigma_{xy}}{xy} = \frac{\sigma_{x/y}}{x/y} = \sqrt{\left(\frac{\sigma_x}{x}\right)^2 + \left(\frac{\sigma_y}{y}\right)^2}$$

$$\sigma_{x^2} = 2x\sigma_x$$

$$\sigma_{\ln x} = \frac{\sigma_x}{\sqrt{x}}$$



- Sometimes data are available as **binned** histogram
 - Most often each bin obeys **Poissonian statistics** (event counting)
- The likelihood function is the product of Poisson PDFs corresponding to each bin having entries n_i
- The expected number of entries n_i depends on some unknown parameters: $\mu_i = \mu_i(\theta_1, \dots, \theta_m)$
- The function to minimize is the following $-2 \ln L$:

$$\begin{aligned} -2 \ln L &= -2 \ln \prod_{i=1}^{n_{\text{bins}}} \text{Poiss}(n_i; \mu_i(\theta_1, \dots, \theta_m)) \\ &= -2 \ln \prod_{i=1}^{n_{\text{bins}}} \frac{e^{-\mu_i(\theta_1, \dots, \theta_m)} \mu_i(\theta_1, \dots, \theta_m)^{n_i}}{n_i!} \end{aligned}$$

- The expected number of entries μ_i is often **approximated** by a **continuous function** $\mu(x)$ evaluated at the center x_i of the bin
- Alternatively, μ_i can be a combination of other histograms (“templates”)
 - E.g.: sum of different **simulated processes** with floating **yields** as fit parameters

- Bin entries can be approximated by Gaussian variables for sufficiently **large number of entries** with standard deviation equal to n_i (**Neyman's χ^2**)
- Maximizing L is equivalent to minimize:

$$\chi^2 = \sum_{i=1}^{n_{\text{bins}}} \frac{(n_i - \mu(x_i; \theta_1, \dots, \theta_m))^2}{n_i}$$

- Sometimes, the denominator n_i is replaced (**Pearson's χ^2**) by:

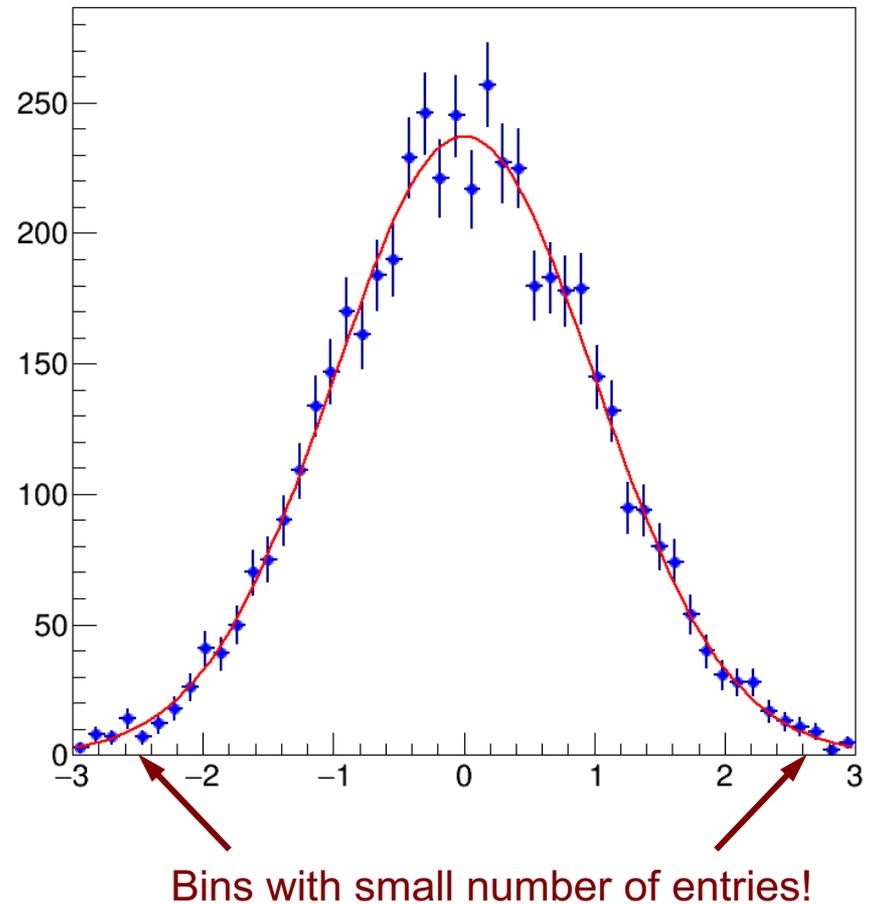
$$\mu_i = \mu(x_i; \theta_1, \dots, \theta_m)$$

in order to avoid cases with zero or small n_i

- Analytic solution exists for linear and other simple problems
 - E.g.: **linear fit model**
- Most of the cases are treated numerically, as for unbinned ML fits

- Binned fits are convenient w.r.t. unbinned fits because the **number of input variables decreases from the number of entries to the number of bins**
 - Usually **simpler and faster** numerically
 - Unbinned fits become unpractical for very large number of entries
- A fraction of the information is lost, hence a possible **loss of precision** may occur for small number of entries
- **Treat correctly bins with small number of entries!**

Gaussian fit (determine yield, μ and σ)



- The maximum value of the likelihood function obtained from the fit doesn't usually give information about the goodness of the fit
- The χ^2 of a fit with a Gaussian underlying model is distributed according to a known PDF

$$P(\chi^2; n) = \frac{2^{-\frac{n}{2}}}{\Gamma(\frac{n}{2})} \chi^{n-2} e^{-\frac{\chi^2}{2}}$$

n is the number of degrees of freedom (n. of bins – n. of params.)

- The cumulative distribution of $P(\chi^2; n)$ follows a uniform distribution between 0 and 1 (p -value)
- If the model deviates from the assumed distribution, the distribution of the p -value will be more peaked around zero
- **Note!** p -values are not the “probability of the fit hypothesis”
 - This would be a Bayesian probability, with a different meaning, and should be computed in a different way

- A better alternative to the (Gaussian-inspired, Neyman and Pearson's) χ^2 has been proposed by Baker and Cousins using the following **likelihood ratio**:

$$\begin{aligned} \chi_\lambda^2 &= -2 \ln \prod_i \frac{L(n_i; \mu_i)}{L(n_i; n_i)} = -2 \ln \prod_i \frac{e^{-\mu_i} \mu_i^{n_i}}{e^{-n_i} n_i^{n_i}} \frac{\cancel{n_i!}}{\cancel{n_i!}} \\ &= 2 \sum_i \left[\mu_i(\theta_1, \dots, \theta_m) - n_i + n_i \ln \left(\frac{n_i}{\mu_i(\theta_1, \dots, \theta_m)} \right) \right] \end{aligned}$$

- Same minimum value as from Poisson likelihood function, since a constant term has been added to the log-likelihood function
- In addition, it **provides goodness-of-fit information**, and asymptotically **obeys chi-squared distribution** with $n - m$ degrees of freedom
(Wilks' theorem, see following slides)

S. Baker, R. Cousins NIM 221 (1984) 437

- Assume two measurements with different **uncorrelated** (Gaussian) errors: $m_1 \pm \sigma_1, m_2 \pm \sigma_2$

- Build the χ^2 :
$$\chi^2 = \frac{(m - m_1)^2}{\sigma_1^2} + \frac{(m - m_2)^2}{\sigma_2^2}$$

- Minimize the χ^2 :
$$0 = \frac{\partial \chi^2}{\partial m} = 2 \frac{(m - m_1)}{\sigma_1^2} + 2 \frac{(m - m_2)}{\sigma_2^2}$$

- Estimate m as:
$$m = \frac{\frac{m_1}{\sigma_1^2} + \frac{m_2}{\sigma_2^2}}{\frac{1}{\sigma_1^2} + \frac{1}{\sigma_2^2}} = \frac{w_1 m_1 + w_2 m_2}{w_1 + w_2}$$

← Weighted average, $w_i = \sigma_i^{-2}$

- Error estimate:
$$\frac{1}{\sigma_m^2} = -\frac{\partial^2 \ln L}{\partial m^2} = \frac{1}{2} \frac{\partial^2 \chi^2}{\partial m^2} = \frac{1}{\sigma_1^2} + \frac{1}{\sigma_2^2}$$

- Correlation coefficient $\rho \neq 0$:

$$m_1 \pm \sigma_1, \quad m_2 \pm \sigma_2$$

- Build χ^2 including correlation terms:

$$\chi^2 = \begin{pmatrix} m - m_1 & m - m_2 \end{pmatrix} \begin{pmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{pmatrix}^{-1} \begin{pmatrix} m - m_1 \\ m - m_2 \end{pmatrix}$$

- The χ^2 minimization gives:

$$m = \frac{m_1(\sigma_2^2 - \rho\sigma_1\sigma_2) + m_2(\sigma_1^2 - \rho\sigma_1\sigma_2)}{\sigma_1^2 - 2\rho\sigma_1\sigma_2 + \sigma_2^2}$$

$$\sigma_m^2 = \frac{\sigma_1^2\sigma_2^2(1 - \rho)^2}{\sigma_1^2 - 2\rho\sigma_1\sigma_2 + \sigma_2^2}$$

a.k.a “**BLUE**”:
Best Linear Unbiased Estimator

- The “common error” σ_C is defined as: $\sigma_C^2 = \rho\sigma_1\sigma_2$
- Using error propagation, this also implies that:

$$\sigma_{m_1 - m_2}^2 = \left(\frac{\partial(m_1 - m_2)}{\partial m_1} \right)^2 \sigma_1^2 + \left(\frac{\partial(m_1 - m_2)}{\partial m_2} \right)^2 \sigma_2^2 + 2 \left(\frac{\partial(m_1 - m_2)}{\partial m_1} \right) \left(\frac{\partial(m_1 - m_2)}{\partial m_2} \right) \rho\sigma_1\sigma_2$$



$$\sigma_{m_1 - m_2}^2 = (\sigma_1^2 - \sigma_C^2) + (\sigma_2^2 - \sigma_C^2)$$

- The previous formulas can be written as a weighted average:

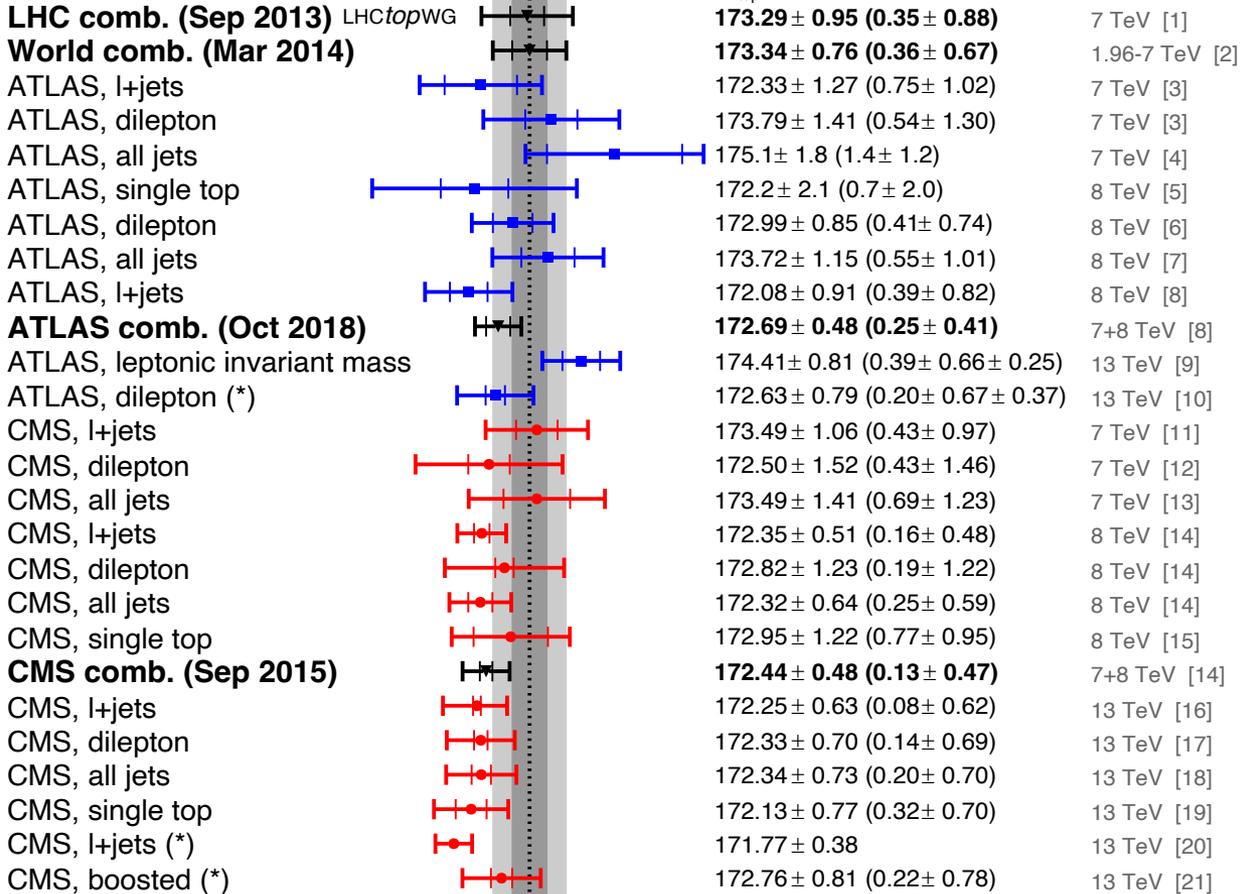
$$m = \frac{\frac{m_1}{\sigma_1^2 - \sigma_C^2} + \frac{m_2}{\sigma_2^2 - \sigma_C^2}}{\frac{1}{\sigma_1^2 - \sigma_C^2} + \frac{1}{\sigma_2^2 - \sigma_C^2}} = \frac{w_1 m_1 + w_2 m_2}{w_1 + w_2}$$

← Note: weights may be negative!

$$\sigma_m^2 = \frac{1}{\frac{1}{\sigma_1^2 - \sigma_C^2} + \frac{1}{\sigma_2^2 - \sigma_C^2}} + \sigma_C^2$$

..... World comb. (Mar 2014) [2]
 ■ stat
 ■ total uncertainty

total stat



* Preliminary

- | | | |
|-------------------------|---------------------------|-------------------------|
| [1] ATLAS-CONF-2013-102 | [8] EPJC 79 (2019) 290 | [15] EPJC 77 (2017) 354 |
| [2] arXiv:1403.4427 | [9] arXiv:2209.00583 | [16] EPJC 78 (2018) 891 |
| [3] EPJC 75 (2015) 330 | [10] ATLAS-CONF-2022-058 | [17] EPJC 79 (2019) 368 |
| [4] EPJC 75 (2015) 158 | [11] JHEP 12 (2012) 105 | [18] EPJC 79 (2019) 313 |
| [5] ATLAS-CONF-2014-055 | [12] EPJC 72 (2012) 2202 | [19] arXiv:2108.10407 |
| [6] PLB 761 (2016) 350 | [13] EPJC 74 (2014) 2758 | [20] CMS-PAS-TOP-20-008 |
| [7] JHEP 09 (2017) 118 | [14] PRD 93 (2016) 072004 | [21] CMS-PAS-TOP-21-012 |

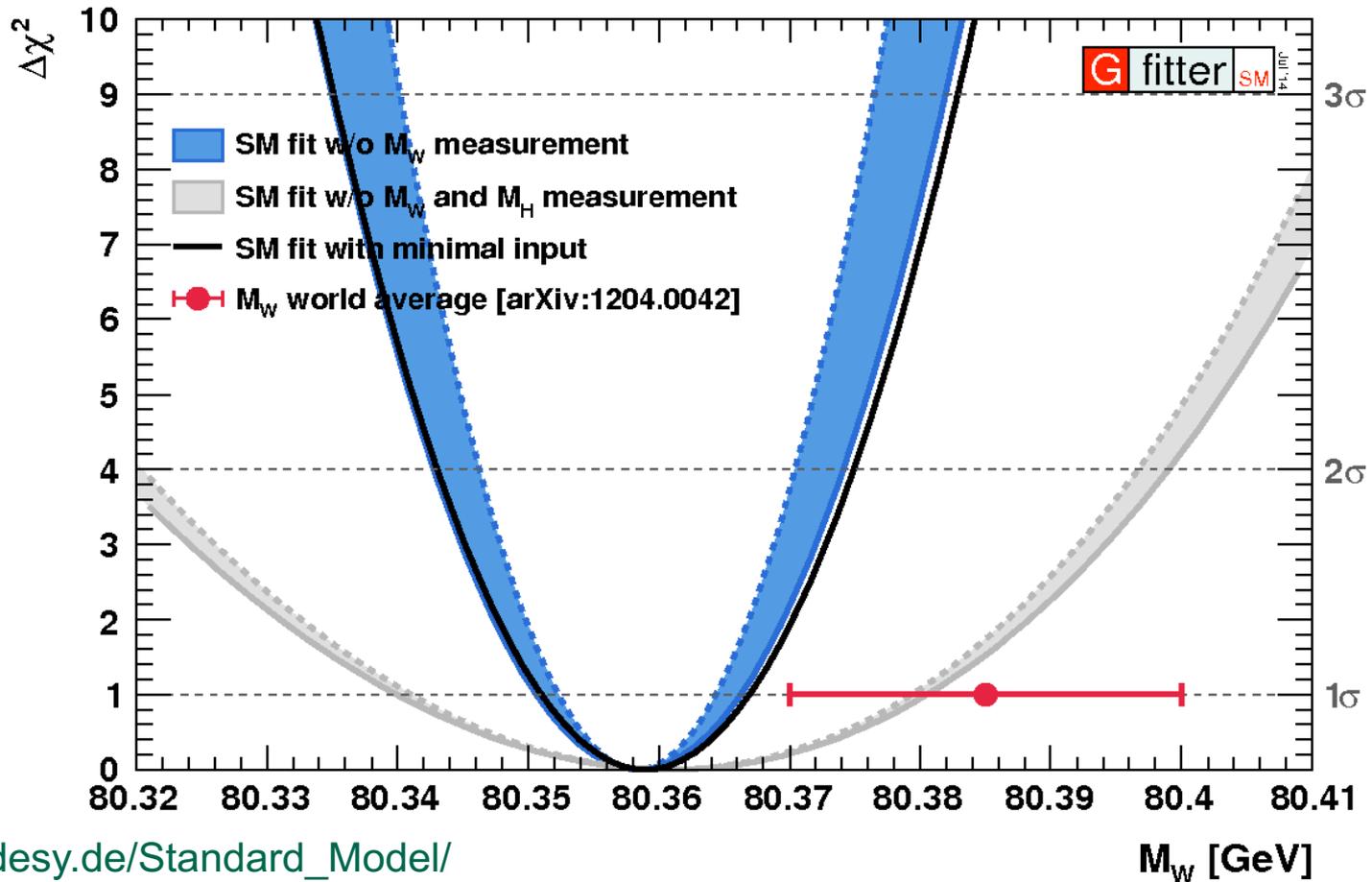
JE
on
and

165 170 175 180 185
 m_{top} [GeV]

- We have n measurements, (m_1, \dots, m_n) with a $n \times n$ covariance matrix (C_{ij})
- Expected values for m_1, \dots, m_n , M_1, \dots, M_n may depend on some theory parameter(s) θ
- The following χ^2 can be minimized to have an estimate of the parameter(s) θ :

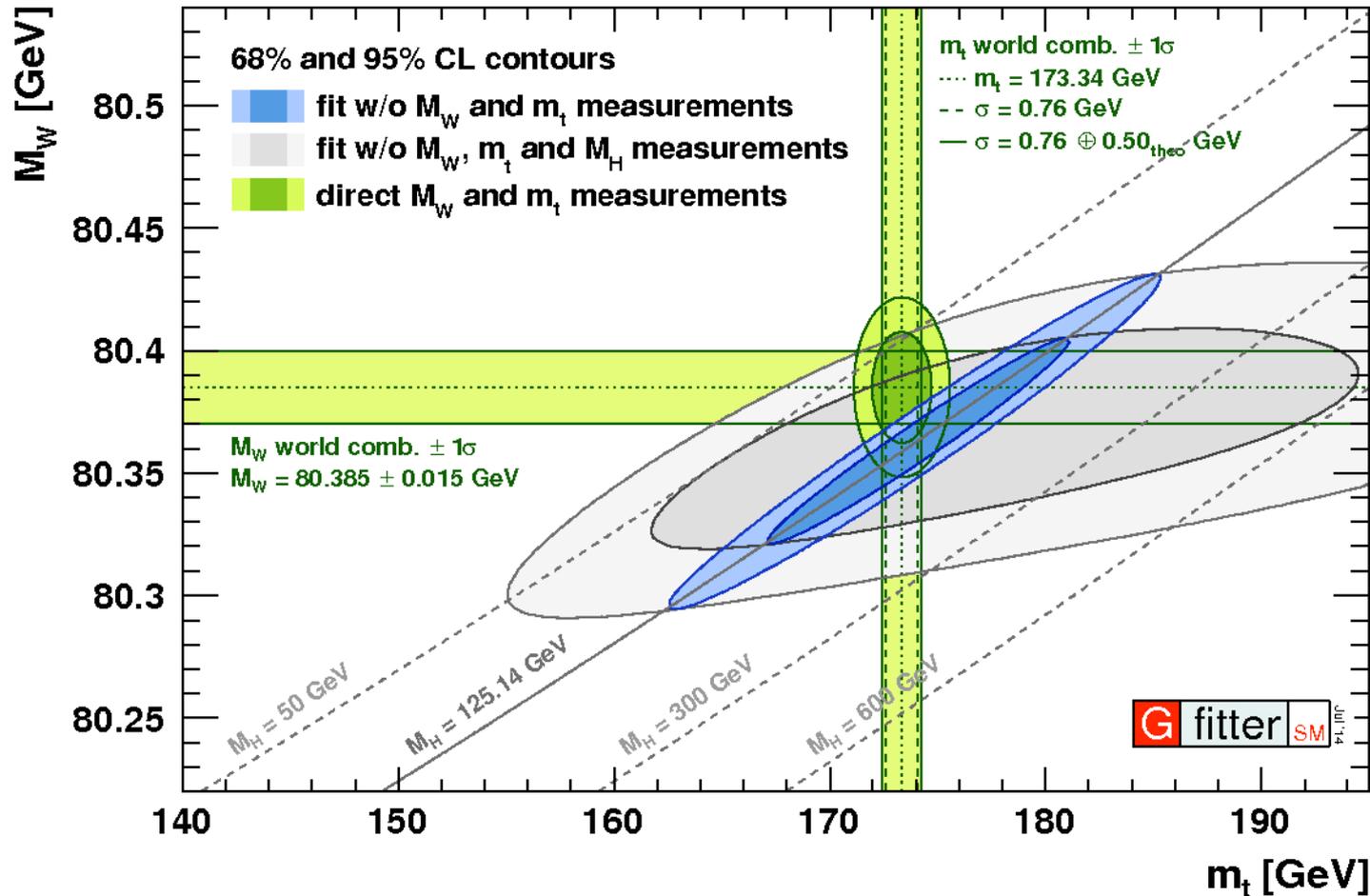
$$\begin{aligned}\chi^2 &= \sum_{i,j=1}^n (m_i - M_i(\theta)) C_{ij}^{-1} (m_j - M_j(\theta)) \\ &= \begin{pmatrix} m_1 - M_1, & \dots, & m_n - M_n \end{pmatrix} \begin{pmatrix} C_{11} & \dots & C_{1n} \\ \vdots & \ddots & \vdots \\ C_{n1} & \dots & C_{nn} \end{pmatrix}^{-1} \begin{pmatrix} m_1 - M_1 \\ \dots \\ m_n - M_n \end{pmatrix} \\ &= (\mathbf{m} - \mathbf{M}(\theta))^T \mathbf{C}^{-1} (\mathbf{m} - \mathbf{M}(\theta))\end{aligned}$$

- A Global χ^2 fit to electroweak measurements predicts the W mass allowing a comparison with direct measurements



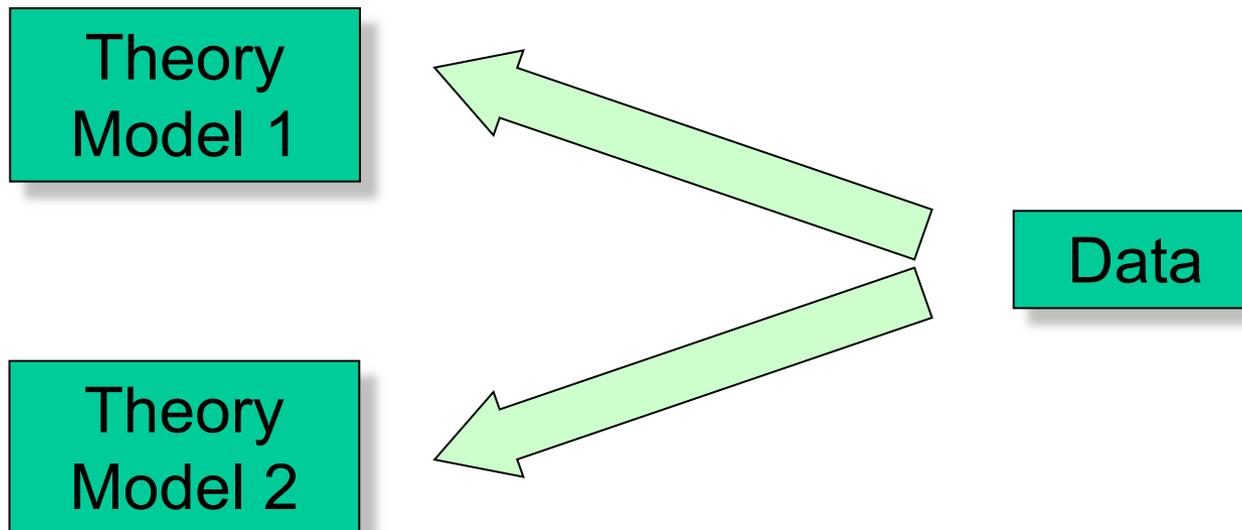
Details on:
http://gfitter.desy.de/Standard_Model/

- W mass vs top-quark mass from global electroweak fit





- Hypothesis testing
 - The ROC curve
 - The Neyman-Pearson lemma
 - Multivariate discrimination
 - Projective likelihood ratio
 - Fisher discriminant
 - Introduction to machine learning
 - Artificial Neural Networks
 - Boosted decision trees
 - The bias-variance tradeoff
 - Issues with machine learning: underfitting, overfitting
-



Which hypothesis is the most consistent with the experimental data?



- Bayesian probability gives meaning to the probability that an hypothesis is true:

$$P(H_1|x) = \frac{P(x|H_1)\pi(H_1)}{P(x)}$$

- The evaluation of $P(x)$ requires the decomposition over all possible hypotheses:

$$P(x) = P(x|H_0)\pi(H_0) + P(x|H_1)\pi(H_1) + \dots$$

- The ratio of probabilities for two hypothesis does not depend on $P(x)$, and can be computed without considering all possible hypotheses:

$$\frac{P(H_1|x)}{P(H_0|x)} = \frac{P(x|H_1)\pi(H_1)}{P(x|H_0)\pi(H_0)}$$



Bayes factors



- It is possible to introduce Bayes factor:

$$\frac{P(H_1|x)}{P(H_0|x)} = B_{1/0}(x) \frac{\pi(H_1)}{\pi(H_0)}$$

- In other words, this defines the posterior odds as a function of prior odds:

$$O_{1/0}(x) = B_{1/0}(x) o_{1/0}$$

- In word:

posterior odds = Bayes factor times prior odds

- Bayes factor can be used to measure how favoured is one hypothesis against another, and, in the simplest cases, it is equal to the likelihood ratio
- Typical range values are:
 - 1-3: very weak evidence
 - 3-20: positive evidence
 - 20-150: strong evidence
 - >150: very strong evidence

$$P(H_1, \theta_1 | x) = \frac{P(x | H_1, \theta_1) \pi(H_1, \theta_1)}{P(x)}$$

$$P(H_1 | x) = \frac{\int P(x | H_1, \theta_1) \pi(H_1, \theta_1) d\theta_1}{P(x)} = \frac{\pi(H_1) \int P(x | H_1, \theta_1) \pi(\theta_1 | H_1) d\theta_1}{P(x)}$$

- Because one has:

$$\begin{aligned} \pi(H_0, \theta_0) &= \pi(\theta_0 | H_0) \pi(H_0) \\ \pi(H_1, \theta_1) &= \pi(\theta_1 | H_1) \pi(H_1) \end{aligned}$$

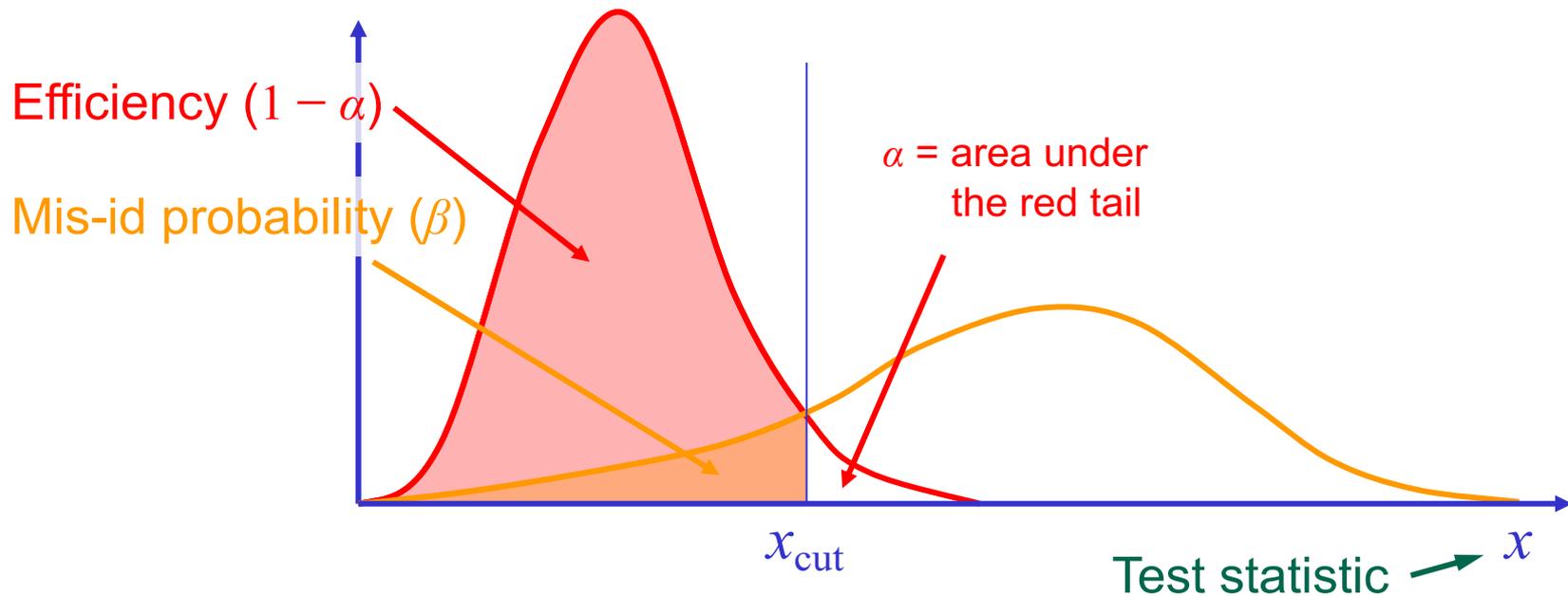
- The Bayes factor should be defined as:

$$B_{1/0}(x) = \frac{P(x | H_1)}{P(x | H_0)}$$

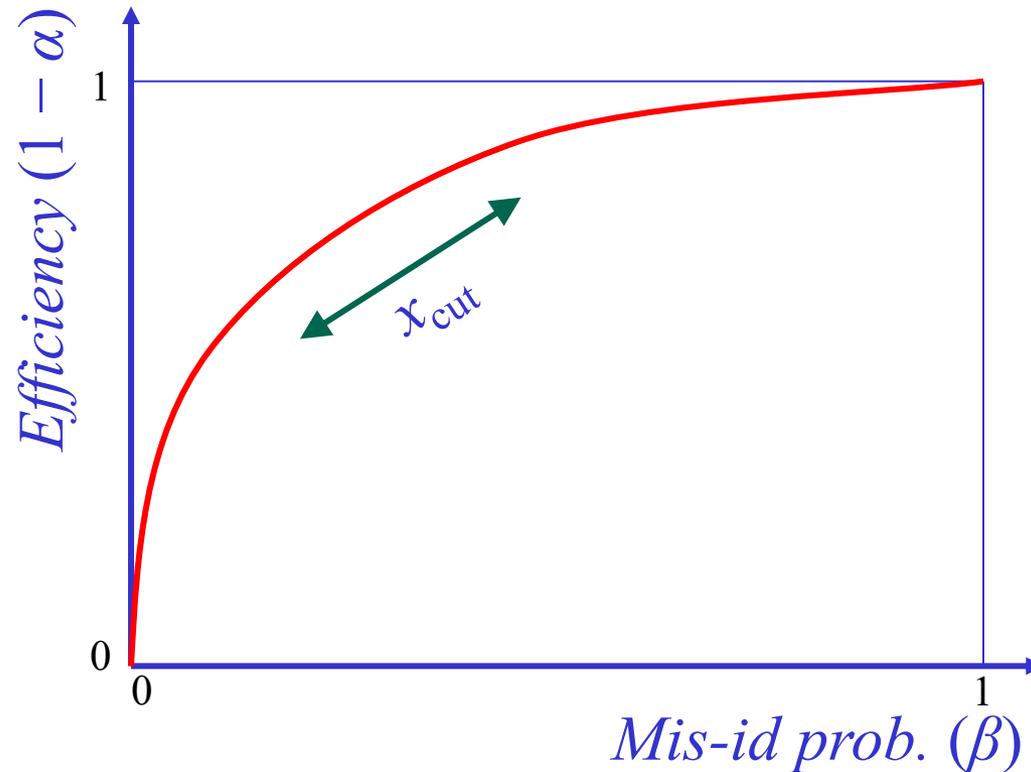
- I.e.: it is the ratio of marginal likelihood for x :

$$B_{1/0}(x) = \frac{P(x | H_1)}{P(x | H_0)} = \frac{\int P(x | H_1, \theta_1) \pi(\theta_1 | H_1) d\theta_1}{\int P(x | H_0, \theta_0) \pi(\theta_0 | H_0) d\theta_0}$$

- Selection (“cut”) on one (or more) variable(s):
 - If $x \leq x_{\text{cut}} \Rightarrow$ **signal**
 - Else, if $x > x_{\text{cut}} \Rightarrow$ **background**

