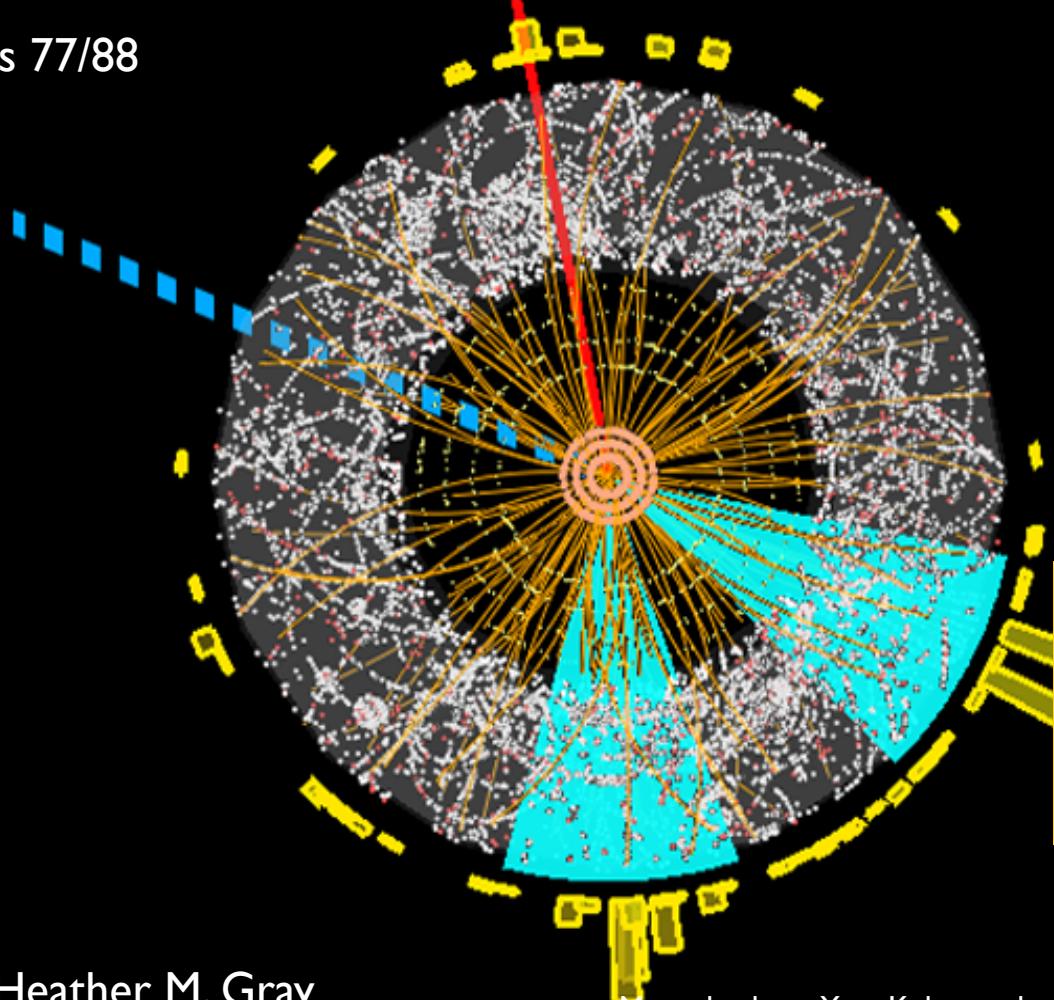


Introduction to Computational Techniques in Physics/Data Science Applications in Physics

Physics 77/88



**Statistics and
Probability,
Interpreting
Measurements**

There are three kinds of lies: lies, damned lies and statistics
- Mark Twain

The Statistics Boot Camp

- Introduction
- Definitions: results of experiments
 - Random variables, probability, PDFs
- Interpreting results
 - Point estimators
 - Max likelihood, least squares fits
- Hypothesis testing, confidence limits
- Simulation (Monte Carlo techniques) — in two weeks

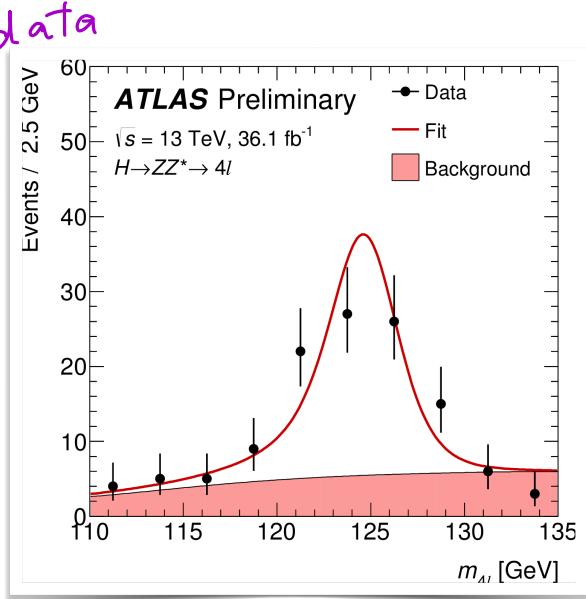
Statistics

- Statistics is a vital tool
 - Physics is an experimental science
 - Requires both qualitative and quantitative understanding
- From observation to a set of laws
 - Make a set of measurements
 - Summarise the result
 - Most conclusions are drawn with some degree of (un)certainty
- Example
 - Gravity exists
 - $F = G_N \frac{m_1 m_2}{r^2}$
 - $G_N = 6.67708 \pm 0.0031 \times 10^{-11} \text{ m}^3 \text{ kg}^{-1} \text{ s}^{-2}$
- Many measurements are a priori uncertain and have to be interpreted in probabilistic terms
- **Classical statistics:** estimate probabilities given a finite amount of data and test whether a given model is consistent with the data

This may seem rather dry, but if you turn out to be an experimentalist it's really important and useful!

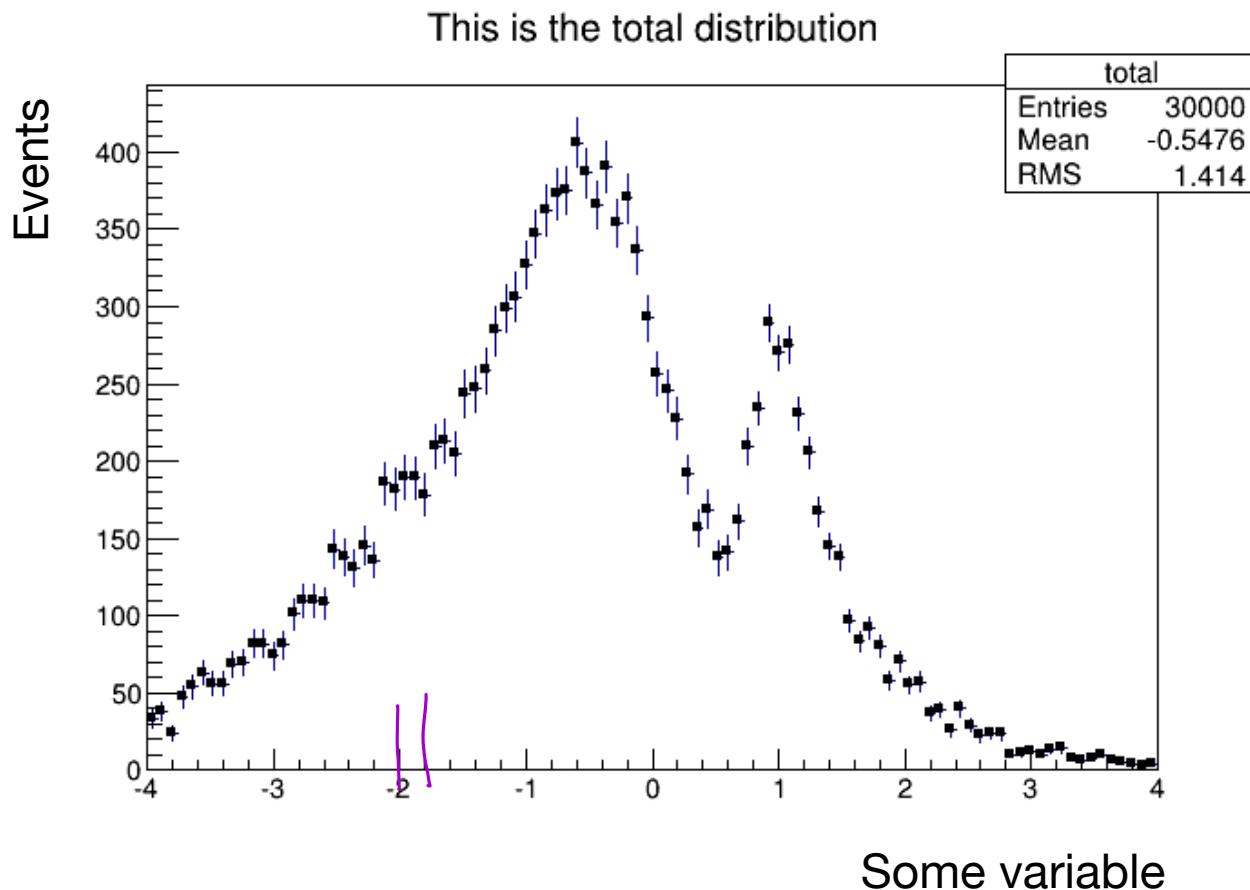
Describing the Data

- **Data:** result of a measurement
 - Physics: typically have quantitative data
 - number
 - Other fields deal with qualitative data
- **Numbers** are easier to handle mathematically; statistics will deal with quantitative measurements
 - discrete data, e.g. integers (counts)
 - continuous data, e.g. energies, momenta
 - Measure with some precision, set by the measuring apparatus or other external conditions

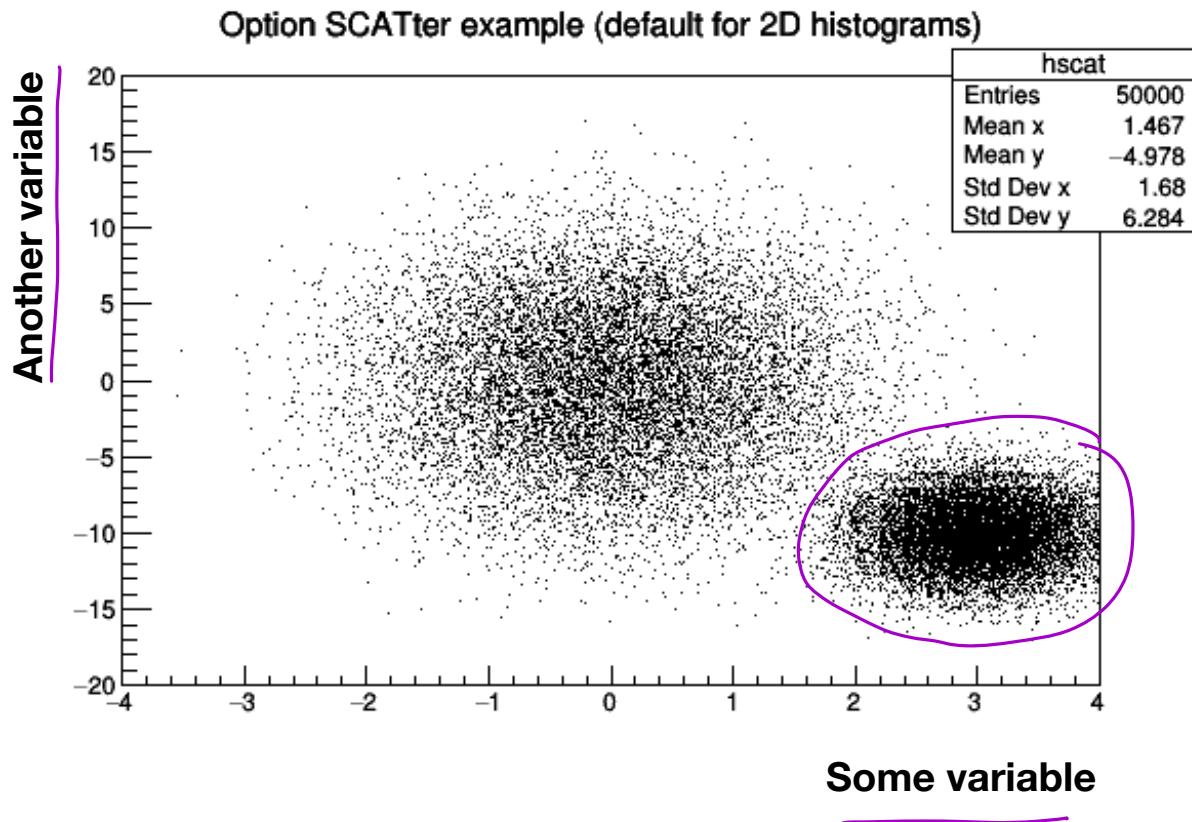


Why are there error bars on the data?

