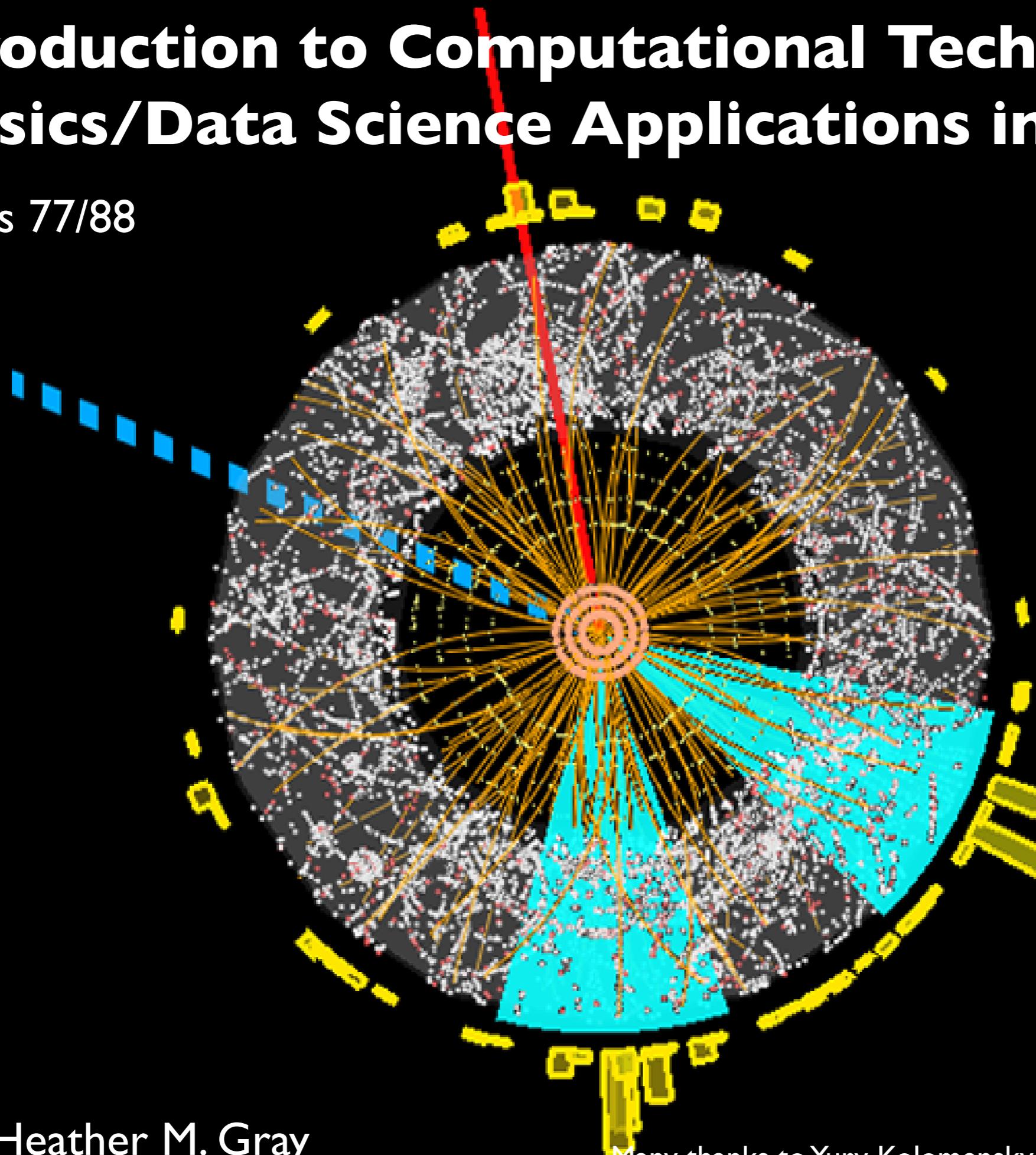


Introduction to Computational Techniques in Physics/Data Science Applications in Physics