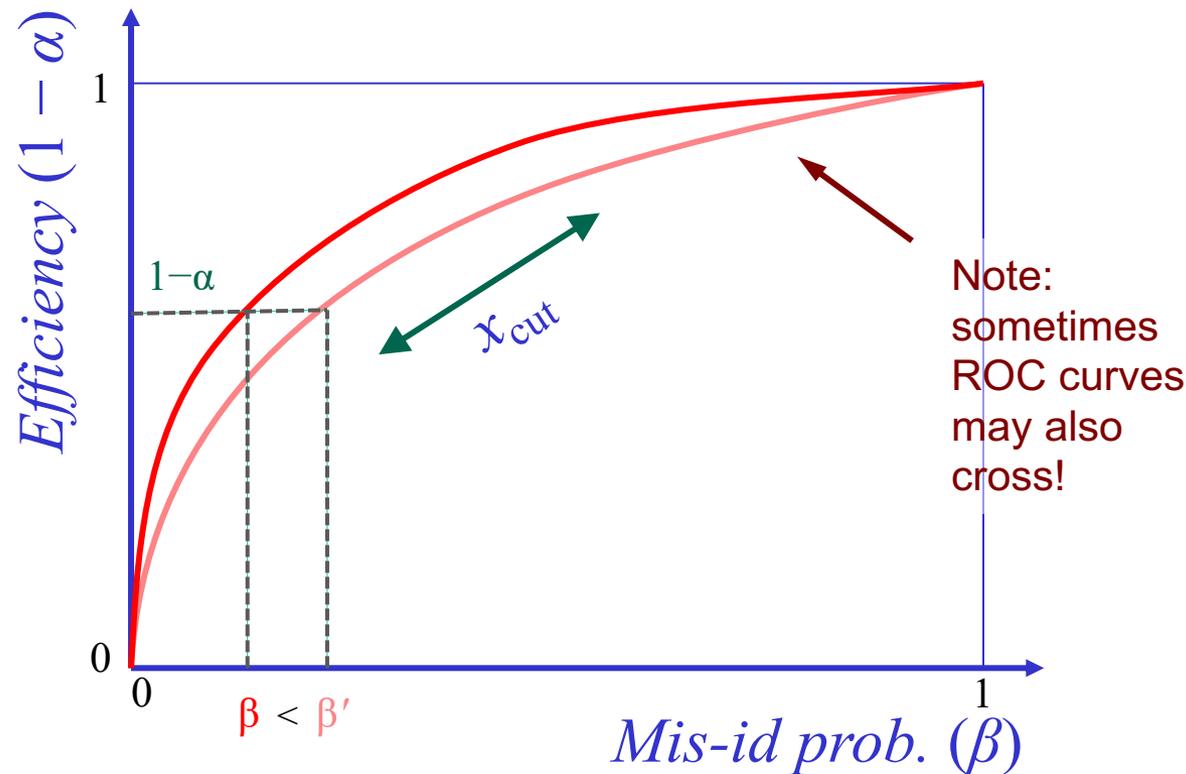


- Varying the applied cut on the **test statistic** both the efficiency and mis-id probability change



Sometimes also referred to as **ROC curve** (*Receiver Operating Characteristic*)

- One test is preferable to another if, for the same level of efficiency ($1 - \alpha$), it has lower mis-id probability (β)

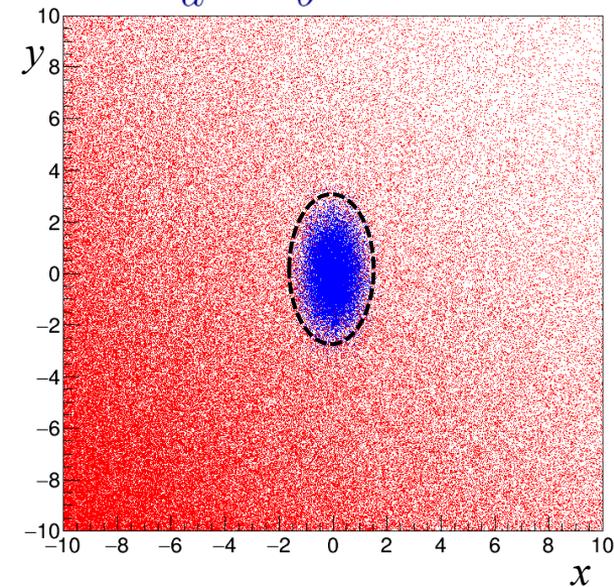
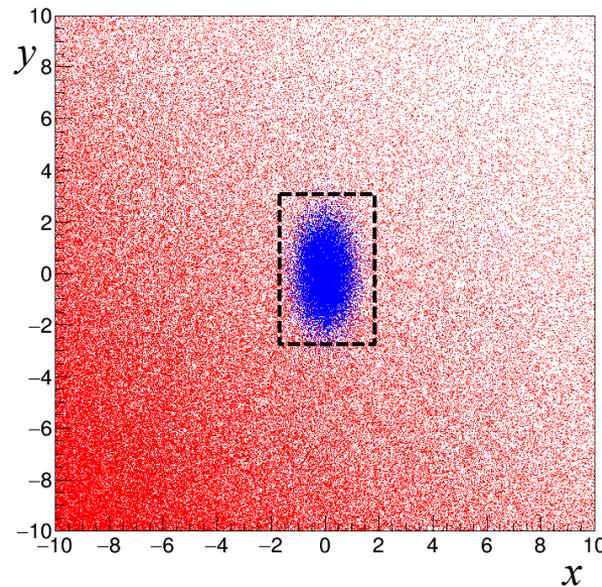
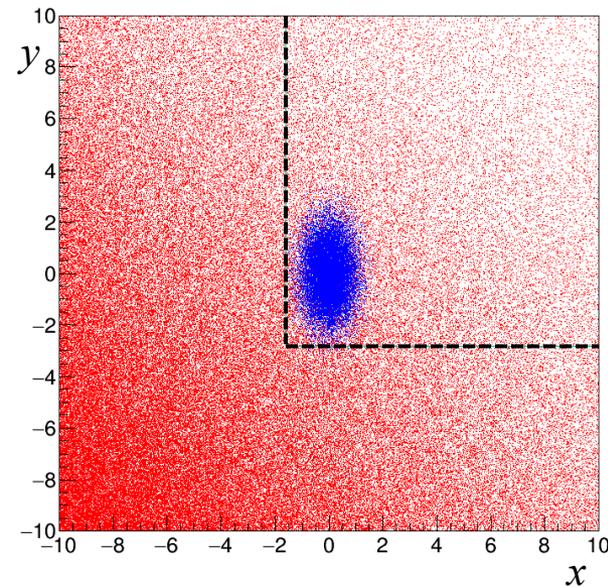


- In case of multiple discriminating variables (**multivariate**), the choice of the optimal selection is not always straightforward

$$x > x_{\text{cut}} \text{ and } y > y_{\text{cut}}$$

$$x_1 < x < x_2 \text{ and } y_1 < y < y_2$$

$$\frac{x^2}{a^2} + \frac{y^2}{b^2} < k^2$$



$$P_s(x, y) = \text{Gauss}(x; 0, \sigma_x) \times \text{Gauss}(y; 0, \sigma_y), \quad P_b(x, y) = \alpha e^{-\alpha x} \times \beta e^{-\beta y}$$

- In many cases it's convenient to find a single variable (**test statistic**) that 'summarizes' all the sample information

- Statisticians' terminology is sometimes not very natural for physics applications, but it has become popular among physicists as well:
- H_0 = **null hypothesis**
 - Ex. 1: *“a sample contains only background”*
 - Ex. 2: *“a particle is a pion”*
- H_1 = **alternative hypothesis**
 - Ex. 1: *“a sample contains background + signal”*
 - Ex. 2: *“a particle is a muon”*
- **Test statistic**: a variable computed from our sample that discriminates between the two hypotheses H_0 and H_1 . Usually a ‘summary’ of the information available in the sample
- α = **significance level**: probability to reject H_1 if H_0 is assumed to be true (error of first kind, false positive)
 - $\alpha = 1 -$ misidentification probability
- β = **misidentification probability**, i.e.: probability to reject H_0 if H_1 is assumed to be true (error of second kind, false negative)
 - $1 - \beta =$ **power of the test** = selection efficiency
- **p-value**: probability, assuming H_0 , of observing a result at least as extreme as the observed **test statistic**

- For a fixed significance level (α) or signal efficiency ($1 - \alpha$), a selection based on the **likelihood ratio** gives the lowest possible mis-id probability (β):

$$\lambda(x) = \frac{L(x|H_1)}{L(x|H_0)} > k_\alpha$$

- **The likelihood function can't always be determined exactly**
- If we can't determine the exact likelihood function, we can choose other discriminators as test statistics that **approximates** the exact likelihood
- **Neural Networks, Boosted Decision Trees** and other **machine-learning** algorithms are example of discriminators that may **closely approximate the performances of the exact likelihood ratio** approaching the Neyman-Pearson limit



- In general, when we consider algorithms that provide a test statistic for samples with multiple variables we talk about **multivariate discriminators**
 - Simple mathematical algorithms exist, as well as complex implementation based on extensive CPU computations
- In general, the algorithms are ‘**trained**’ using input samples whose nature is known (**training samples**)
 - I.e.: where H_0 or H_1 is known to be true
 - Example: use data samples simulated with computer algorithms (**Monte Carlo**)
- Some of the most common problems:
 - **The size of samples is finite**, hence the true distributions for the considered hypotheses can't be determined exactly
 - The distribution of the input samples **does not reproduce exactly the true distribution** of real data (e.g.: systematic uncertainties)



- The likelihood function is **approximated** by the product of projective PDF in each variable

$$\lambda(x) = \frac{L(x_1, \dots, x_n | H_1)}{L(x_1, \dots, x_n | H_0)} \simeq \frac{\prod_{i=1}^n f_i(x_i | H_1)}{\prod_{i=1}^n f_i(x_i | H_0)}$$

x_1, \dots, x_n approximately
considered independent
variables

- Exact only in case of **independent variables**
- The approximation may be improved if the variables are first rotated in order to **eliminate correlation** (principal component analysis)
 - Find eigenvectors of the covariance matrix
 - Note:** uncorrelated variables are not necessarily independent

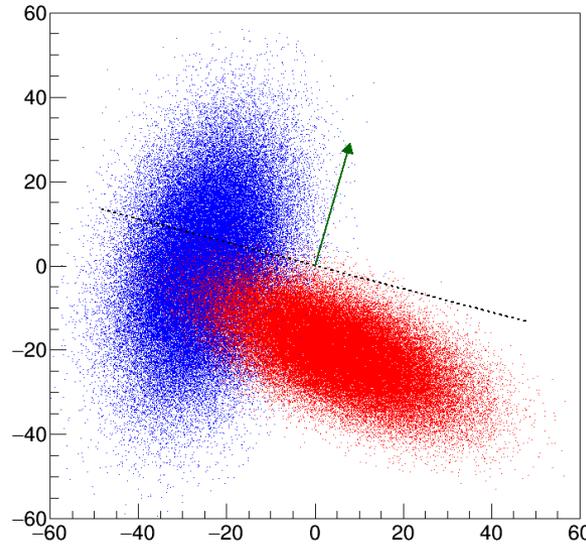
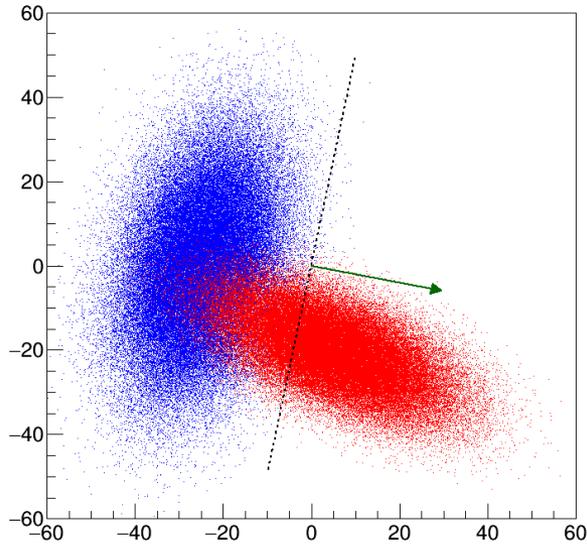
- Linear combination of input variables that maximizes the distance of the means of the two classes while minimizing the variance projected along a direction \mathbf{w} :

$$J(\mathbf{w}) = \frac{|\mu_0 - \mu_1|^2}{\sigma_0^2 + \sigma_1^2} = \frac{\mathbf{w}^T \cdot (\mathbf{m}_0 - \mathbf{m}_1)}{\mathbf{w}^T (\boldsymbol{\Sigma}_0 + \boldsymbol{\Sigma}_1) \mathbf{w}}$$



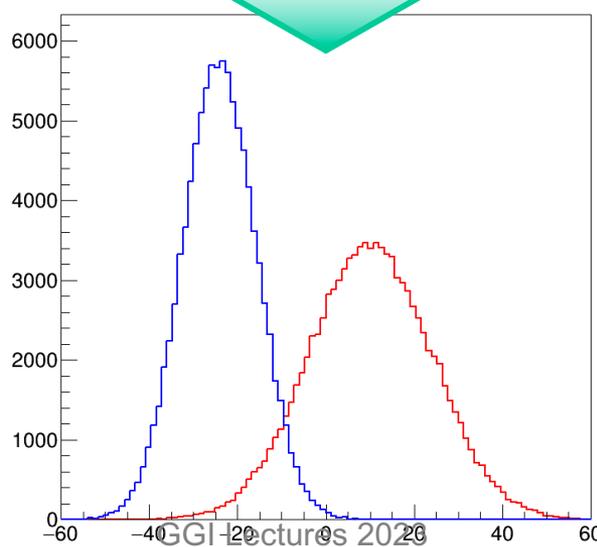
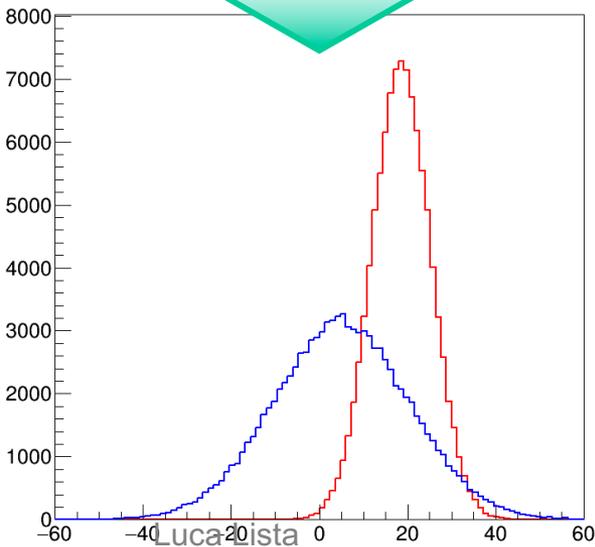
Sir Ronald Aylmer Fisher
(1890-1962)

- The selection is achieved by requiring $J(\mathbf{w}) > J_{\text{cut}}$, which determines an **hyperplane perpendicular to \mathbf{w}** that separates the two samples
- The maximization problem can be **solved analytically** using linear algebra

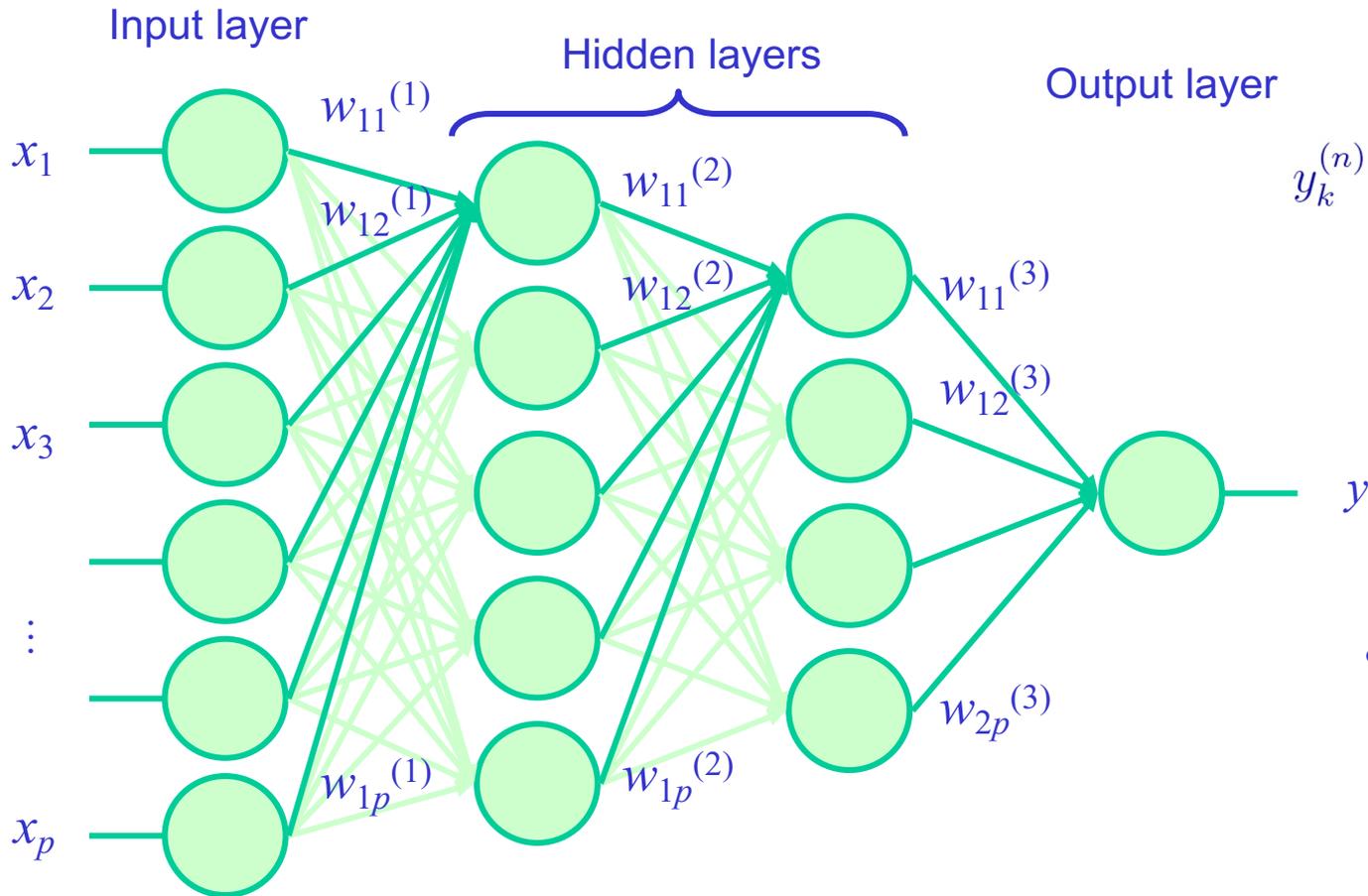


Projections along different directions achieve different overlap level

Maximum separation achieved by maximizing Fisher discriminant



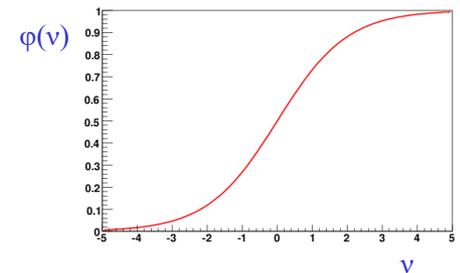
- Artificial simplified model of how neuron cells work: multilayer perceptron



$$y_k^{(n)}(\vec{x}) = \varphi \left(\sum_{j=1}^p w_{kj}^{(n)} x_j \right)$$

$\varphi(v)$ = Activation function

$$\varphi(v) = \frac{1}{1 + e^{-\lambda v}}$$



- Find the optimal set of network parameters $w_{ij}^{(n)}$ that minimize the “**loss function**” defined on a set of N training events:

$$L(w) = \sum_{i=1}^N (y_i^{true} - y(\vec{x}_i))^2$$

- Variation of loss function are also used, like cross-entropy, adopted more for binomial models with parameter $p_i = y_i^{true}$:

$$\begin{aligned} L(w) &= -\log \left(\prod_{i=1}^N y(\vec{x}_i)^{y_i^{true}} (1 - y(\vec{x}_i))^{1-y_i^{true}} \right) = \\ &= -\sum_{i=1}^N [y_i^{true} \log y(\vec{x}_i) - (1 - y_i^{true}) \log(1 - y(\vec{x}_i))] \end{aligned}$$

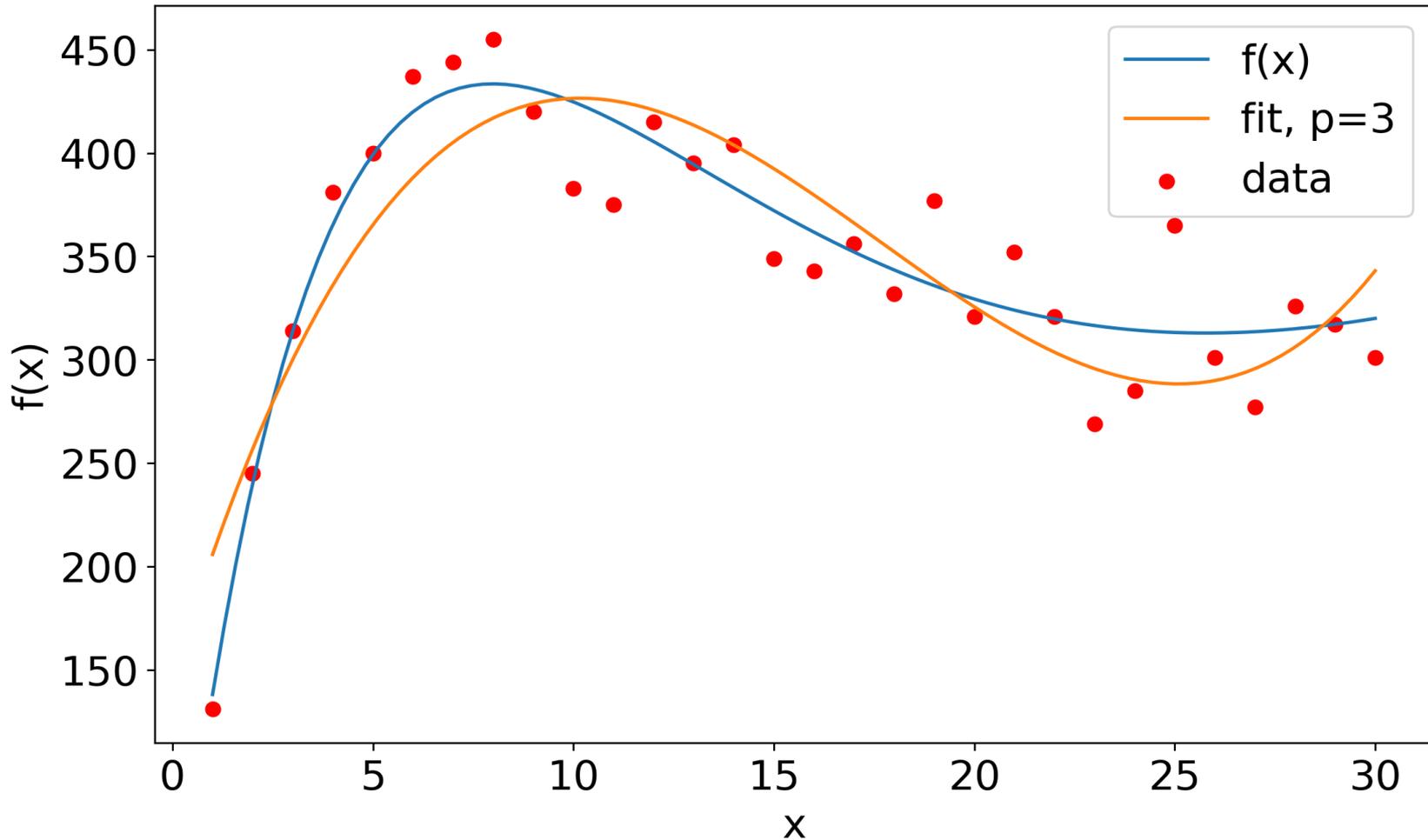
- Where $y_i^{true} = 1$ for signal (H_1), 0 for background (H_0)
- Usually achieved with **stochastic gradient descent**: weights are modified for each training event (**back propagation**):

$$w_{ij} \rightarrow w_{ij} - \eta \frac{\partial L(w)}{\partial w_{ij}}$$

- Artificial neural network with a **single** hidden layer may **approximate any analytical** function within a given approximation if the number of neurons is sufficiently high
- Demonstration in:
 - H. N. Mhaskar, Neural Computation, Vol. 8, No. 1, Pages 164-177 (1996), *Neural Networks for Optimal Approximation of Smooth and Analytic Functions*:
“We prove that neural networks with a single hidden layer are capable of providing an optimal order of approximation for functions assumed to possess a given number of derivatives, if the activation function evaluated by each principal element satisfies certain technical conditions”
- Anyway, the finite number of layer may lead to reaching this goal only approximately
- **Deep learning**: networks with several hidden layers
 - Can manage **complex variables combinations**, e.g.: exploiting invariant mass distributions using four-vectors as input!
 - Almost untreatable in the past, a lot of progress has been done recently, with better training algorithms and more easily available CPU power



- Let us take a polynomial fit of a points whose true model is not known
- The degree of the polynomial, i.e.: the number of parameter, can be chosen arbitrarily and determines the complexity of the model, and how it can flexibly adapt to the data
- By varying the number of parameter, one may improve the agreement of the fit curve to data at the cost of introducing a larger variance

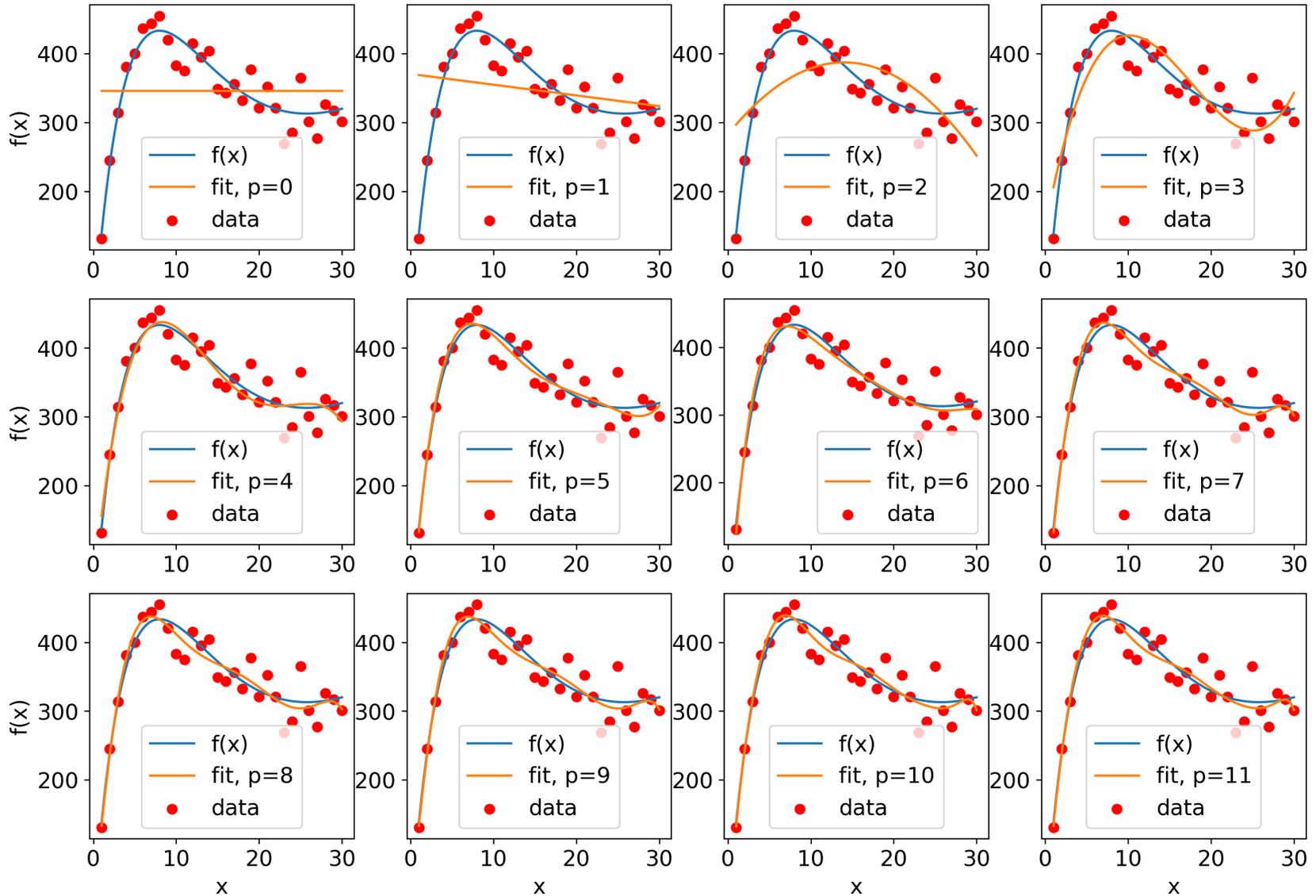




- Mean squared error:

$$\varepsilon^2 = \frac{1}{N} \sum_{i=1}^N (y_i^* - \hat{f}(x_i))^2$$

- By increasing the number of parameters, the mean squared error decreases for the fit sample, but may increase for independently extracted samples
- The fit follows the fluctuations in data more closely than the original distribution, that is not known (overfitting)





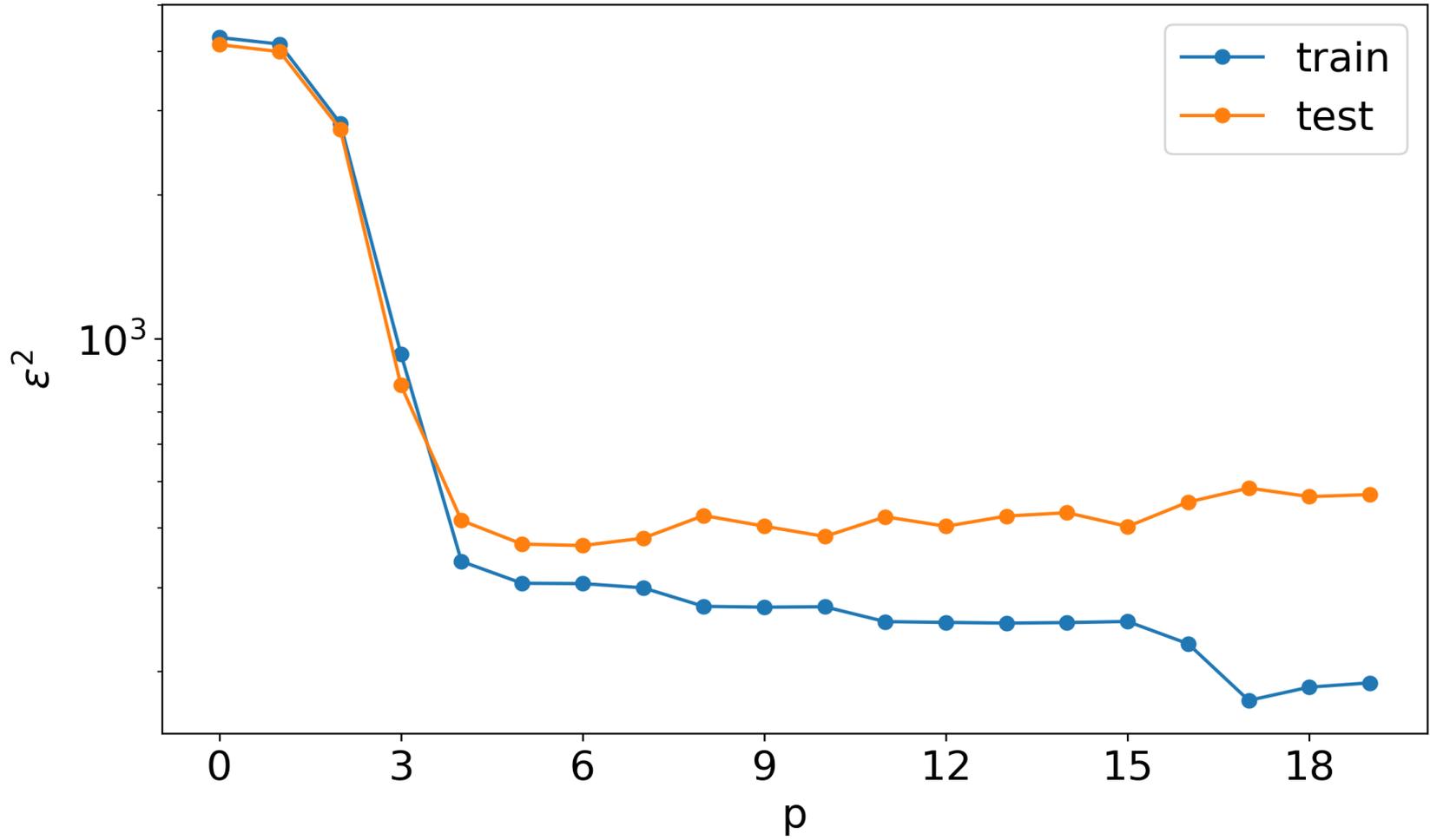
- Expected mean squared error:

$$\langle \varepsilon^2 \rangle = \mathbb{E} \left[\frac{1}{N} \sum_{i=1}^N (y_i^* - \hat{f}(x_i))^2 \right]$$

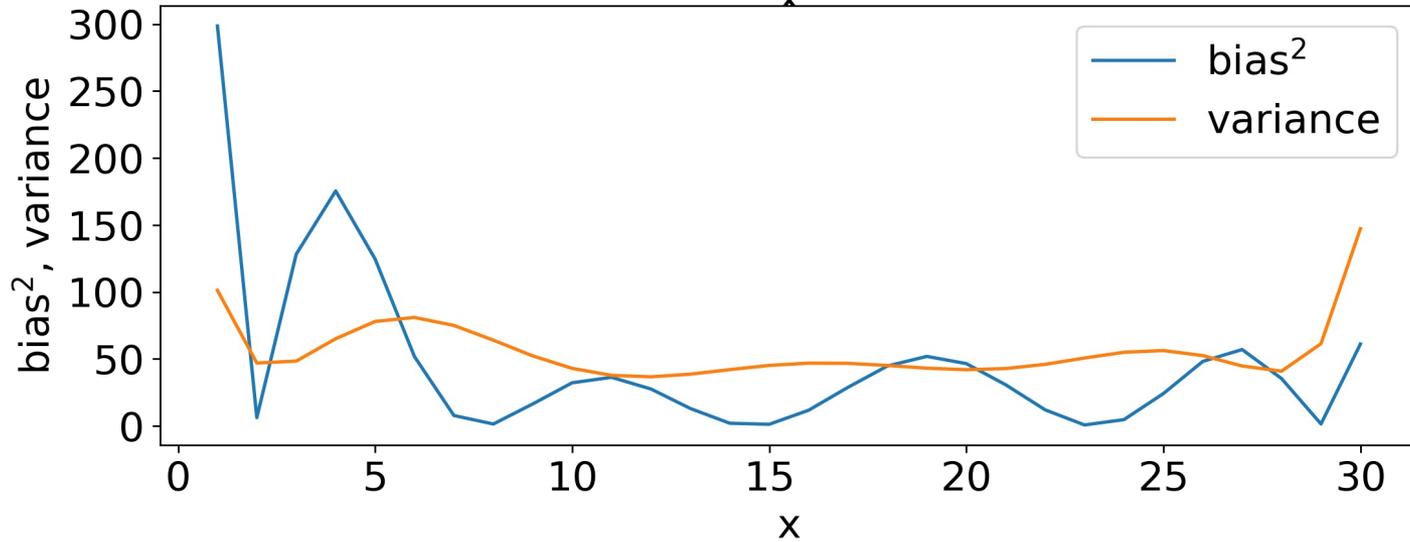
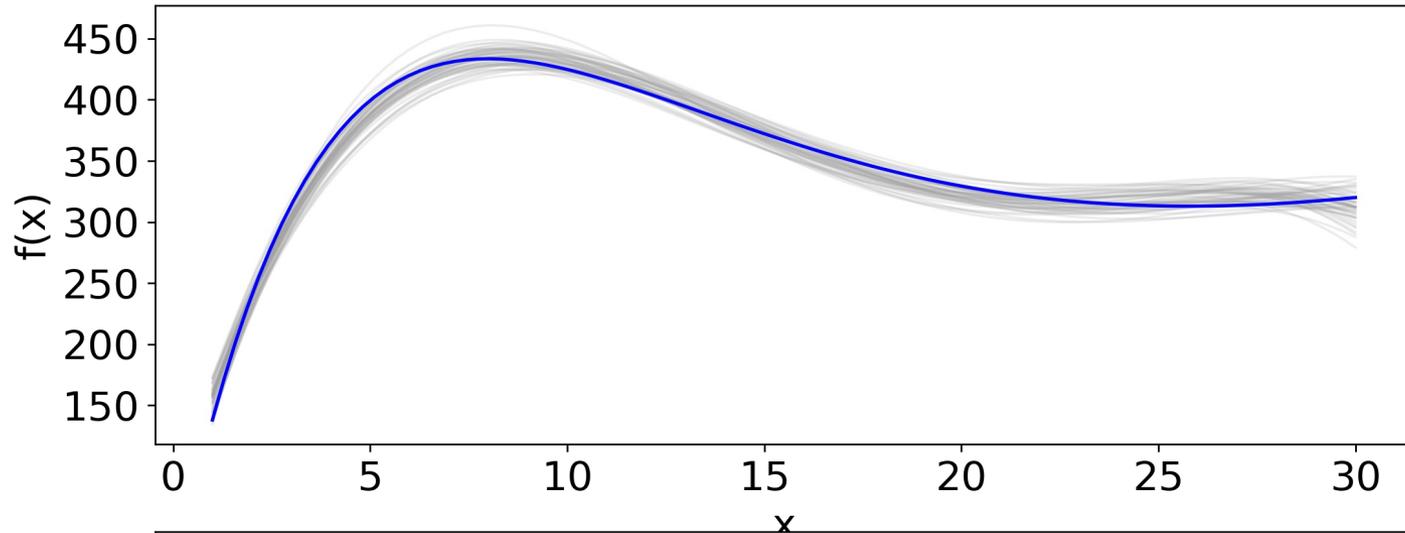
- It is possible to demonstrate that the following decomposition holds, if we use the notation $\hat{\theta} = \hat{f}(x_i)$, $\theta^* = y_i^*$:

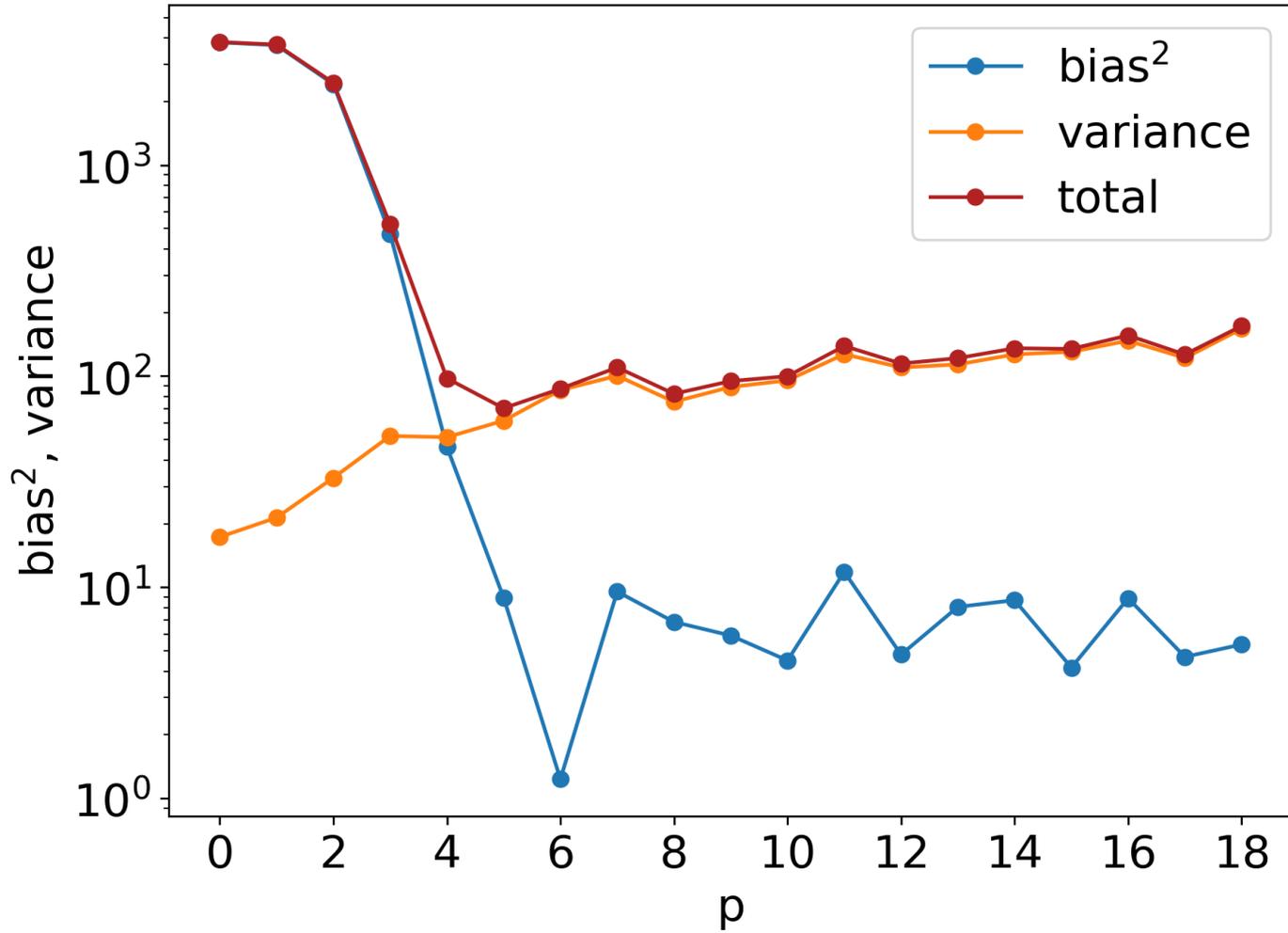
$$\langle \varepsilon^2 \rangle = \text{Var}[\hat{\theta}] + \text{Bias}[\hat{\theta}]^2 + \text{Var}[\theta^*]$$

- While the intrinsic noise of the data $\text{Var}[\theta^*]$ can't be improved with the fit, $\text{Var}[\hat{\theta}]$ and $\text{Bias}[\hat{\theta}]$ depend on the fit model, i.e.: on the number of parameters
- One cannot achieve at the same time optimum variance and optimum bias, and a trade off must be chosen

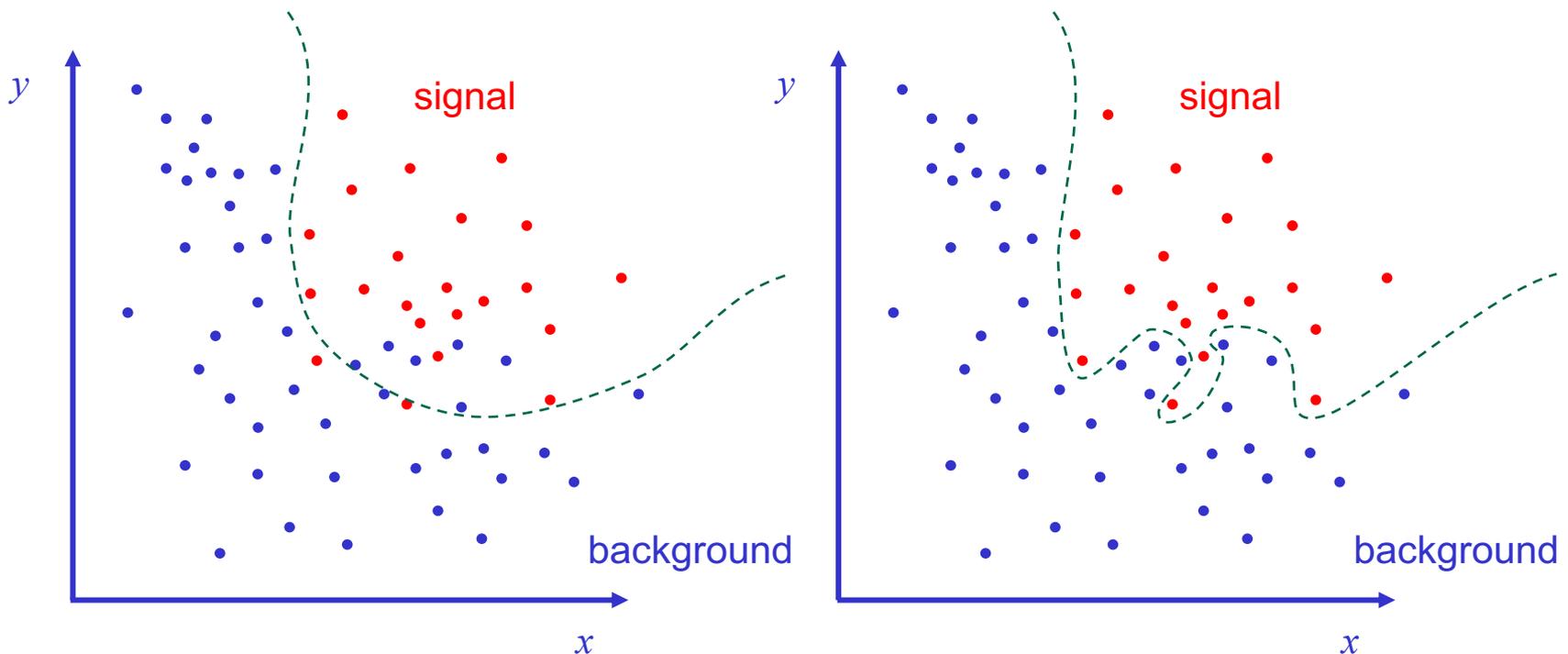


$p = 4$





- Algorithms may learn too much from the training sample, exploiting features that are only due to **random fluctuations**
- Check for **overtraining** comparing the discriminator's distributions for the **training sample** and for an independent **test sample** (consistent distributions = no overtraining)



INFN Decision Trees



- Cuts are applied sequentially
- Each cut splits the sample into **nodes**
 - Nodes where signal or background is largely dominant are classified as **leaves**
 - Alternatively: stop splitting if too few events per node, total number of nodes too high, ...
 - One branch = one sequence of cuts
- Cuts can be optimized to achieve the best split level: maximize for each node the **gain of Gini index**:

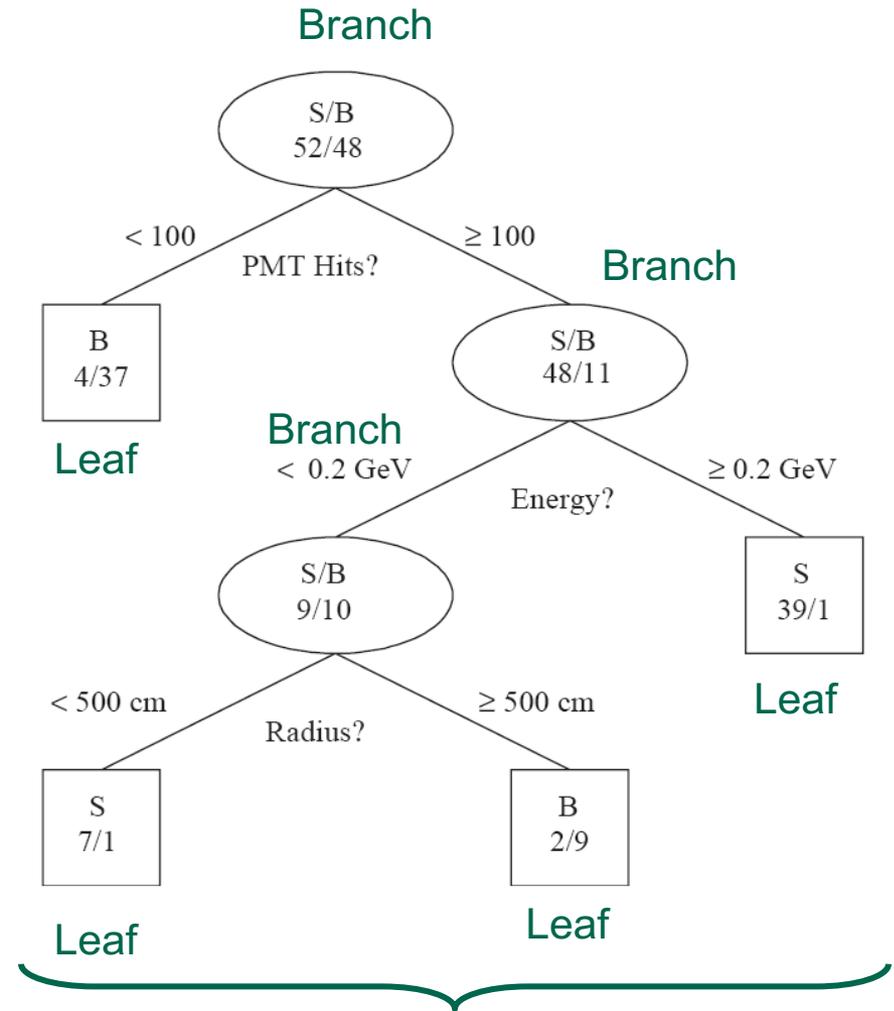
$$G = P(1 - P)$$

P = node purity

$G = 0$ for nodes with all S or all B events

$$\text{Gain} = N_{\text{parent}} G_{\text{parent}} - N_{\text{left ch.}} G_{\text{left ch.}} - N_{\text{right ch.}} G_{\text{right ch.}}$$

- Alternative metrics exist
E.g.: **cross entropy** = $-(P \ln P + (1-P) \ln(1-P))$, ...



Decision tree

- A single tree tends to adapt very closely to data, and performances are usually not very satisfactory
- It may easily provide overtraining if the tree is too deep
- It means: it has large variance, but in general low bias
- Combining many independent trees may reduce the variance
- Random forests are algorithms that combine many different trees and provide as output the combination of individual results



- Combine a large number of decision trees (**forest**) using different weights. Usually: $\mathcal{O}(1000)$ trees used
 - More performant and stable than a single optimized tree
 - **Boosting** achieved by **iteratively** reweighting training sample according to classifier output in previous iteration
1. **Reweight events** using previous iteration's classifier result
 2. **Build and optimize a new tree** with reweighted events
 3. Give a **score** to each tree
 4. The final **BDT classifier** result is the weighted average over all trees, using the given scores as weights:

$$y(\vec{x}) = \sum_{k=1}^{N_c} w_i C^{(i)}(\vec{x})$$

- Misclassified events are reweight according to the fraction of classification errors of the previous tree:

$$\frac{1 - f}{f}, \quad f = \frac{N_{\text{misclassified}}}{N_{\text{tot}}}$$

- Also use (log of) misclassification fraction as score for each tree:

$$y(\vec{x}) = \sum_{k=1}^{N_c} \log \left(\frac{1 - f^{(i)}}{f^{(i)}} \right) C^{(i)}(\vec{x})$$

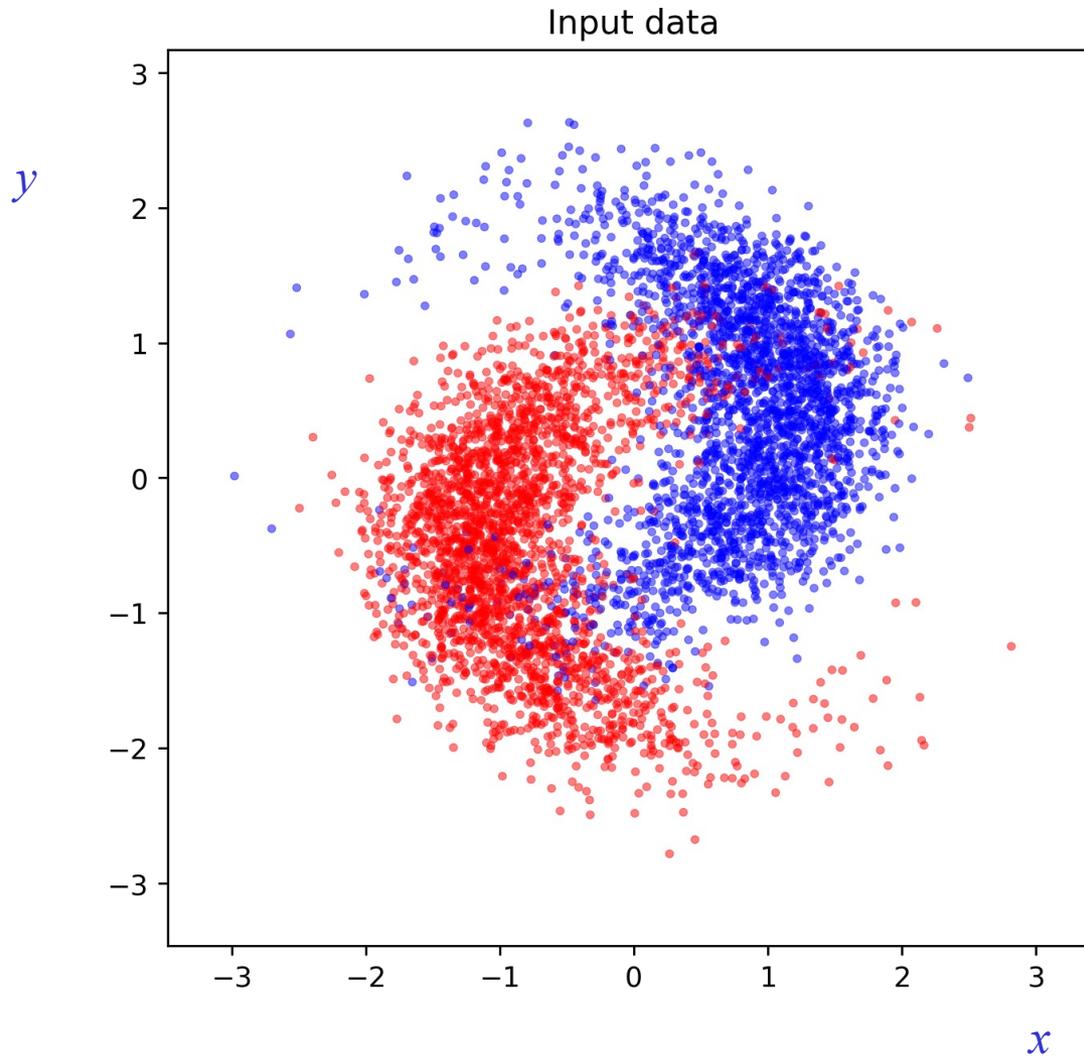
- Next iteration will better perform on events poorly classified in the previous iteration
- Further variations and more algorithms available

- Introducing a loss function for decision trees allows to use stochastic gradient descent algorithms
- Moreover, it also allows to use decision trees for regression, not only for classification
- The output of a tree can be taken as the index corresponding to one of the possible leaves, for a classification problem: $i = q(\vec{x})$
- For regression problems, we can assume that the output of a tree is a function of the index: $g = w_{q(\vec{x})}$ so that for many trees the output is the combination of their outputs:

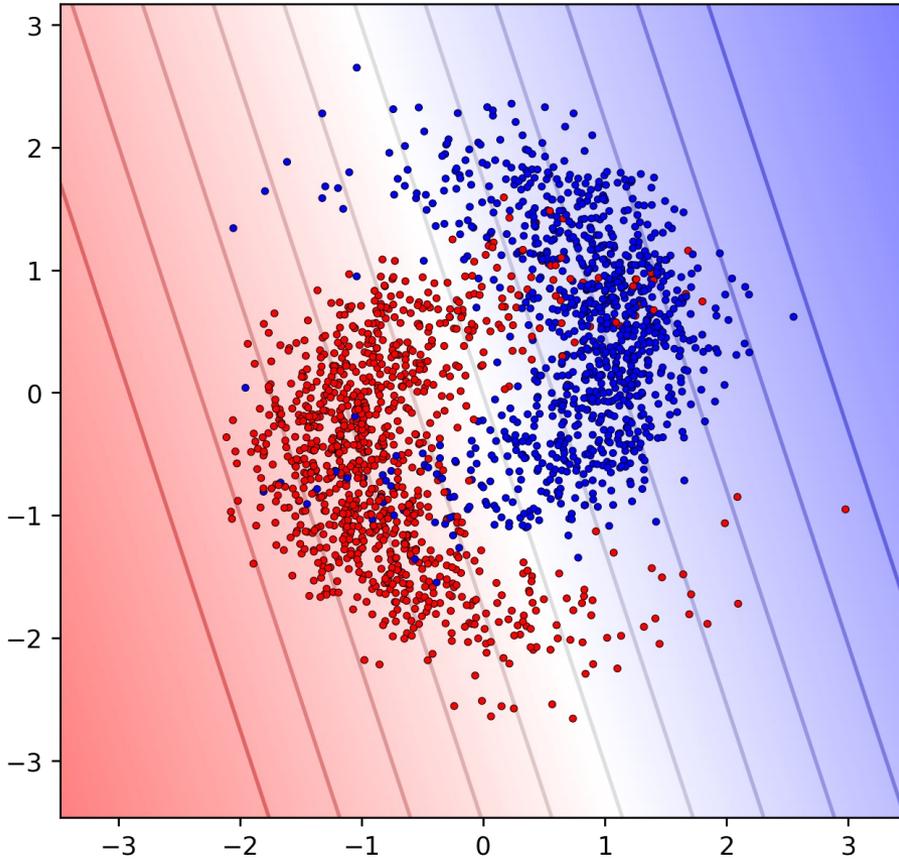
$$\hat{y}(\vec{x}) = \sum_{j=1}^M g^{(j)}(\vec{x}) = \sum_{j=1}^M w_{q(\vec{x})}^{(j)}$$