Histogram

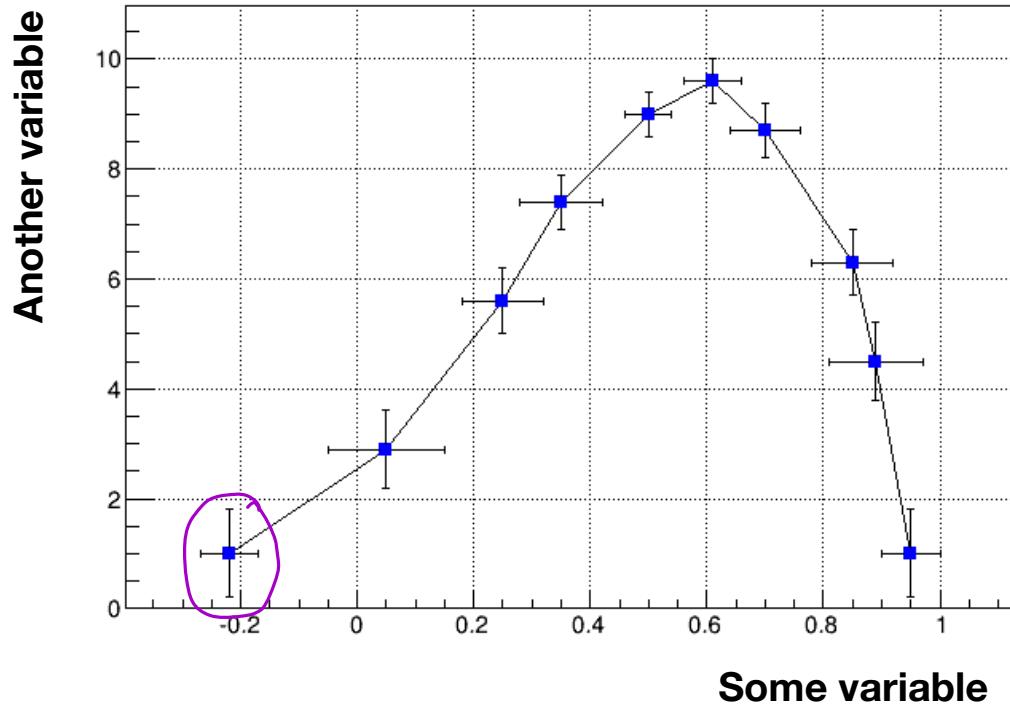


Scatterplots



Graph

TGraphErrors Example



Uncertainty and Error

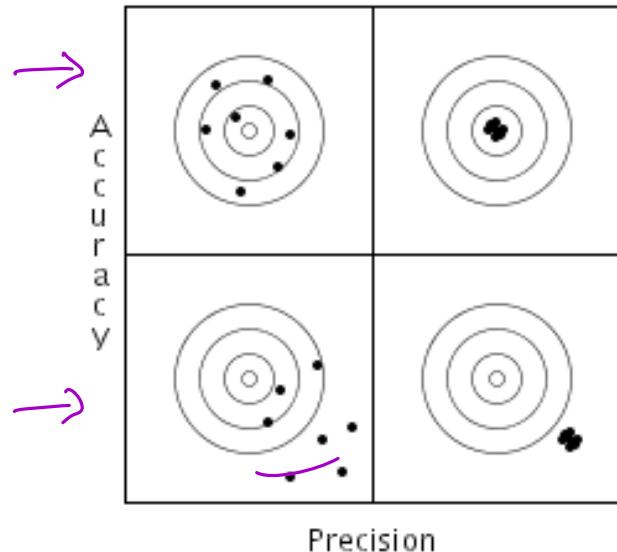
- In physics, the words **uncertainty** and **error** are used interchangeably to describe how far a **particular measurement** is expected to deviate from the **true value** — typically
 - Use symbol σ for the error
 - Formal definition is probabilistic: **68% chance** to find the **experimental result** within $\pm 1\sigma$ of the **true value**
(frequentist interpretation)
 - Often interpreted as a **range of possible true values**
(Bayesian interpretation)
 - We'll come back to the difference between **Bayesian** and **frequentist**

Uncertainty and Error

- How do we define what is typical ?
 - Underlying assumption: our experiment is one sample of a population of similar measurements
 - Derive the value of σ from the properties of the population
 - Implicit assumption: our experiment is mistake free
 - i.e. all similar experiments would return similar results

Precision vs Accuracy

- Precision: *spread* of the data around the *average* value
 - Typically associated with *statistical uncertainty*
- Accuracy: *deviation* of the *average* value from *true* value
 - *Bias*
 - Typically associated with *systematic uncertainty*
- Bad data: *outliers*
 - Data *inconsistent* with distribution (e.g. *noise, mistakes*)



http://anomaly.org/wade/blog/2006/01/accuracy_and_precision.html

Golden Rules

- When reporting results of a measurement, always report its uncertainty
- Round off values to 1-2 digits of uncertainty
 - Rule of thumb: 1 digit if the last digit is > 4, 2 digits otherwise
 - $x = \underline{3.142} \pm \underline{0.024}$
 - $y = \underline{3.1} \pm \underline{0.6}$
- Uncertainty can come from the spread in the data and/or precision of the instrument
 - Rule of thumb: "half of the last digit"
 - Statistically correct: $\sigma_{\text{inst}} = \frac{\text{last digit}}{\sqrt{12}}$

Probability: Definitions

- For numerical data, probabilistic definition is often most convenient (and quantitative)
 - Let's define probability now
 - Formally, it is a quantity that is defined by Kolmogorov axioms
-
- 1. For every subset A in S $P(A) \geq 0$
 - 2. For disjoint subsets $(A \cap B = \emptyset)$
$$P(A \cup B) = P(A) + P(B)$$
 - 3. $P(S) = 1$

Two Interpretations

- “Frequentist” interpretation:
 - Probability is a limiting frequency of a given outcome when experiments are repeated an infinite number of times
 - Measurable parameters are represented by estimator with assigned confidence levels (CL)
 - CL measures a probability an estimator would fall in a certain range, given a true value of a parameter
 - No probability is assigned to constants of nature
- “Bayesian” interpretation:
 - More general: define probability as a degree of belief that a given statement is true
 - E.g. that the true value of parameter α is in interval $[a, b]$
 - This is somewhat subjective, but follows how most people think

Frequentist Probability

- Definitions
 - Let S be set of all possible outcomes of a measurement
 - Any subset A with only one element (single outcome) is the elementary outcome
 - Define $P(A) = \lim_{N \rightarrow \infty} (\# \text{occ of } A \text{ in } N \text{ trials})$
- Assume outcomes are (in principle) repeatable
- Confidence in a measurement grows with N
- Frequentist statistics is appropriate (and often argued for) in situations where measurements can be reproducibly repeated so that validity of approach can be tested (e.g. particle physics)



John von Neumann



Jerzy Neyman

Bayes Theorem

- Conditional probability of A given B
- $P(A|B) = \frac{P(A \cap B)}{P(B)}$
- Interpreted within Bayesian statistics as
 - $P(\text{theory}|\text{data}) \propto P(\text{data}|\text{theory}) \underbrace{\frac{P(\text{theory})}{\text{prior probability}}}$
 - Posterior probability likelihood/result
- Allows one to interpret a single experiment as a measure of (subjective) probability that a given hypothesis is correct (e.g. that some fundamental const. is in some range)
- Requires assigning some probability to prior knowledge
 - That's where subjectivity comes in



Thomas Bayes

Probability: Random Variables and PDFs

- For a continuous variable, x , we define the probability density function (pdf)

- $f(x, \theta)$ = probability that x lies between x and $x + dx$
- θ = parameters of the function

- Integrate to obtain the probability

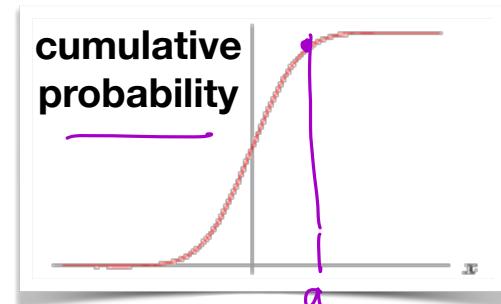
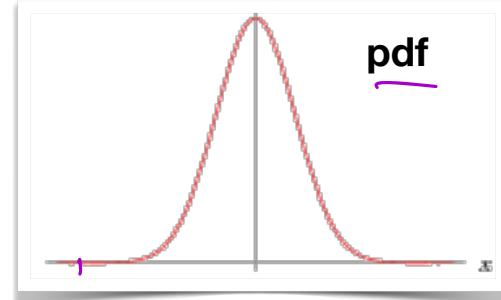
$$\underline{F(a)} = \int_{-\infty}^a f(x) d(x)$$

- Probability that $x < a$

- Discrete variables: integral \rightarrow sum $\sum_{x_i}^a f(x_i)$

- Expectation value

$$\underline{E[u(x)]} = \langle u(x) \rangle = \int_{-\infty}^{\infty} u(x) f(x) d(x)$$



Expectation Values

- Expectation value of a function

$$\bullet E[u(x)] = \int_{-\infty}^{\infty} u(x) f(x) dx$$

- Moments of a function

$$\bullet \alpha_n \equiv E(x^n) = \int_{-\infty}^{\infty} x^n f(x) dx \quad n\text{th moment}$$

$$\bullet m_n \equiv E((x-\alpha)^n) = \int_{-\infty}^{\infty} (x-\alpha)^n f(x) dx$$

n-th central
moment

Mean and Variance

- Mean

$$\bullet \mu = \int_{-\infty}^{\infty} x f(x) dx$$

- Variance

$$\bullet \sigma^2 = \text{Var}(x) = \int_{-\infty}^{\infty} x^2 f(x) dx - \mu^2$$

- σ is the standard deviation

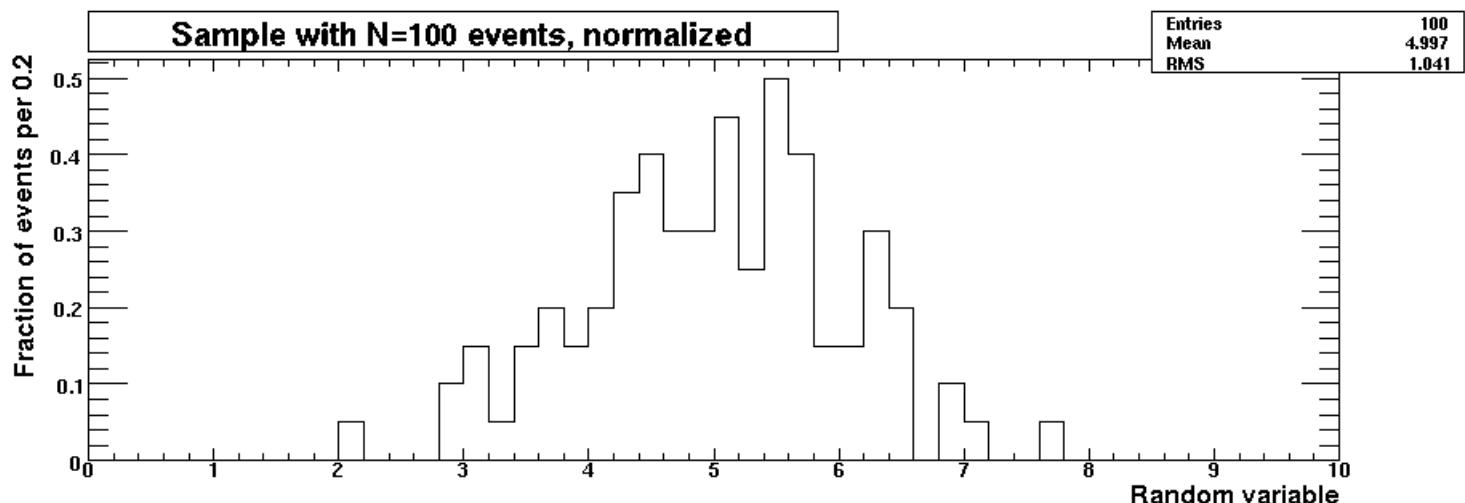
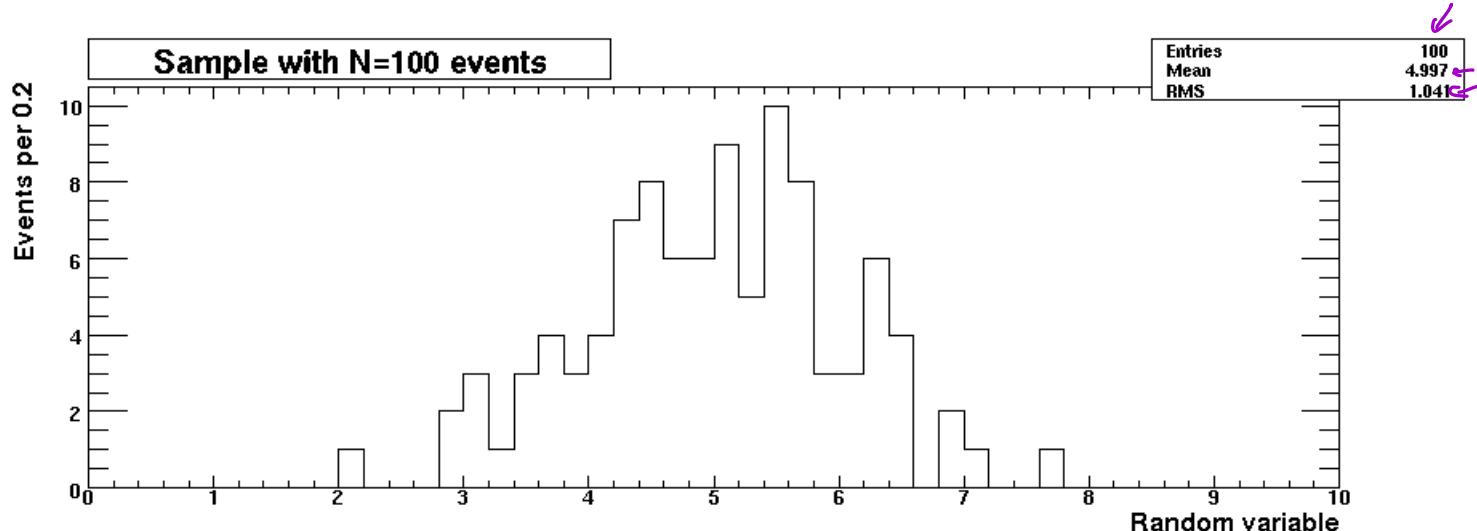
✓