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

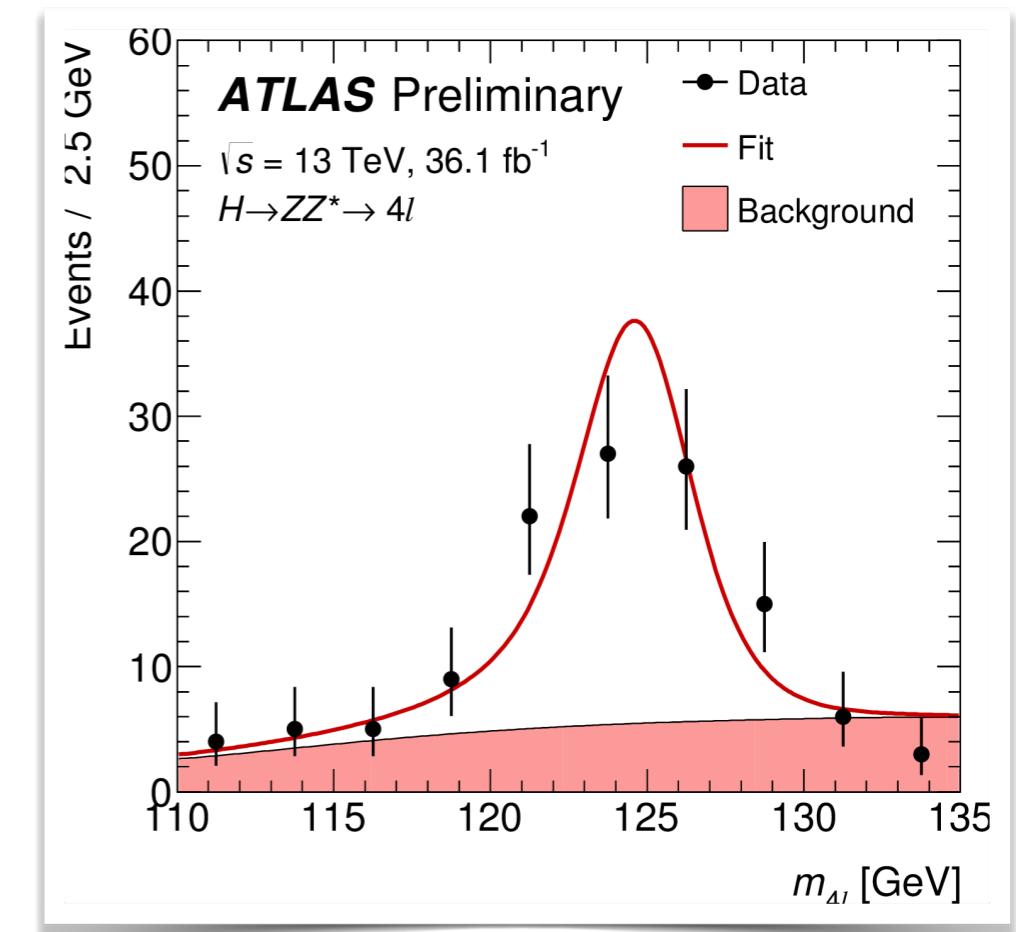
- Introduction
- Definitions: results of experiments
 - Random variables, probability, PDFs
- Interpreting results
 - Point estimators
 - Max likelihood, least squares fits
- Hypothesis testing, confidence limits
- Simulation (Monte Carlo techniques) — in two weeks

Statistics

- Statistics is a vital tool
 - Physics is an experimental science
 - Requires both and understanding
 - From to
 - Make a set of
 - Summarise the
 - Most conclusions are drawn with some degree of
 - Example
 -
 -
 -
 - Many measurements are a priori uncertain and have to be interpreted in probabilistic terms
 - **Classical statistics:** estimate given a and test whether a given is 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