- The loss function may be the sum of

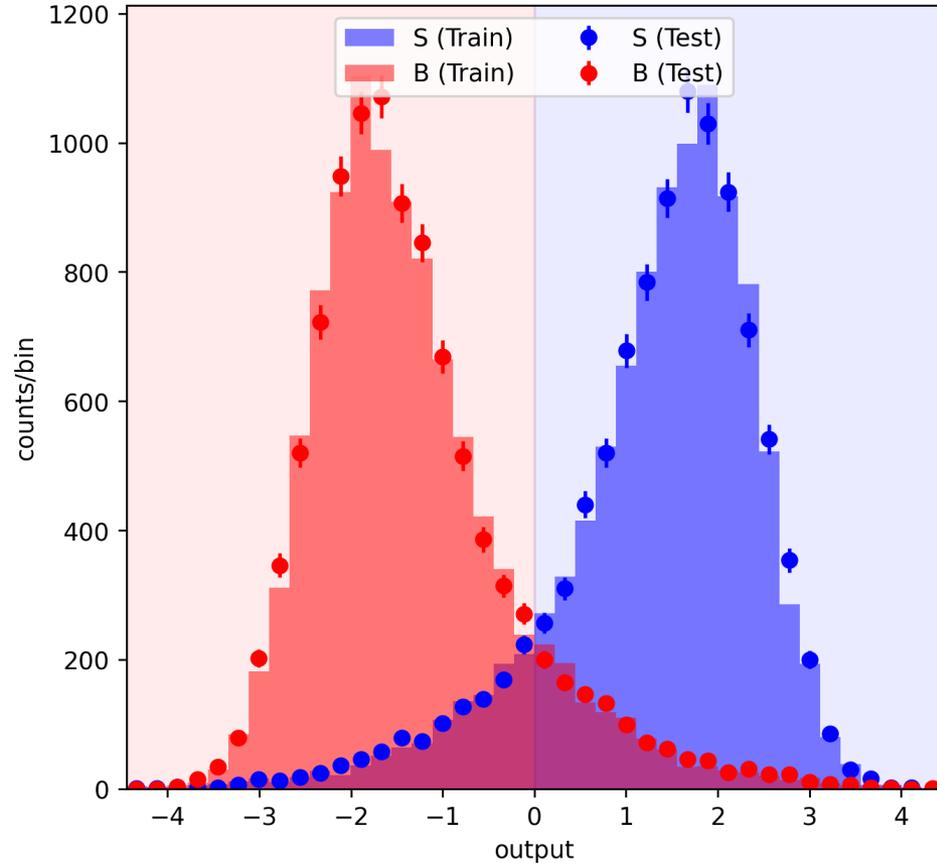
Strongly nonlinear example



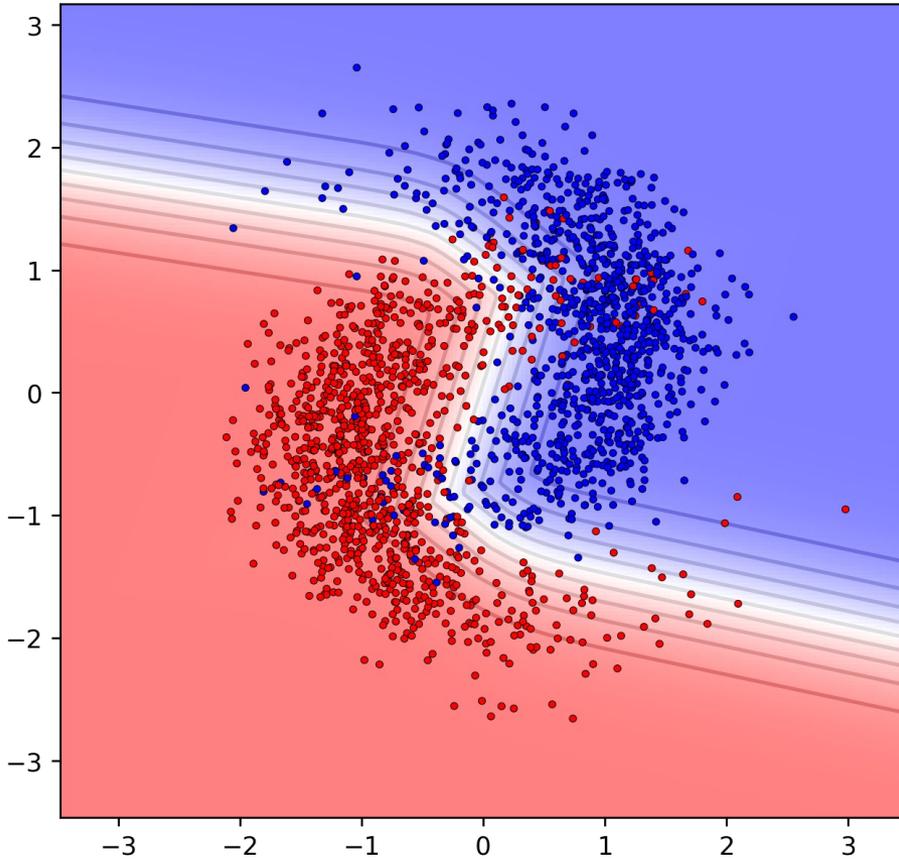
Linear classifier



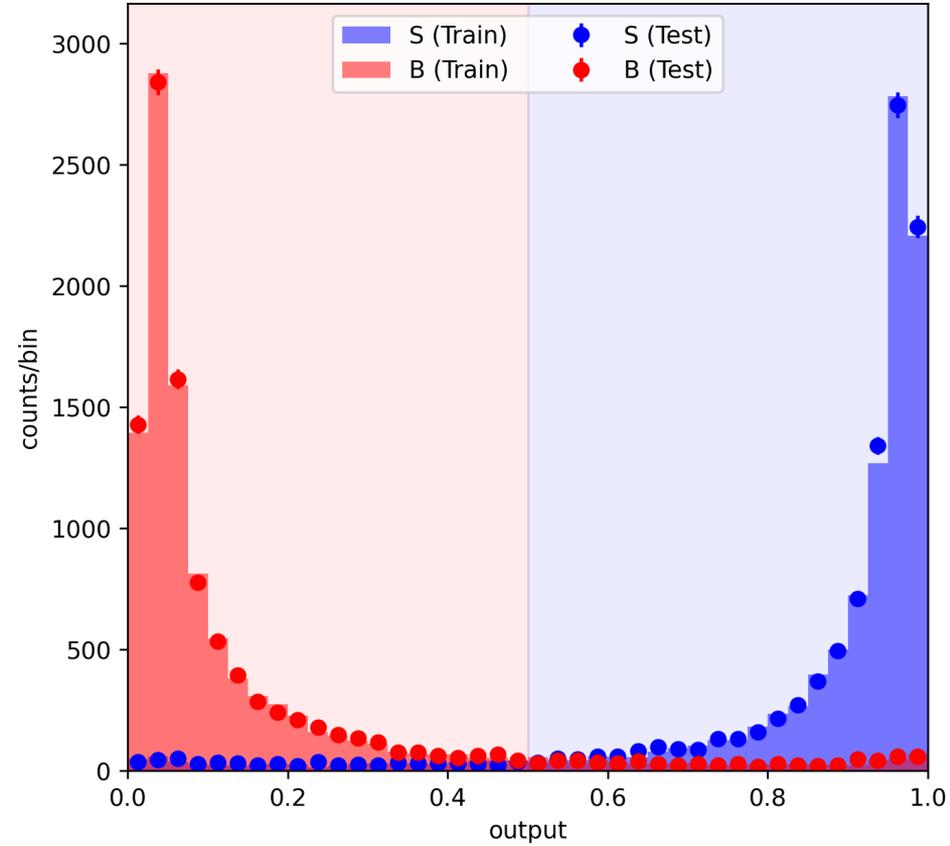
Linear classifier



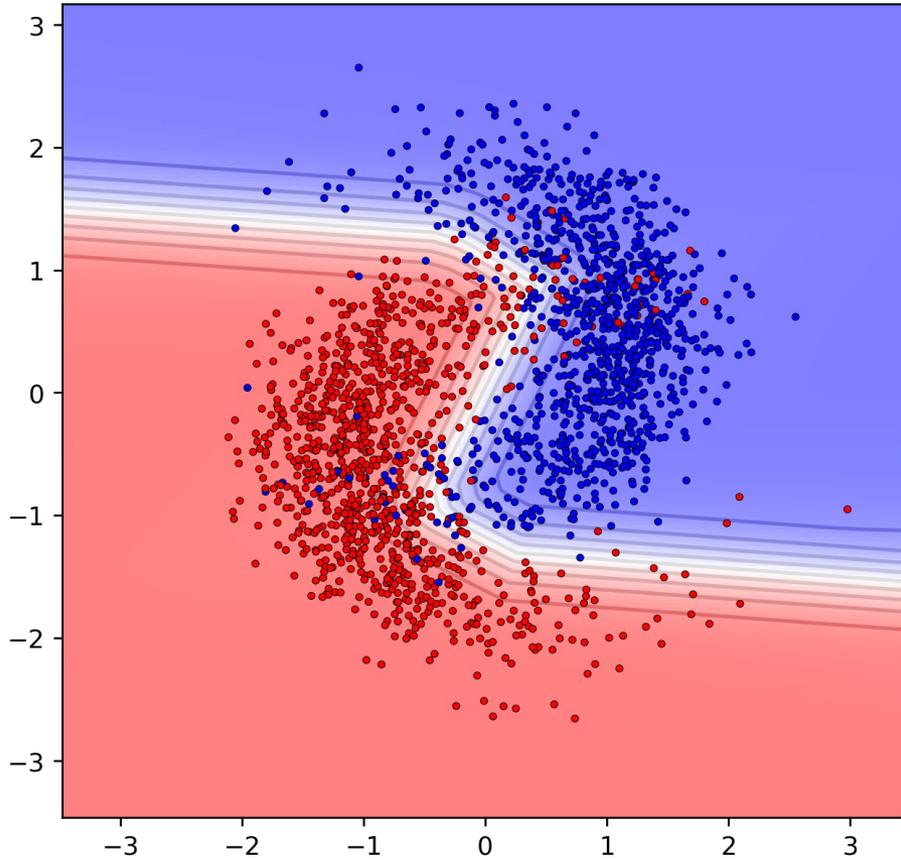
Neural network, shallow



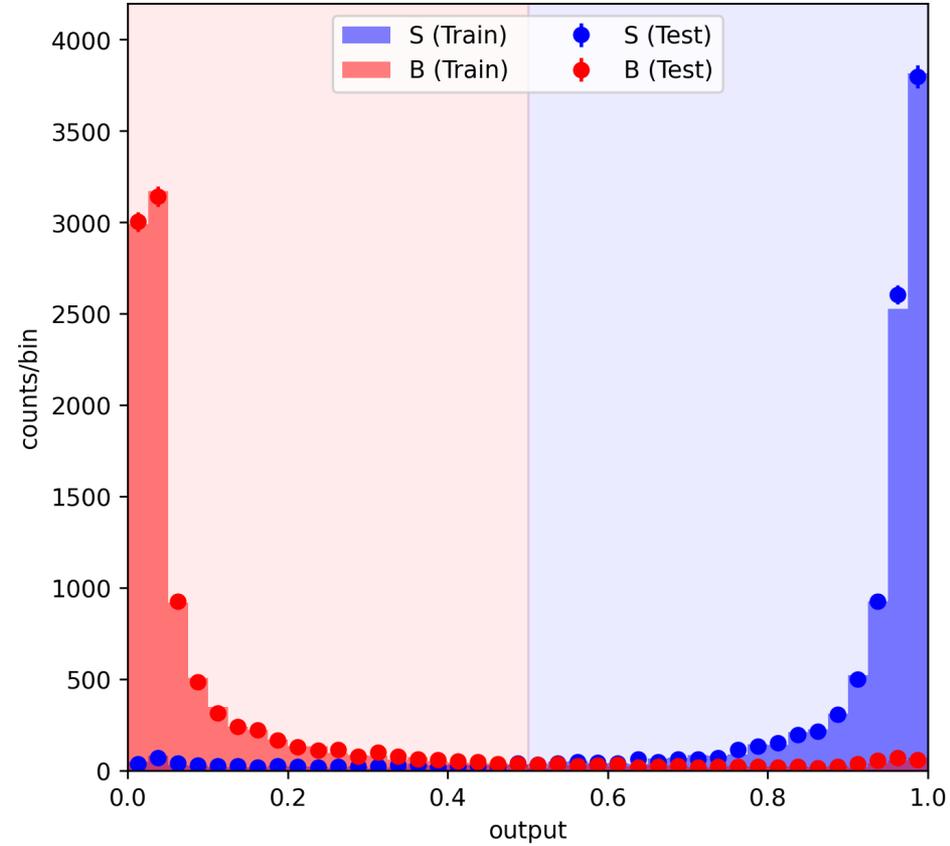
Neural network, shallow



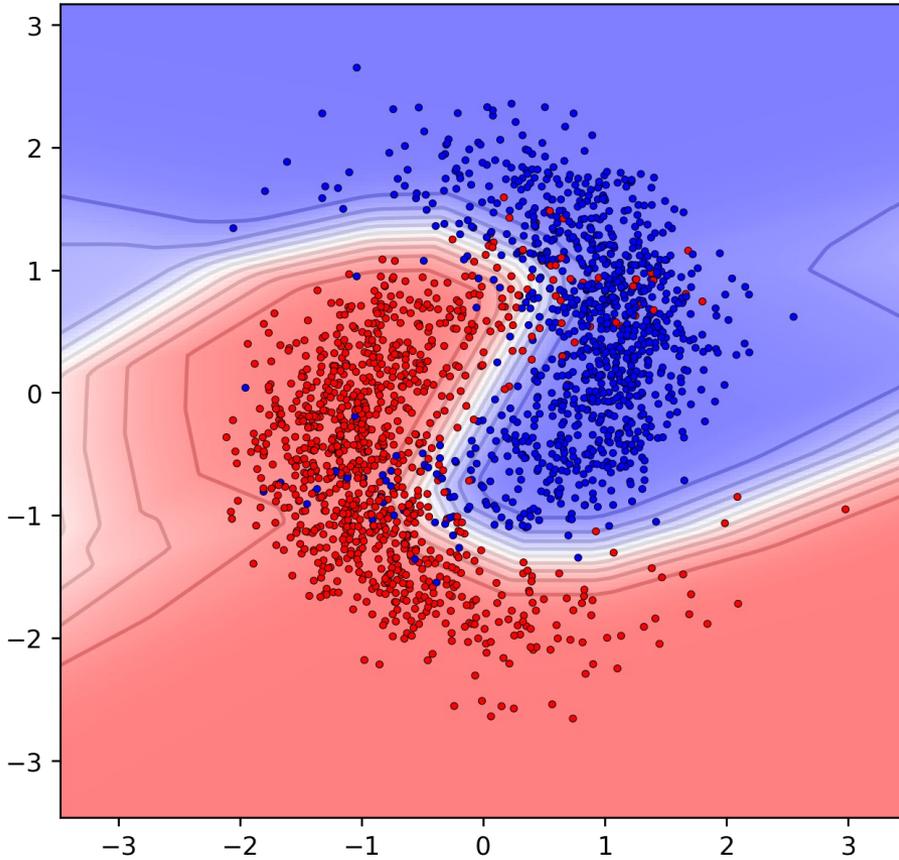
Neural network, middle



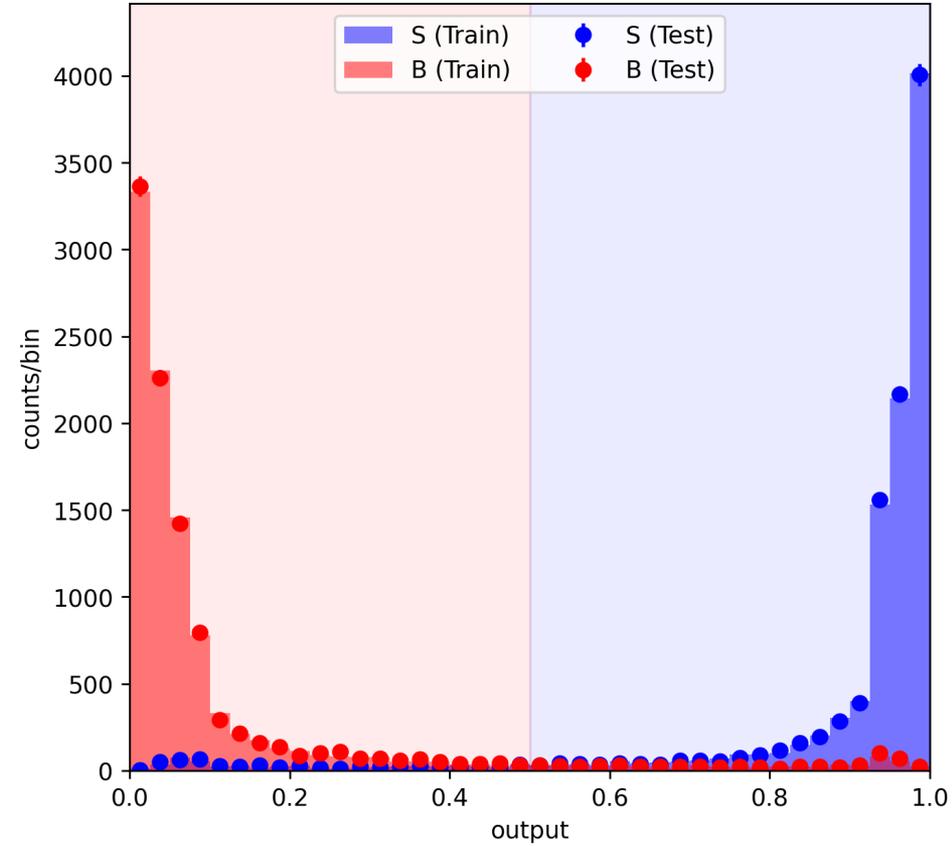
Neural network, middle



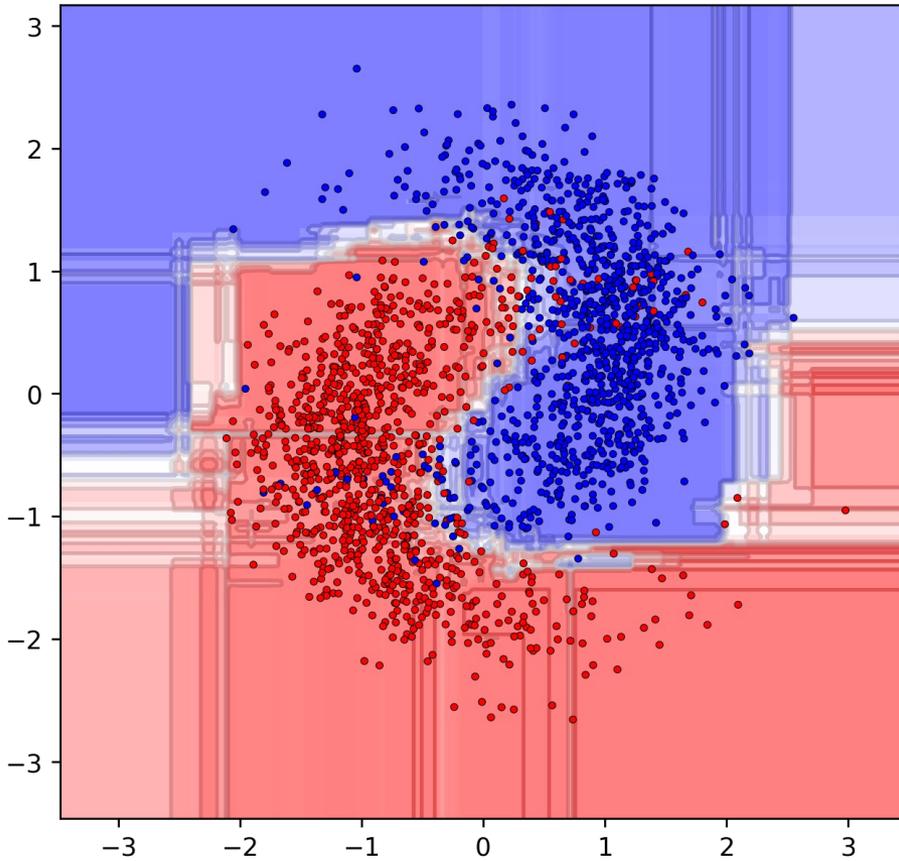
Neural network, deep



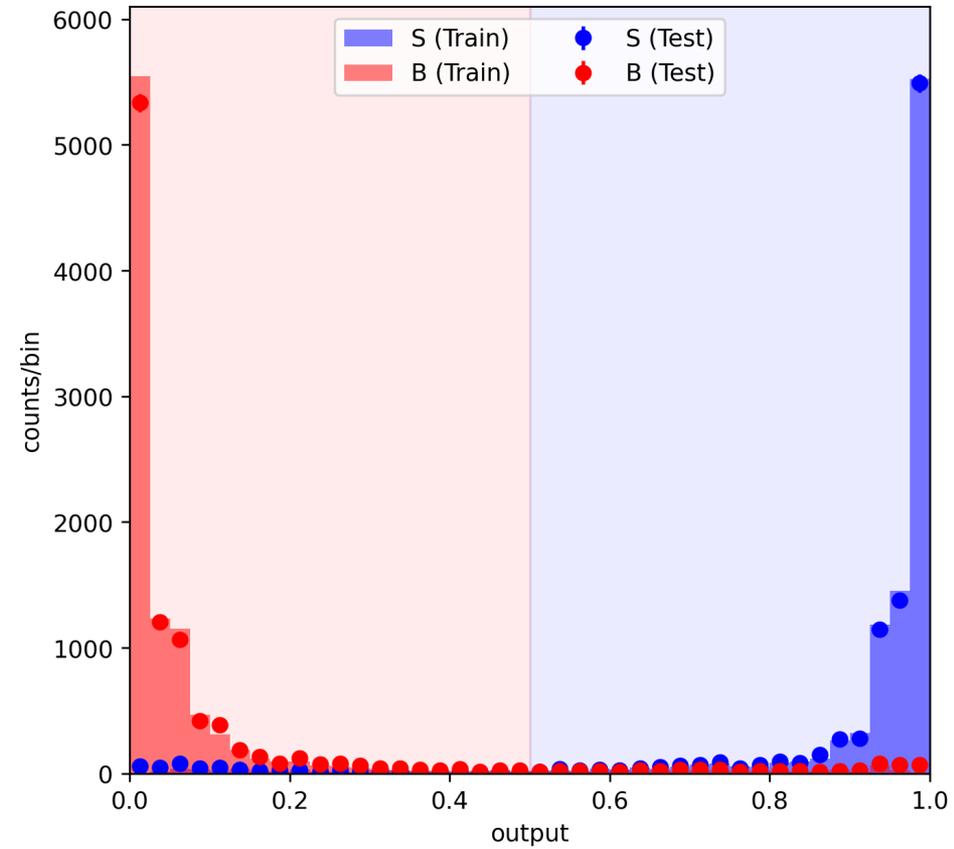
Neural network, deep



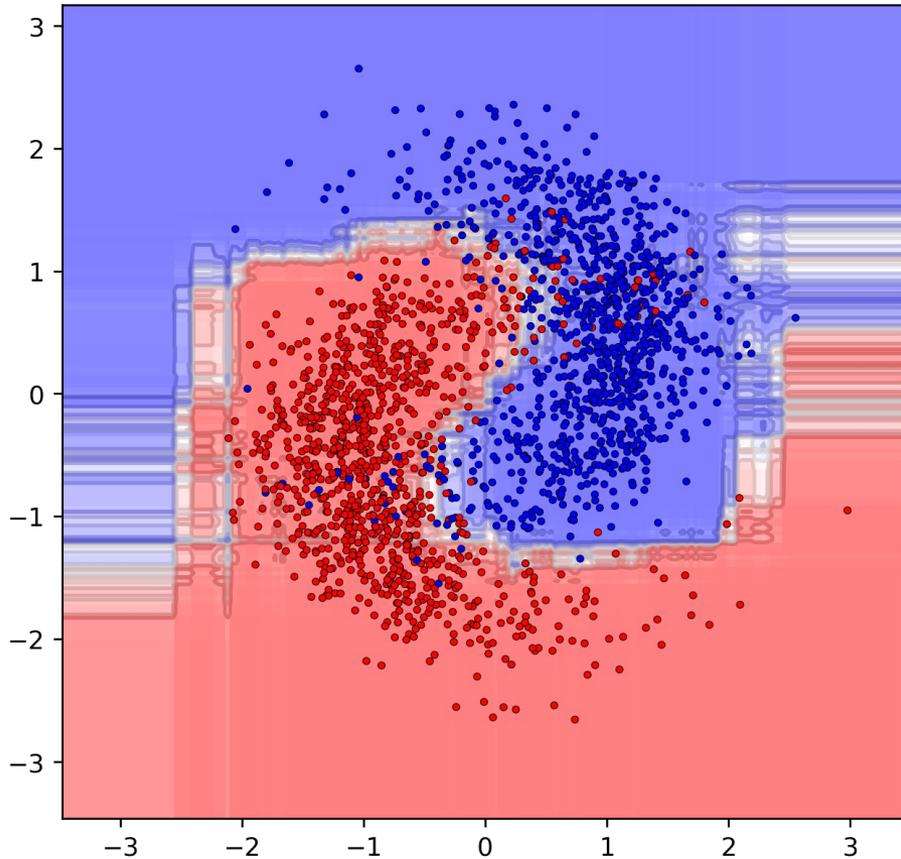
Random forest



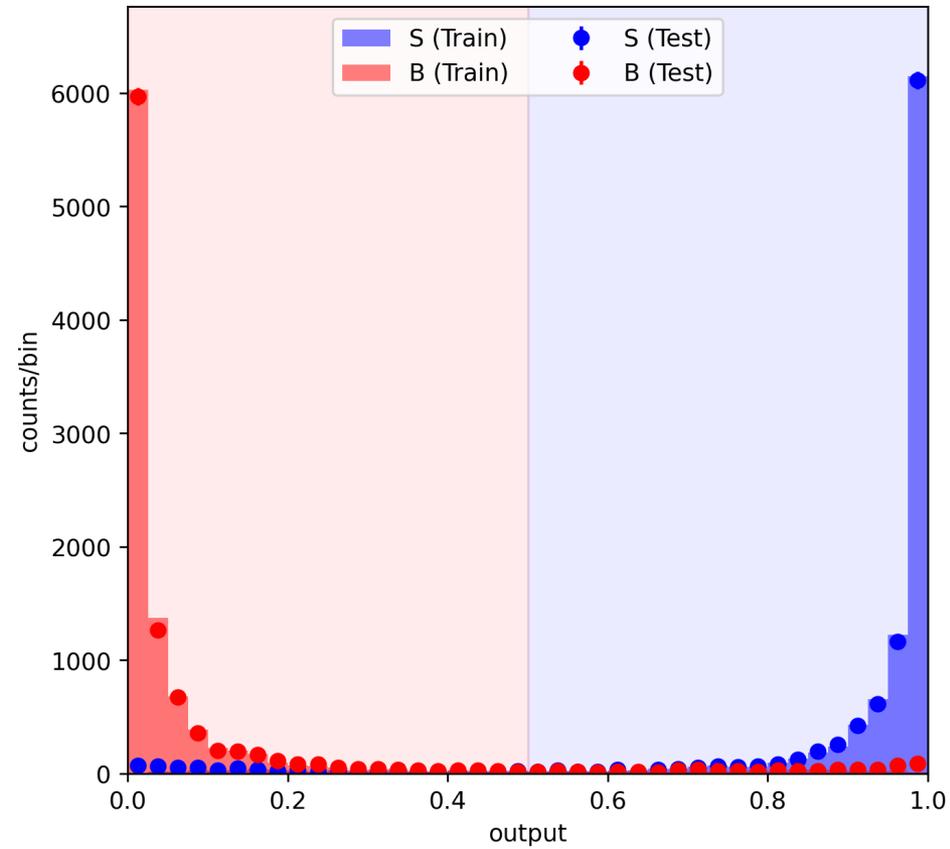
Random forest



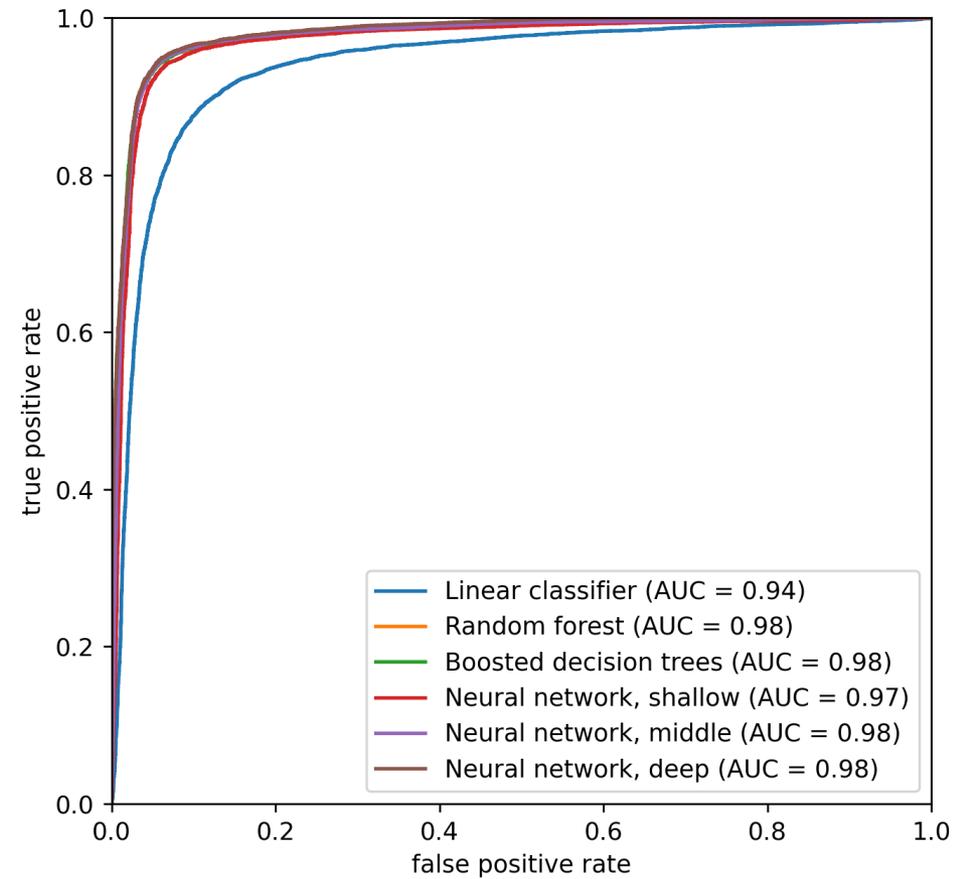
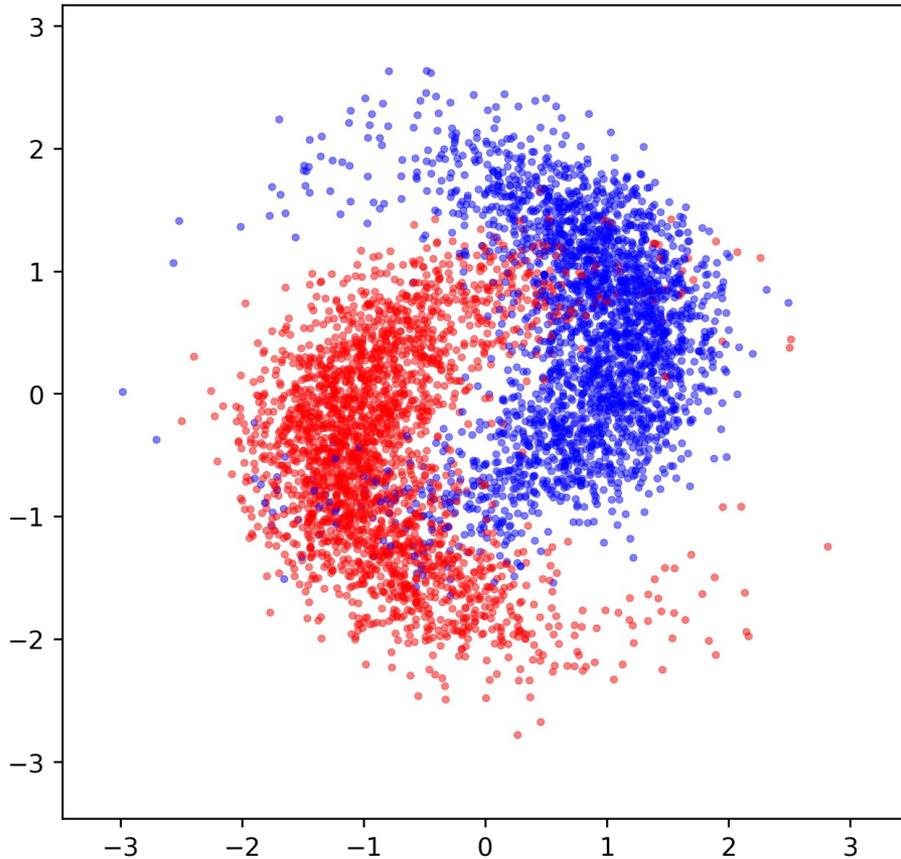
Boosted decision trees