If you only remember one thing today, remember this

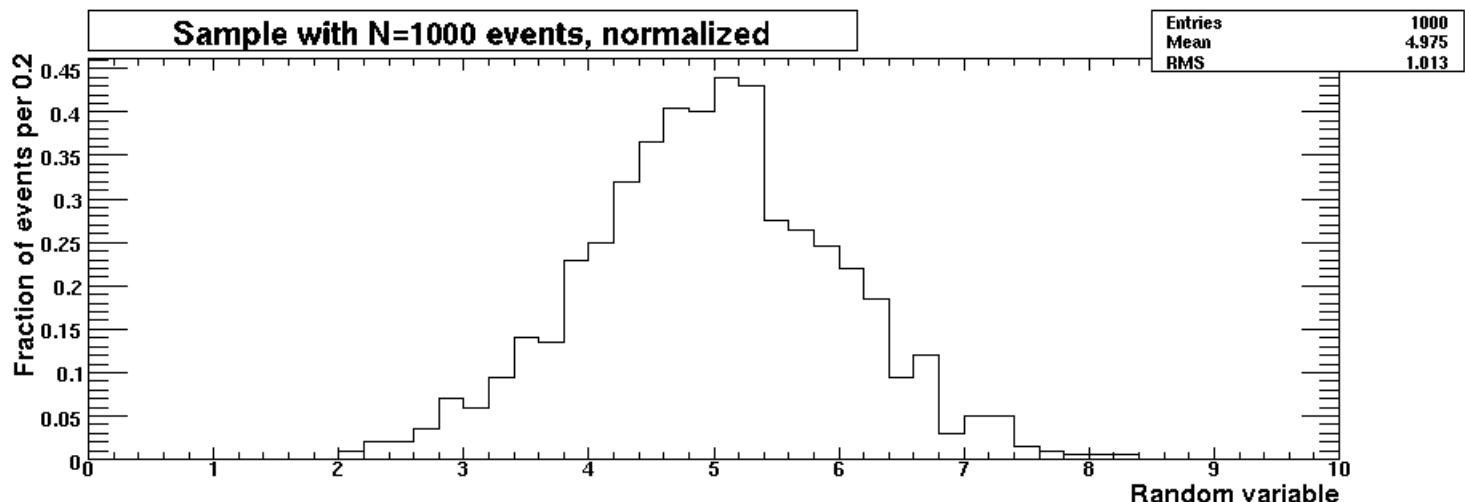
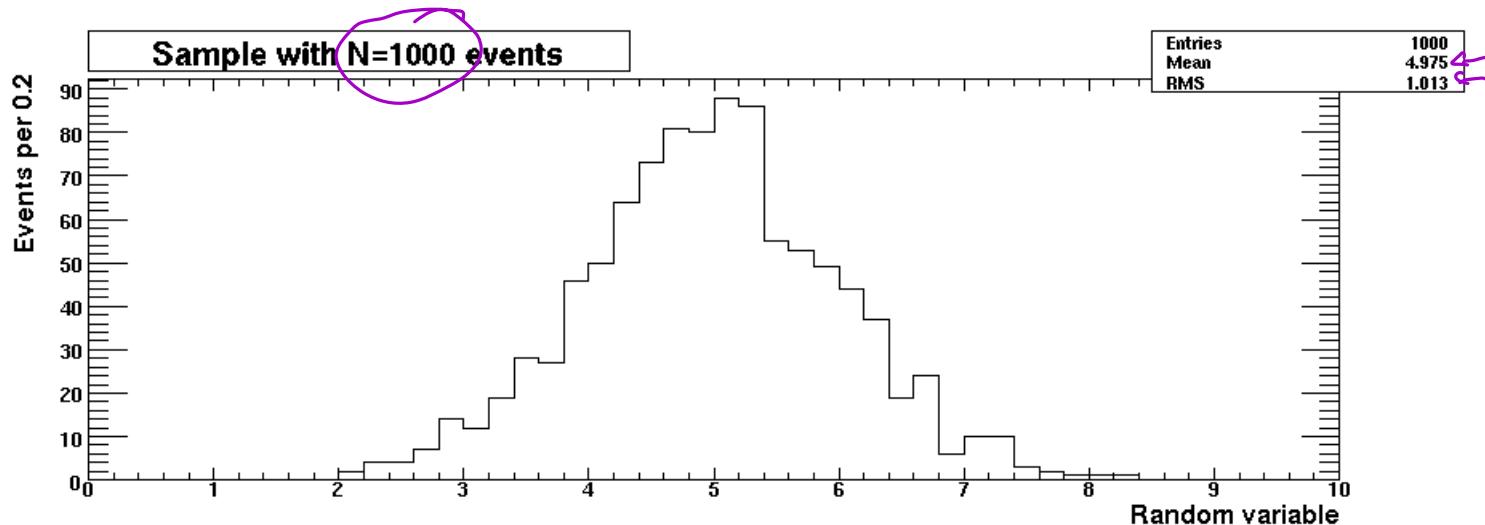
If we know the PDF, we know how to determine μ and σ

Example in jupyter notebook

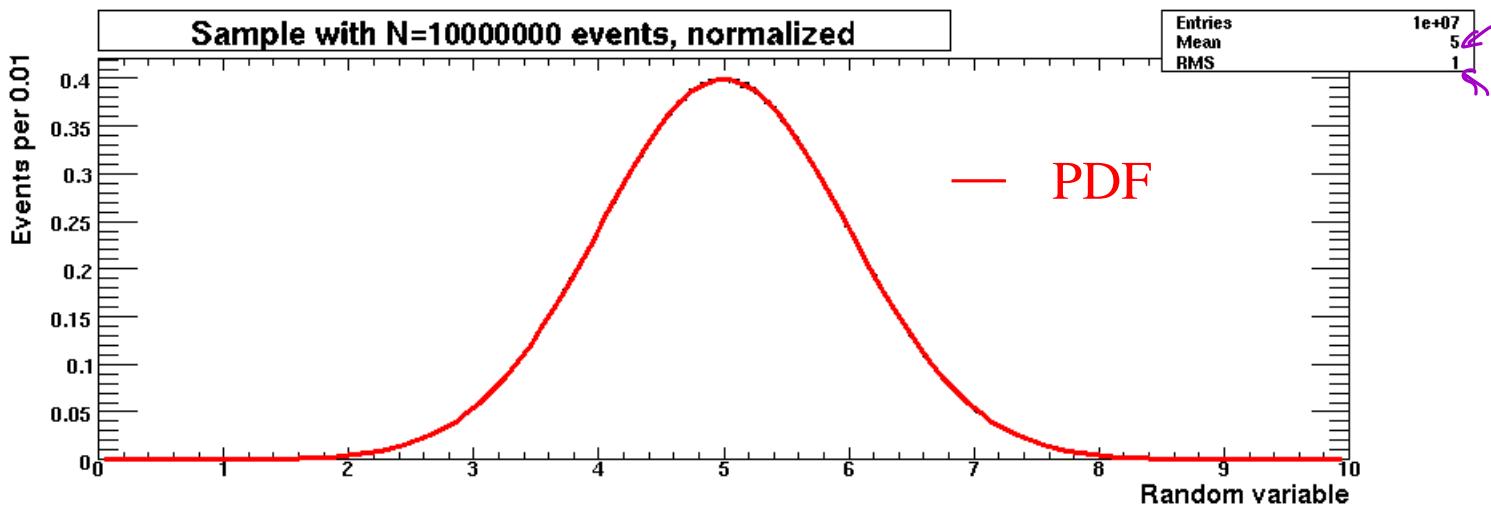
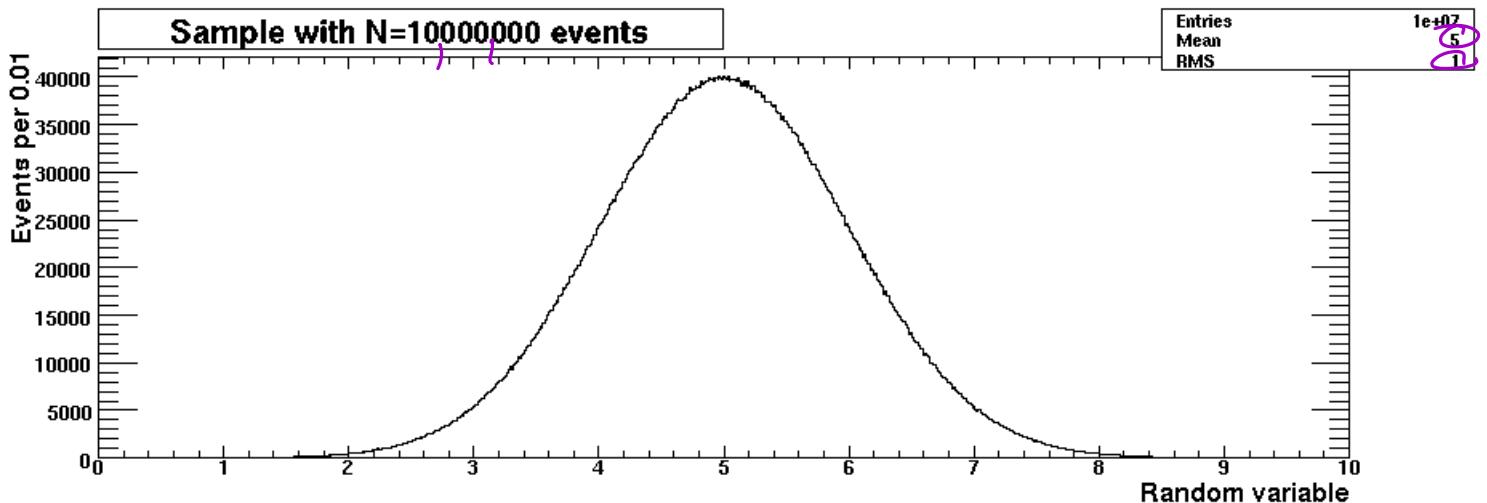
Sample from a Continuous PDF



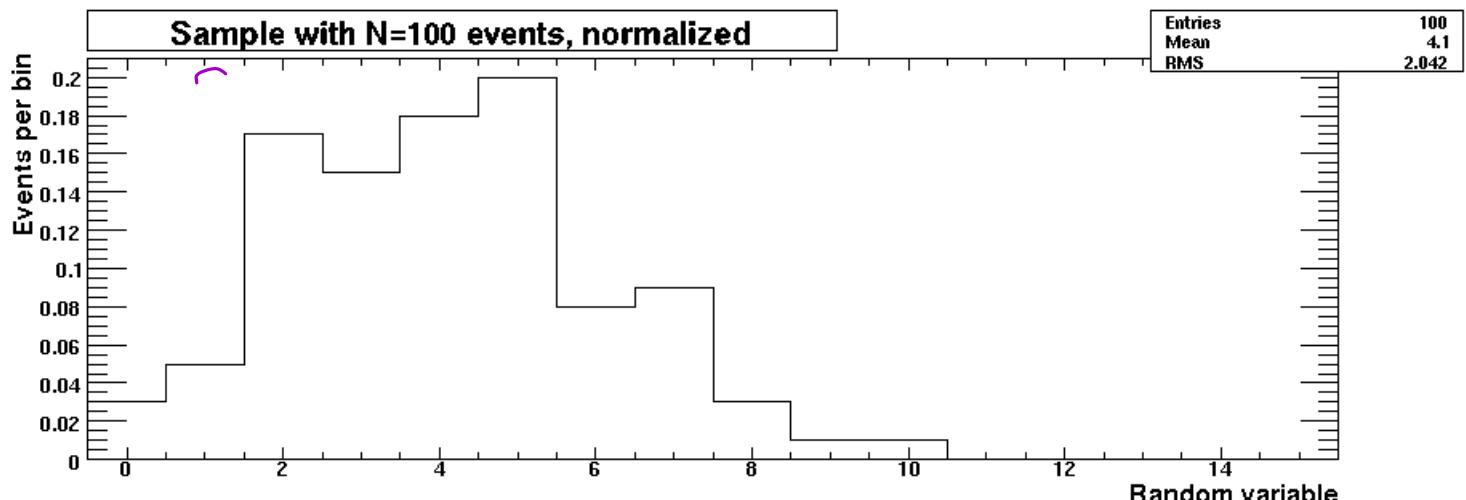
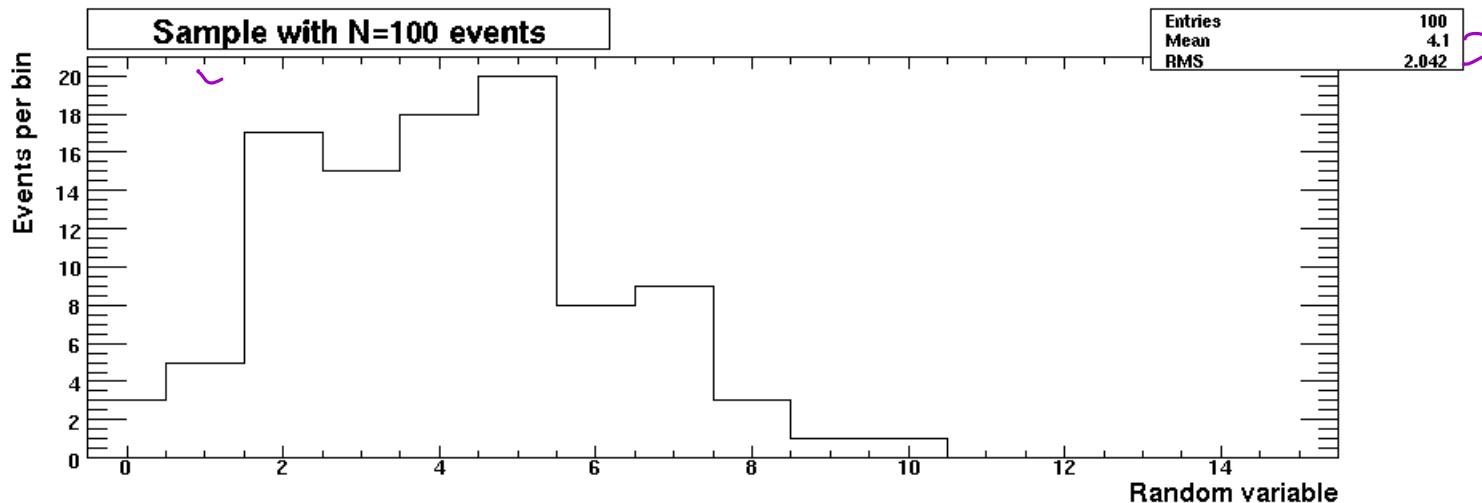
Sample From a Continuous PDF