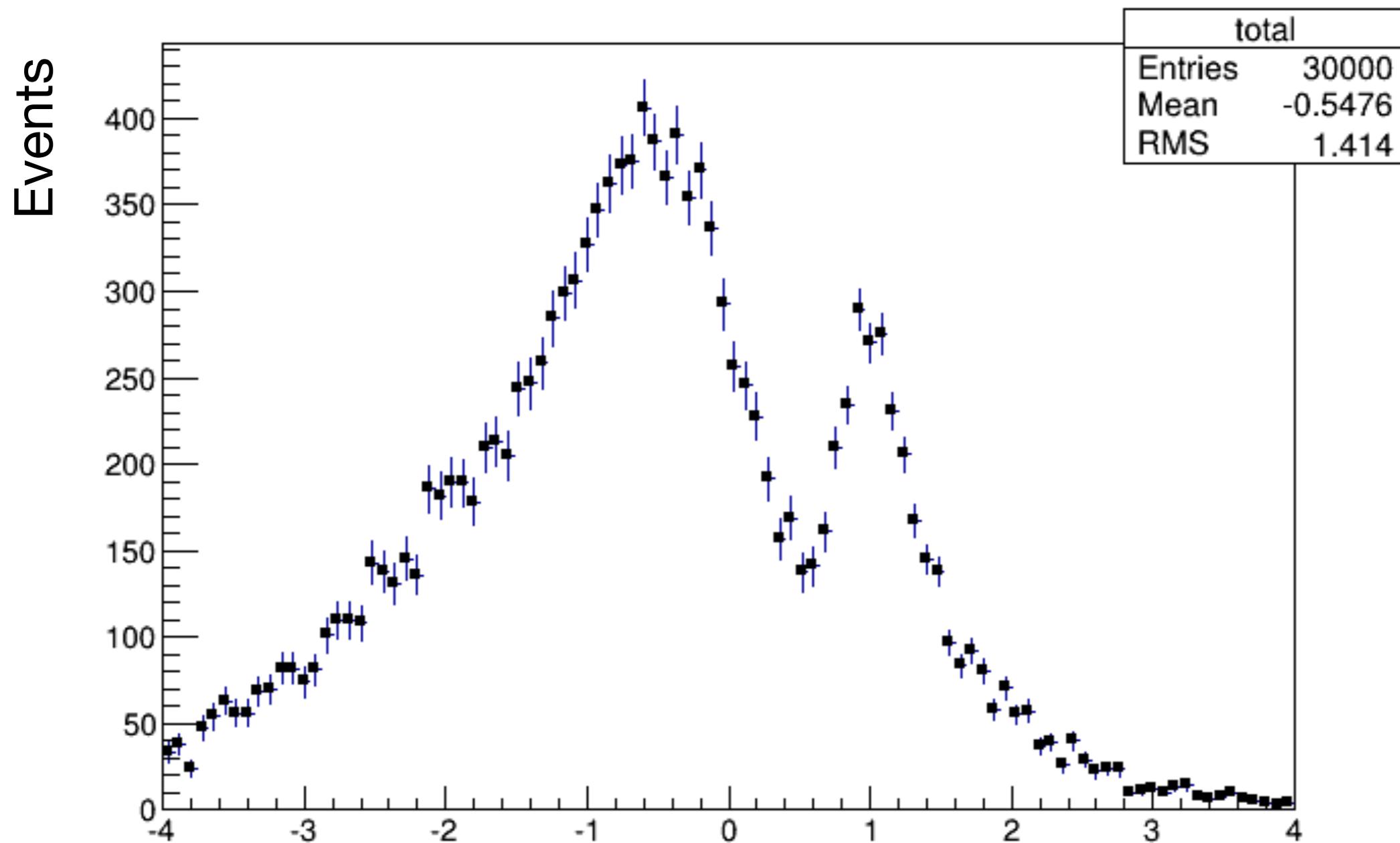
- **Data:**
 - Physics:
 - Other fields deal with
 - **Numbers** are easier to handle mathematically; statistics will deal with quantitative measurements
 - data, e.g. integers (counts)
 - data, e.g. energies, momenta
 - Measure with some , set by the measuring apparatus or other external conditions



Why are there error bars on the data?