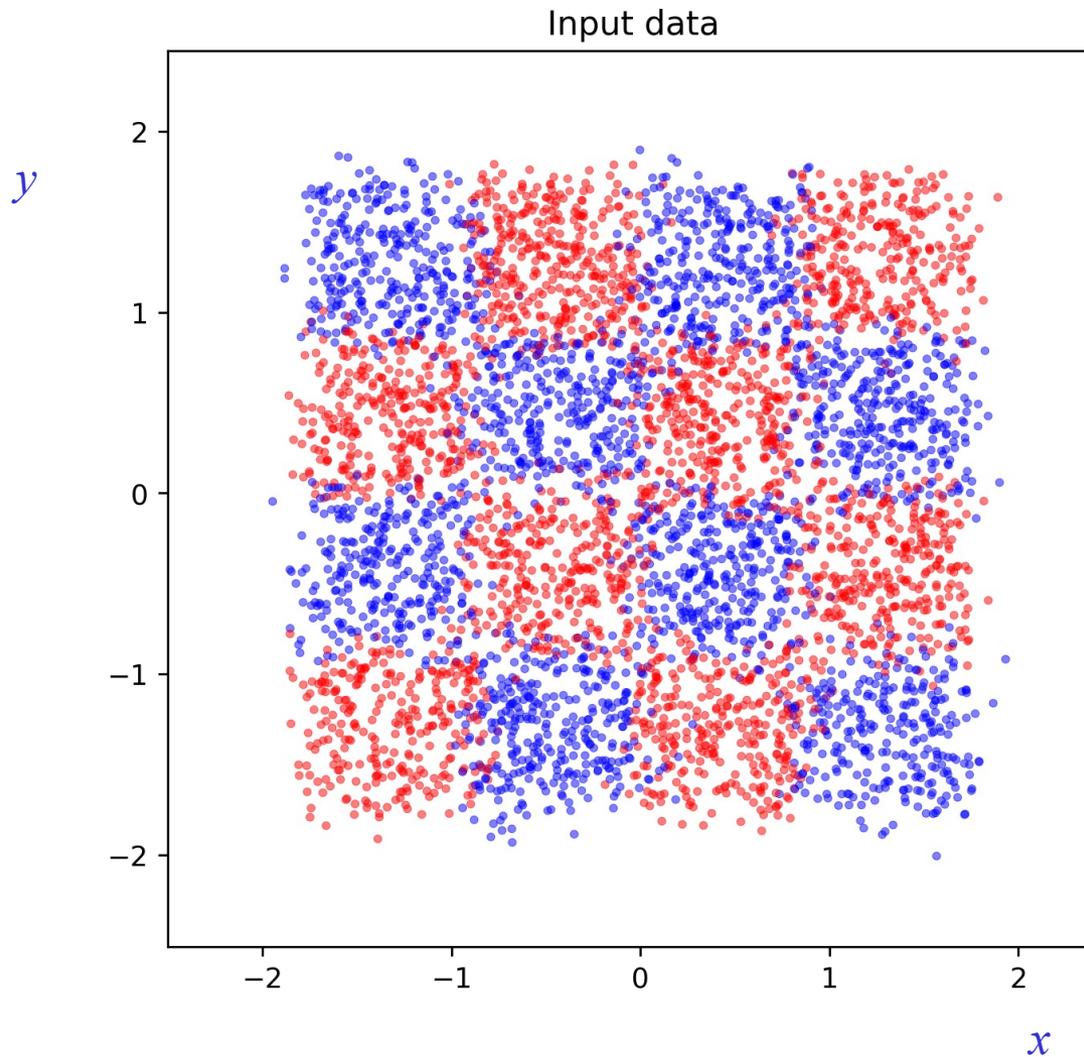


Boosted decision trees

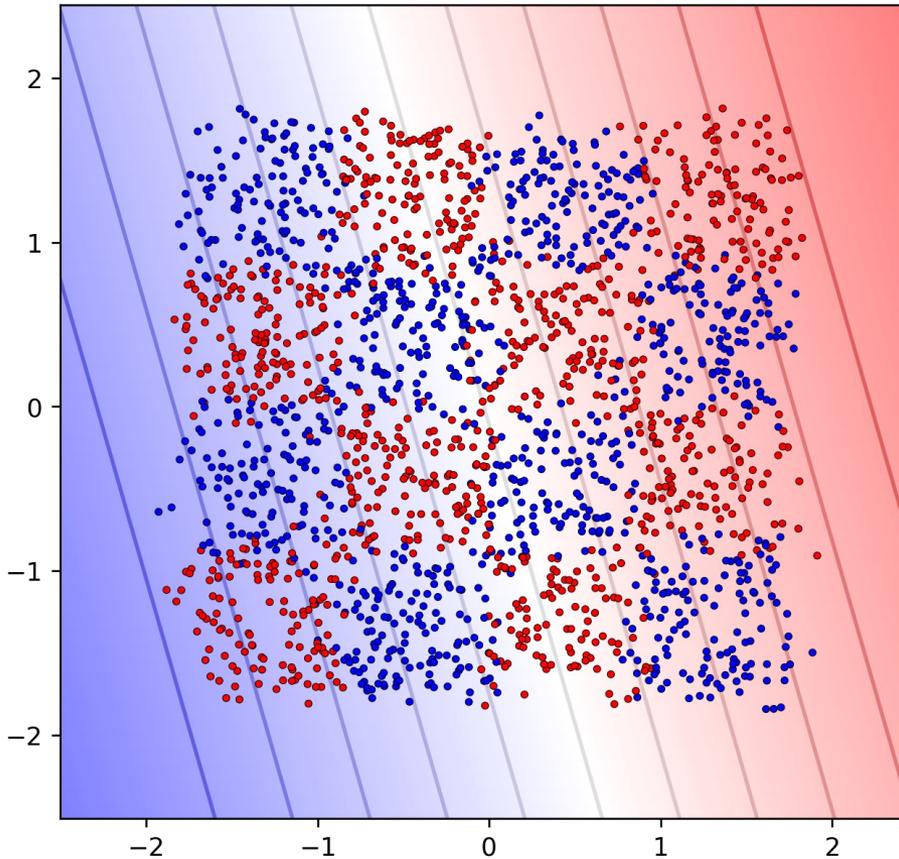


Input data

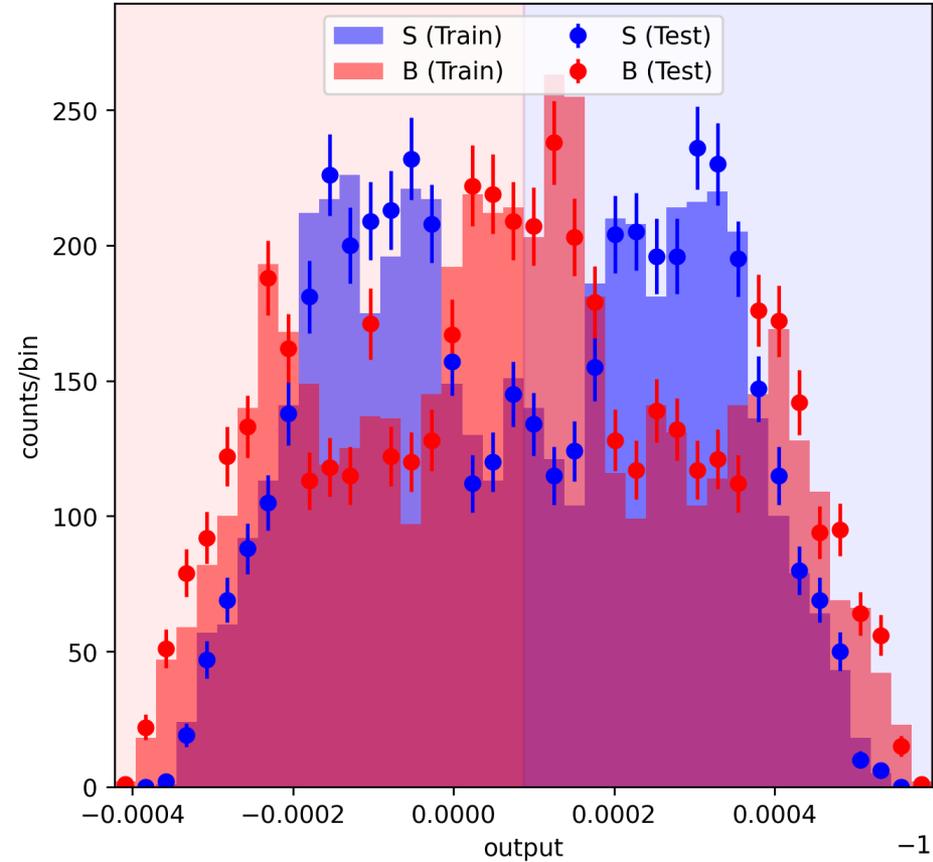




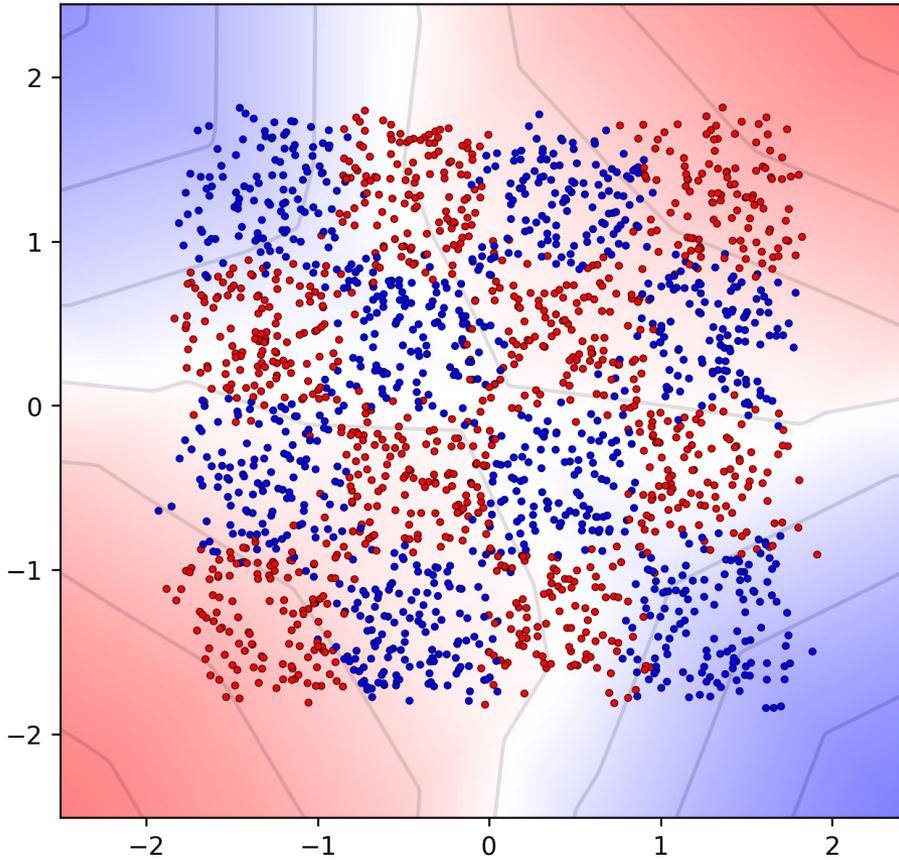
Linear classifier



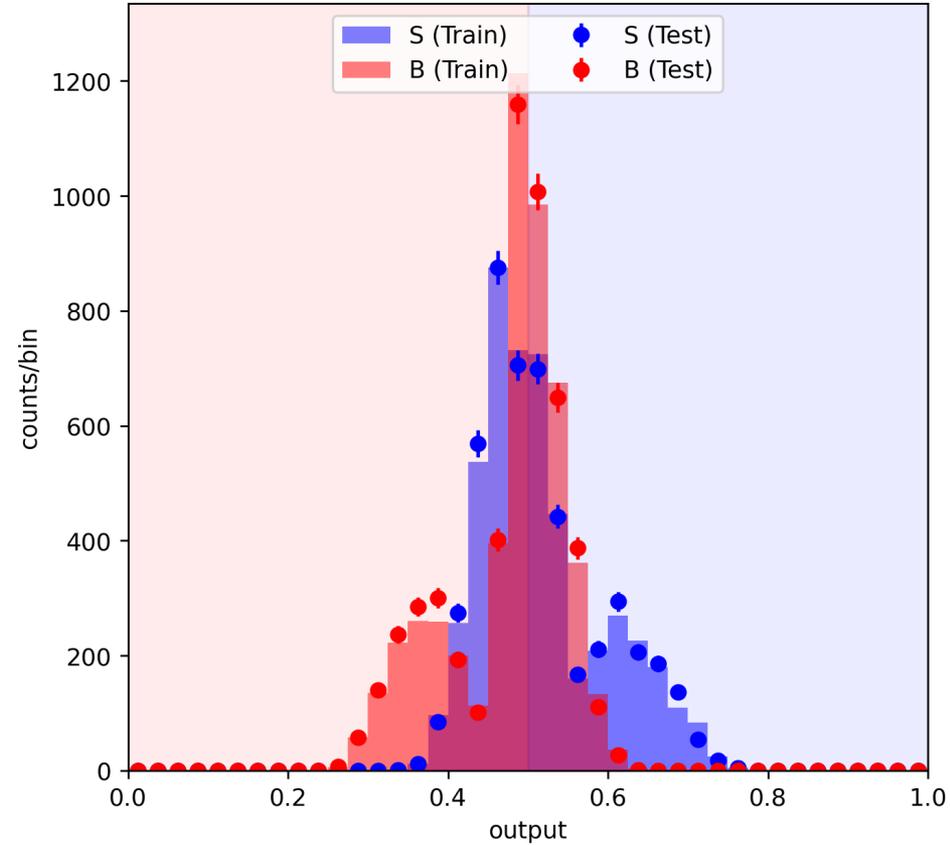
Linear classifier



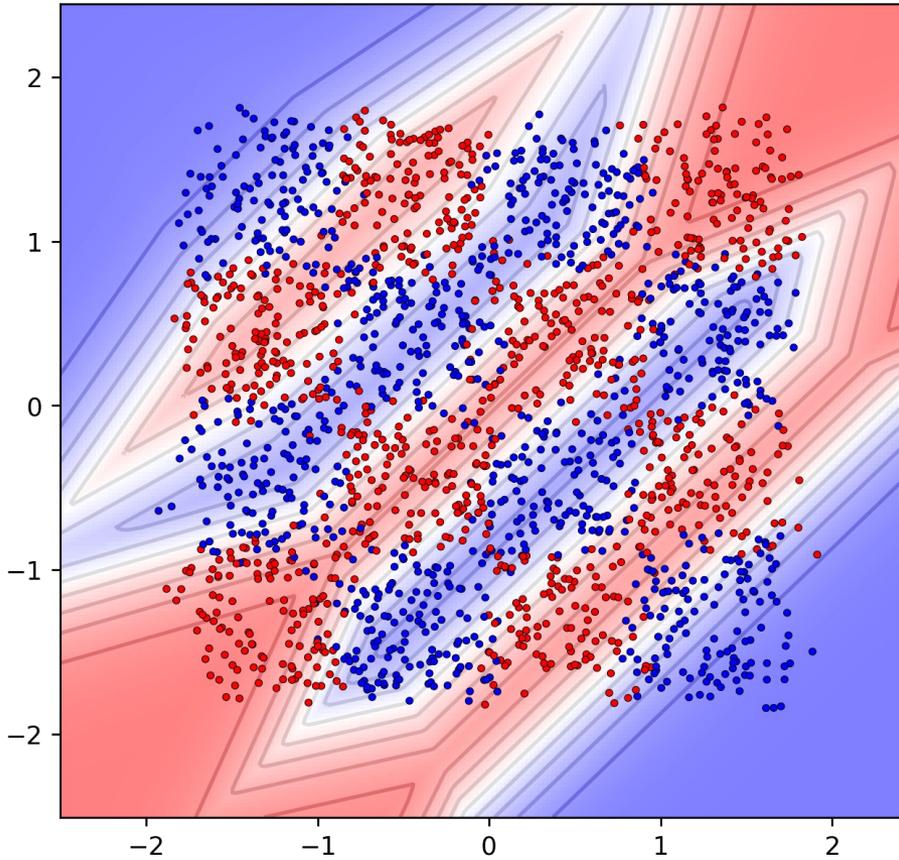
Neural network, shallow



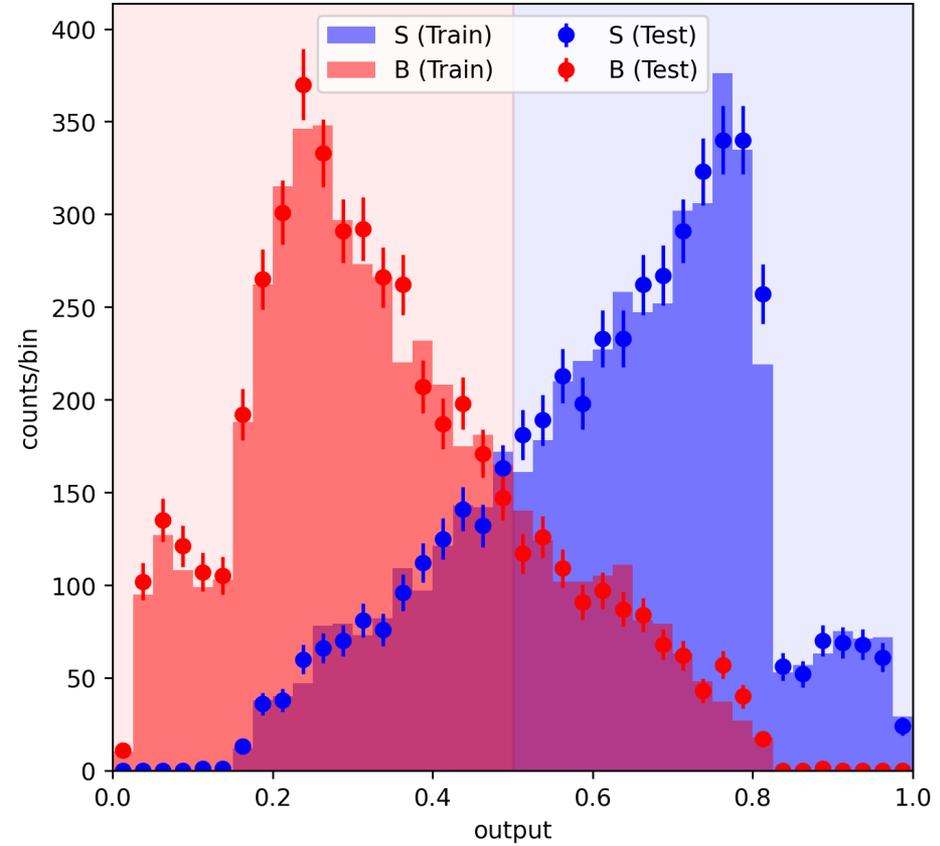
Neural network, shallow



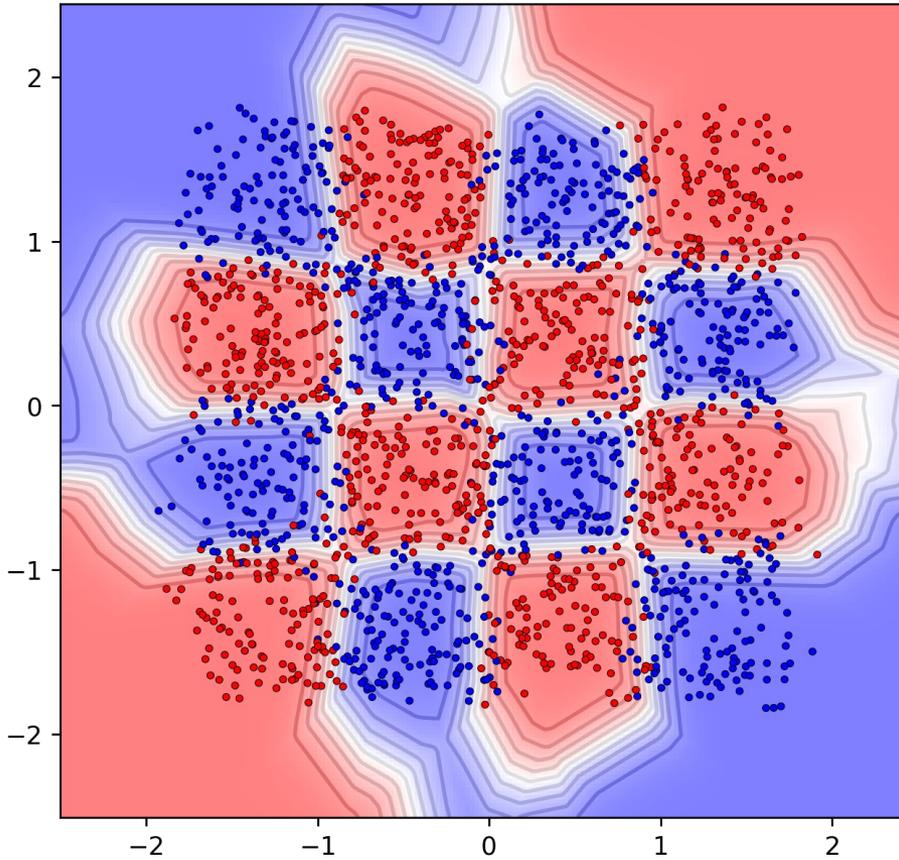
Neural network, middle



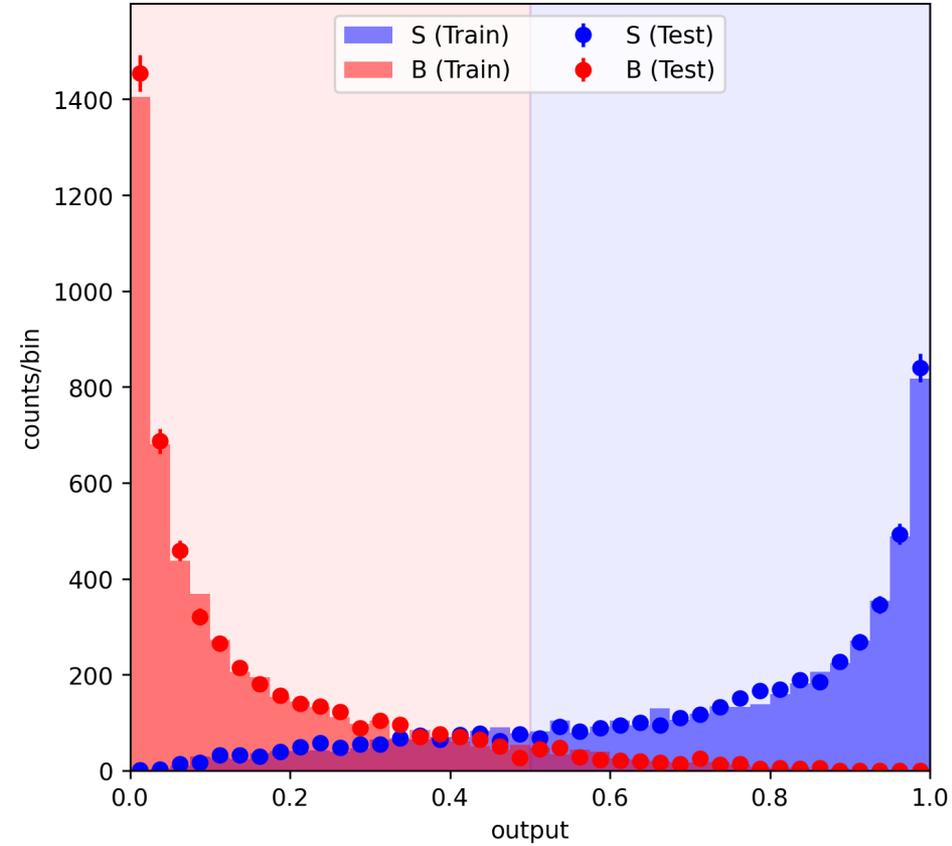
Neural network, middle



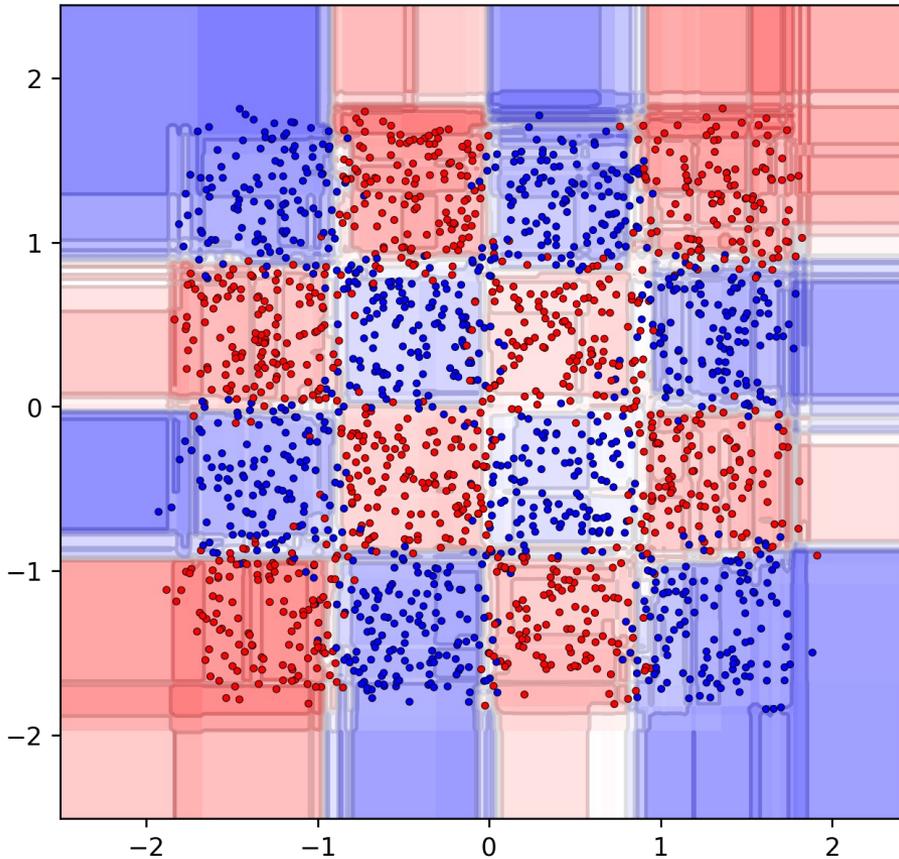
Neural network, deep



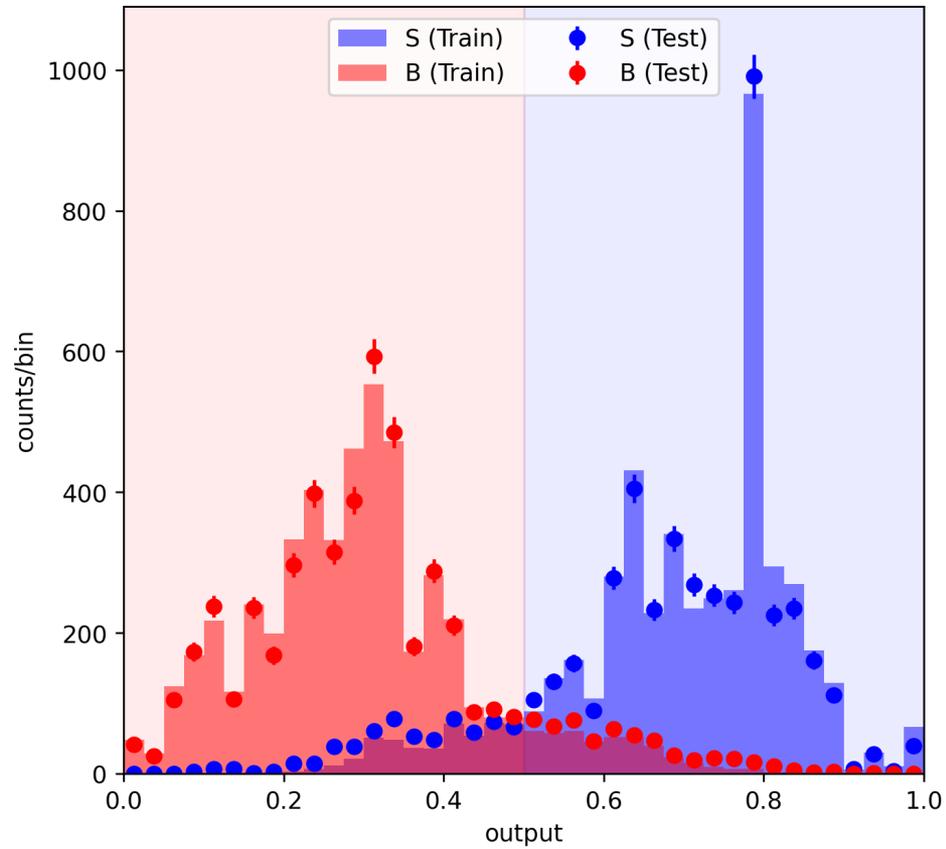
Neural network, deep



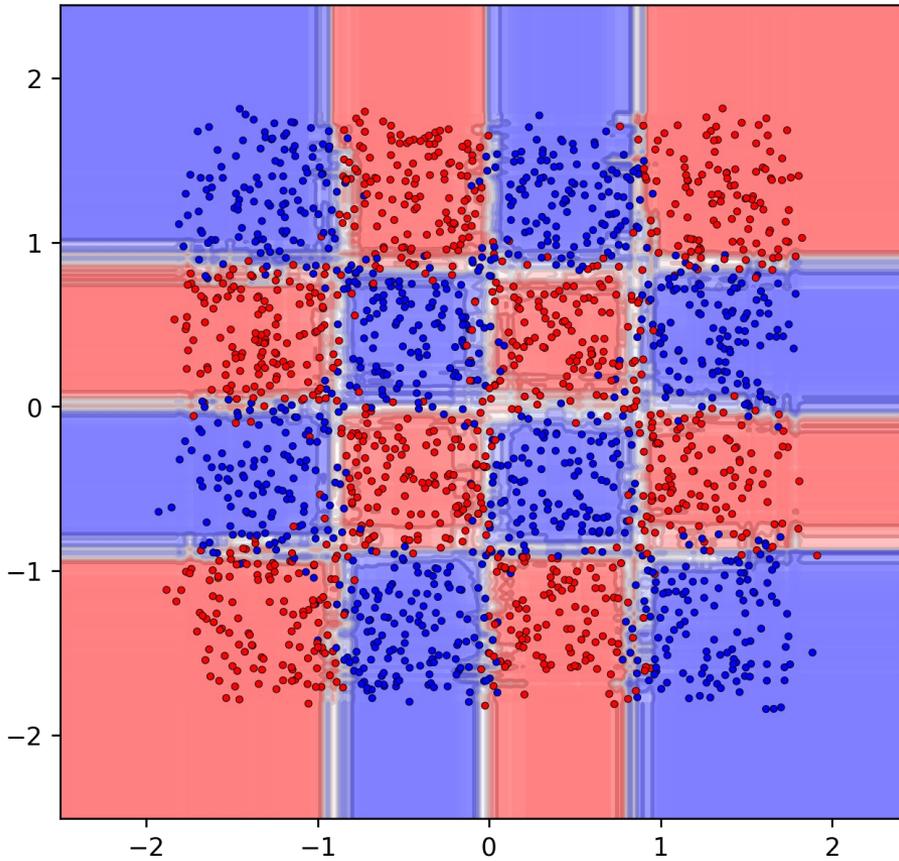
Random forest



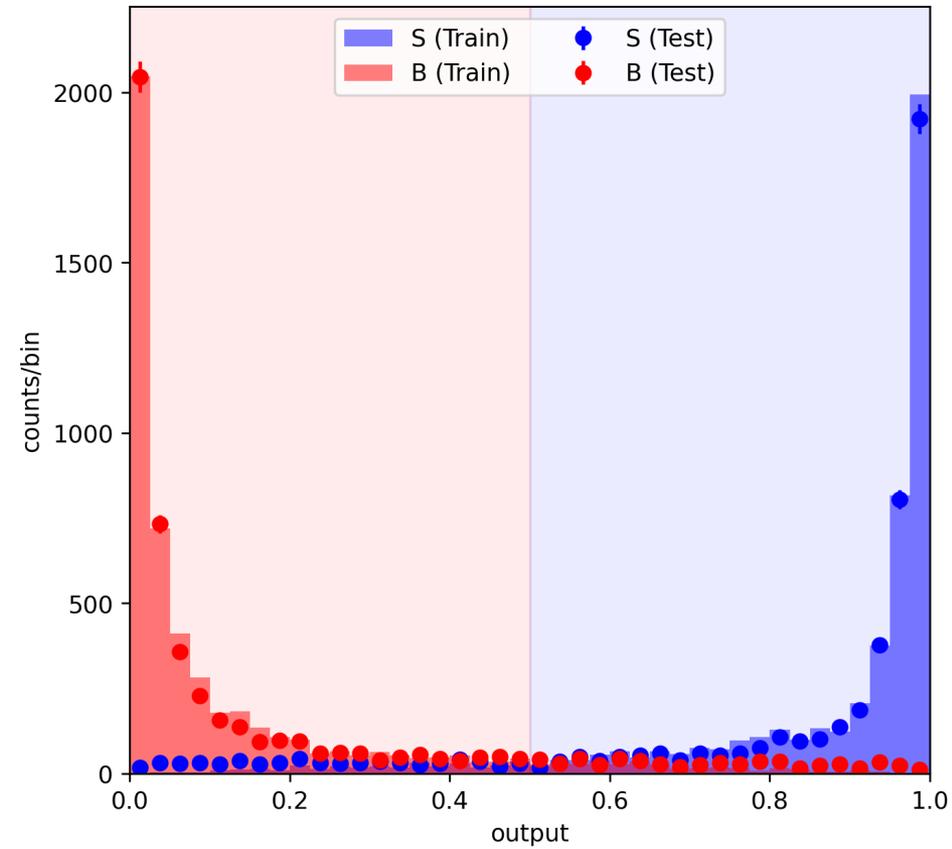
Random forest



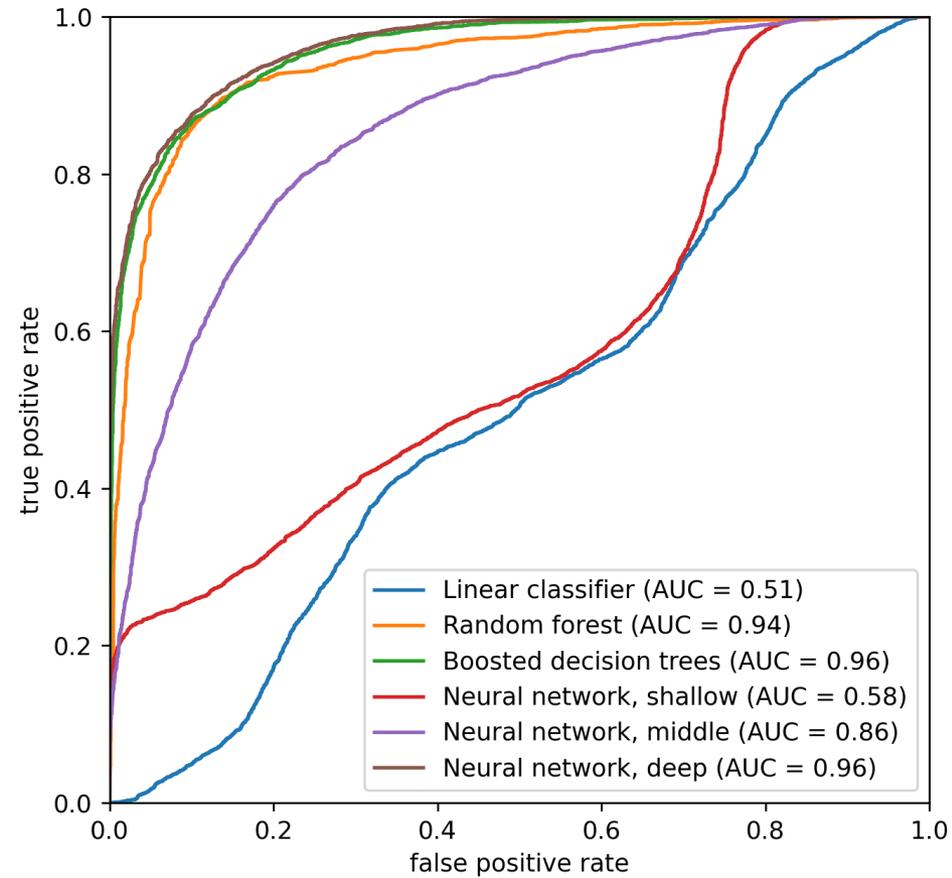
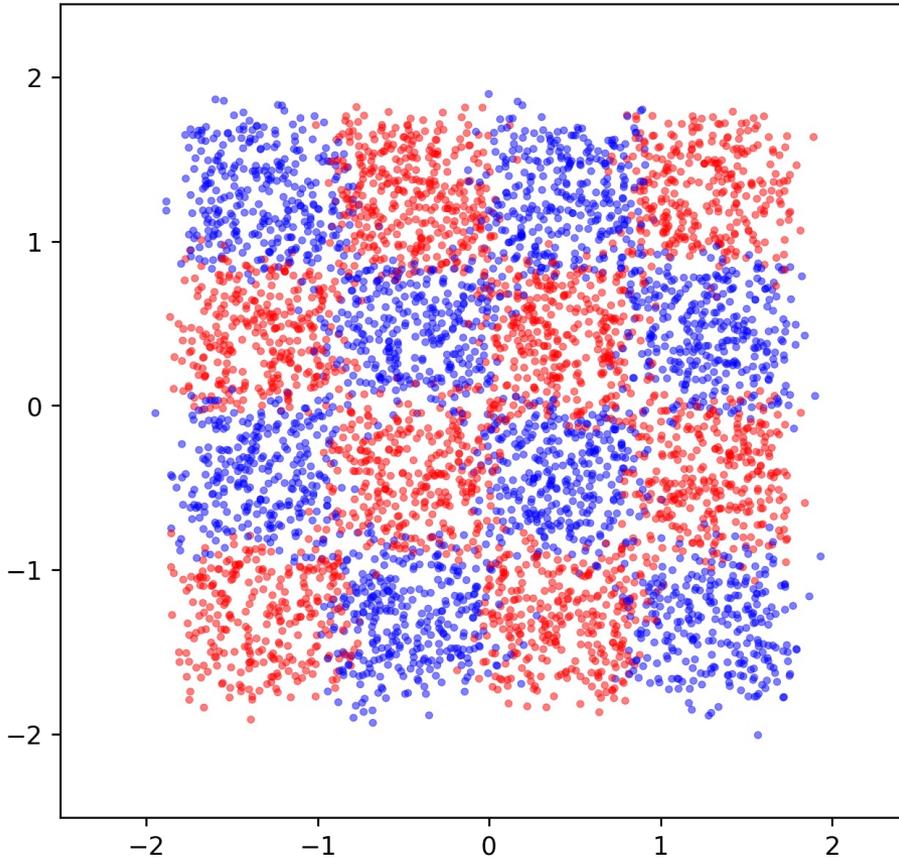
Boosted decision trees