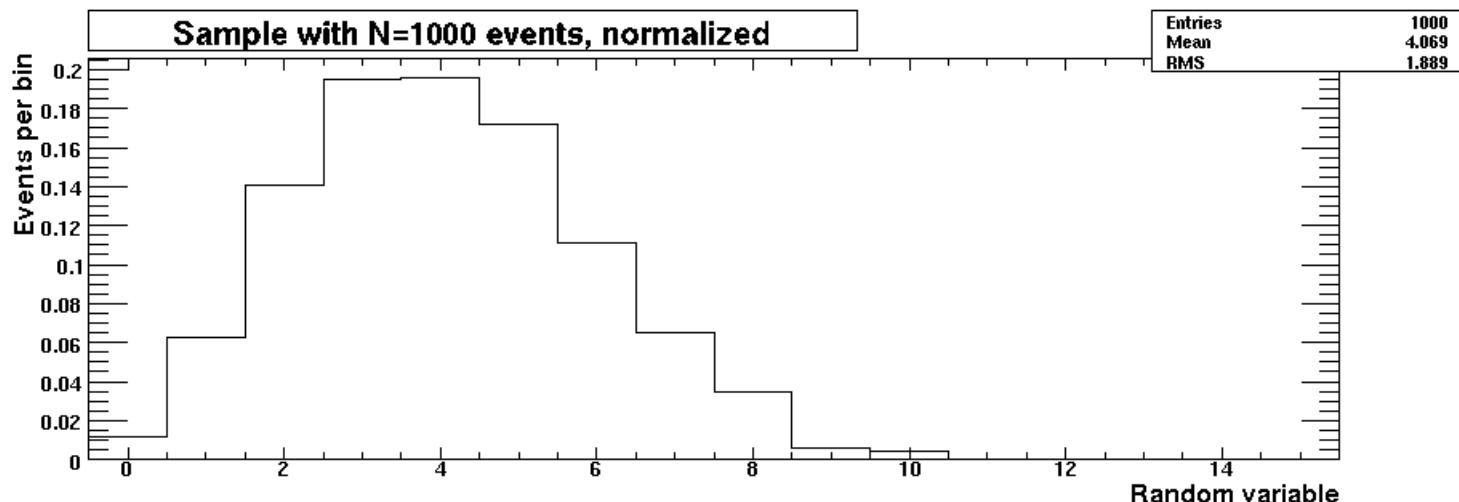
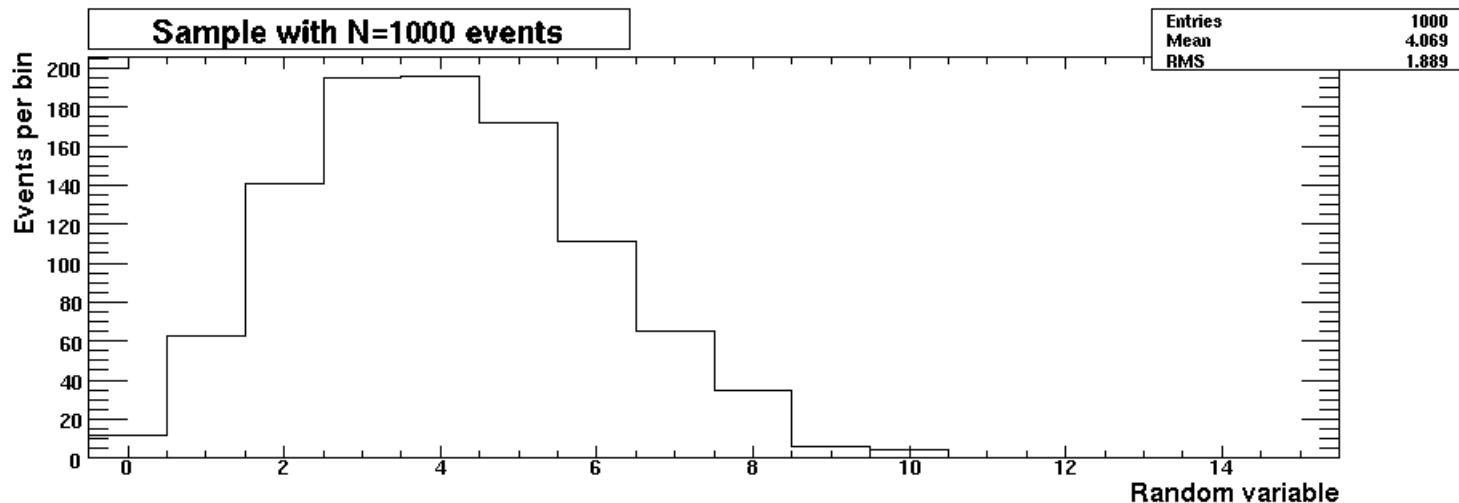
Sample From a Continuous PDF



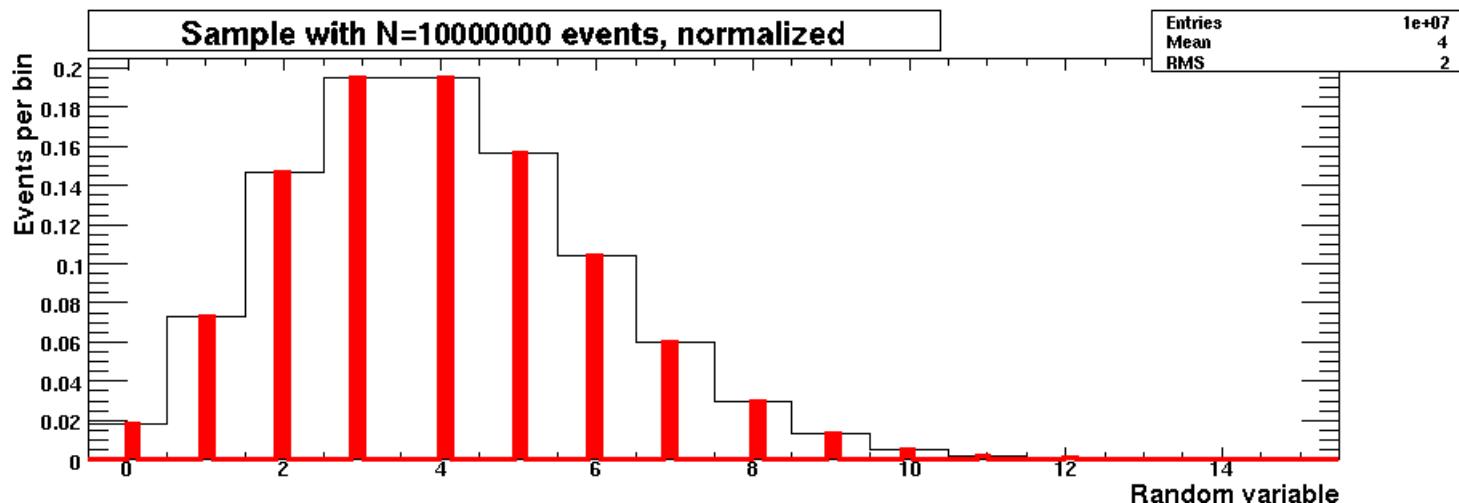
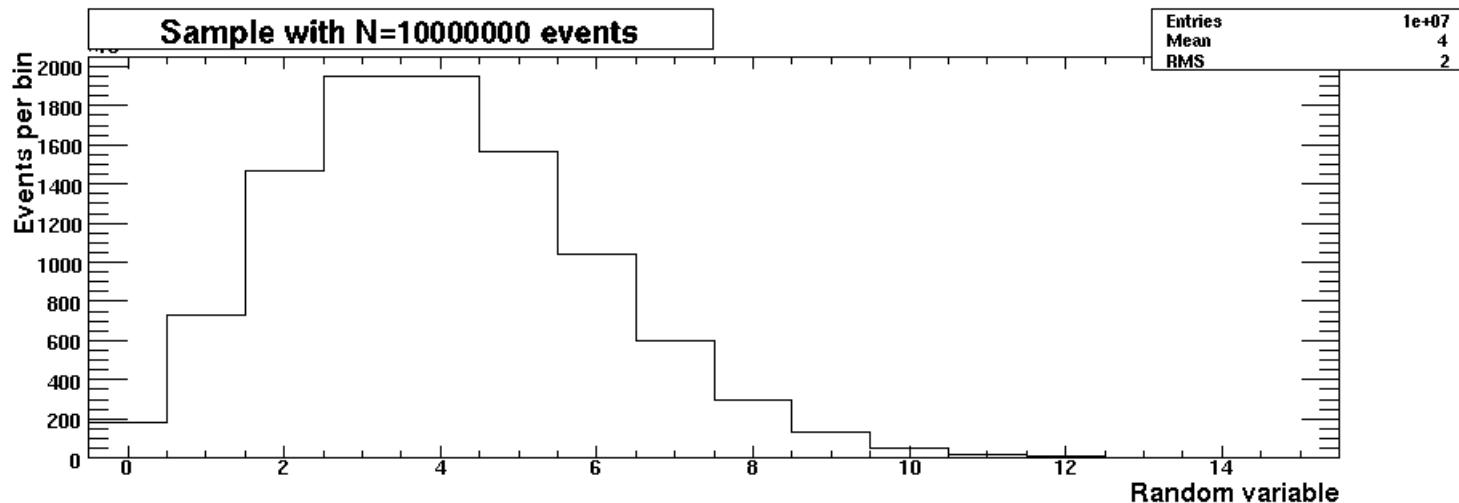
Sample From a Discrete PDF