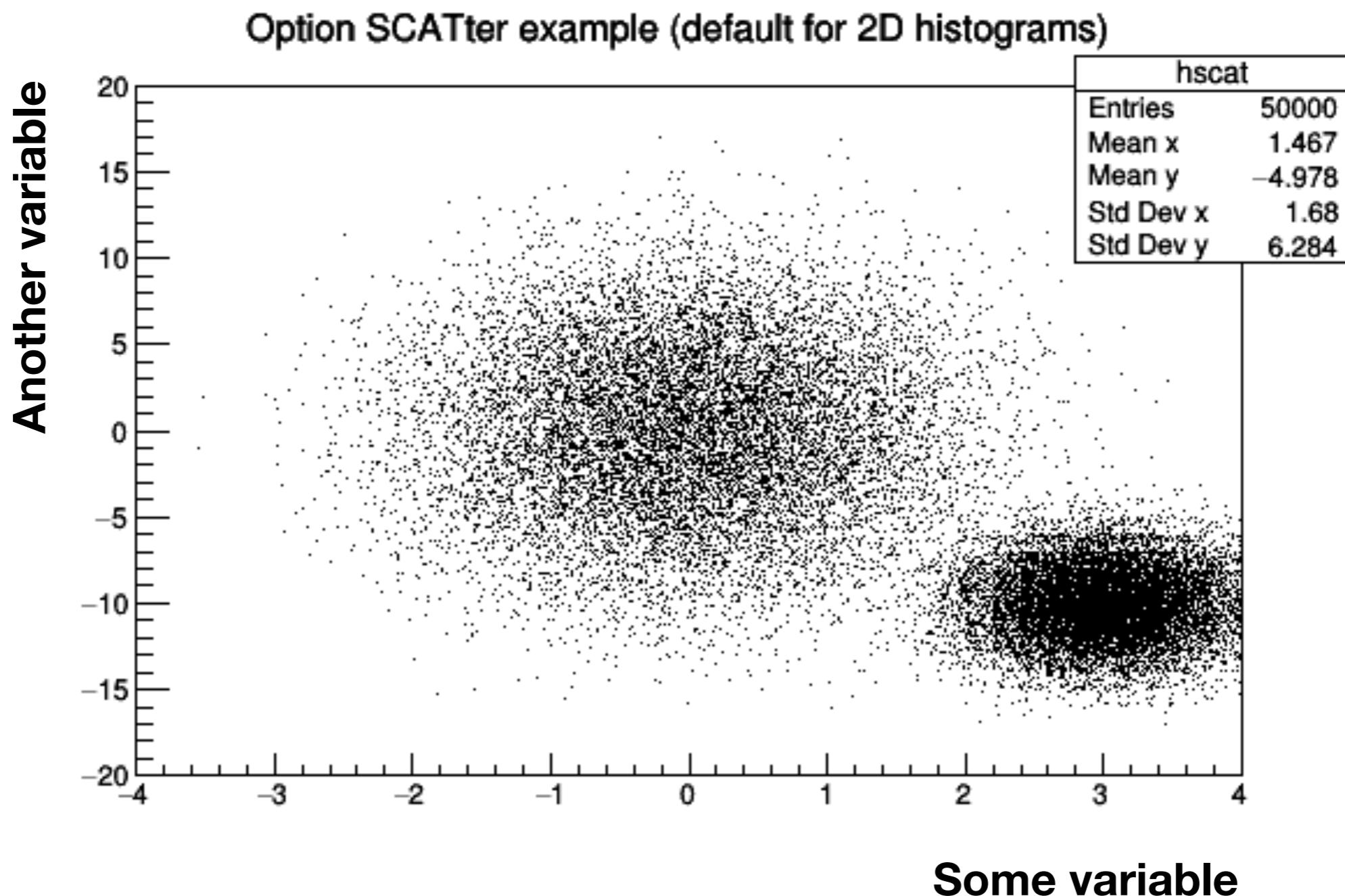
Histogram

This is the total distribution



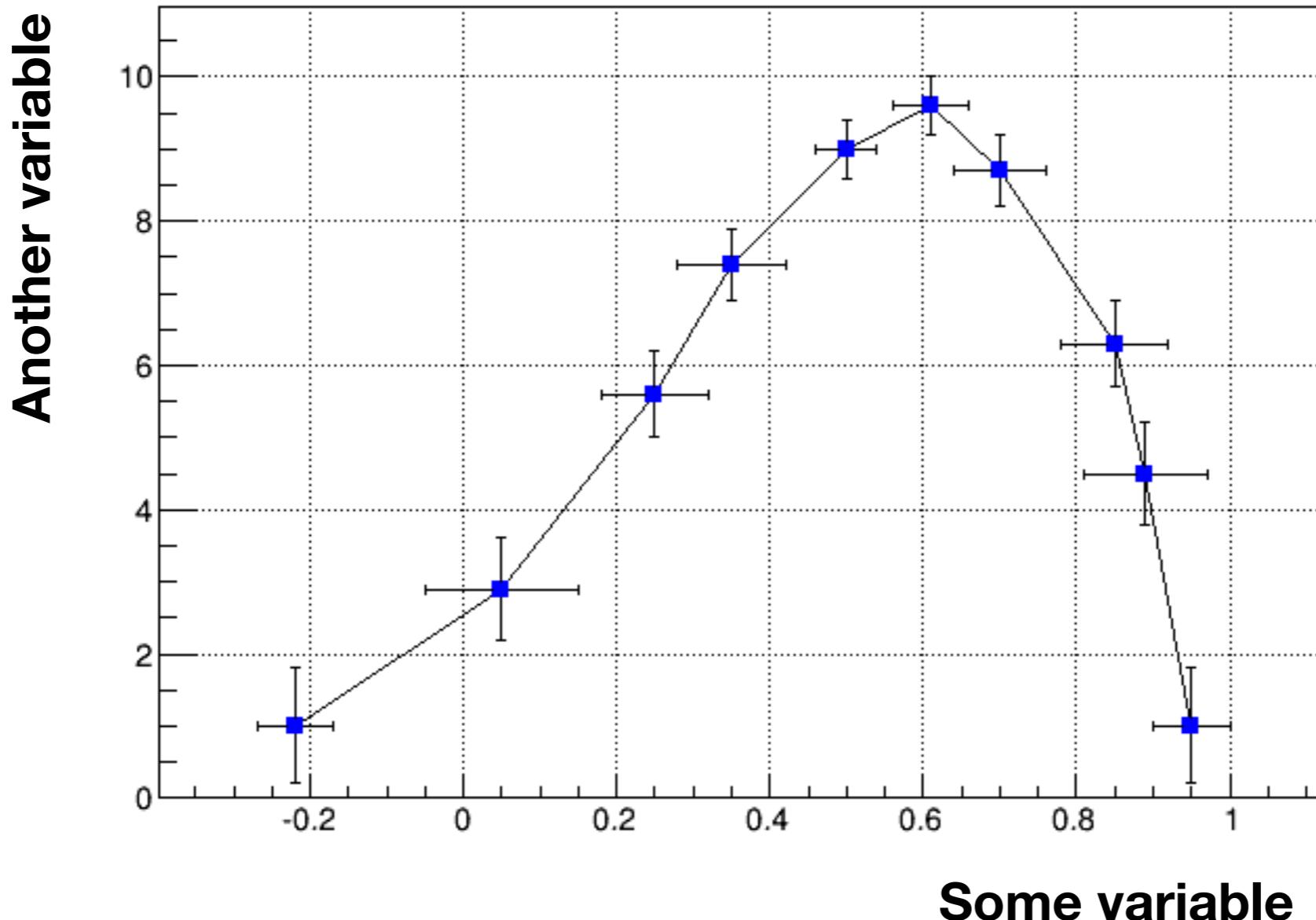
Some variable

Scatterplots



Graph

TGraphErrors Example



Uncertainty and Error

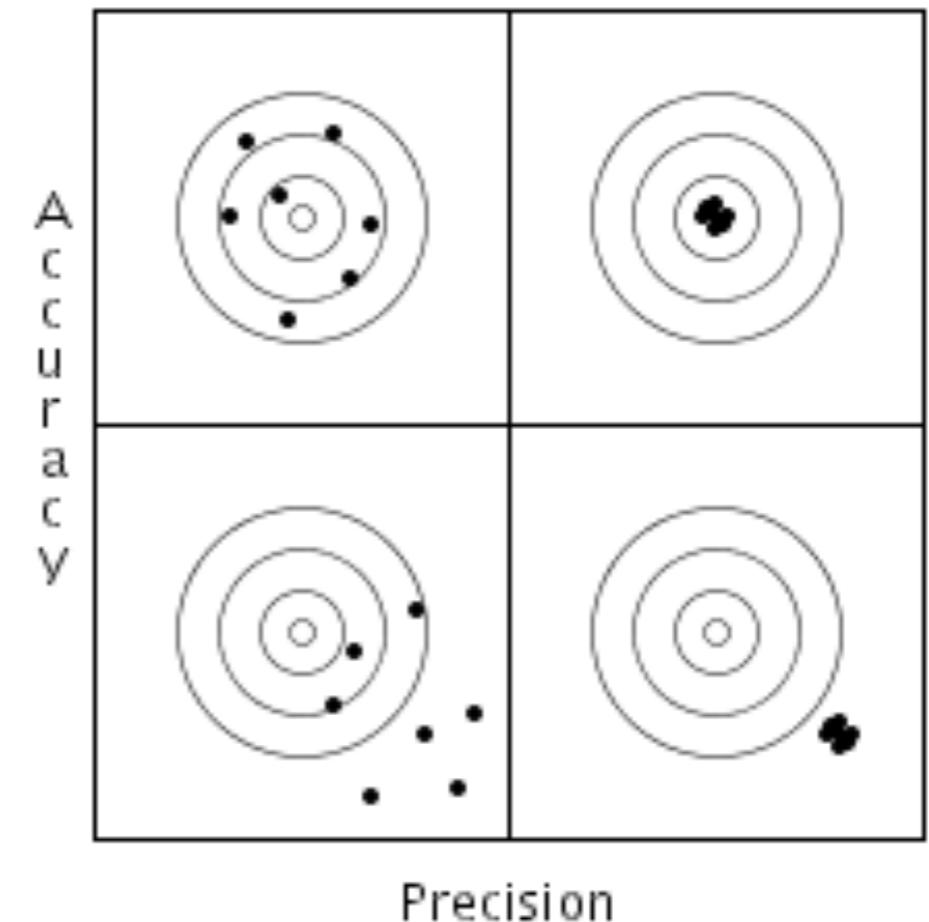
- In physics, the words **uncertainty** and **error** are used interchangeably to describe how far a measurement — typically
 - Use symbol Δ for the uncertainty
 - Formal definition is probabilistic: Δ is the range within which the true value is expected to fall to find the best fit
- Often interpreted as a measure of the spread of the data around the true values
- We'll come back to the difference between uncertainty and error

Uncertainty and Error

- How do we what is ?
 - Underlying assumption: our experiment is of a of
 - Derive the from the of the
 - Implicit assumption: our experiment is
 - i.e. all experiments would return results

Precision vs Accuracy

- Precision: of the data around the value
 - Typically associated with uncertainty
- Accuracy: of the value from the value
 - Typically associated with uncertainty
- Bad data:
 - Data with distribution (e.g.)



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

Golden Rules

- When reporting **of a measurement**, report its **uncertainty**
- Round off values to **3 significant digits**
 - Rule of thumb: **if the last digit is 5 or greater, round up**, otherwise **round down**
 - $x = 12.345 \pm 0.024$
 - $y = 12.3456 \pm 0.6$
- Uncertainty can come from the **accuracy** and/or **precision** of the measurement
 - Rule of thumb: **use the smaller uncertainty**
 - Statistically correct:

Probability: Definitions

- For numerical data, is often most convenient (and quantitative)
- Let's define probability now
 - Formally, it is a that is defined by
 - 1.
 - 2.
 - 3.