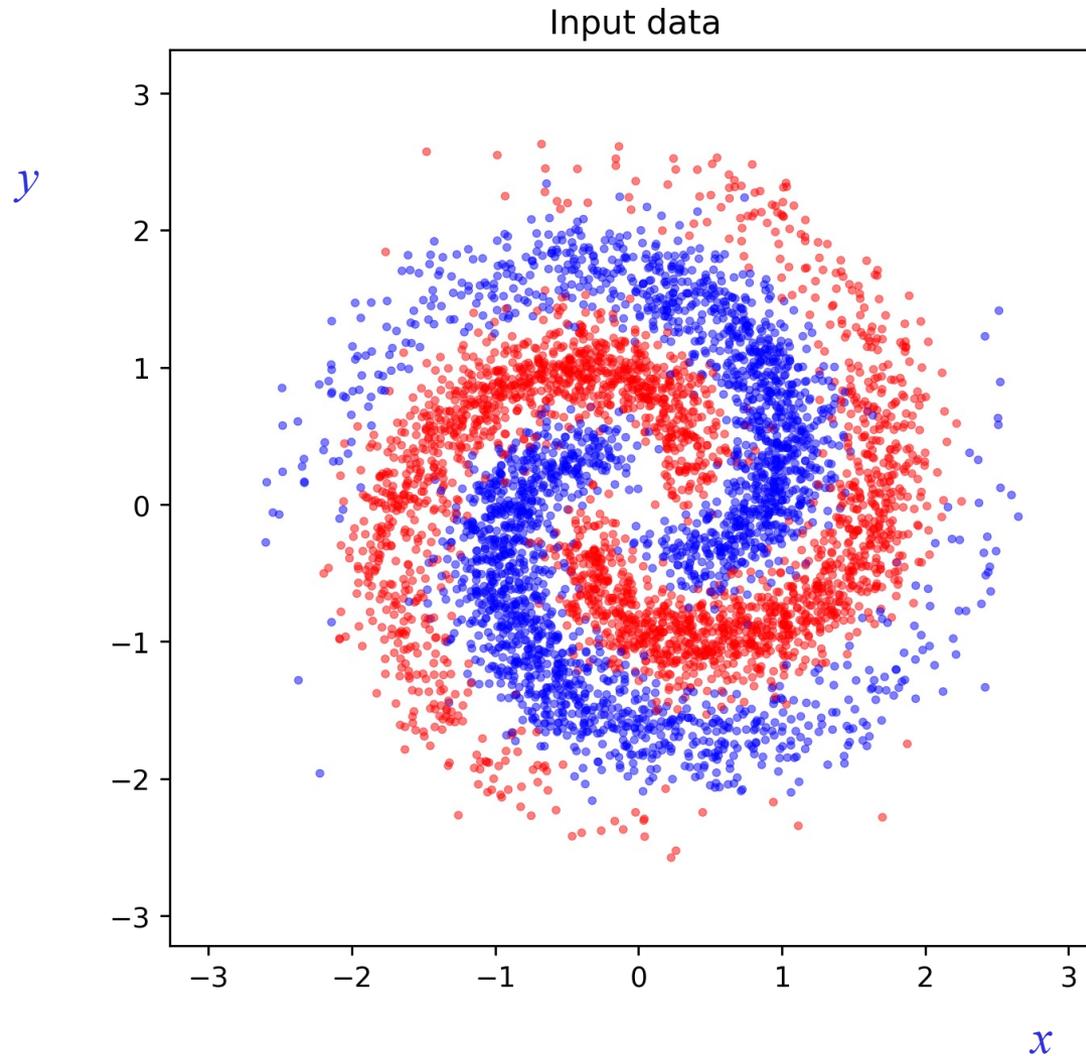


Boosted decision trees

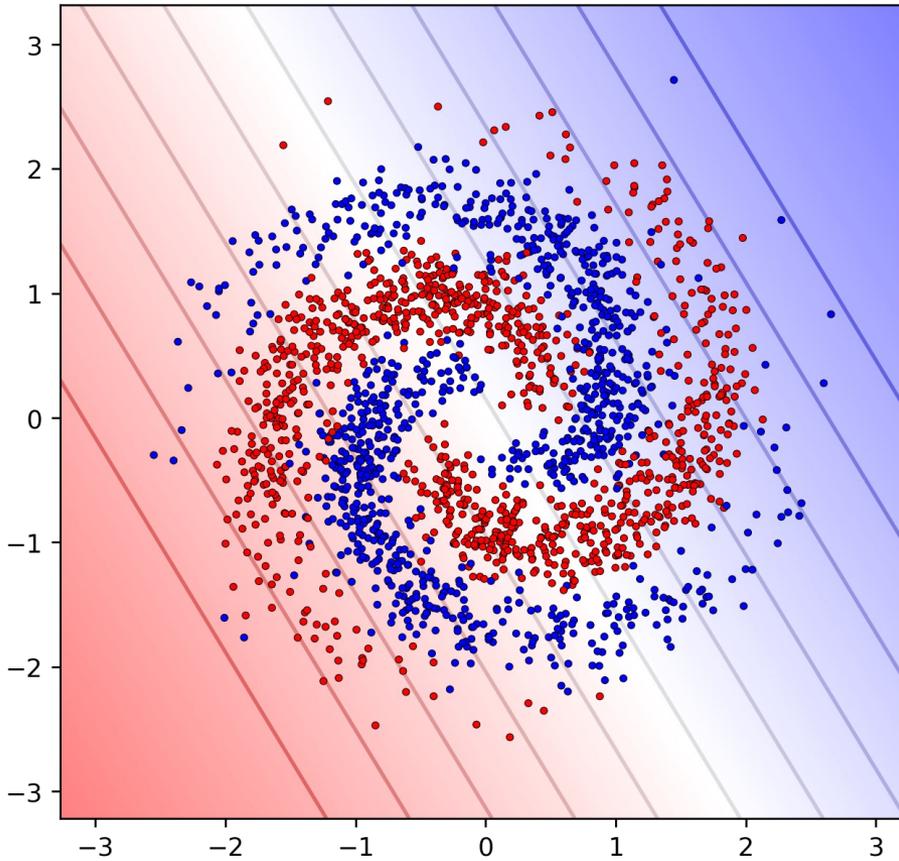


Input data

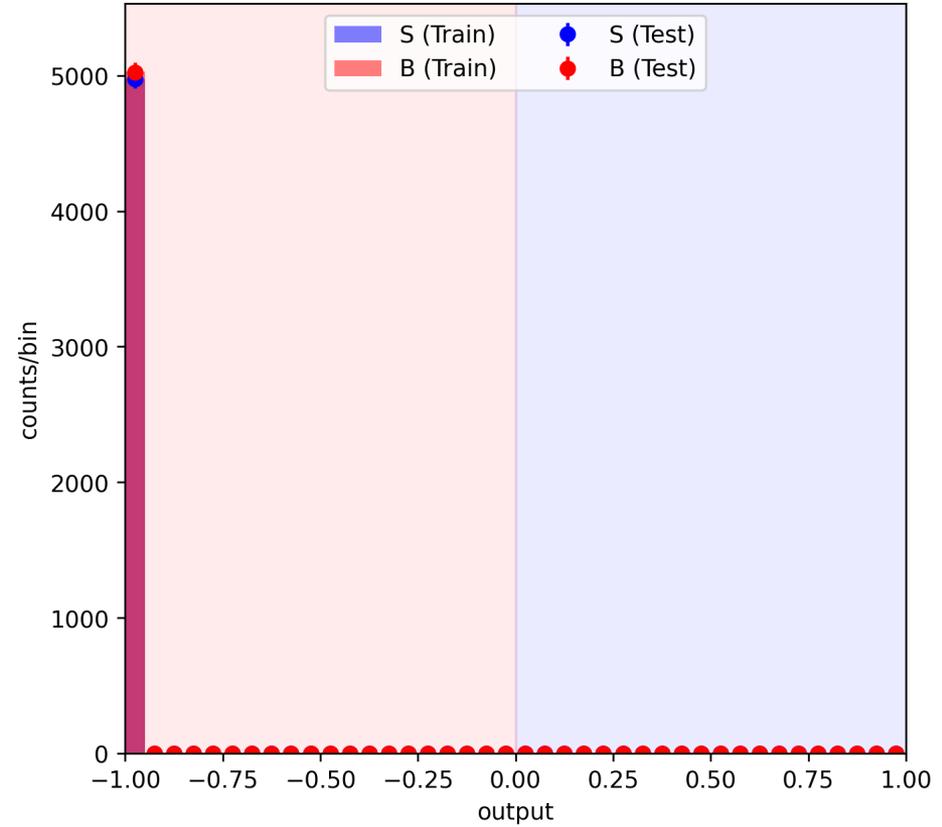




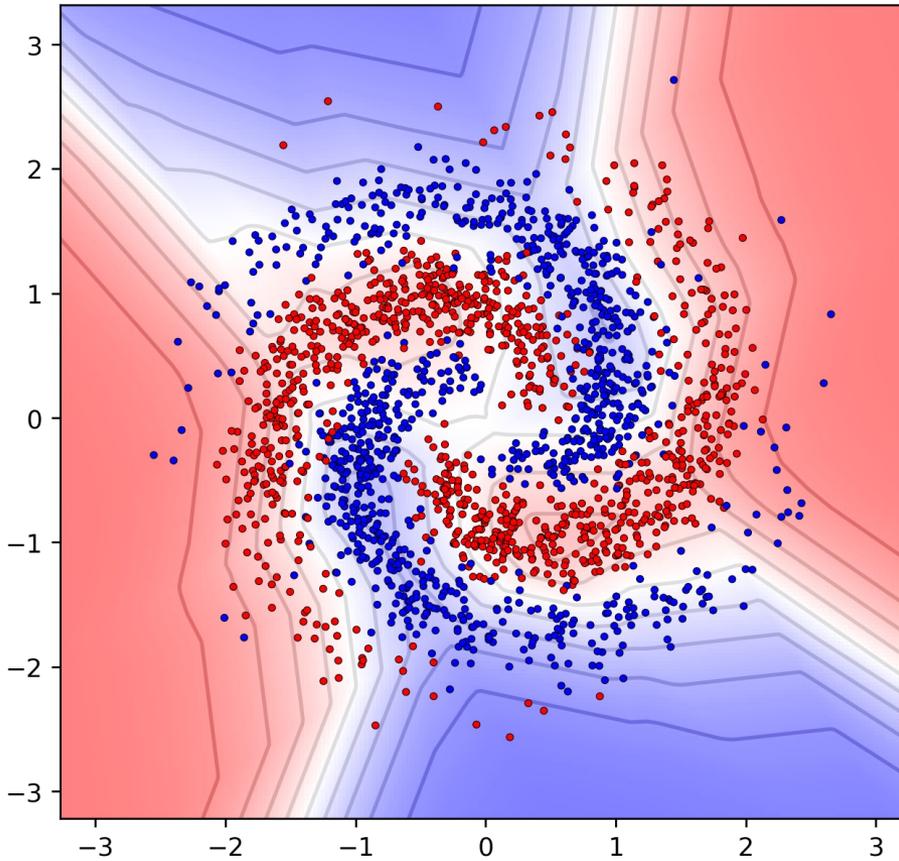
Linear classifier



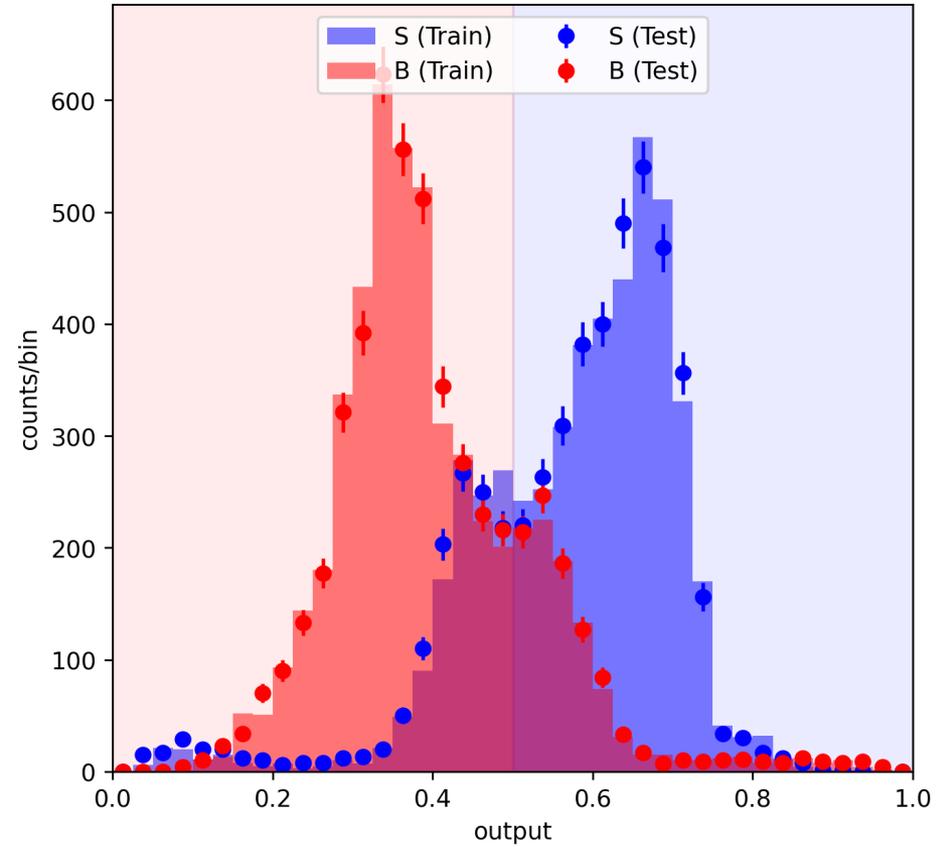
Linear classifier



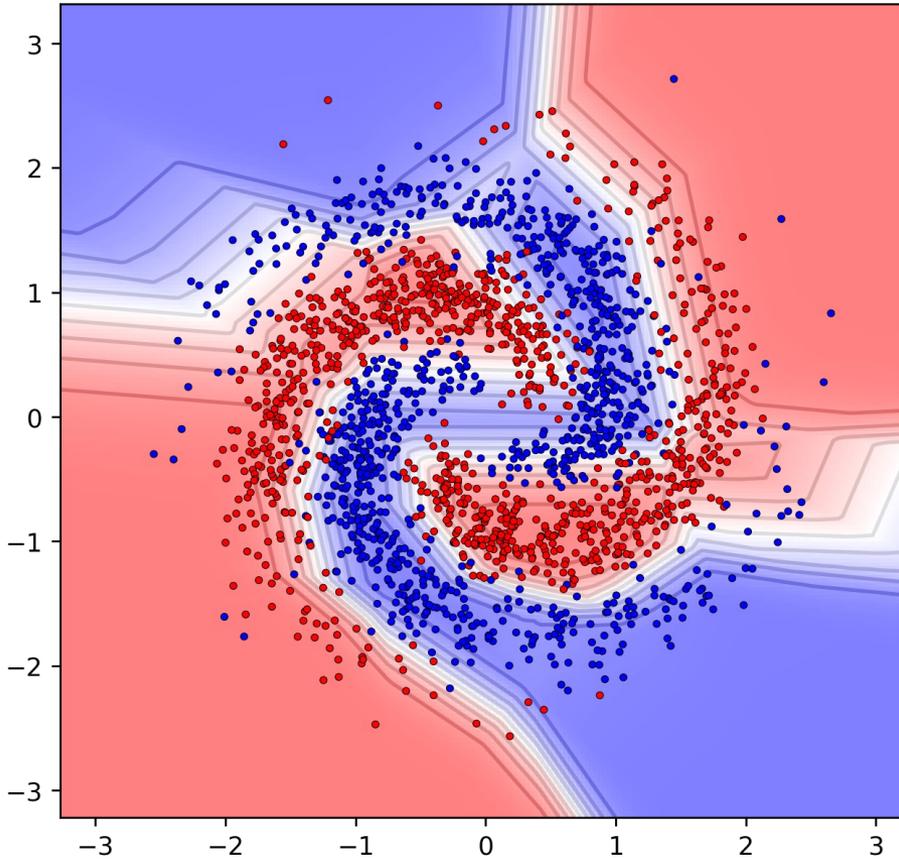
Neural network, shallow



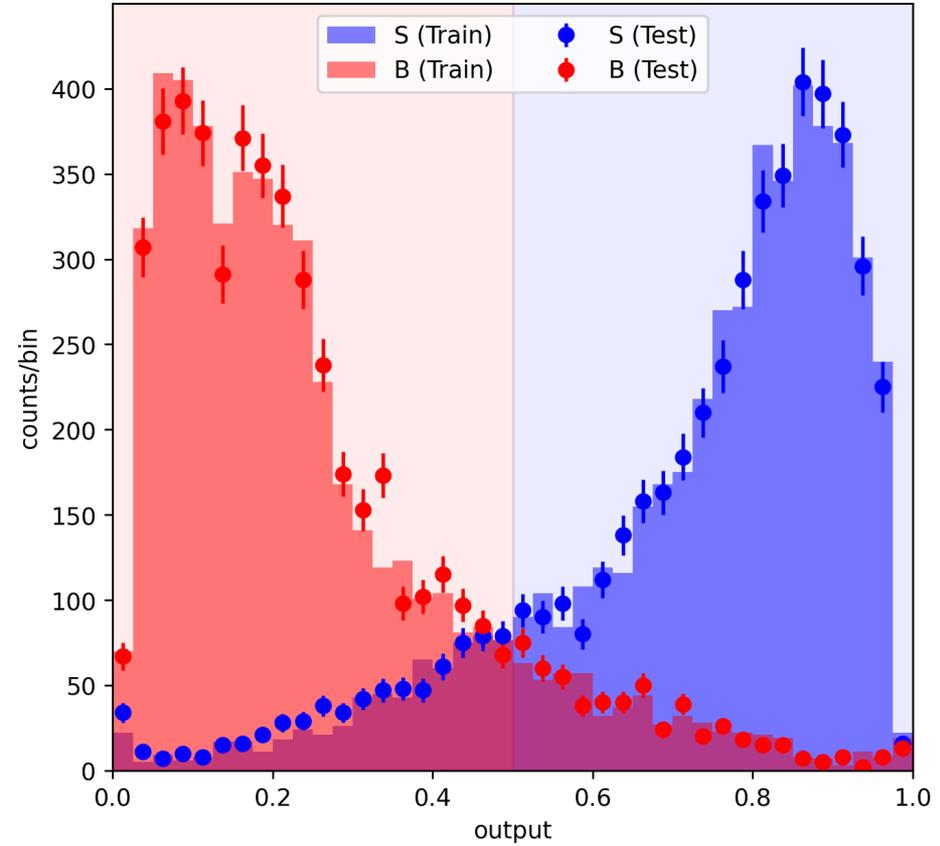
Neural network, shallow



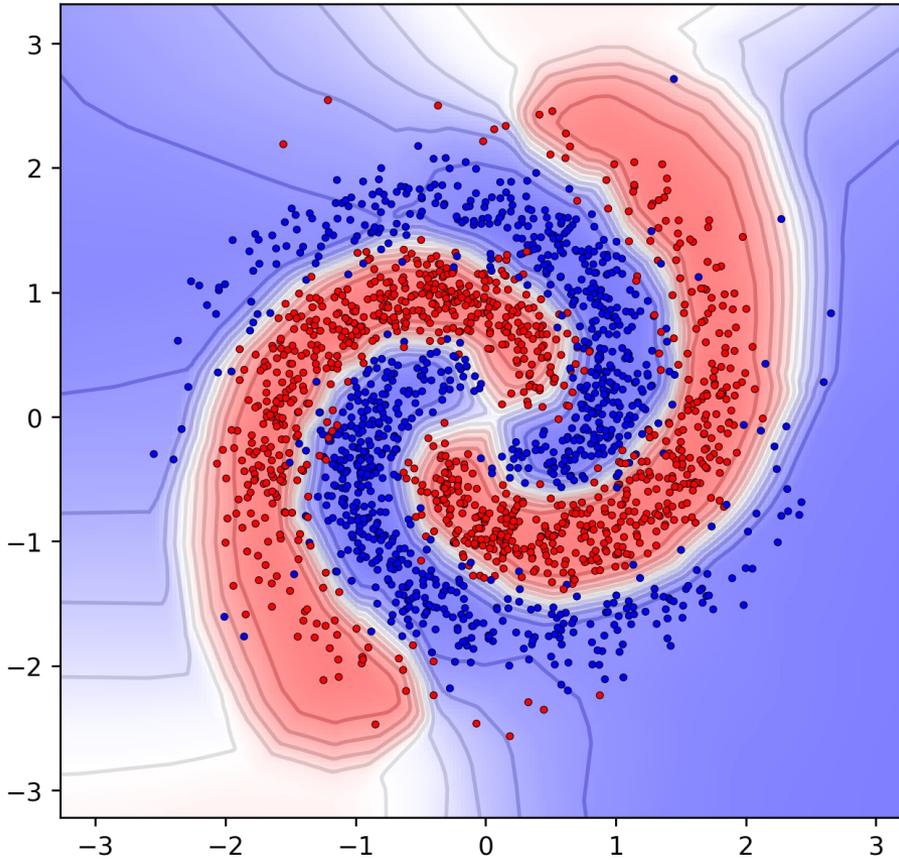
Neural network, middle



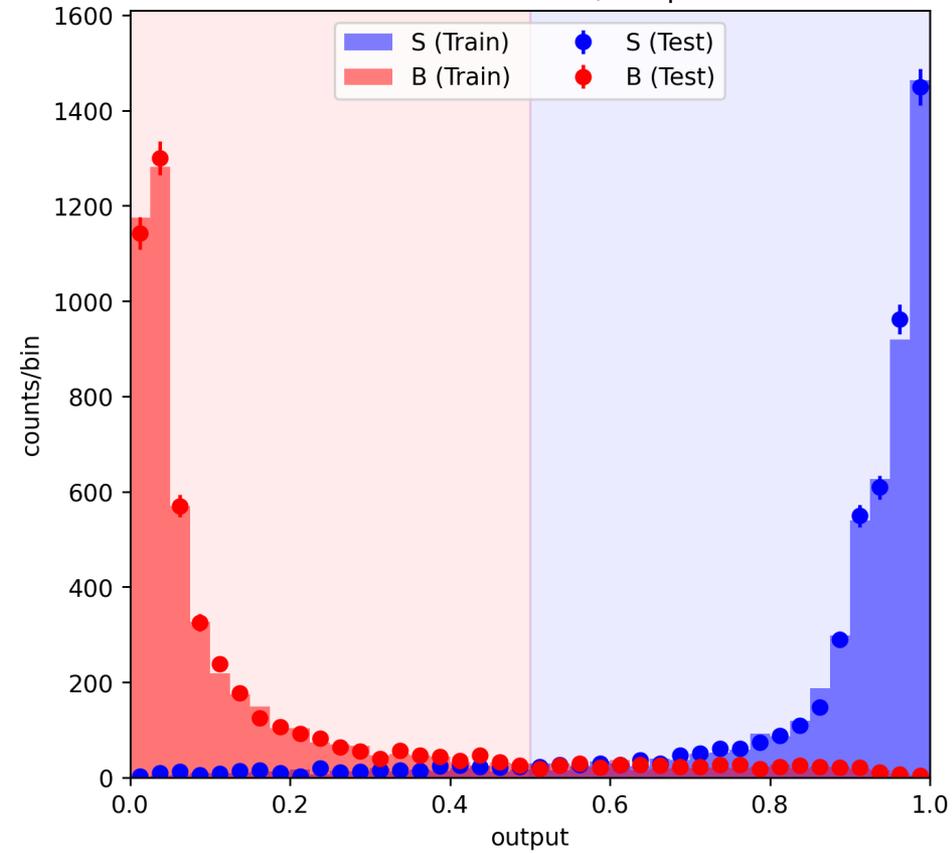
Neural network, middle



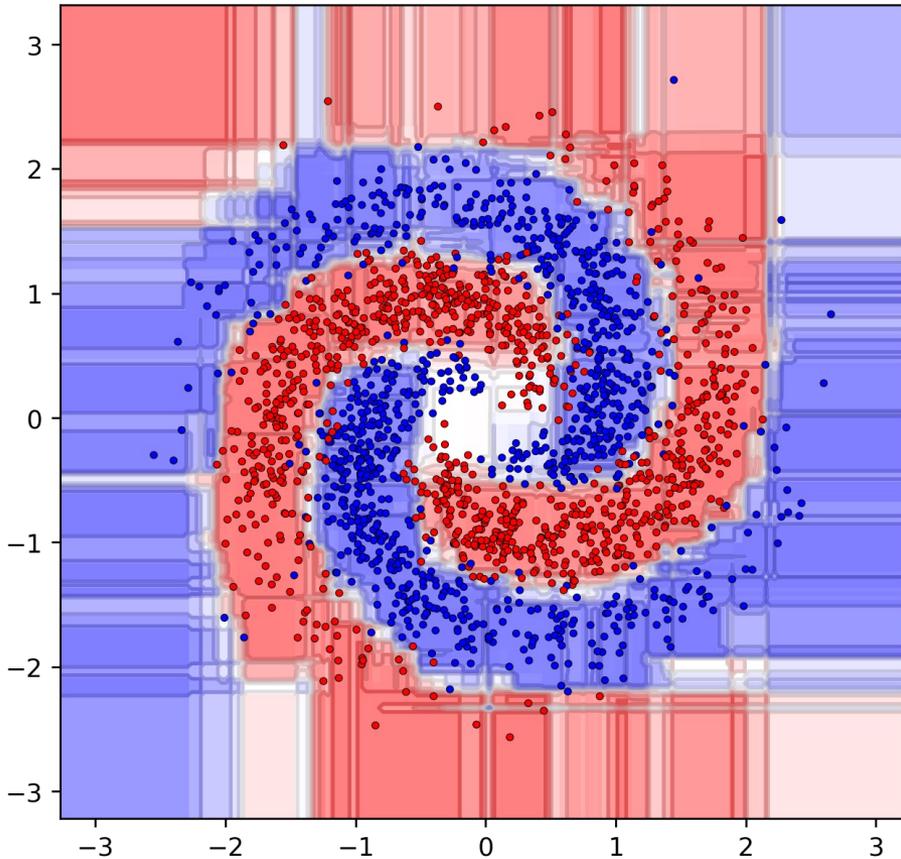
Neural network, deep



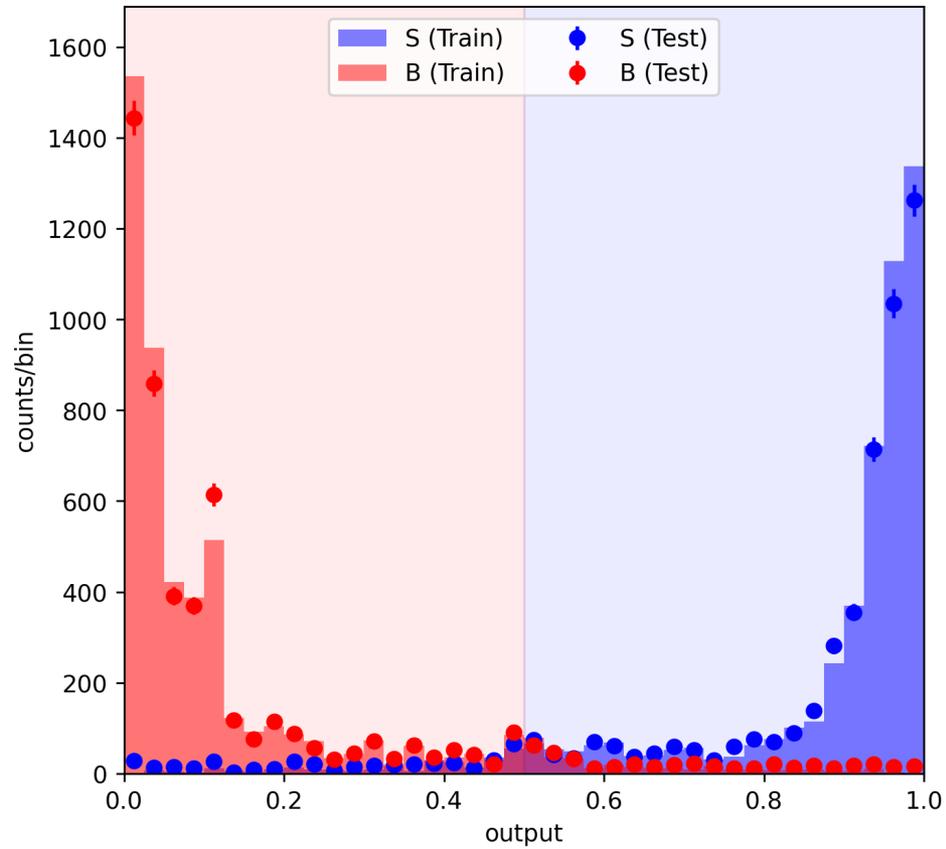
Neural network, deep



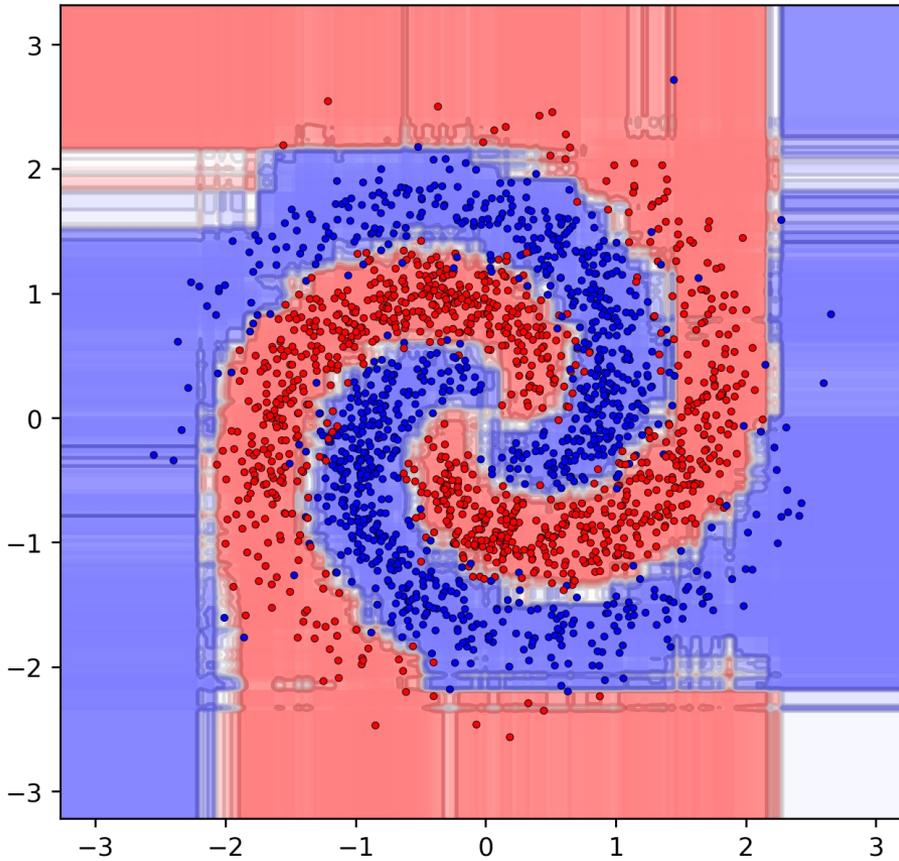
Random forest



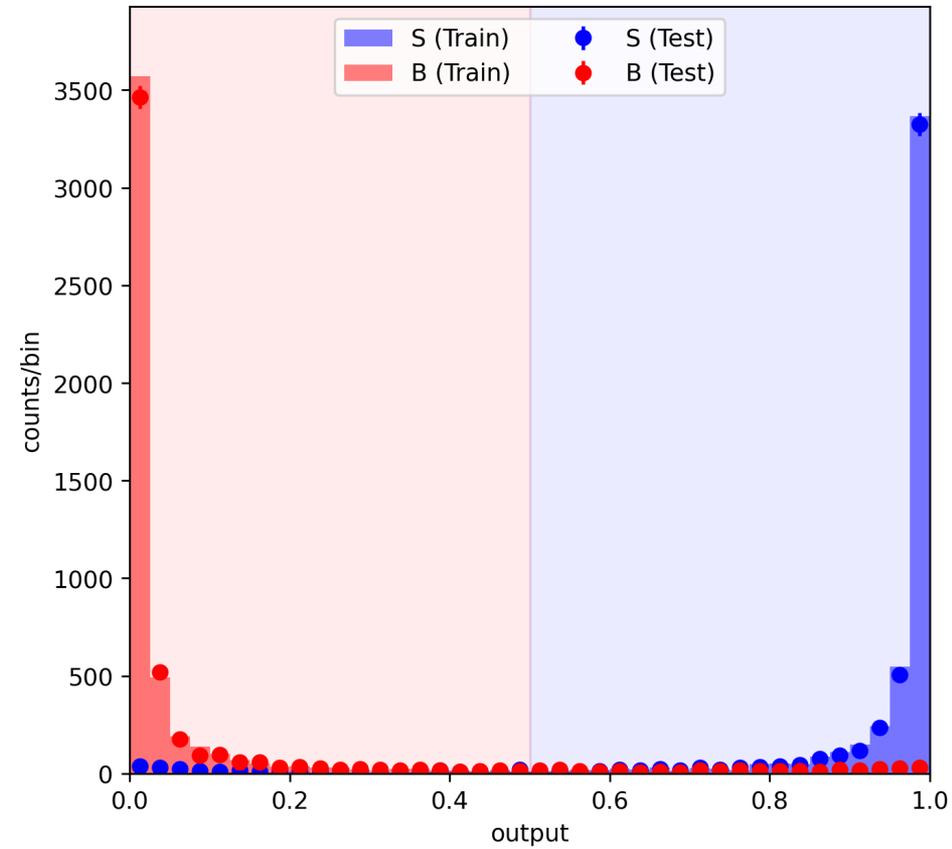
Random forest



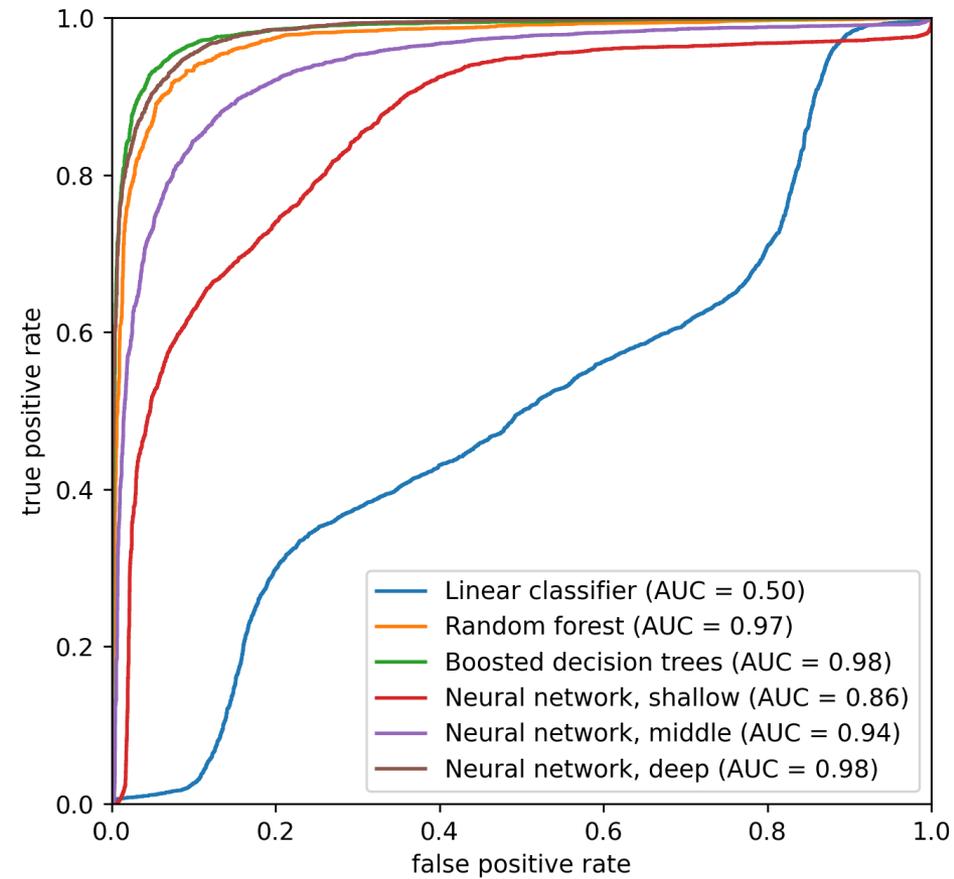
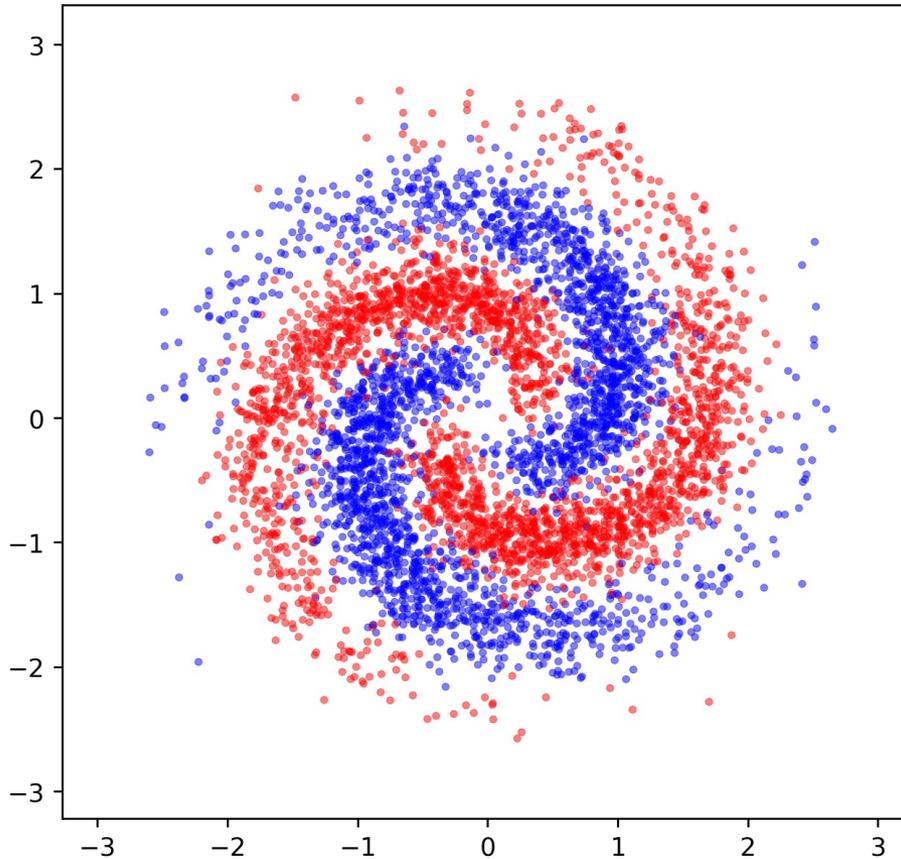
Boosted decision trees



Boosted decision trees



Input data





- Convolutional neural networks
- Unsupervised learning
 - Clustering Anomaly detection
- Autoencoders
- Generative Networks

Tensor flow playground



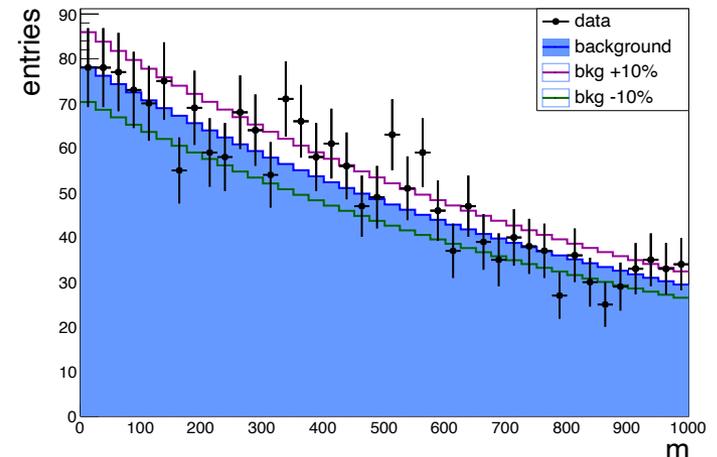
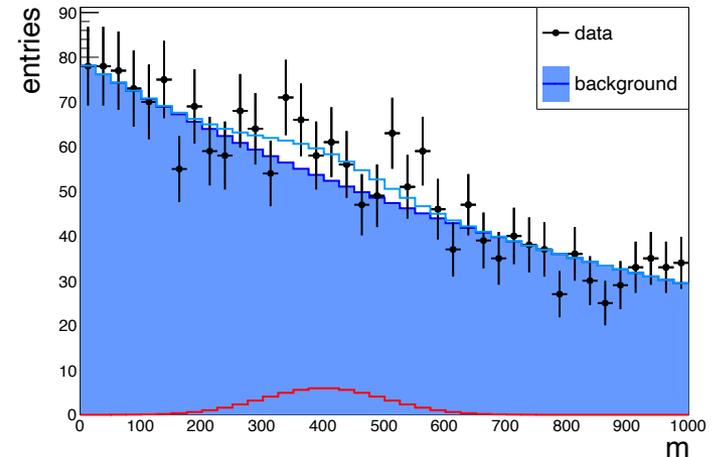
- <https://playground.tensorflow.org/>



- Discoveries and significance level
 - Upper limits
 - The flip-flopping issue and the Feldman-Cousins approach
 - Bayesian upper limits for event-counting problems
 - Modified frequentist approach: the CLs method
 - Treatment of nuisance parameters and systematic uncertainties
 - The profile likelihood method
 - Applications of Wilks' theorem
 - Different test statistics
 - Asimptotic approximations and the Asimov dataset
 - The look-elsewhere effect
 - Understanding the "Brazilian" exclusion plots
-

- We want to test our data sample against two hypotheses about the theoretical underlying model:
 - H_0 : the data are described by a model that contains background only
 - H_1 : the data are described by a model that contains signal plus background
- Our discrimination is based on a test statistic λ whose distribution is known under the two hypotheses
 - Let's assume λ tends to have (conventionally) large values if H_1 is true and small values if H_0 is true
 - This convention is consistent with λ being the likelihood ratio $L(x|H_1)/L(x|H_0)$
- Under the frequentist approach, compute the p -value as the probability that λ is greater or equal to than the value λ_{obs} we observed

Are data below more consistent with a background fluctuation or with a peaking excess?



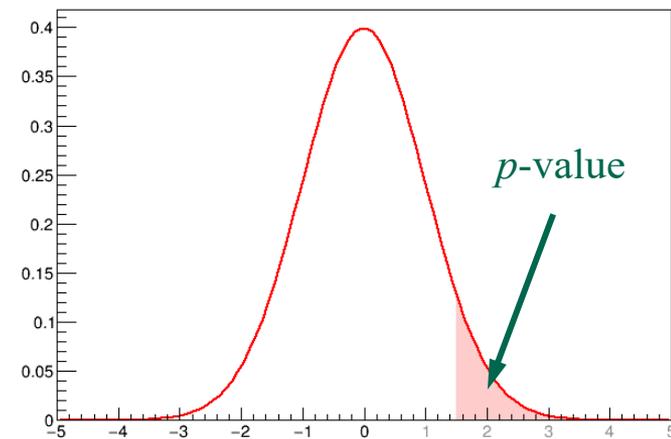
Significance

- The *p-value* is usually converted into an equivalent area of a Gaussian tail:

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

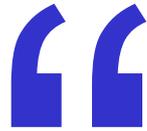
$$Z = \Phi^{-1}(1 - p)$$

Φ = cumulative of a normal distribution



$Z =$
significance level

- In literature we find, by convention:
 - If the significance is $Z > 3$ (“ 3σ ”) one claims “*evidence of*”
 - Probability that background fluctuation will produce a test statistic at least as extreme as the observed value : $p < 1.349 \times 10^{-3}$
 - If the significance is $Z > 5$ (“ 5σ ”) one claims “*observation*” (**discovery!**)
 - $p < 2.87 \times 10^{-7}$
- Note:** the probability that background produces a large test statistic is not equal to probability of the null hypothesis (background only), which has only a Bayesian sense



The p-value was never intended to be a substitute for scientific reasoning. Well-reasoned statistical arguments contain much more than the value of a single number and whether that number exceeds an arbitrary threshold. The ASA statement is intended to steer research into a ‘post $p < 0.05$ era’.

- 1. p-values can indicate how incompatible the data are with a specified statistical model.*
- 2. p-values do not measure the probability that the studied hypothesis is true, or the probability that the data were produced by random chance alone.*
- 3. Scientific conclusions and business or policy decisions should not be based only on whether a p-value passes a specific threshold.*
- 4. Proper inference requires full reporting and transparency.*
- 5. A p-value, or statistical significance, does not measure the size of an effect or the importance of a result.*
- 6. By itself, a p-value does not provide a good measure of evidence regarding a model or hypothesis.*

Ronald L. Wasserstein · Nicole A. Lazar

The ASA's statement on p-values: context, process, and purpose

DOI:10.1080/00031305.2016.1154108

<http://amstat.tandfonline.com/doi/abs/10.1080/00031305.2016.1154108>



- From Cowan *et al.*, EPJC 71 (2011) 1554:

“ *It should be emphasized that in an actual scientific context, rejecting the background-only hypothesis in a statistical sense is only part of discovering a new phenomenon. One’s **degree of belief** that a new process is present will depend in general on other factors as well, such as the **plausibility of the new signal hypothesis** and the **degree to which it can describe the data.***

Here, however, we only consider the task of determining the p -value of the background-only hypothesis; if it is found below a specified threshold, we regard this as “discovery”.