Sample From a Discrete PDF

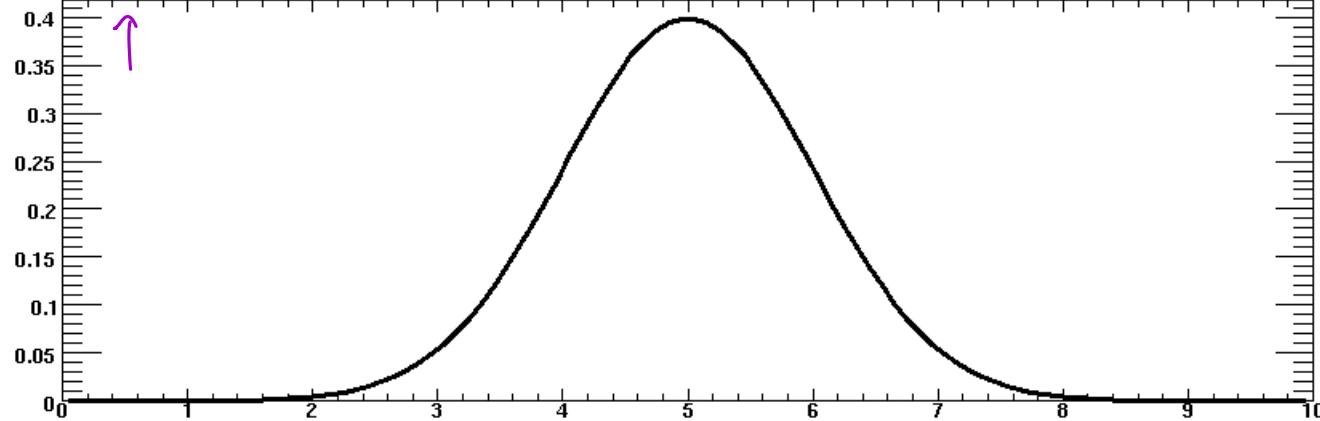


Sample From a Discrete PDF

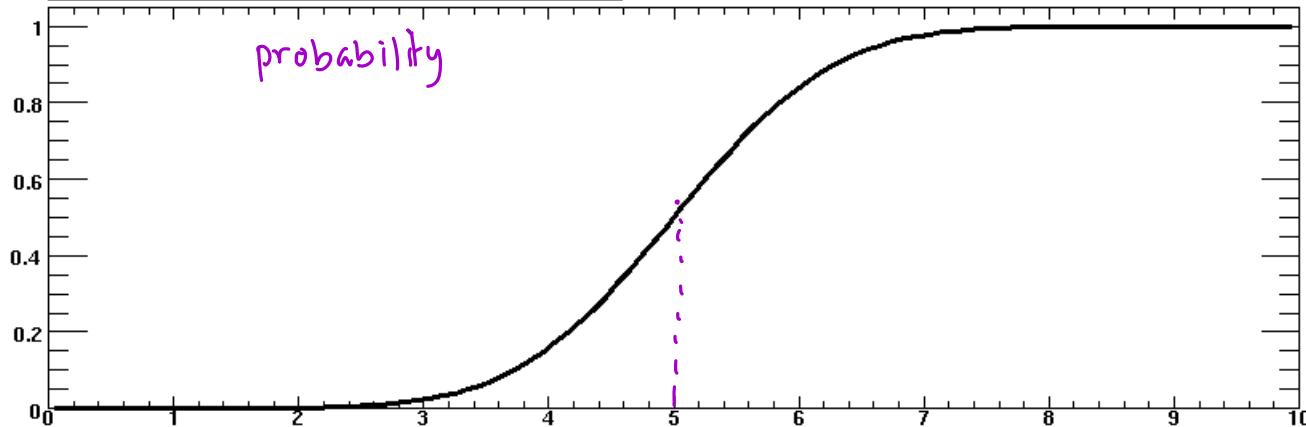


Cumulative Distribution

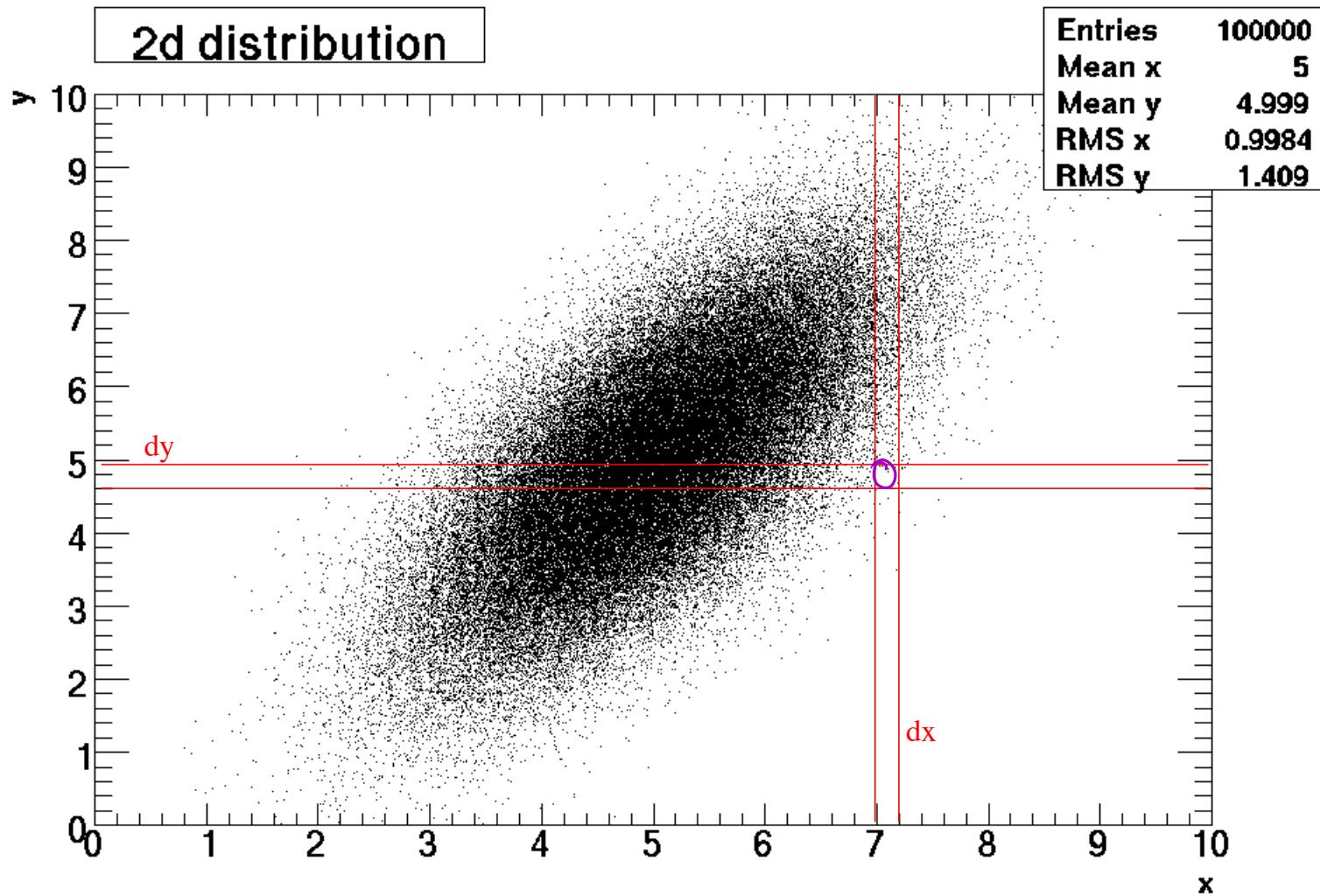
PDF



Cumulative Distribution



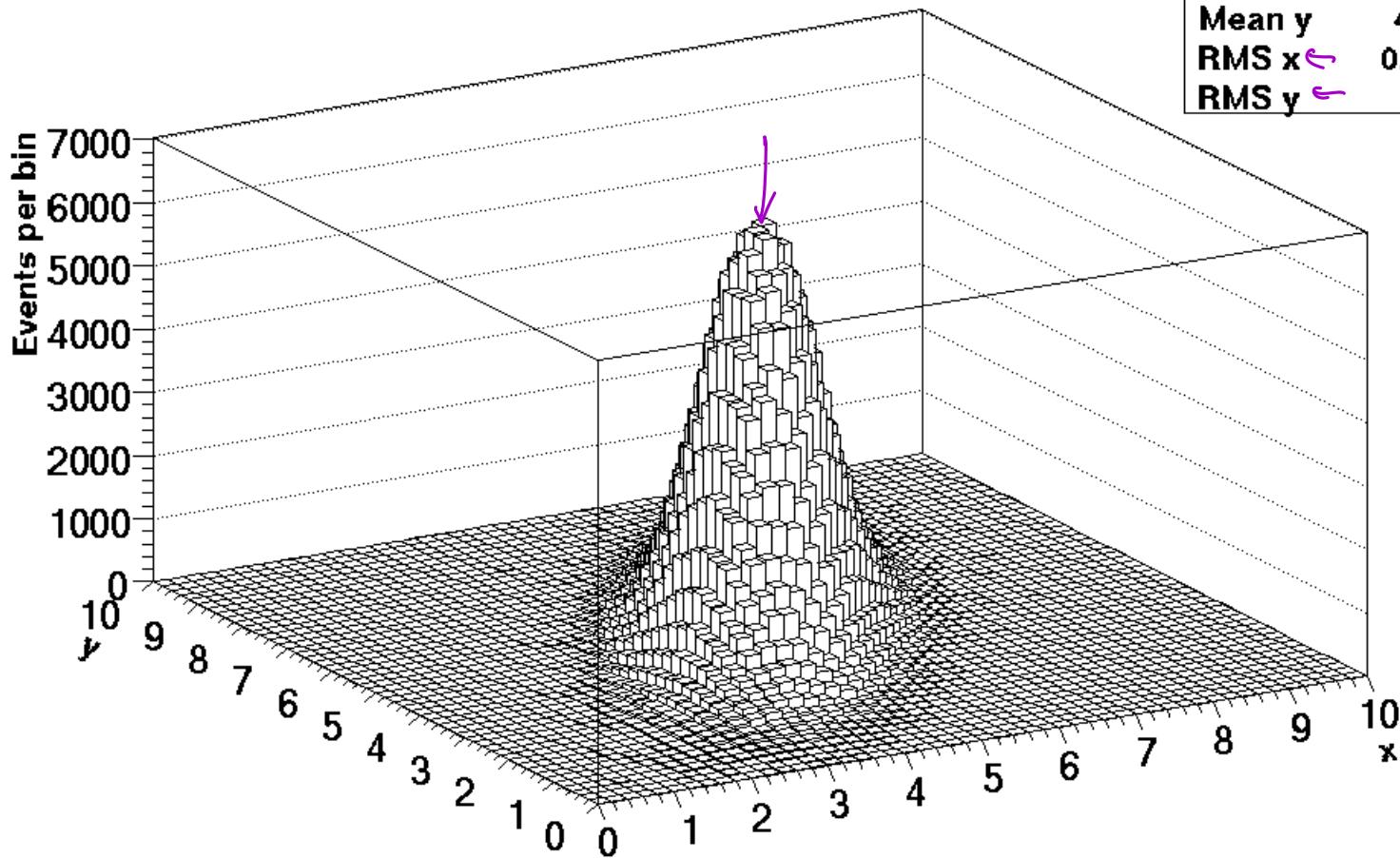
2d Distribution



2d Distribution

2d distribution

Entries	1000000
Mean x	4.999
Mean y	4.999
RMS x	0.9974
RMS y	1.41



Most important distributions

- Gaussian or normal distribution
 - central limit theorem → use all the time
- Poisson distribution
 - Number of events in data, e.g. rare event
 - Geiger counter click
- Binomial distribution
 - Efficiencies, # hits in detectors
- There are many others
 - uniform
 - exponential
 - Gamma
 - Chi² or χ^2

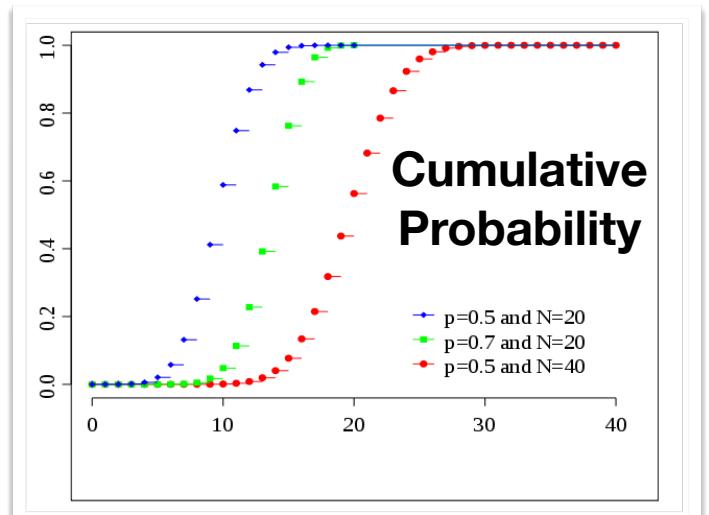
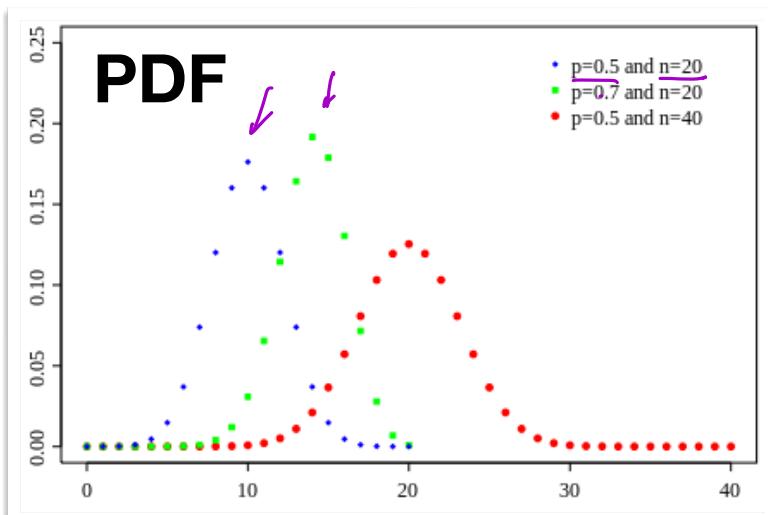
Binomial Distribution

- Random process with two possible outcomes
- p = Prob of outcome 1, q = Prob of outcome 2, $q=1-p$
- After n trials, prob of getting outcome #1 exactly k times is

$$\underline{f(k,p) = \binom{n}{k} p^k q^{n-k}}$$

where $\binom{n}{k} = \frac{n!}{k!(n-k)!}$

$$\mu = np \quad \sigma = np(1-p)$$



Poisson Distribution

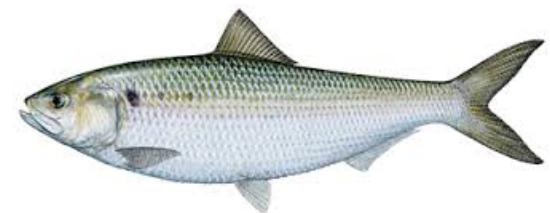
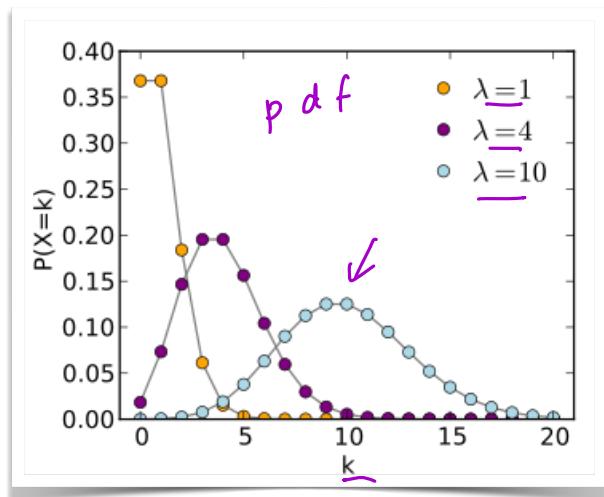
- Prob of finding exactly k events in the interval $[x, x+dx]$ if the

events occur with an rate of λ

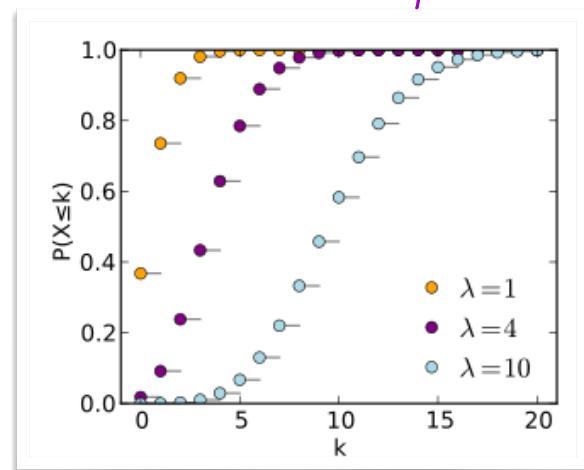
$$\frac{x^k e^{-\lambda}}{k!}$$

- For large λ approaches a Gaussian

$$\mu = \lambda \quad \sigma = \sqrt{\lambda}$$



cumulative probability



Example: Measure Efficiency

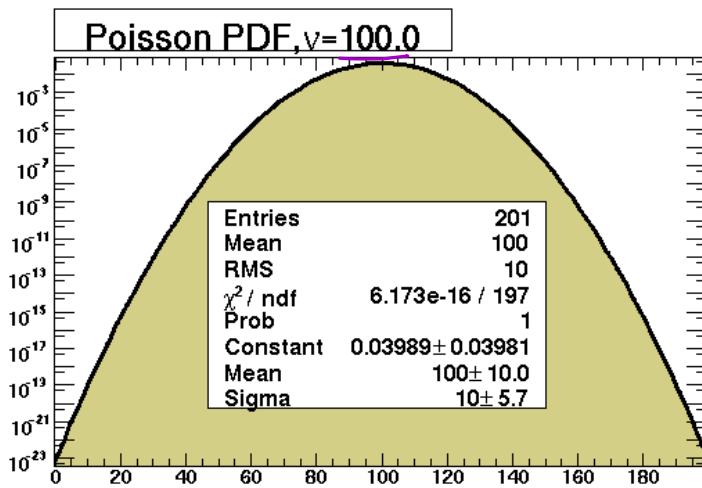
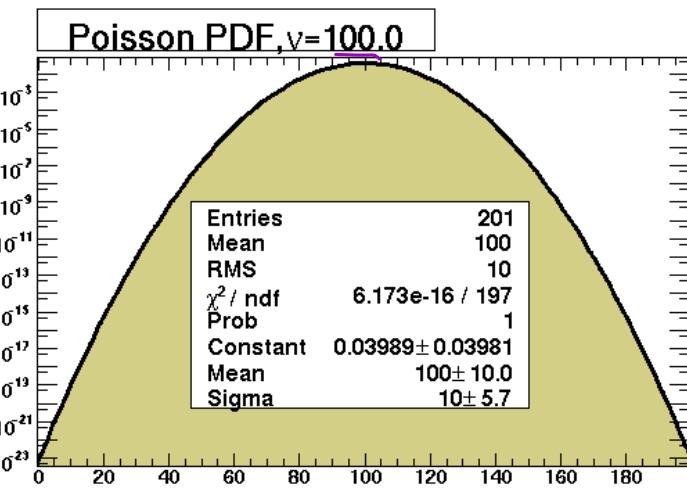
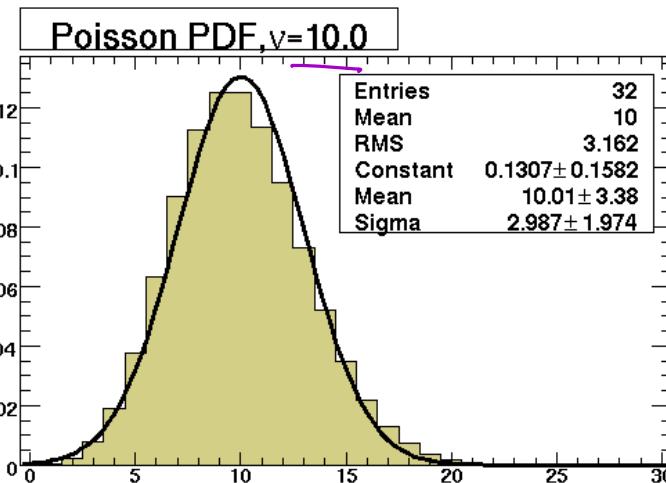
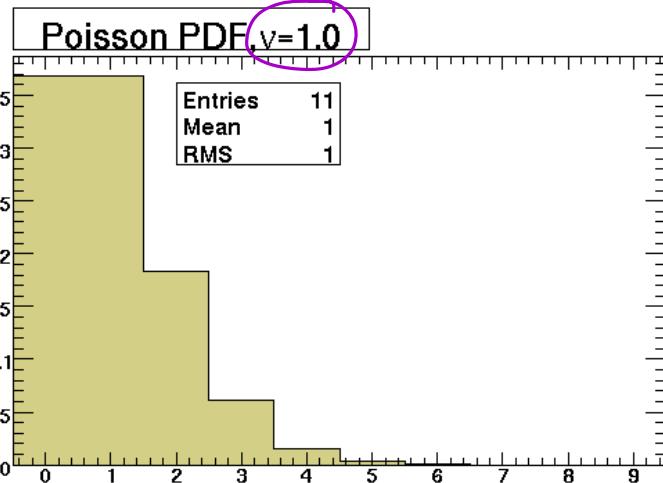
- Generate a sample of N events
- Apply selection; suppose n_{pass} events passed
- Estimate

$$\hat{\epsilon} = \frac{n_{\text{pass}}}{N}$$

$$\sigma(\hat{\epsilon}) = \sqrt{V(\hat{\epsilon})} = \sqrt{\frac{V(n_{\text{pass}})}{N^2}} = \sqrt{\frac{\hat{\epsilon}(1-\hat{\epsilon})}{N}}$$

$$\sigma(\hat{\epsilon}) \neq \frac{\sqrt{n_{\text{pass}}}}{N} = \sqrt{\frac{\hat{\epsilon}}{N}}$$