Two Interpretations

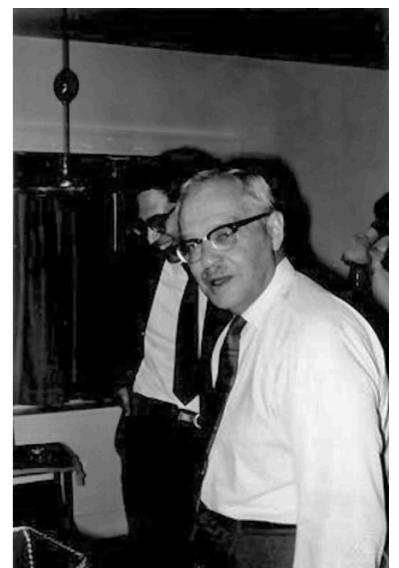
- “Frequentist” interpretation:
 - Probability is a measure of how often a given parameter value is reported when experiments are repeated many times
 - Measurable parameters are represented by random variables assigned to the outcomes of experiments
 - CL measures a confidence interval, given a sample size, such that an estimator would fall in a certain range a certain percentage of the time
 - No probability is assigned to individual parameter values
- “Bayesian” interpretation:
 - More general: define probability as a measure of the degree of belief that a given statement is true
 - E.g. that the true value of a parameter is in a certain interval
 - This is somewhat subjective, but follows how people actually reason

Frequentist Probability

- Definitions
 - Let Ω be the sample space of all possible outcomes of an experiment.
 - Any subset $A \subseteq \Omega$ with only finitely many elements is called an event.
 - Define $P(A) = \frac{|A|}{|\Omega|}$
- Assume that the probabilities are (in principle) known.
- Confidence in a hypothesis grows with the number of observations.
- Frequentist statistics is often used in situations where it can be assumed that the underlying process is random (e.g. particle physics).



John von Neumann



Jerzy Neyman

Bayes Theorem

- probability of given

-

- Interpreted within as

- $P(\text{theory} \mid \text{data}) \propto$



Thomas Bayes

- Allows one to a as a measure
of that a given is correct
(e.g. that some is in some)

- Requires assigning some to prior knowledge

- That's where comes in

Probability: Random Variables and PDFs

- For a continuous variable, x , we define the

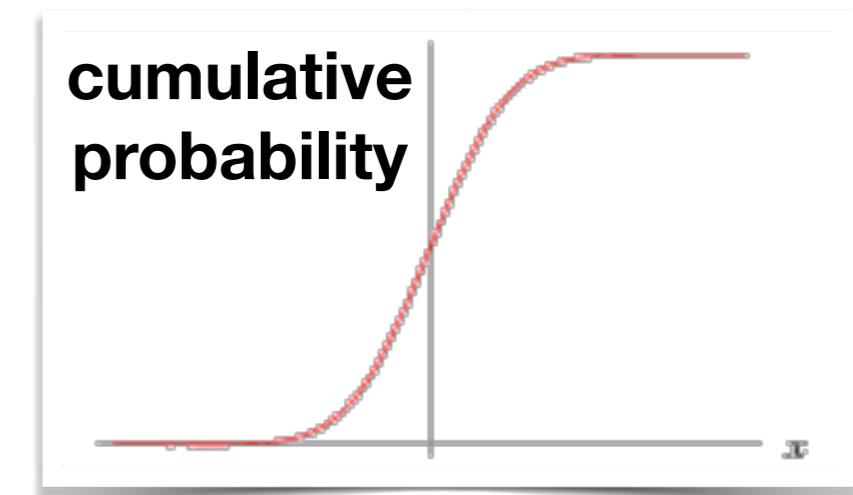
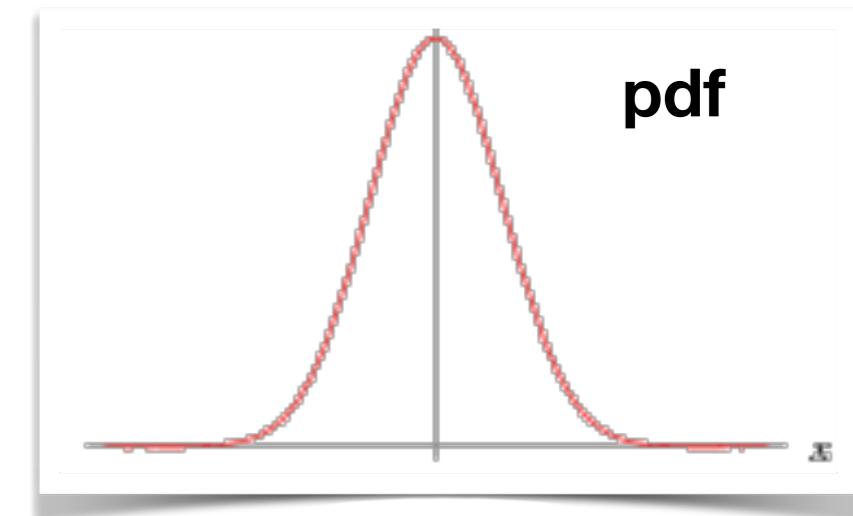
- $f(x, \theta) =$