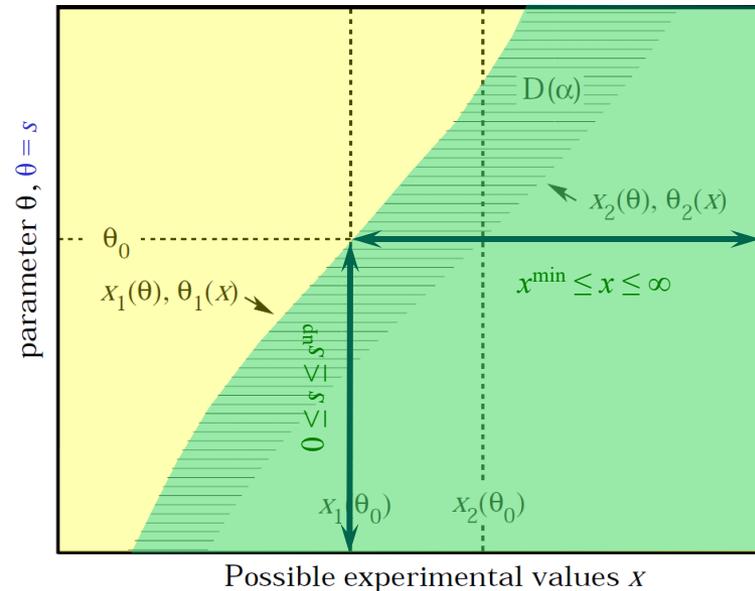
”

Complementary role of Frequentist and Bayesian approaches ☺

- Measure the amount of excluded region resulting from our (negative) search for a new signal
- Building a **fully asymmetric Neyman confidence belt** based on the considered test statistic x
- Invert the belt, find the allowed interval:

$$s \in [s_1, s_2] \Rightarrow s \in [0, s^{\text{up}}]$$

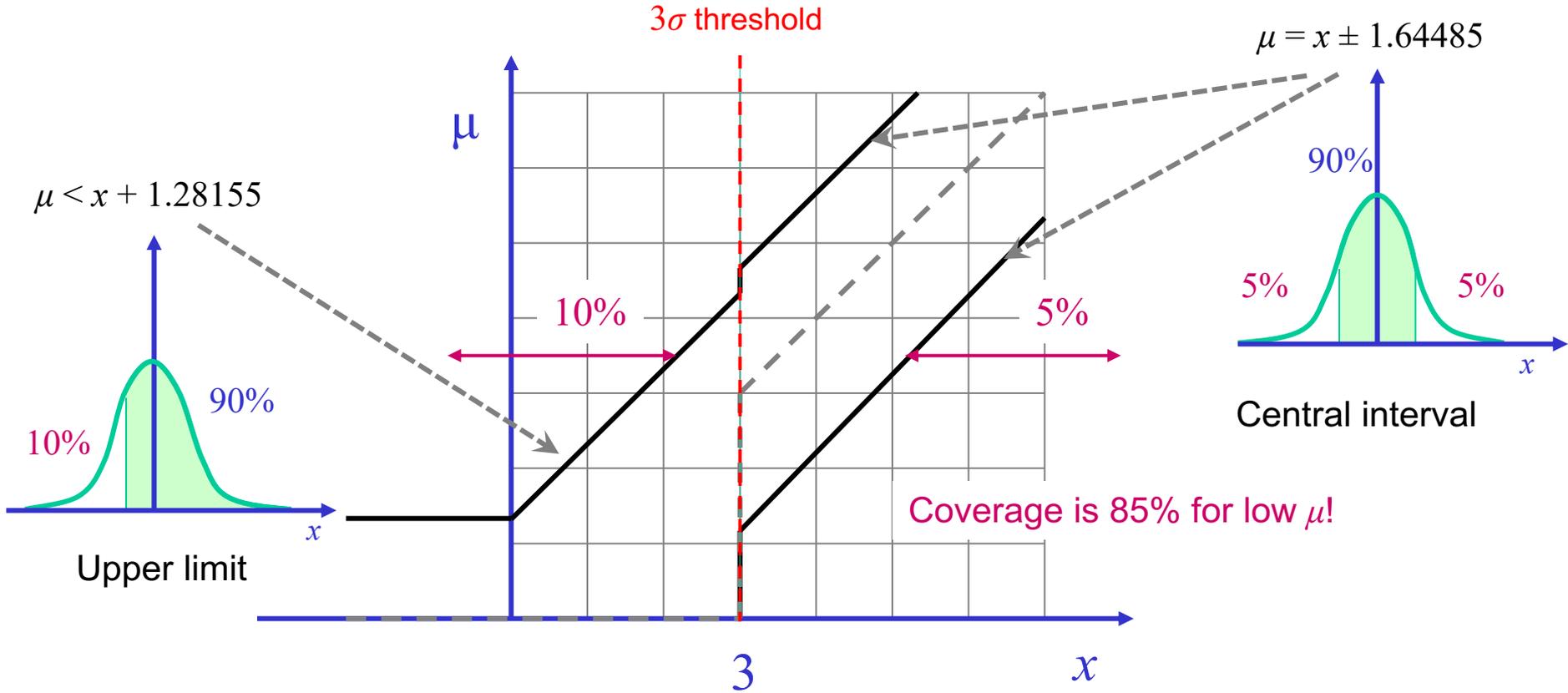
- **Upper limit** = upper extreme of the asymmetric interval $[0, s^{\text{up}}]$
- In case the observable x is **discrete** (e.g.: the number of events n in a counting experiments), **the coverage may not be exact**



- When to quote a **central value** or **upper limit**?
- A popular choice was:
 - *“Quote a 90% CL upper limit of the measurement if the significance is below 3σ ; quote a central value otherwise”*
 - Upper limit \leftrightarrow central interval decided according to observed data
- **This produces an incorrect coverage!**



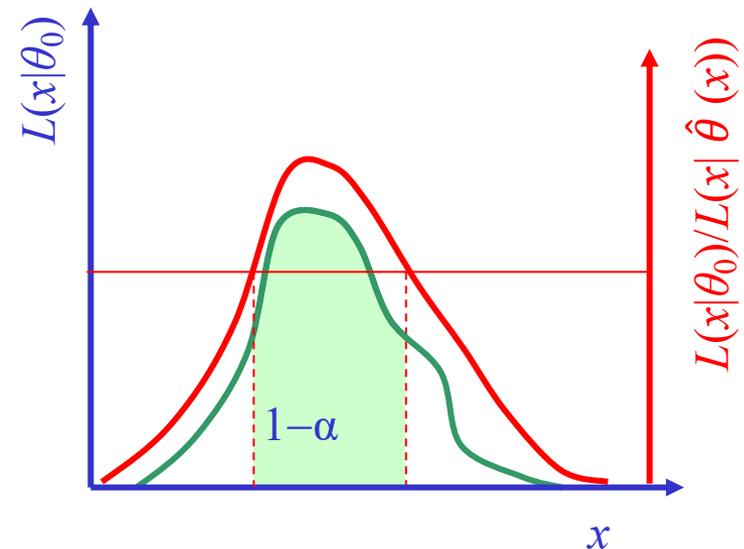
- Assume a Gaussian with a fixed width: $\sigma = 1$



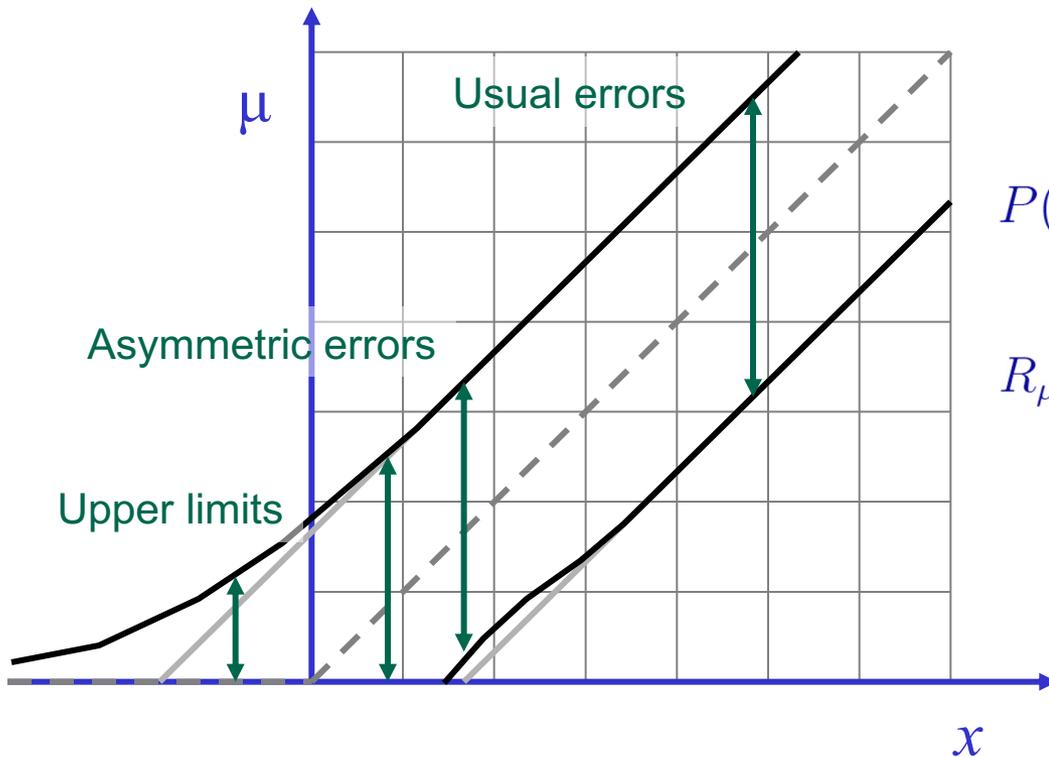
- Feldman and Cousins proposed a criterion to define the Neyman belt based in a likelihood ratio test:

$$R_\mu = \{x : L(x|\theta_0) / L(x|\hat{\theta}) > k_\alpha\}$$

- The value k_α depends on the desired significance level α
- $H_0: \theta = \hat{\theta}$, the best-fit value
- $H_1: \theta = \theta_0$, the specific value considered for the Neyman belt construction



- Application to the Gaussian case:



$$\hat{\mu} = \max(x, 0)$$

$$P(x|\hat{\mu}) = \begin{cases} \frac{1}{\sqrt{2\pi}}, & x \geq 0, \\ \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}, & x < 0. \end{cases}$$

$\hat{\mu} = x$ for $x \geq 0$

$$R_{\mu}(x) = \frac{P(x|\mu)}{P(x|\hat{\mu})} = \begin{cases} e^{-\frac{(x-\mu)^2}{2}}, & x \geq 0, \\ e^{-\frac{x\mu - \mu^2}{2}}, & x < 0. \end{cases}$$

Confidence intervals must be computed numerically, even for this simple Gaussian case!



- The simplest search for a new signal consists of counting the number of events passing a specified selection
- The number of selected events n is distributed according to a Poissonian distribution
- Expected n for signal + background (H_1): $s + b$
- Expected n for background only (H_0): b

- We measure n events, we want to compare with the two hypotheses H_1 and H_0 .
- Simplest case: b is known with negligible uncertainty
 - If not, uncertainty on its estimate must be taken into account

- Let's assume the background b is known with no uncertainty:

$$L(n; s) = \frac{(s + b)^n}{n!} e^{-(s+b)}$$

- A uniform prior, $\pi(s) = 1$ simplifies, as usual, the computation:

$$1 - \alpha = \int_0^{s^{\text{up}}} P(s|n) ds = \frac{\int_0^{s^{\text{up}}} L(n; s) \pi(s) ds}{\int_0^{\infty} L(n; s) \pi(s) ds}$$

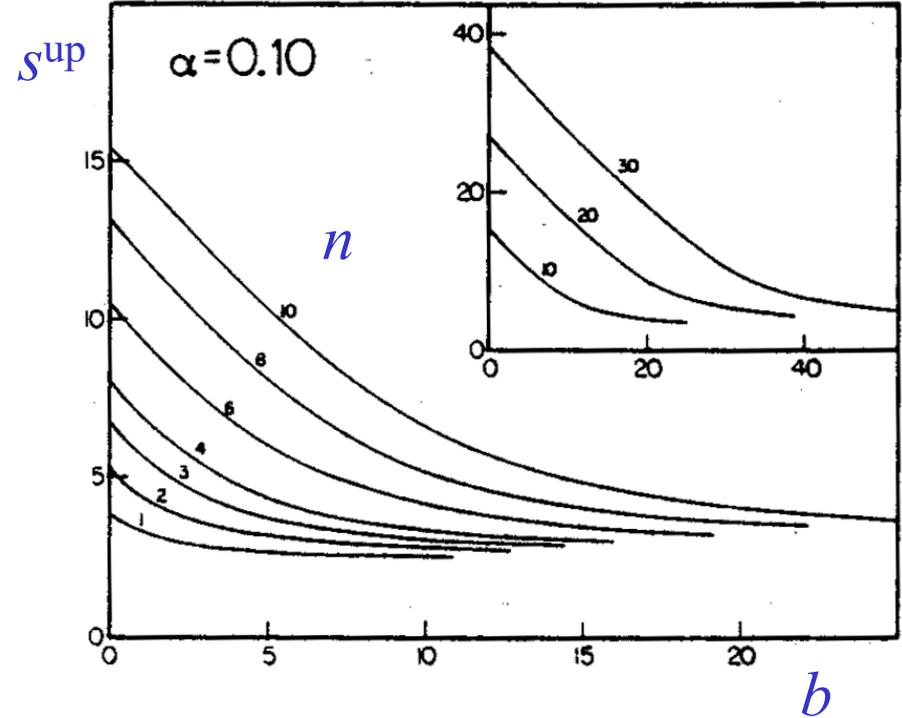
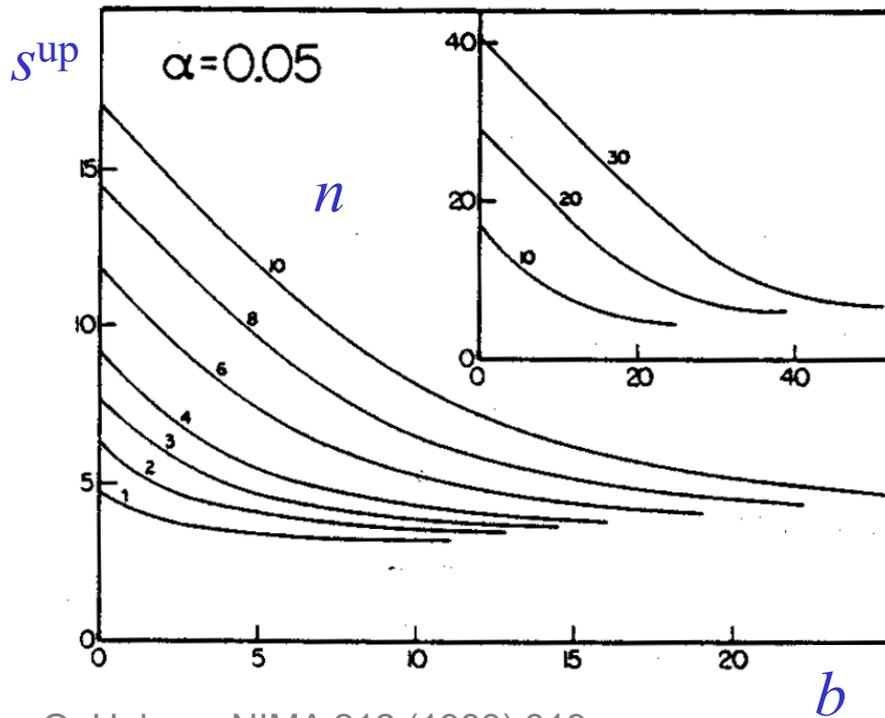

- Inverting the equation gives the upper limit s^{up}
- For $n = 0$ s^{up} does not depend on b :

$$\alpha = e^{-s^{\text{up}}}$$

- $s < 2.303$ (90% CL) $\leftarrow \alpha = 0.1$
- $s < 2.996$ (95% CL) $\leftarrow \alpha = 0.05$

$$\alpha = e^{-s^{\text{up}}} \frac{\sum_{m=0}^n \frac{(s^{\text{up}} + b)^m}{m!}}{\sum_{m=0}^n \frac{b^m}{m!}}$$

- Upper limits decrease as b increases and increase as n increases
- For $n = 0$, upper limits are not sensitive on b (given in prev. slide)



O. Helene. NIMA 212 (1983) 319



- Assume we have negligible background ($b = 0$) and we measure zero events ($n = 0$)
- The likelihood function simplifies as:

$$L(n = 0; s) = \text{Pois}(0; s) = e^{-s}$$

- The (fully asymmetric) Neyman belt inversion is pretty simple:

$$P(n \leq 0; s^{\text{up}}) = \alpha \rightarrow s^{\text{up}} = -\ln \alpha$$

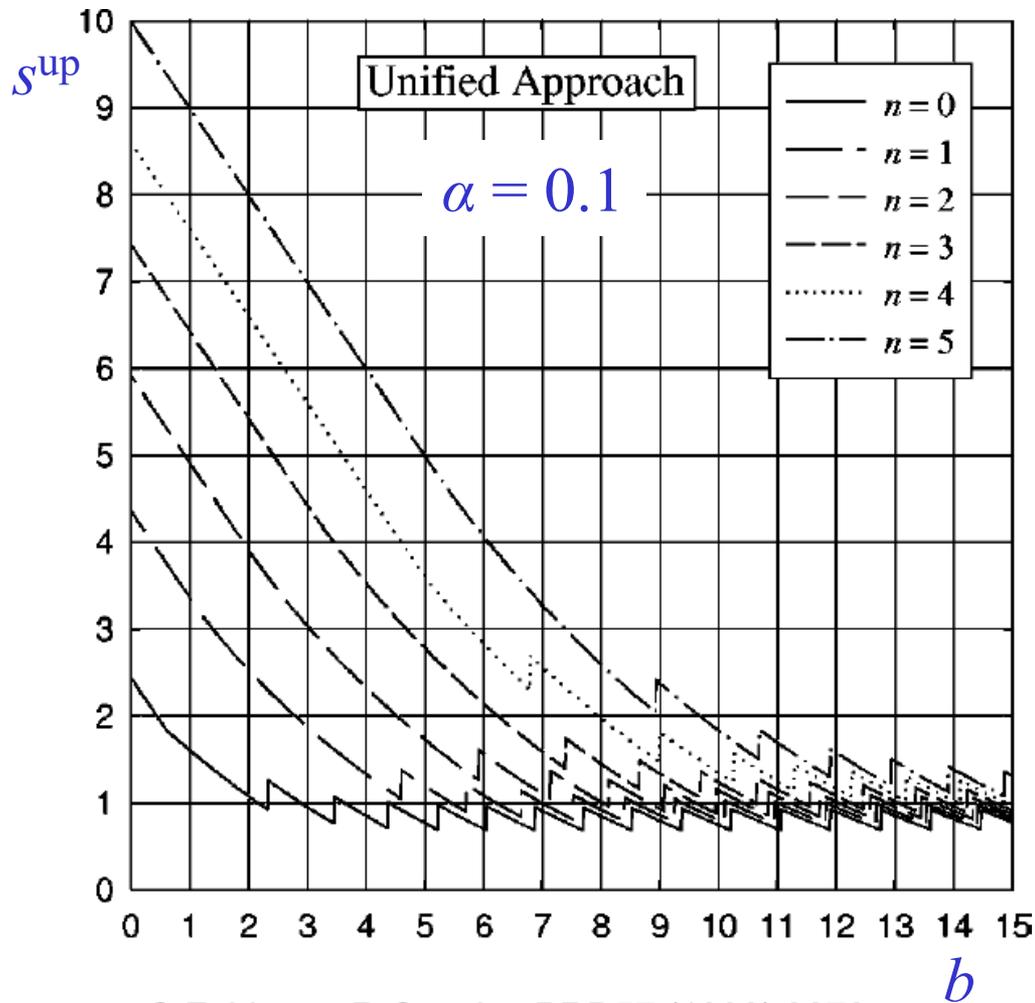
- The results are by chance identical to the Bayesian computation:

$$s < 2.303 \text{ (90\% CL)} \leftarrow \alpha = 0.1$$

$$s < 2.996 \text{ (95\% CL)} \leftarrow \alpha = 0.05$$

- In spite of the numerical coincidence, the interpretation of frequentist and Bayesian upper limits remain very different!
- **Warning:** this evaluation suffer from the “flip-flopping” problem, so the coverage is spoiled if you decide to switch from upper limit to a central value depending on the observed significance!

- F&C intervals cure the flip-flopping issue and ensure the correct coverage
 - May overcover for discrete variables
- The “ripple” structure is due to the discrete nature of Poissonian counting
- Note that even for $n = 0$ the upper limit decrease as b increases (apart from ripple effects)
- If two experiment are designed for an expected background of –say– 0.5 and 0.01, the “worse” one has the best expected upper limit



G.Feldman, R.Cousins PRD57 (1998) 3873
 C. Giunti, PRD59 (1999), 053001



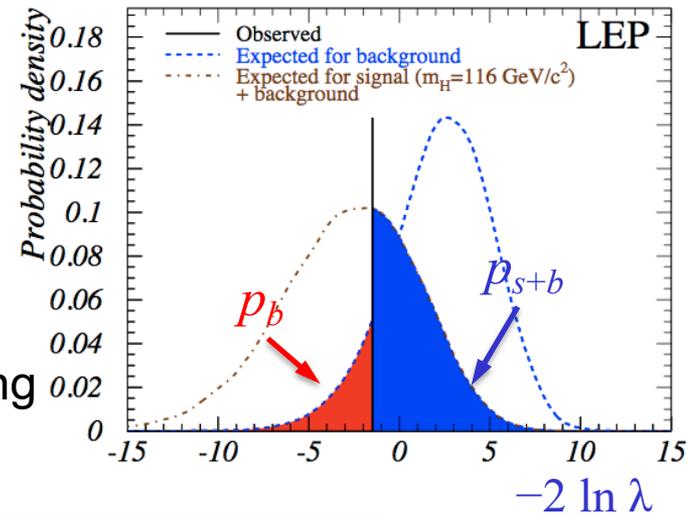
- *“The intervals constructed according to the unified procedure [FC] for a Poisson variable n consisting of signal and background have the property that for $n = 0$ observed events, the upper limit decreases for increasing expected background. This is counter-intuitive, since it is known that if $n = 0$ for the experiment in question, then no background was observed, and therefore one may argue that the expected background should not be relevant. The extent to which one should regard this feature as a drawback is a subject of some controversy”*

- A **modified approach** was proposed for the first time when combining the limits on the Higgs boson search from the four LEP experiments, ALEPH, DELPHI, L3 and OPAL
- Given a test statistic $\lambda(x)$, determine its distribution for the two hypotheses $H_1(s + b)$ and $H_0(b)$, and compute:

$$\left\{ \begin{array}{l} p_{s+b} = P(\lambda(x|H_1) \leq \lambda^{\text{obs}}) \\ p_b = P(\lambda(x|H_0) \geq \lambda^{\text{obs}}) \end{array} \right.$$

- The upper limit is computed, instead of requiring $p_{s+b} \leq \alpha$, on the modified statistic $CL_s \leq \alpha$:

- Since $1 - p_b \leq 1$, $CL_s \geq p_{s+b}$, hence upper limits computed with the CL_s method are always **conservative**



$$CL_s = \frac{p_{s+b}}{1 - p_b}$$

Note: $\lambda \leq \lambda^{\text{obs}}$ implies $-2 \ln \lambda \geq \lambda^{\text{obs}}$

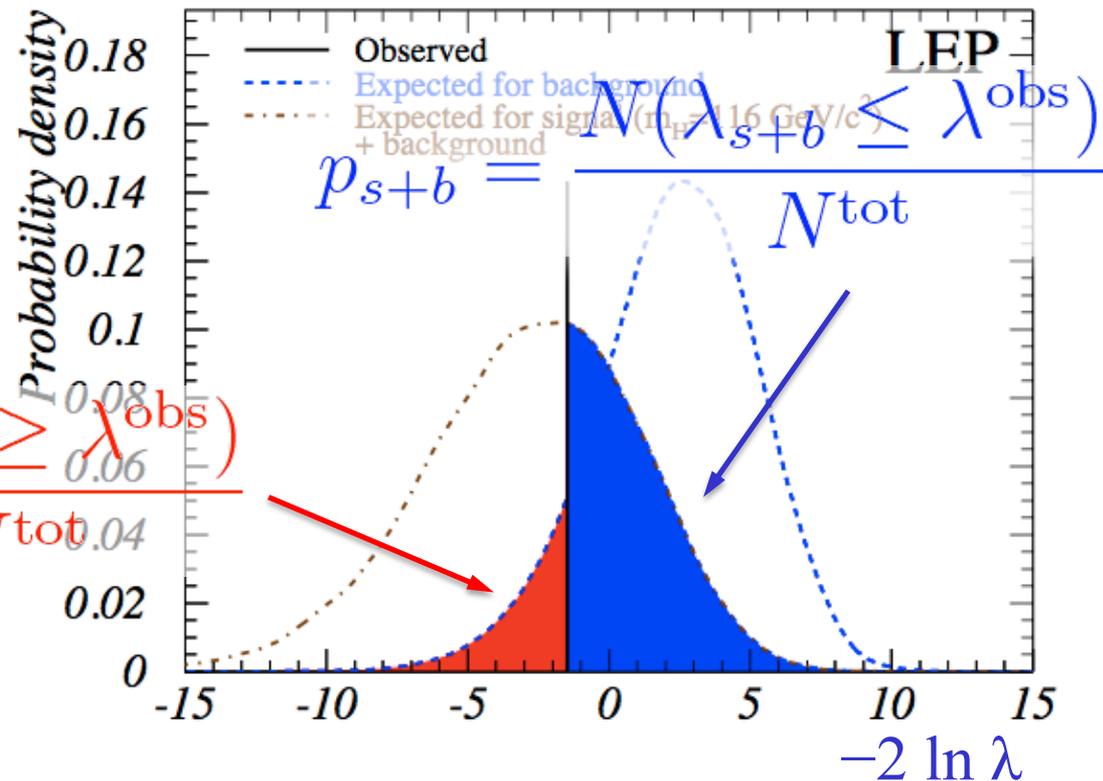
CL_s with toy experiments

- In practice, p_b and p_{s+b} are computed in from simulated pseudo-experiments (“toy Monte Carlo”)

$$CL_s = \frac{N(\lambda_{s+b} \leq \lambda^{obs})}{N(\lambda_b \leq \lambda^{obs})}$$

$$p_b = \frac{N(\lambda_b \geq \lambda^{obs})}{N^{tot}}$$

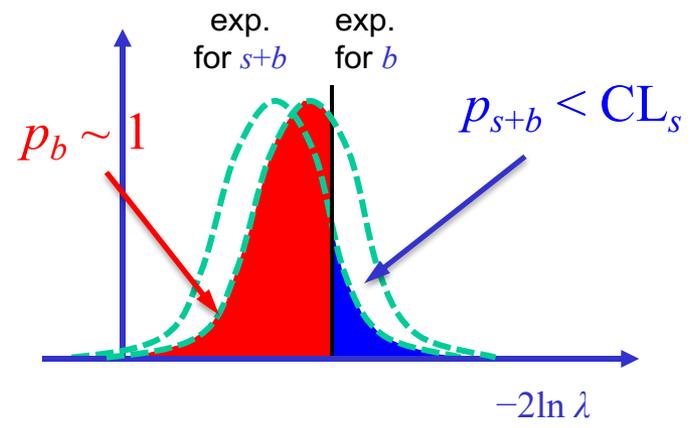
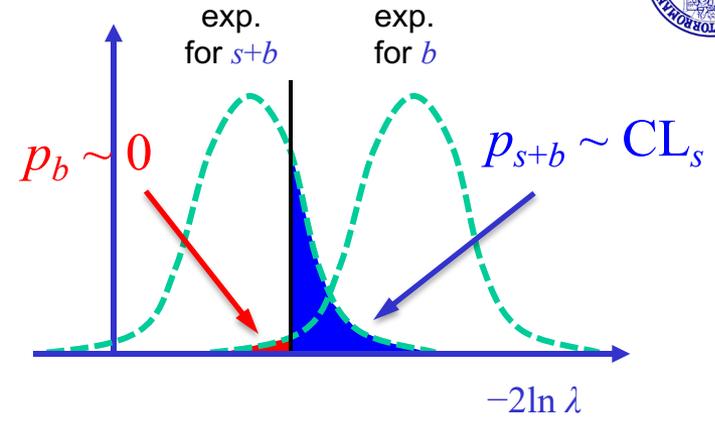
Plot from LEP Higgs combination paper



INFN Main CL_s features



- p_{s+b} : probability to obtain a result which is **less compatible** with the signal than the observed result, **assuming the signal hypothesis**
- p_b : probability to obtain a result **less compatible** with **the background-only hypothesis** than the observed one
- If the two distributions are **very well separated** and H_1 is true, then p_b will be very small $\Rightarrow 1 - p_b \sim 1$ and $CL_s \sim p_{s+b}$, i.e: the ordinary p -value of the $s+b$ hypothesis
- If the two distributions **largely overlap**, then if p_b will be large $\Rightarrow 1 - p_b$ **small**, preventing CL_s to become very small
- $CL_s < 1 - \alpha$ prevents rejecting cases where the experiment has little sensitivity



$$CL_s = \frac{p_{s+b}}{1 - p_b} = \frac{P(\lambda_{s+b} \leq \lambda^{\text{obs}})}{P(\lambda_b \leq \lambda^{\text{obs}})}$$

- Let's consider the previous event counting experiment, using $n = n^{\text{obs}}$ as test statistic
- In this case CL_s can be written as:

$$CL_s = \frac{P(n \leq n^{\text{obs}} | s + b)}{P(n \leq n^{\text{obs}} | b)}$$

- Explicitating the Poisson distribution, the computation gives the same result as for the Bayesian case with a uniform prior
- In many cases the CL_s upper limits give results that are very close, numerically, to Bayesian computations done assuming a uniform prior
- **But the interpretation is very different from Bayesian limits!**

$$\alpha = e^{-s^{\text{up}}} \frac{\sum_{m=0}^n \frac{(s^{\text{up}} + b)^m}{m!}}{\sum_{m=0}^n \frac{b^m}{m!}}$$



- *“A specific modification of a purely classical statistical analysis is used to **avoid excluding or discovering signals which the search is in fact not sensitive to**”*
- *“The use of CL_s is a conscious decision not to insist on the frequentist concept of full coverage (to guarantee that the confidence interval doesn't include the true value of the parameter in a fixed fraction of experiments).”*
- *“**confidence intervals obtained in this manner do not have the same interpretation as traditional frequentist confidence intervals nor as Bayesian credible intervals**”*

A. L. Read, Modified frequentist analysis of search results
(the CL_s method), 1st Workshop on Confidence Limits, CERN, 2000

- Usually, signal extraction procedures (fits, upper limits setting) determine, together with parameters of interest, also nuisance parameters that model effects not strictly related to our final measurement

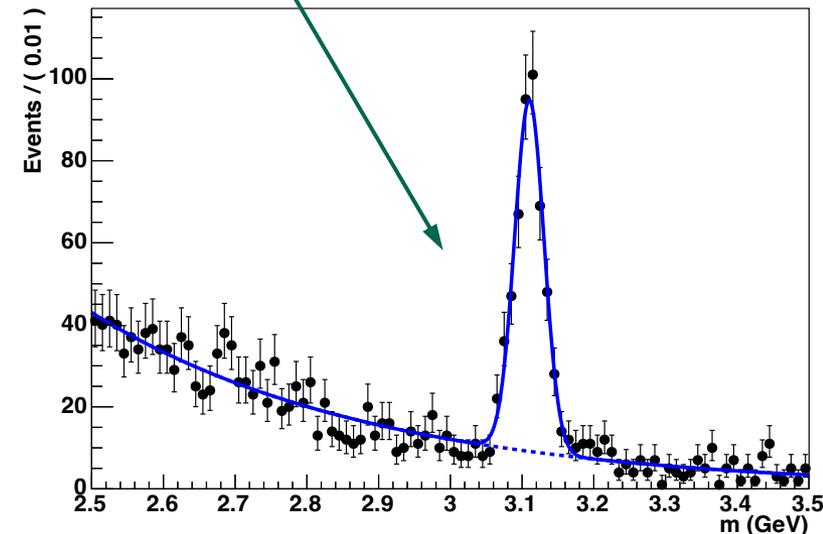
- Background yield and shape parameters
- Detector resolution
- ...

$$L(m; s, b, \mu, \sigma, \lambda) = \frac{e^{-(s+b)}}{n!} \left(s \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(m-\mu)^2}{2\sigma^2}} + b\lambda e^{-\lambda m} \right)$$

- Nuisance parameters are also used to model sources of **systematic uncertainties**

- Often referred to nominal values

- Examples: cross section × int. lumi
- $b = \beta \sigma_b L_{\text{int}}$ with $\beta^{\text{nominal}} = 1$
- $b = e^\beta \sigma_b L_{\text{int}}$ with $\beta^{\text{nominal}} = 0$
(negative yields not allowed!)





- Notation below: μ = parameter(s) of interest, θ = nuisance parameter(s)
- No special treatment:

$$P(\mu, \theta|x) = \frac{L(x; \mu, \theta)\pi(\mu, \theta)}{\int L(x; \mu', \theta')\pi(\mu', \theta')d\mu'd\theta'}$$

- $P(\mu|x)$ obtained as marginal PDF of μ obtained integrating on θ :

$$P(\mu|x) = \int P(\mu, \theta|x)d\theta = \frac{\int L(x; \mu, \theta)\pi(\mu, \theta)d\theta}{\int L(x; \mu', \theta)\pi(\mu', \theta)d\mu'd\theta}$$



- Introduce a complementary dataset to constrain the nuisance parameters θ (e.g.: calibration data, background estimates from control sample...)
- Formulate the statistical problem in terms of both the main data sample (x) and the control sample (y)

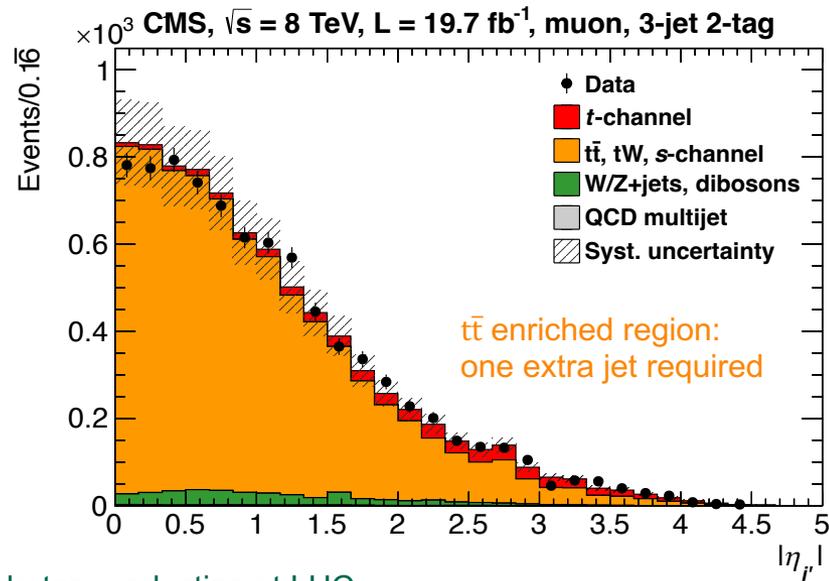
$$L(x, y; \mu, \theta) = L(x; \mu, \theta)L(y; \theta)$$

- Not always the control sample data are available
 - E.g.: calibration from test beam, stored in different formats, control samples analyzed with different software framework...
 - In some cases may be complex and CPU intensive
- Simplest case; assume known PDF for “nominal” value of θ^{nom} (e.g.: estimate with Gaussian uncertainty)

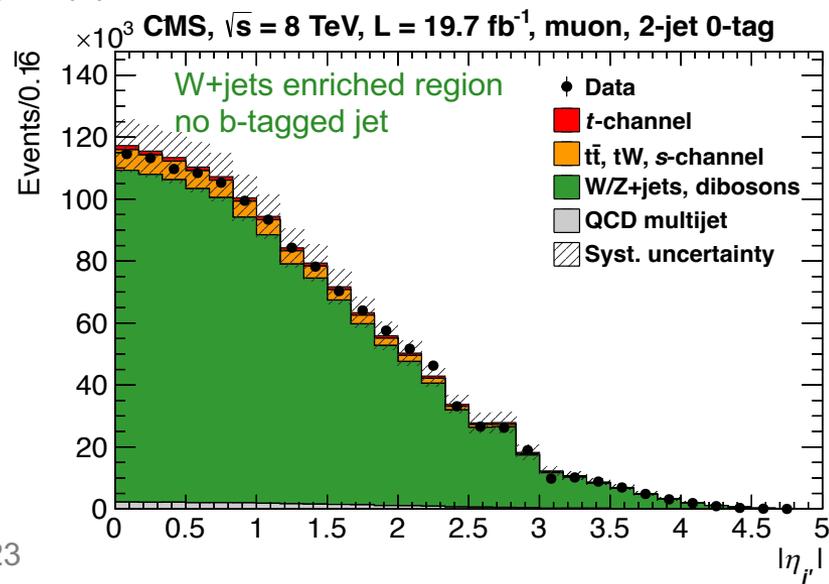
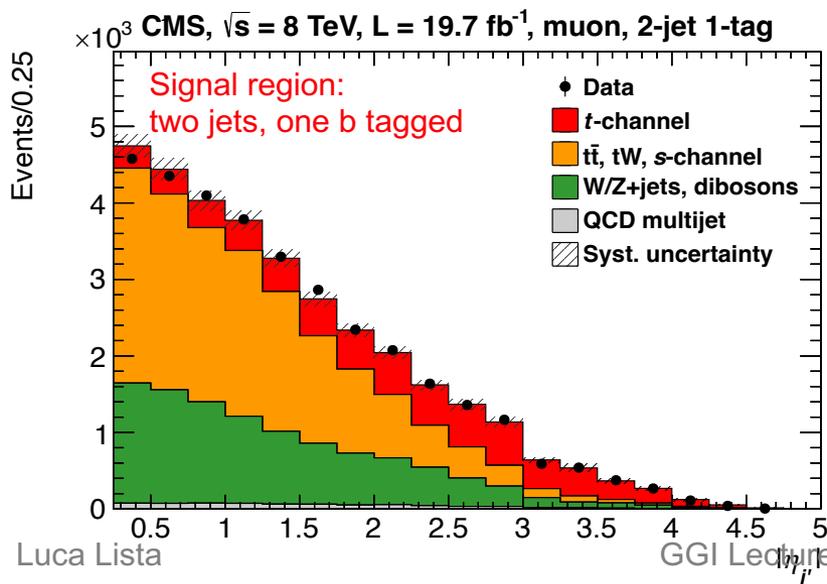
$$L(x, \theta^{\text{nom}}; \mu, \theta) = L(x; \mu, \theta)L(\theta^{\text{nom}}; \theta)$$

Fitting control regions

- In some cases, background parameters can be constrained from statistically independent **control samples**
 - Consider possible signal contamination!
- Background yield can be measured in **background-enriched regions** and extrapolated to **signal regions** applying scale factors predicted by simulation
- Complete likelihood function = product of likelihood functions in each considered regions, sharing common nuisance parameters
 - Typically: **background rates**



Measurement of single-top production at LHC





- Method proposed by Cousins and Highland
 - Add posterior from another experiment into the likelihood definition
 - Integrate the likelihood function over the nuisance parameters

$$L^{\text{hybrid}}(x; \mu) = \int L(x; \mu, \theta) L(\theta^{\text{nom}}; \theta) d\theta$$

- Also called “hybrid” approach, because a partial Bayesian approach is implicit in the integration
 - Bayesian integration of PDF, then likelihood used in a frequentist way
- **Not guaranteed to provide exact frequentist coverage!**
- Numerical studies with pseudo experiments showed that the **hybrid CL_s upper limits** gives very similar results to **Bayesian limit** assuming a uniform prior