Poisson Distribution



Gaussian Distribution

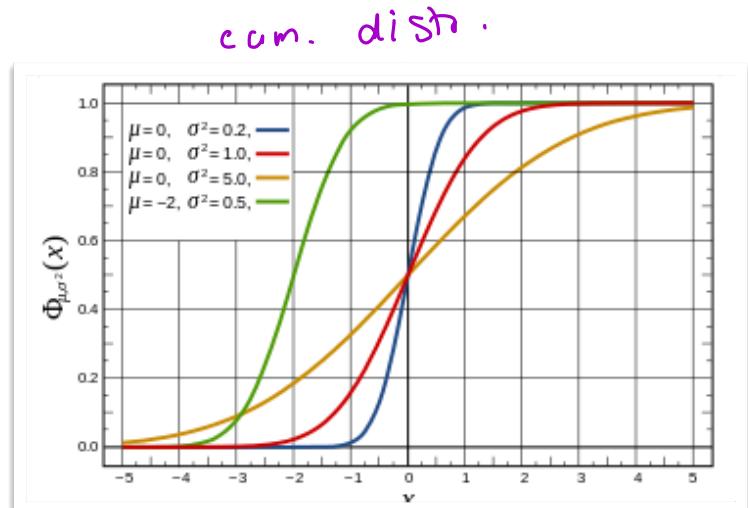
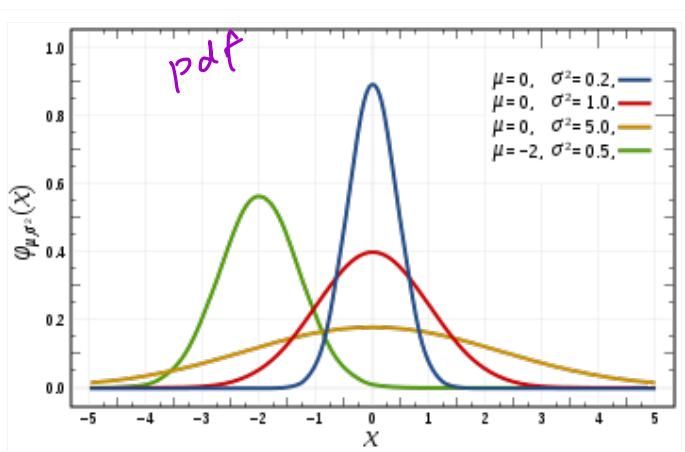
- **Central Limit Theorem**

- Given a random sample (x_1, x_2, \dots, x_n) drawn from a pdf with mean μ and variance σ^2 , if the mean is

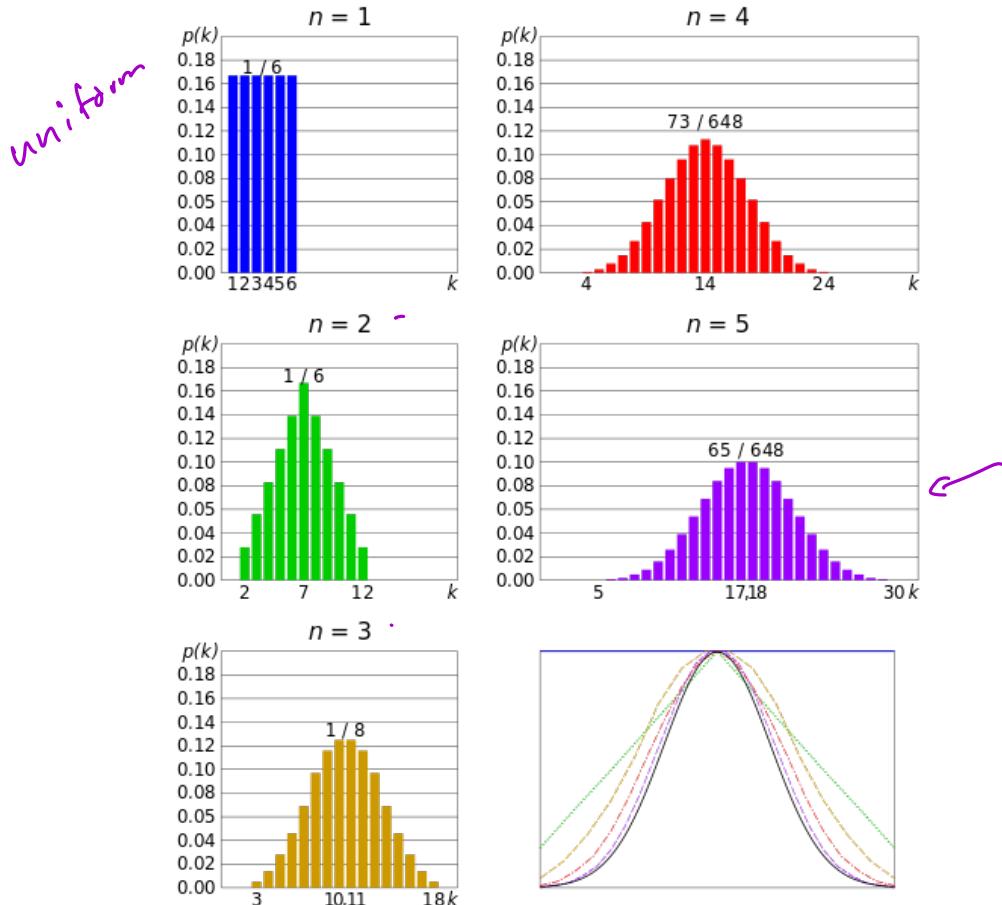
$$\cdot S/n = \frac{1}{n} \sum_{i=1}^n x_i$$

- the distribution of S/n approaches the normal distribution as $n \rightarrow \infty$

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



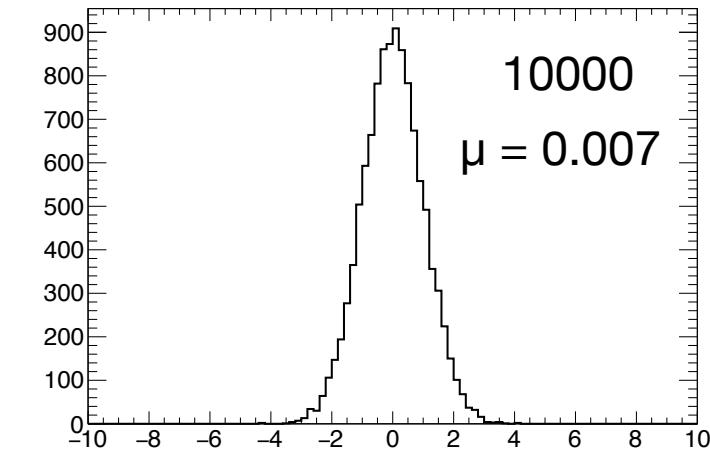
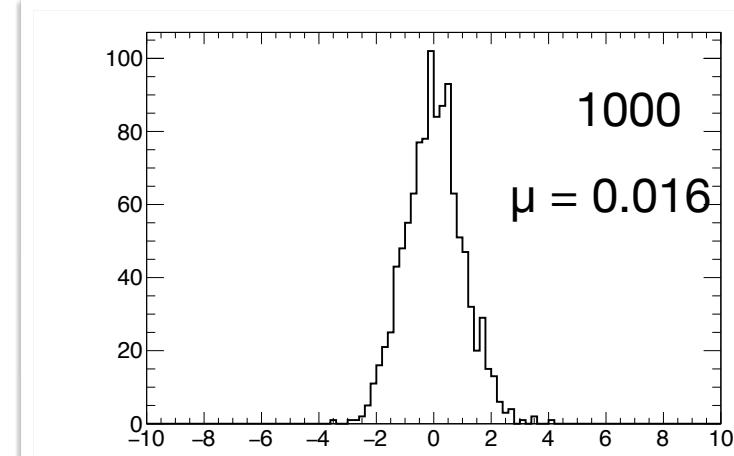
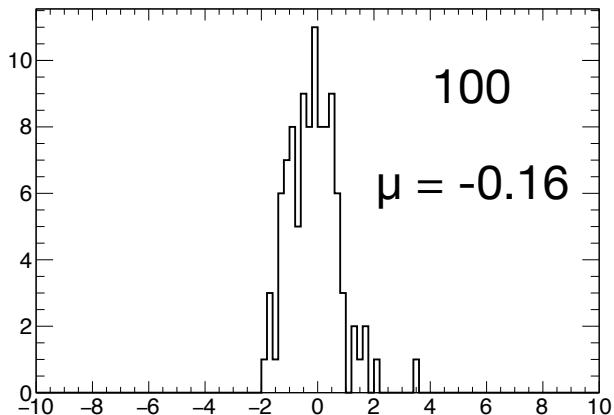
Central Limit Theorem



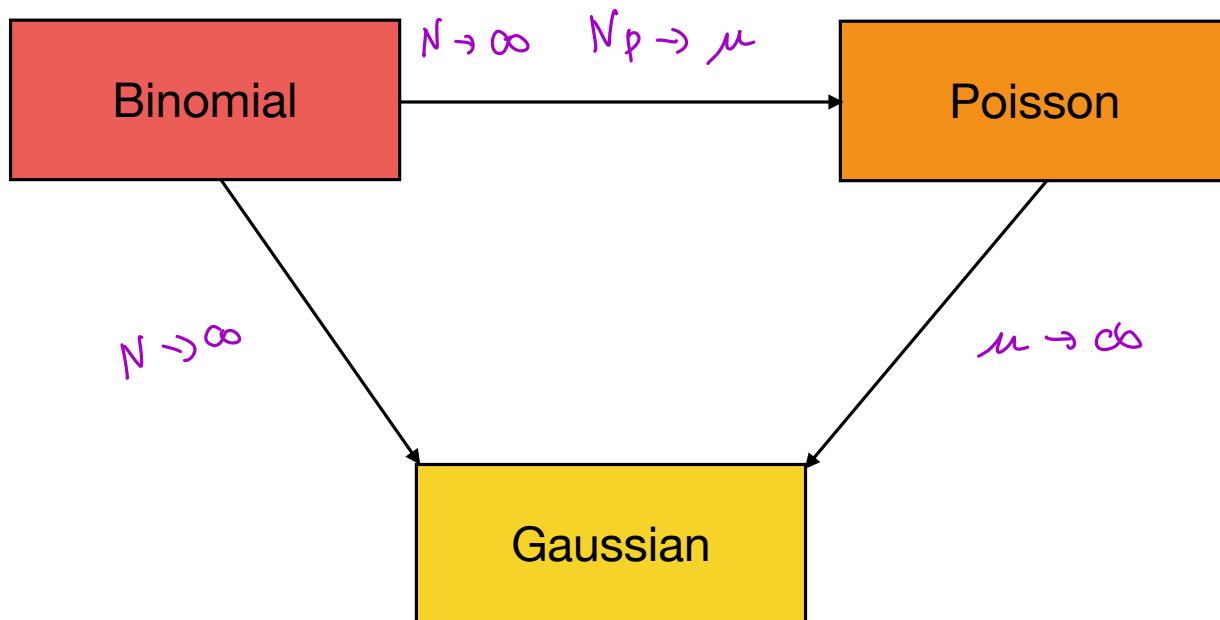
http://en.wikipedia.org/wiki/File:Dice_sum_central_limit_theorem.svg

Example

Example in jupyter notebook



Relationships between pdfs



Point Estimation

- Standard problem: set of values described by x_1, x_2, \dots, x_n

$$\bullet f(x) \equiv f(x_n; \theta)$$

data parameters

- Point estimation:

$$\bullet \hat{\theta} = \underline{\theta}(x_1, x_2, \dots, x_n)$$

Estimator of θ

Estimators

- Typical goal: estimate **true values of one or more parameters** from experimental data and understand the uncertainty on that measurement
- **Characteristics** of an estimator
 - consistency
→ approaches true value asymptotically for infinite data
 - bias
→ difference wrt true value for finite dataset
 - efficiency
→ variance of the estimator
 - sufficiency
→ depends on true value
 - robustness
→ sensitivity to data, e.g. outliers
- **Uncertainty**: how far the **true parameter** might be from our estimate due to statistical fluctuations in the ensemble of the measurements

Basic Estimators

- Estimators for mean and variance
- Shape of the PDF (fitting)
 - Maximum likelihood
 - Most efficient , but may be biased
 - Goodness-of-fit is not readily available
 - Least Chi-squared
 - Maximum likelihood for Gaussian distributed data
 - Convenient for binned data , analytic solve for linear functions
 - Automatic goodness-of-fit measure
 - Be careful of Gaussian approximation
 - (e.g. when Poisson becomes Gaussian)

Mean and Variance from a Sample

- Estimators (equally weighted data)

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N x_i \quad N > 0$$

$$\hat{\sigma}^2 = \frac{1}{N-1} \sum_{i=1}^N (x_i - \hat{\mu})^2 \quad N > 0$$

- Variances of these

$$\underline{V[\hat{\mu}] = \frac{\sigma^2}{N}} \quad \sigma(\hat{\mu}) = \frac{\sigma}{\sqrt{N}}$$

$$\underline{V[\hat{\sigma}^2] = \frac{1}{N} \left(m_4 - \frac{N-3}{N-1} \sigma^4 \right)}$$

$$\underline{\sigma(\hat{\sigma}) = \frac{\sigma}{\sqrt{2N}}} \quad \text{Gaussian data of } x \text{ and large } N$$