- Define a test statistic based on a likelihood ratio:

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

← Fix μ , fit θ

← Fit both μ and θ

- μ is usually the “signal strength” (i.e.: $\sigma/\sigma_{\text{th}}$) in case of a search for a new signal
- Different ‘flavors’ of test statistics
 - E.g.: deal with unphysical $\mu < 0$, ...
- The distribution of $q_\mu = -2 \ln \lambda(\mu)$ may be asymptotically approximated to the distribution of a χ^2 with one degree of freedom (one parameter of interest = μ) due to the **Wilks’ theorem**
(→ next slide)

- Consider a likelihood function from N measurements:

$$\prod_{i=1}^N L(x_1^i, \dots, x_n^i; \theta_1, \dots, \theta_m) = \prod_{i=1}^N L(\vec{x}_i; \vec{\theta})$$

- Assume that H_0 and H_1 are two **nested hypotheses**, i.e.: they can be expressed as:

$$\vec{\theta} \in \Theta_0 \quad \vec{\theta} \in \Theta_1$$

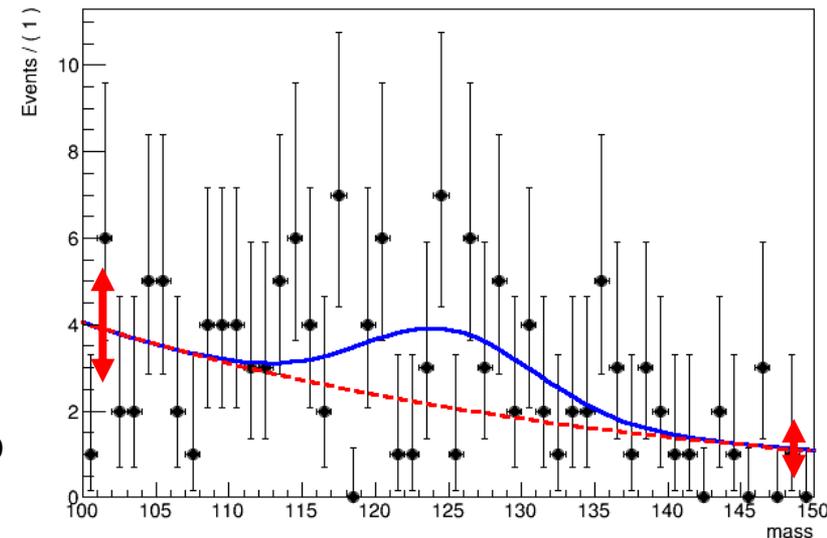
- Where $\Theta_0 \subseteq \Theta_1$. Then, the following quantity for $N \rightarrow \infty$ is distributed as a χ^2 with n.d.o.f. equal to the difference of Θ_0 and Θ_1 dimensionality:

$$\chi_r^2 = -2 \ln \frac{\sup_{\vec{\theta} \in \Theta_0} \prod_{i=1}^N L(\vec{x}_i; \vec{\theta})}{\sup_{\vec{\theta} \in \Theta_1} \prod_{i=1}^N L(\vec{x}_i; \vec{\theta})}$$

- E.g.: searching for a signal with strength μ , $H_0: \mu = 0$, $H_1: \mu \geq 0$ we have the profile likelihood (**supremum = best fit value**):

$$\chi_r^2(\mu) = -2 \ln \frac{\sup_{\vec{\theta}} \prod_{i=1}^N L(\vec{x}_i; \mu, \vec{\theta})}{\sup_{\mu', \vec{\theta}} \prod_{i=1}^N L(\vec{x}_i; \mu', \vec{\theta})}$$

- Gaussian signal over an exponential background
- Fix all parameters from theory prediction, fit only the signal yield
- Assume a –say– 30% uncertainty on the background yield
- A log normal model may be assumed to avoid unphysical negative yields



b_0 = true
(unknown)
value



– $b_0 = b e^\beta$, where our estimate β is known
with a Gaussian uncertainty $\sigma_\beta = 0.3$

b = our
estimate

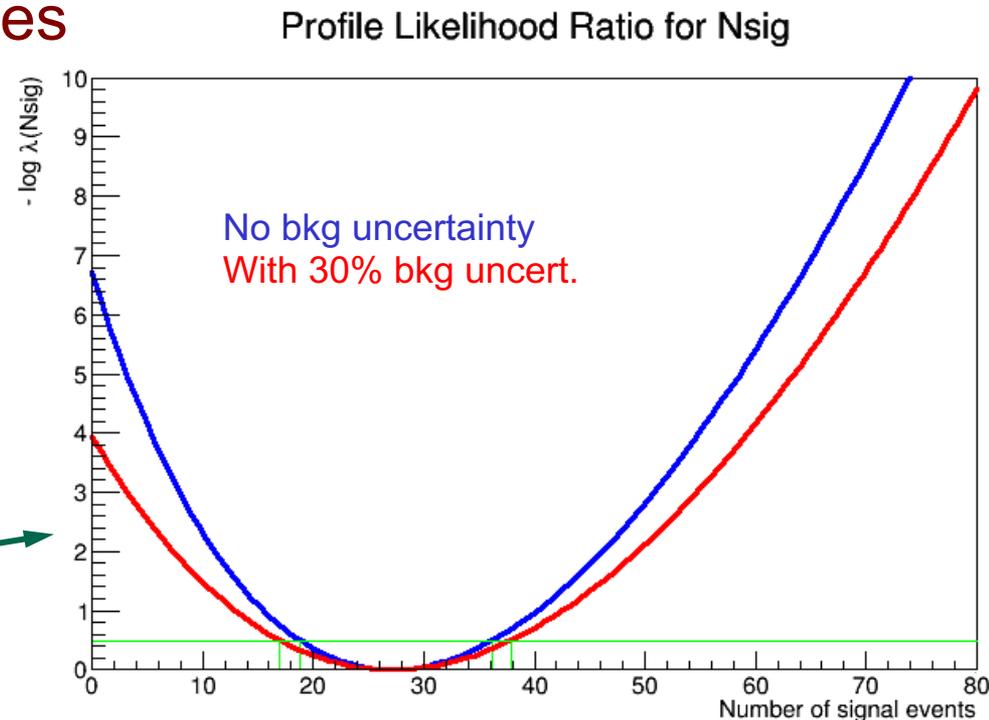
$$L(m; s, \beta) = L_0(m; s, b_0 = b e^\beta) P(\beta; \sigma_\beta)$$

$$L_0(m; s, b_0) = \frac{e^{-(s+b_0)}}{n!} \left(s \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(m-\mu)^2}{2\sigma^2}} + b_0 \lambda e^{-\lambda m} \right)$$

$$P(\beta; \sigma_\beta) = \frac{1}{\sqrt{2\pi}\sigma_\beta} e^{-\frac{\beta^2}{2\sigma_\beta^2}}$$

- The profile likelihood shape is broadened, with respect to the usual likelihood function, due to the presence of nuisance parameter β (loss of information) that model systematic uncertainties
- Uncertainty on s increases
- Significance for discovery using s as test statistic decreases

This implementation is based on RooStats, a package, released as optional library with ROOT <http://root.cern.ch>



Significance evaluation



- Assume $\mu = 0$, if $q_0 = -2 \ln \lambda(0)$ can be approximated by a χ^2 with one d.o.f., then the significance is approximately equal to:

$$Z \cong \sqrt{q_0}$$

- The level of approximation can be verified with a computation done using pseudo experiments:
- Generate a large number of toy samples with zero background and determine the distribution of $q_0 = -2 \ln \lambda(0)$, then count the fraction of cases with values greater than the measured value (*p-value*), and convert it to Z :

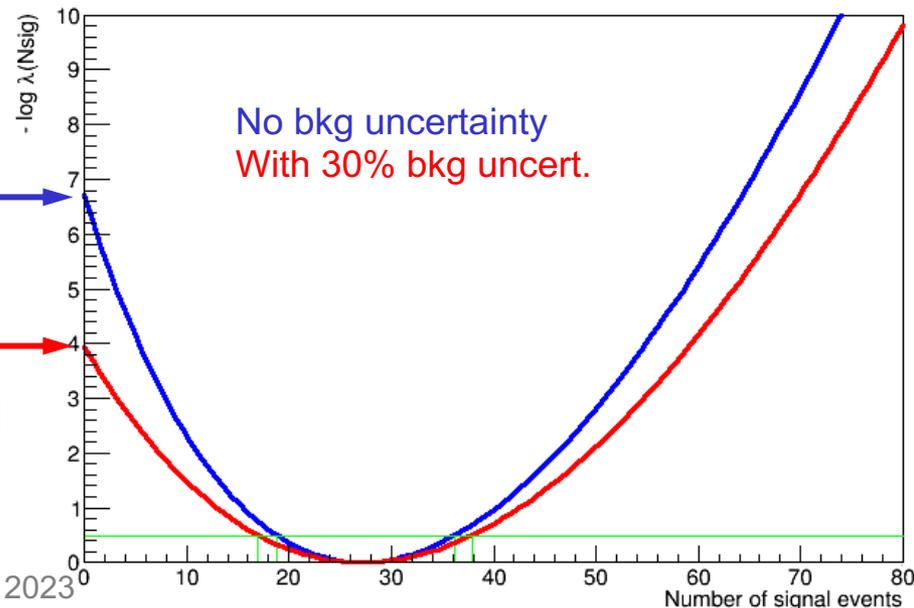
$$Z \cong \sqrt{2 \times 6.66} = 3.66$$

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

$$Z \cong \sqrt{2 \times 3.93} = 2.81$$

- Toy samples may be unpractical for very large Z

Profile Likelihood Ratio for Nsig





- Asymptotic approximate formulae exist for most of adopted estimators
- If we want to test μ and we suppose data are distributed according to μ' , we can write:

$$-2 \ln \lambda(\mu) = \frac{(\mu - \hat{\mu})^2}{\sigma_{\hat{\mu}}^2} + \mathcal{O}(1/\sqrt{N})$$

- where $\hat{\mu}$ is distributed according to a Gaussian with average μ' and standard deviation σ (A. Wald, 1943)
- The covariance matrix can be asymptotically approximated by:

$$V_{ij}^{-1} = - \left\langle \frac{\partial^2 \ln L}{\partial \theta_i \partial \theta_j} \right\rangle$$

where μ' is assumed as signal strength value

- Case by case, the estimate of σ (from the inversion of V_{ij}^{-1}) can be determined

A. Wald, Trans. of AMS 54 n.3 (1943) 426-482

G. Cowan et al., EPJ C71 (2011) 1554



- Under the true hypothesis μ , $\hat{\mu}$ is distributed around μ and the test statistic, neglecting the $\mathcal{O}(1/\sqrt{N})$ term, is distributed according to a χ^2 with one degree of freedom (Wilks' theorem):

$$-2 \ln \lambda(\mu) = \frac{(\mu - \hat{\mu})^2}{\sigma_{\hat{\mu}}^2}$$

- If $\hat{\mu}$ is distributed around $\mu' \neq \mu$ the distribution of the test statistic is a **non-central χ^2** .

- Test statistic for **discovery**:

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

- In case of a negative estimate of μ , set the test statistic to zero: consider only positive μ as evidence against the background-only hypothesis. Approximately: $Z \cong \sqrt{q_0}$.

- Test statistic for **upper limits**:

$$q_\mu = \begin{cases} -2 \ln \lambda(\mu), & \hat{\mu} \leq \mu, \\ 0, & \hat{\mu} > \mu. \end{cases}$$

- If the estimate is larger than the assumed μ , an upward fluctuation occurred. Don't exclude μ in those cases, hence set the statistic to zero

- **Higgs** test statistic:

$$\tilde{q}_\mu = \begin{cases} -2 \ln \frac{L(\vec{x}|\mu, \hat{\theta}(\mu))}{L(\vec{x}|0, \hat{\theta}(0))}, & \hat{\mu} < 0, \quad \leftarrow \text{Protect for unphysical } \mu < 0 \\ -2 \ln \frac{L(\vec{x}|\mu, \hat{\theta}(\mu))}{L(\vec{x}|\hat{\mu}, \hat{\theta})}, & 0 \leq \hat{\mu} \leq \mu, \\ 0, & \hat{\mu} > \mu. \quad \leftarrow \text{As for upper limits statistic} \end{cases}$$

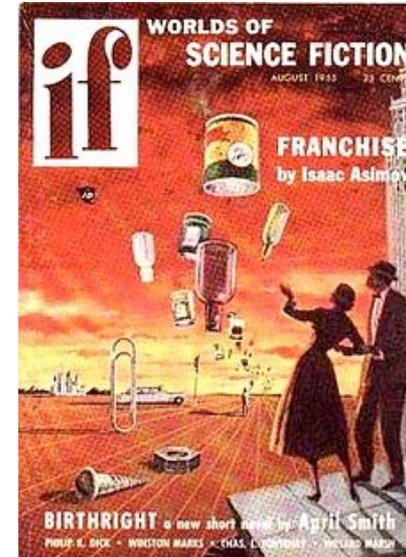
	Test statistic	Profiled?	Test statistic sampling
LEP	$q_\mu = -2 \ln \frac{\mathcal{L}(data \mu, \tilde{\theta})}{\mathcal{L}(data 0, \tilde{\theta})}$	no	Bayesian-frequentist hybrid
Tevatron	$q_\mu = -2 \ln \frac{\mathcal{L}(data \mu, \hat{\theta}_\mu)}{\mathcal{L}(data 0, \hat{\theta}_0)}$	yes	Bayesian-frequentist hybrid
LHC	$\tilde{q}_\mu = -2 \ln \frac{\mathcal{L}(data \mu, \hat{\theta}_\mu)}{\mathcal{L}(data \hat{\mu}, \hat{\theta})}$	yes $(0 \leq \hat{\mu} \leq \mu)$	frequentist

- Convenient to compute approximate values:
*“We define the **Asimov data set** such that when one uses it to evaluate the estimators for all parameters, one obtains the true parameter values”*
- Imagine that our only parameter is μ . We would like to have a dataset where the fit value is the true value μ' . This can be obtained using as number of counts the (non-integer) value $n = \mu's + b$
- In this case, we have:

$$-2 \ln \lambda(\mu) \cong \frac{(\mu - \mu')^2}{\sigma_{\hat{\mu}}^2}$$

- Therefore we can estimate the variance of $\hat{\mu}$ to be used in Wald's approximation of the test statistic:

$$\sigma_{\hat{\mu}}^2 \cong - \frac{(\mu - \mu')^2}{2 \ln \lambda(\mu)}$$



- In practice: all observables are replaced with their expected value
- Yields expected values are possibly non integer:

$$\lambda_A(\mu) = \frac{L_A(\mu, \hat{\theta})}{L_A(\hat{\mu}, \hat{\theta})} = \frac{L_A(\mu, \hat{\theta})}{L_A(\mu', \theta)}$$

- The variance of the test statistic, in Wald's approximation, is estimated as:

$$\sigma_{\hat{\mu}}^2 \cong \frac{(\mu - \mu')^2}{-2 \ln \lambda_A}$$

- Median significance for discovery or exclusion (and their $\pm 1\sigma$ bands) can be obtained using the Asimov dataset

$$\text{med}[Z_0|\mu'] = \sqrt{q_{0,A}} \quad \leftarrow \text{For discovery using } q_0$$

$$\text{med}[Z_\mu|0] = \sqrt{q_{\mu,A}} \quad \leftarrow \text{For upper limit using } q_\mu$$

$$\text{med}[Z_\mu|0] = \sqrt{\tilde{q}_{\mu,A}} \quad \leftarrow \text{Upper limits using } \tilde{q}_\mu$$

In practice: all the interesting formulae are implemented in RooStats package, released as optional library in ROOT

- Consider a search for a **signal peak** over a background distribution that is smoothly distributed over a wide range
- You could either:
 - Know which mass to look at, e.g.: search for a rare decay with a known particle, like $B_s \rightarrow \mu\mu$
 - Search for a peak at an **unknown mass value**, like for the Higgs boson
- In the former case it's easy to compute the peak significance:
 - Evaluate the test statistics for $\mu = 0$ (background only) at your observed data sample
 - Evaluate the **p -value** according to the expected distribution of your test statistic q **under the background-only hypothesis**, convert it to the equivalent area of a Gaussian tail to obtain the significance level:

$$p = \int_{q^{\text{obs}}}^{\infty} f(q|\mu = 0) dq \qquad Z = \Phi^{-1}(1 - p)$$

- In case you search for a peak at an unknown mass, the previous p -value has only a **local** meaning:
 - Probability to find a background fluctuation as large as your signal or more at a fixed mass value m :

$$p(m) = \int_{q^{\text{obs}}(m)}^{\infty} f(q|\mu = 0) dq$$
 - We need the probability to find a background fluctuation at least as large as your signal at **any** mass value (**global**)
 - local p -value would be an overestimate of the global p -value
- The chance that an over-fluctuation occurs on **at least one mass value** increases with the searched range
- Magnitude of the effect:**
 - Roughly proportional to the **ratio of resolution over the search range**, also depending on the significance of the peak
 - Better resolution = less chance to have more events compatible with the same mass value
- Possible approach: let also m fluctuate in the test statistics fit:

$$\hat{q}_0 = -2 \ln \frac{L(\mu = 0)}{L(\hat{\mu}; \hat{m})} \leftarrow \begin{array}{l} \text{Note: for } \mu=0 \\ L \text{ doesn't depend on } m \\ \text{Wilks' theorem doesn't apply} \end{array} \quad p^{\text{glob}} = \int_{\hat{q}_0^{\text{obs}}}^{\infty} f(\hat{q}_0|\mu = 0) d\hat{q}_0$$

- The effect can be evaluated with brute-force Toy Monte Carlo:

- Run N experiments with background-only
- Find the maximum \hat{q} of the test statistic q in the entire search range
- Determine its distribution, hence compute the observed global p -value
- Requires very large toy Monte Carlo samples (5σ : $p = 2.87 \times 10^{-7}$)

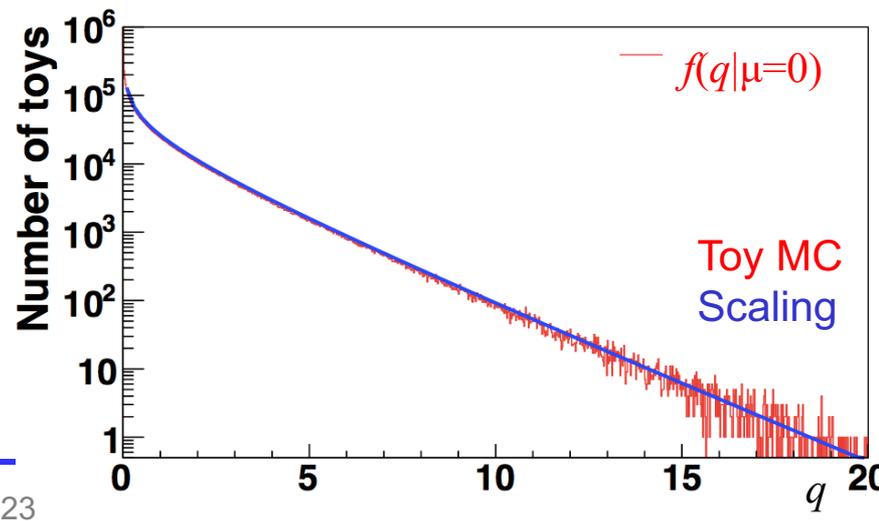
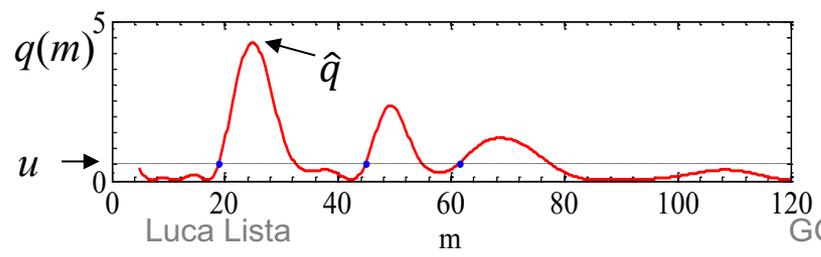
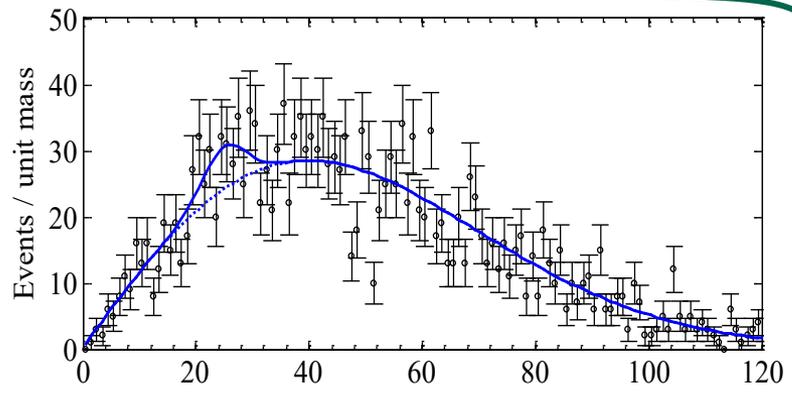
$$\hat{q} = \max_m q(m)$$

- Approximate evaluation based on local p -value, times correction factors (“trial factors”, Gross and Vitells, EPJC 70:525-530,2010)

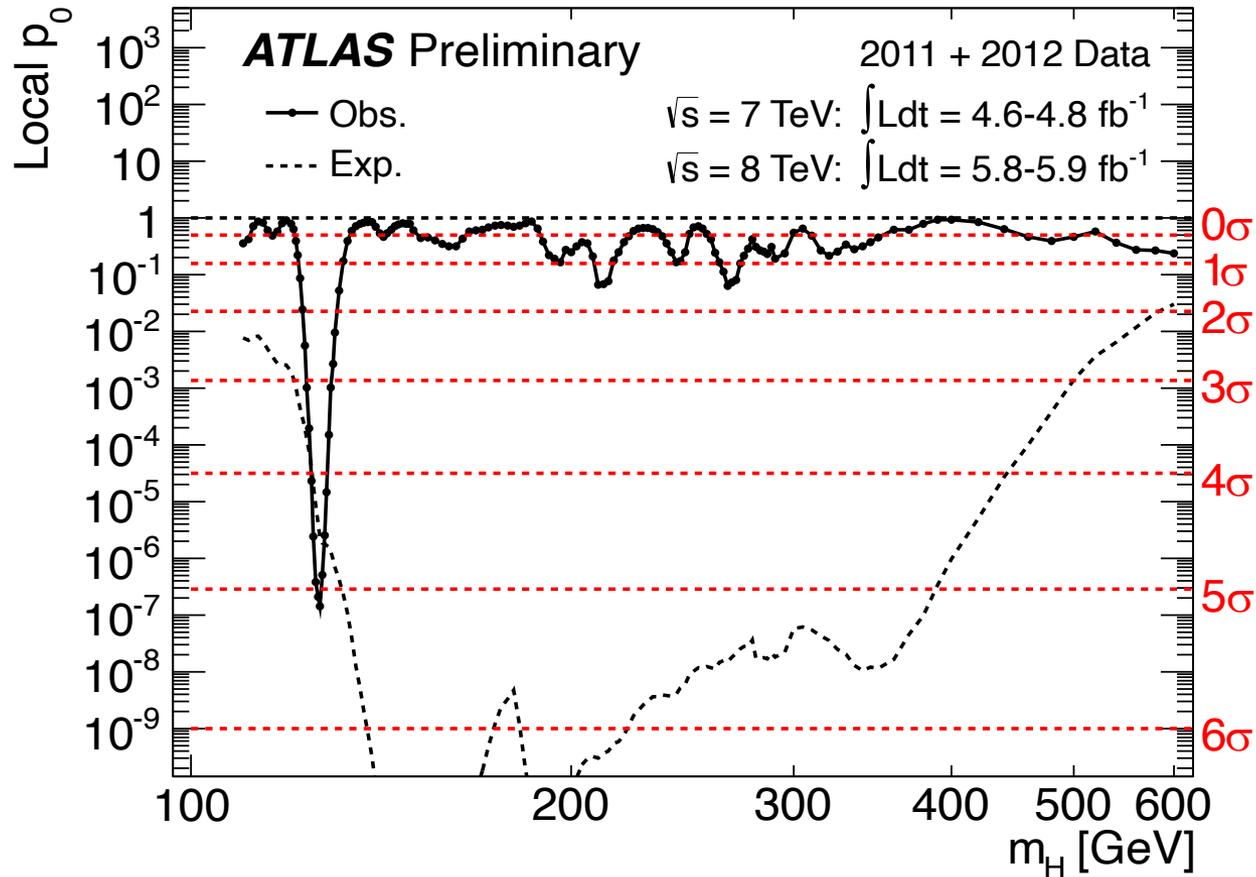
$$p^{\text{glob}} = P(\hat{q} > u) \leq \langle N_u \rangle + \frac{1}{2} P(\chi^2 > u)$$

$\langle N_u \rangle$ is the average number of up-crossings of the test statistic, can be evaluated at some lower reference level (toy MC) and scaled by:

$$\langle N_u \rangle = \langle N_{u_0} \rangle e^{-\frac{u-u_0}{2}}$$



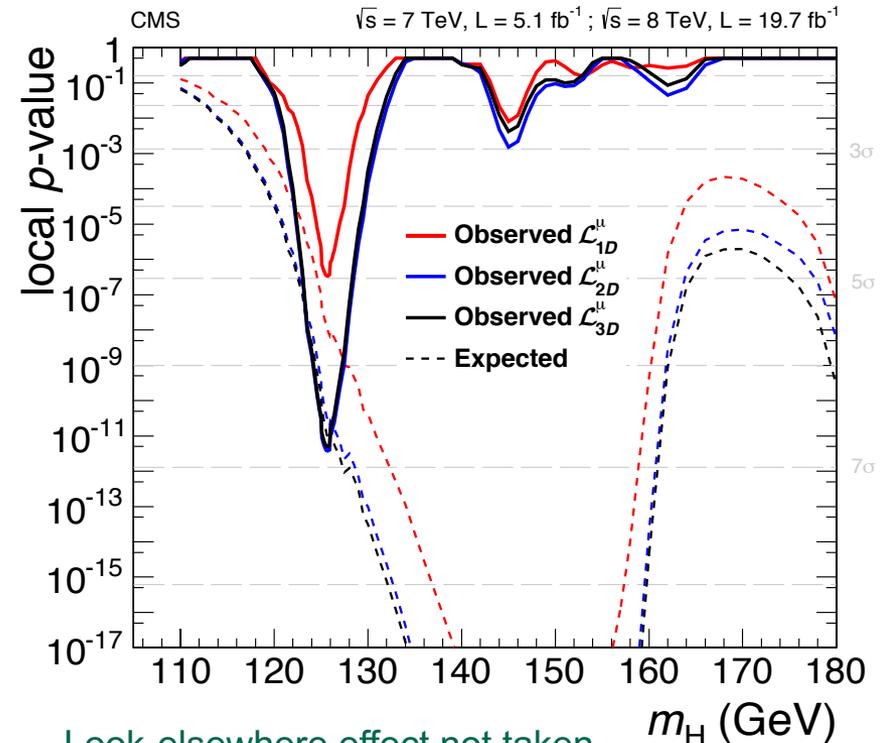
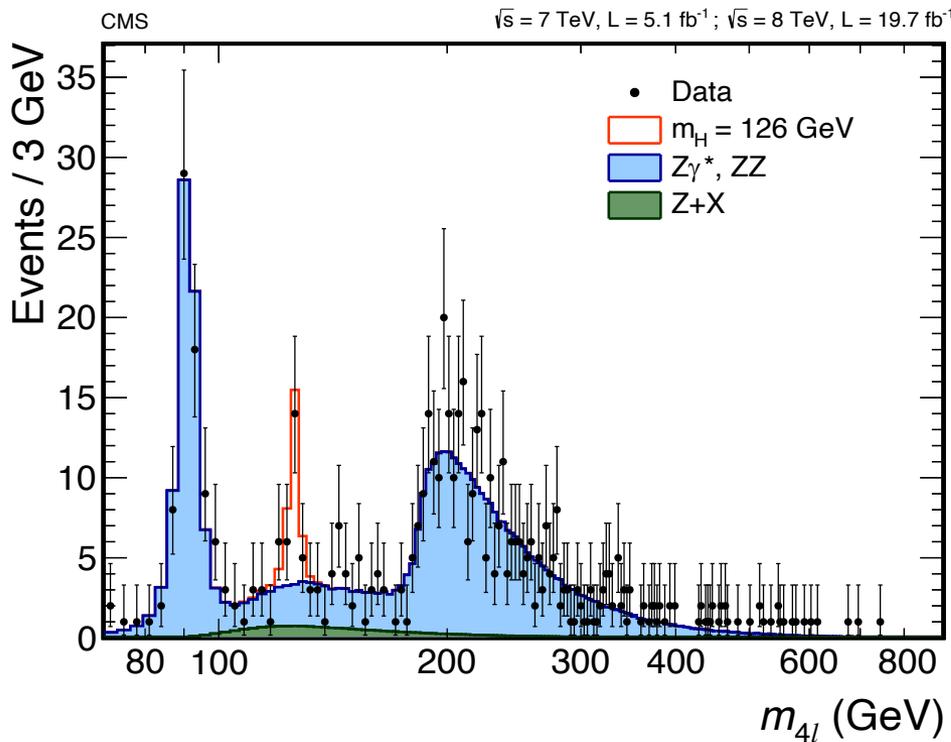
- Higgs search at ATLAS



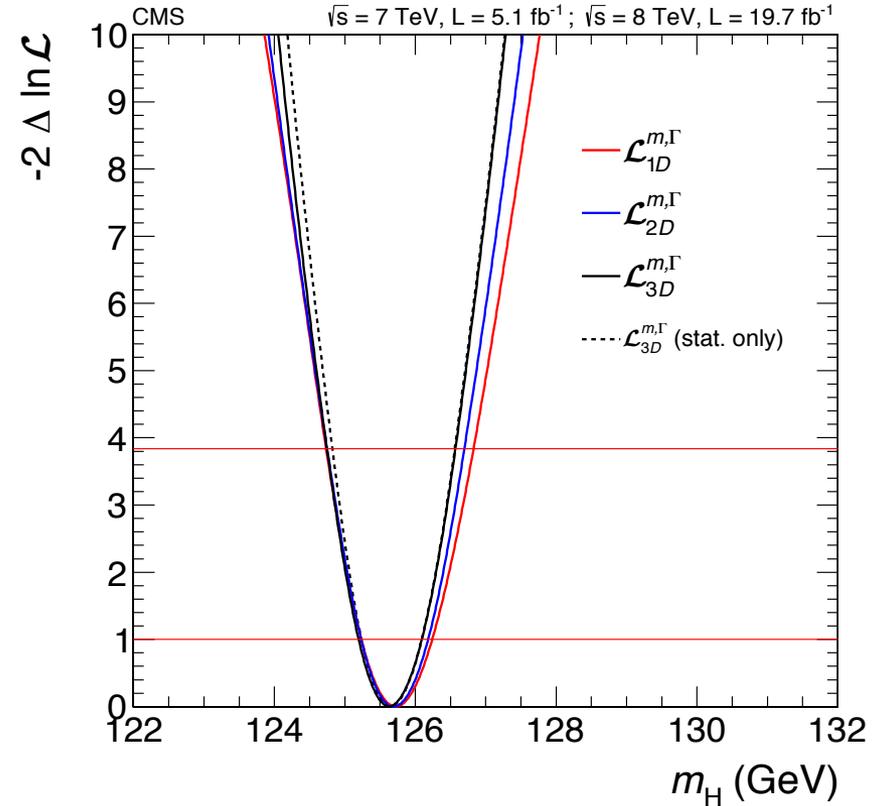
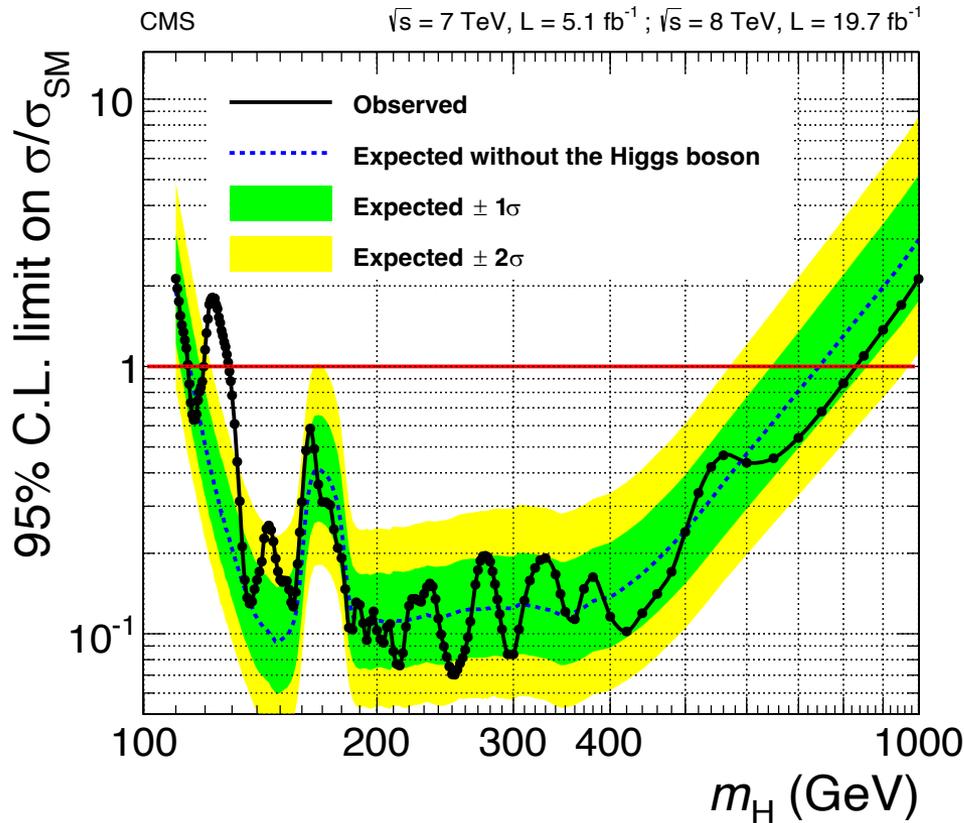
- Use the number of σ s Z as test statistic: $u = Z^2$ behaves as a chi-square
- Use the 0σ level ($p = 0.5$) as level u^0 , then extrapolate to the minimum p value, where $Z \cong 5$, i.e.: $u = Z^2 \cong 5^2 = 25$
- The number of upcrossing can be counted from the plot, and is equal to $N_0 = 9$, which allows us to estimate: $\langle N_0 \rangle = 9 \pm 3$
- Estimate the global p-value as:
 - $p^{glob} \cong \langle N_u \rangle + \frac{1}{2} P(\chi^2 > u) \cong \langle N_u \rangle + 3 \times 10^{-7}$
 - $\langle N_u \rangle \cong \langle N_0 \rangle e^{-(5^2 - 0^2)/2}$
 - $\langle N_u \rangle \cong (9 \pm 3) e^{-25/2} \cong (3 \pm 1) \times 10^{-5}$
 - $p^{glob} \cong 3 \times 10^{-5} + 3 \times 10^{-7} \cong 3 \times 10^{-5} \Rightarrow Z \cong 4\sigma$ instead of 5σ

Putting all together

- Search for Higgs boson in $H \rightarrow 4l$ at LHC
- 1D, 2D, 3D: different test statistics using $4l$ invariant mass plus other discriminating variables based on the event kinematics



Look-elsewhere effect not taken into account here



“The modified frequentist construction CLs is adopted as the primary method for reporting limits. As a complementary method to the frequentist construction, a Bayesian approach yields consistent results.”

Agreed statistical procedure described in:
 ATLAS and CMS Collaborations,
 LHC Higgs Combination Group
 ATL-PHYS-PUB 2011-11/CMS NOTE
 2011/005, 2011.