Likelihood Function

- Likelihood $\underline{\mathcal{L}}(x; \theta)$: probability that a measurement of \underline{x} will yield a specific value for a given theory
 - Need to specify both the theory and the value for any parameters in that theory
- With an ensemble of measurements, overall likelihood is obtained from the product of likelihoods of the measurements

$$\underline{\mathcal{L}}(x; \theta) = \prod_{i=1}^n \underline{\mathcal{L}_i}$$

- Here θ represents one or more parameters

Log Likelihood

- To estimate the parameter(s), Θ maximise the likelihood

- Set derivative to zero

- Typically easier to maximise the

$$\ln \mathcal{L}$$

$$\frac{\partial \mathcal{L}}{\partial \theta} = \frac{\partial}{\partial \theta} \ln \prod_{i=1}^n f_i$$

$$\begin{aligned} \ln(ab) \\ = \ln a + \ln b \end{aligned}$$

$$= \frac{\partial}{\partial \theta} \sum_{i=1}^n \ln f_i$$

$$\approx 0$$

- If there are several Θ_i we can minimise with respect to each of them

$$\frac{\partial \mathcal{L}}{\partial \theta_1} = 0 \quad \frac{\partial \mathcal{L}}{\partial \theta_2} = 0 \quad \dots \quad \frac{\partial \mathcal{L}}{\partial \theta_i} = 0$$

Likelihood Example: Poisson

- N independent trials with results n_i
- Likelihood function for observing n_i if true mean is μ

$$\mathcal{L}(n_i; \mu) = \frac{e^{-\mu} \mu^{n_i}}{n_i!}$$

- Product over N measurements

$$\mathcal{L}(\text{data}, \mu) = \prod_{i=1}^N \frac{e^{-\mu} \mu^{n_i}}{n_i!}$$

$$\begin{aligned} \Rightarrow \ln \mathcal{L} &= \sum_{i=1}^N (-\mu + n_i \ln \mu - \ln(n_i!)) \\ &= -N\mu + \sum_{i=1}^N n_i \ln \mu + \text{const} \end{aligned}$$

$$\frac{\partial \ln \mathcal{L}}{\partial \mu} = -N + \sum_{i=1}^N \frac{n_i}{\mu} = 0$$

$$\begin{aligned} N &= \frac{\sum_{i=1}^N n_i}{\mu} \\ \Rightarrow \mu &= \frac{1}{N} \sum_{i=1}^N n_i \end{aligned}$$

Best estimator is
the mean value

Likelihood Example: Gaussian

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

- Take the derivative of the log likelihood

$$\begin{aligned} \bullet \frac{\partial}{\partial \mu} (\ln \mathcal{L})|_{\hat{\mu}=\mu} &= \frac{\partial}{\partial \mu} \left(-\sum_i \frac{(x_i - \mu)^2}{2\sigma^2} + \text{const} \right) \\ &= \sum_i \frac{x_i - \mu}{\sigma^2} \\ \Rightarrow \hat{\mu} &= \frac{1}{N} \sum_i x_i \end{aligned}$$

- The unbiased estimator for σ is

$$\bullet \hat{\sigma} = \frac{1}{N-1} \sum_i (x_i - \mu)^2$$

Binned vs unbinned likelihood functions

- Likelihood formalism works for any well-behaved probability density function (pdf)
- Product of the likelihood is a product over measurements
- Example measurement:** Measure the lifetime of a particle of a given species for an ensemble of such particles produced at $t=0$ such that the # p at time t :

$$\bullet f(t) = \frac{1}{\tau} e^{-t/\tau}$$

- Consider two ways to construct the likelihood
 - For decay i measure t_i and take the product of all measured times \mathcal{L} (unbinned likelihood)
 - Make a histogram of the number of decays in bins of time
 - Measurement is the number of decays in each time bin (binned likelihood)

The Likelihood and χ^2

- If the data is Gaussian, we have

- $\ln \mathcal{L} = - \sum_i \frac{(x_i - \mu)^2}{2\sigma^2} + \text{const.}$

- Compare to

- $\underline{\chi^2} = \sum_i \frac{(x_i - \mu)^2}{\sigma^2}$

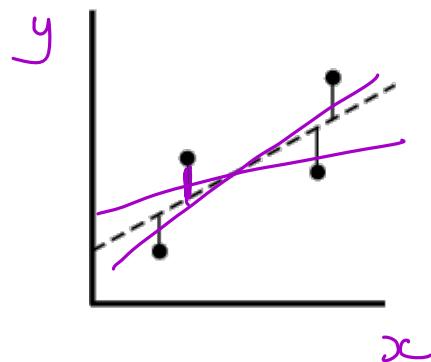
- By inspection

- $\chi^2 = \cdot - 2 \ln \mathcal{L}$

- Note: the likelihood formulation works for all pdfs not just Gaussians

Method of Least Squares

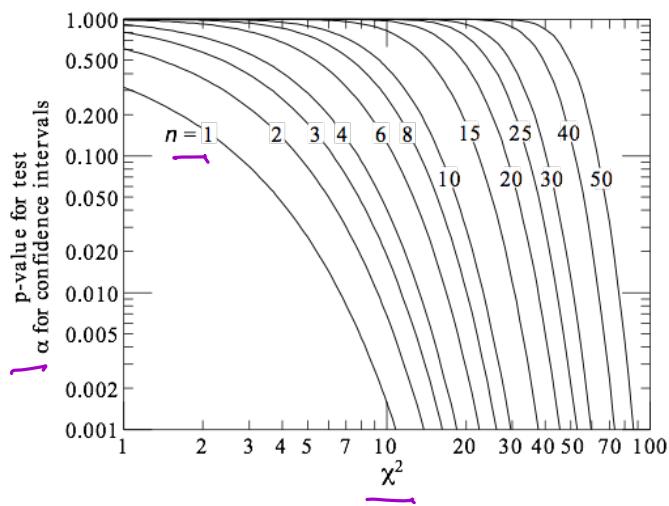
- Assume we have enough statistics for our measurements such that we can assume we are in the Gaussian regime
- Goal: Find the best estimates of parameter of a function that describes the data
- Minimize the distance of the data from the fit
 - Account for the uncertainties
- Scatter defined by χ^2
- $\chi^2 = \sum_{i=1}^N \frac{(x_i - \mu)^2}{\sigma^2}$
- Can write the χ^2 in terms of our observables
- $\chi^2 = \sum_{i=1}^N \frac{(y_i - F(x_i, \theta))^2}{\sigma^2}$
- Minimise χ^2 with respect to θ
- Useful when minimising $\ln \mathcal{L}$ is slow (high statistics samples)



Example: chi-squared p-values

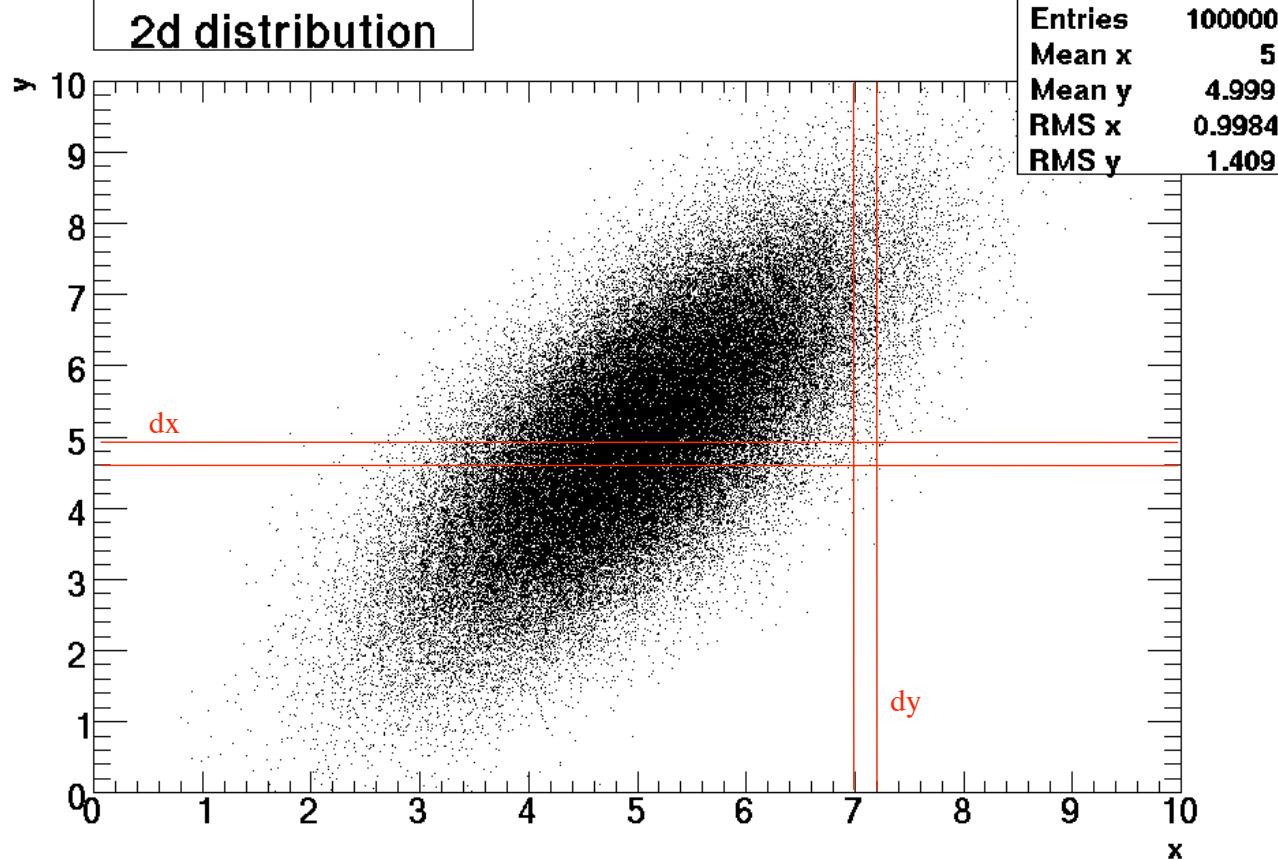
- One advantage of a χ^2 is that the value of the $\min \chi^2$ can be interpreted as a *goodness-of-fit*
 - iff *errors* on each data point are *known* and
 - the *noise* (distribution of *data* around their *expected value*) are *Gaussian*
- In the plot below
 - $n = \text{number of degrees of freedom} = n_{\text{dof}} = N_{\text{data}} - N_{\text{parameters}}$
 - For a *good fit*, expect χ^2 close to be close to *dof*

$$\chi^2 / \text{dof} \sim 1$$

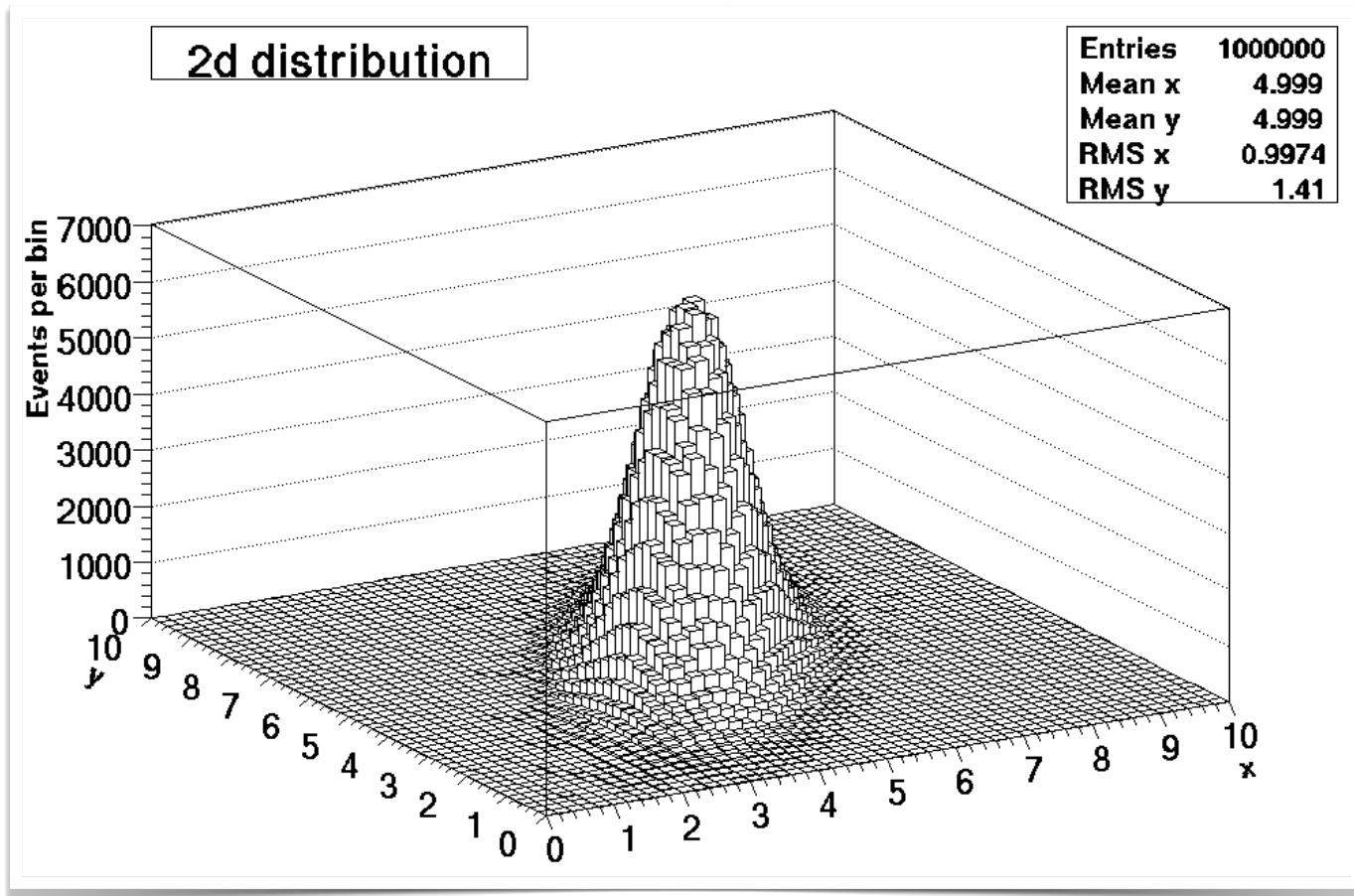


Example in jupyter notebook

2D distribution



2d distribution



Covariance and Correlation

Covariance Matrix

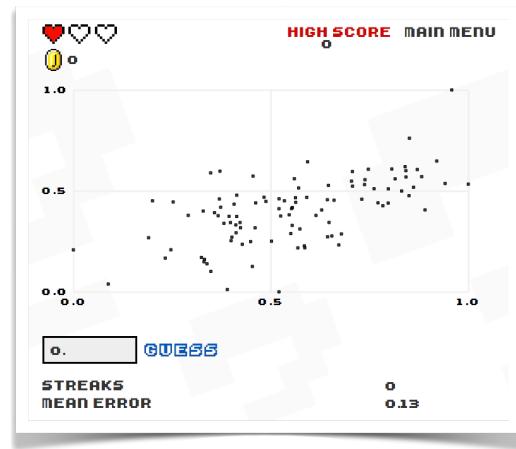
$$\text{cov}[x, y] = \underline{\mathbb{E}[(x - \mu_x)(y - \mu_y)]} = E[xy] - \mu_x \mu_y$$

- A representation of the N-dimensional parameter space as a covariance matrix
 - Diagonal elements: variance
 - Off-diagonal: covariance

Correlation (normalised covariance)

If two variables are uncorrelated, independent variables, then $\text{cov}[x, y] = 0$ for $x \neq y$

$$\rho_{xy} = \frac{\text{cov}[x, y]}{\sigma_x \sigma_y}$$



<http://guessthecorrelation.com/>

Covariance Matrix for a Gaussian

- If x and y are independent variables

$$\bullet G(x, y | \mu_x, \sigma_x, \mu_y, \sigma_y) = \underbrace{\frac{1}{\sqrt{2\pi} \sigma_x} e^{-\frac{(x-\mu_x)^2}{2\sigma_x^2}}}_{\frac{1}{\sqrt{2\pi} \sigma_y} e^{-\frac{(y-\mu_y)^2}{2\sigma_y^2}}}$$

$$\bullet \frac{\partial}{\partial \mu_x} \ln \mathcal{L} = \frac{1}{\sigma_x^2}$$

- Now, assume that x and y are correlated

- Covariance matrix is defined by

$$\bullet \langle \hat{V}^{-1} \rangle_{ij} = - \frac{\partial^2 \ln \mathcal{L}}{\partial \mu_i \partial \mu_j}$$

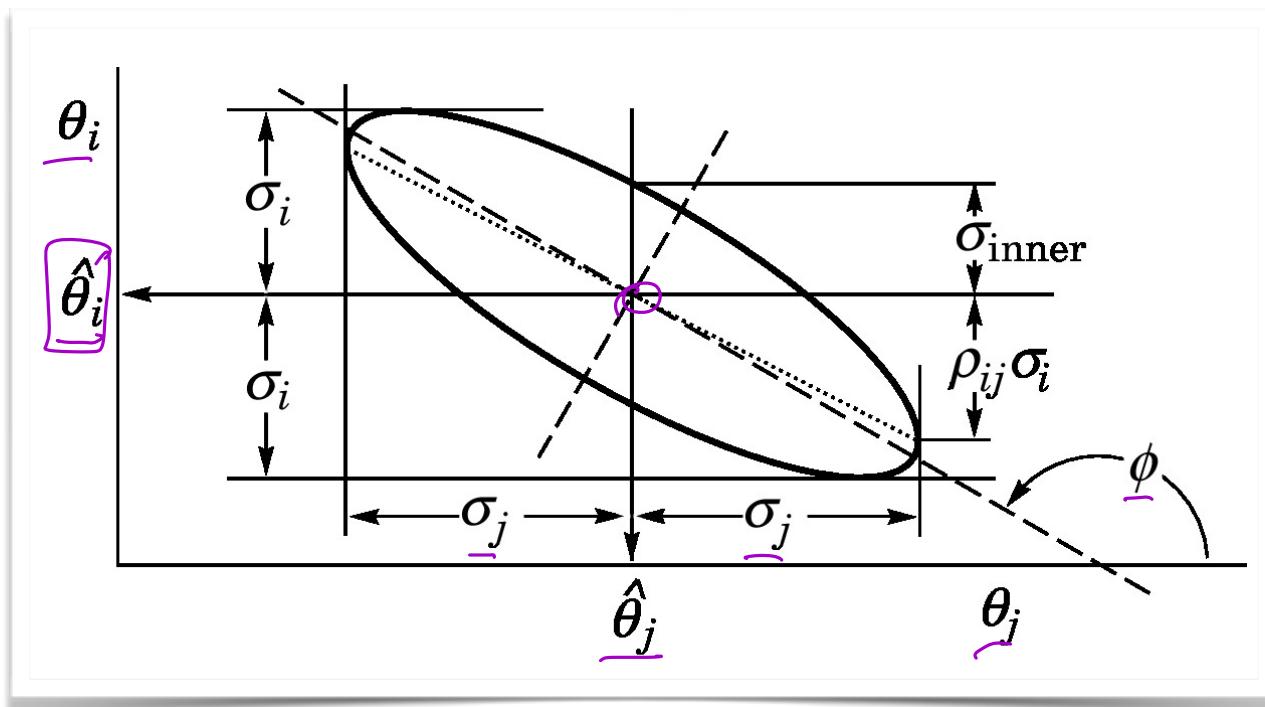
- For a binned likelihood, where N is large and the likelihood can be reduced

to a χ^2

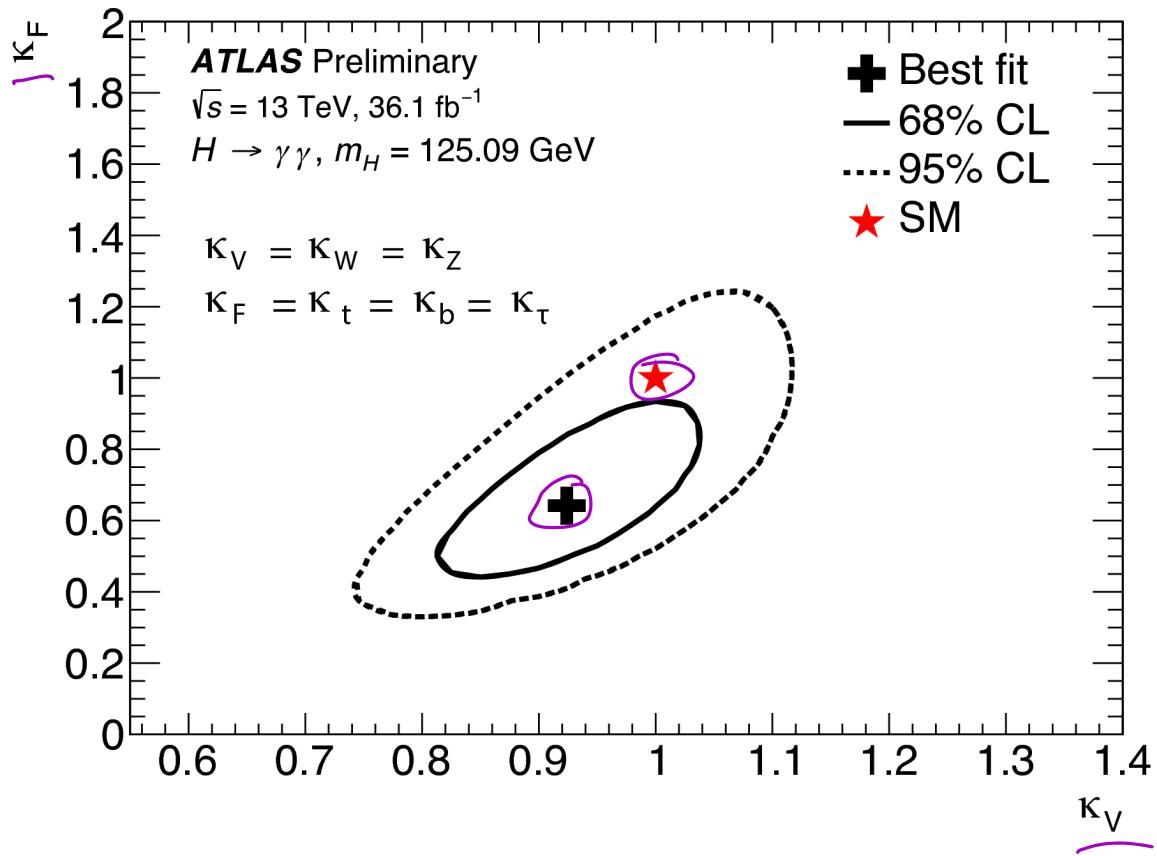
$$\bullet \langle \hat{V}^{-1} \rangle_{ij} = - \frac{2\chi^2}{\partial \mu_i \partial \mu_j}$$

Correlated Uncertainties

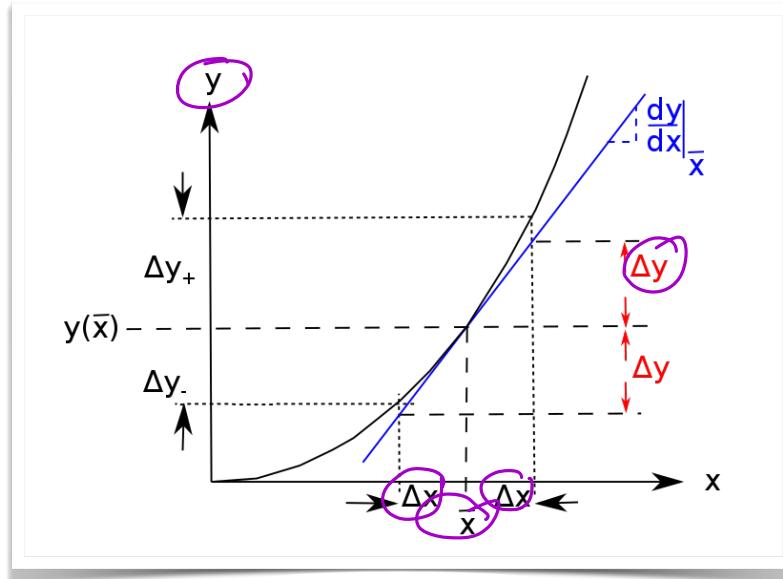
- Standard error ellipse for two parameters with a negative correlation
- Slope related to correlation coefficient $\frac{\partial \theta_i}{\partial \theta_j}$
- Correlation matrix typically determined from data numerically during fitting procedure



Correlation Example from Higgs



Propagation of Errors



- Determine error on final measurement from known errors on input measurements

$$\sigma_f^2 = \left(\frac{\partial f}{\partial \alpha} \right)^2 + \left(\frac{\partial f}{\partial \beta} \right)^2 + 2 \underbrace{\frac{\partial f}{\partial \alpha} \frac{\partial f}{\partial \beta}}_{\text{cov } \alpha \beta}$$

- More dimensions are usually expressed as a matrix
- Useful reference: https://en.wikipedia.org/wiki/Propagation_of_uncertainty

Confidence Interval

$\alpha = 5\%$

95%

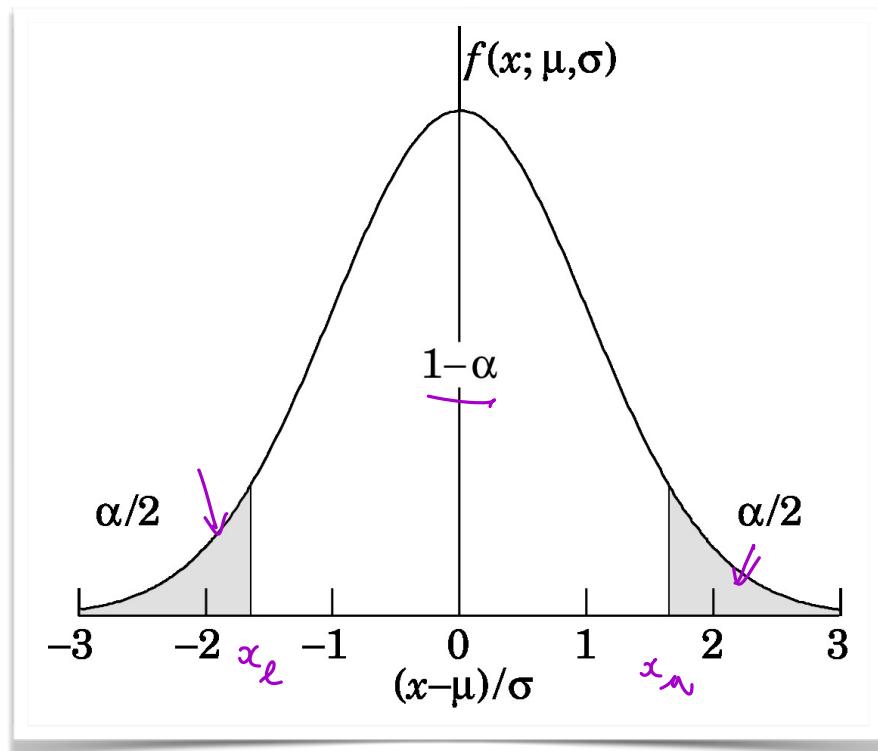
- Fraction of the result not between x_L and x_U is α

$\alpha = 10\%$

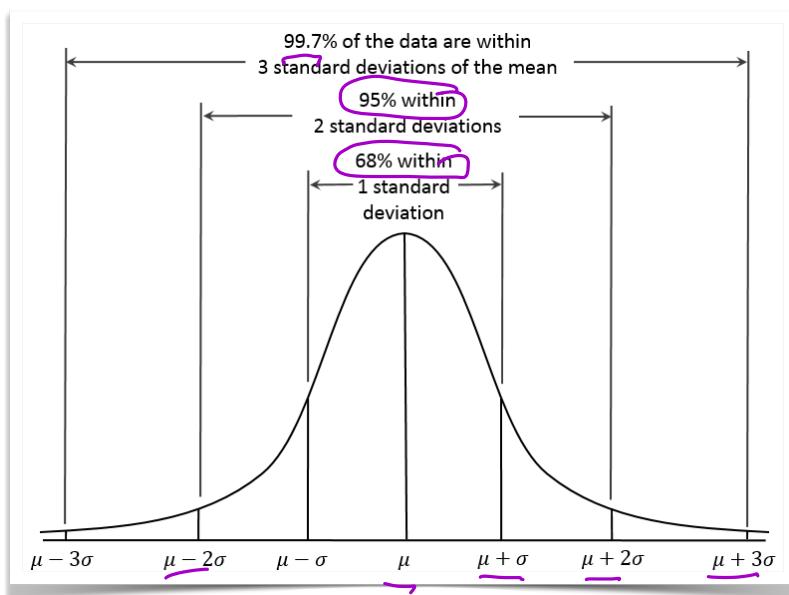
90%

$$\underline{1 - \alpha} = \int_{x_L}^{x_U} f(x; \theta) dx$$

\nwarrow PDF



Confidence Levels for a Gaussian



p-value

α	δ
0.3173	1σ
4.55×10^{-2}	2σ
2.7×10^{-3}	3σ
5.7×10^{-7}	4σ
2.0×10^{-9}	5σ

Confidence Levels: Higgs