- θ

- Integrate to obtain the

- $F(a) = \int$

- Probability that $x < a$
- Discrete variables:
- Expectation value



- $E[u(x)] =$

Expectation Values

- Expectation value of a

- $E[u(x)] =$

- Moments of a

- $\alpha_n \equiv$

- $m_n \equiv$

Mean and Variance

- Mean

- $\mu = \int$

- Variance

- $\sigma^2 =$

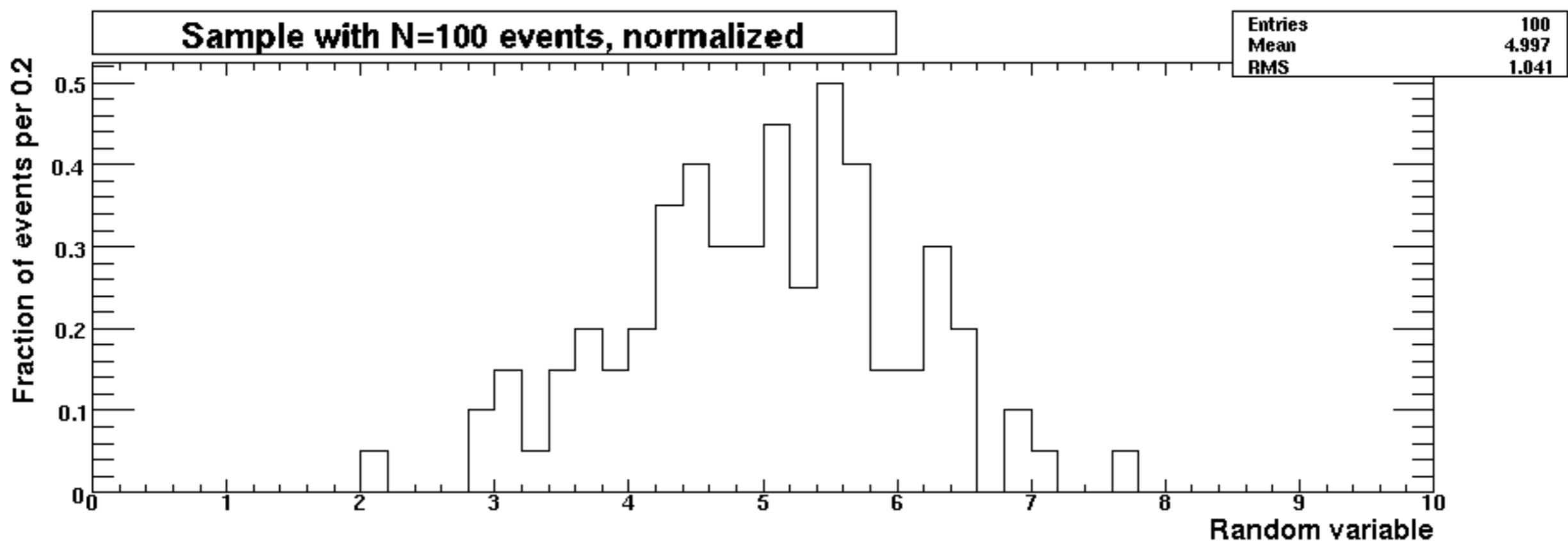
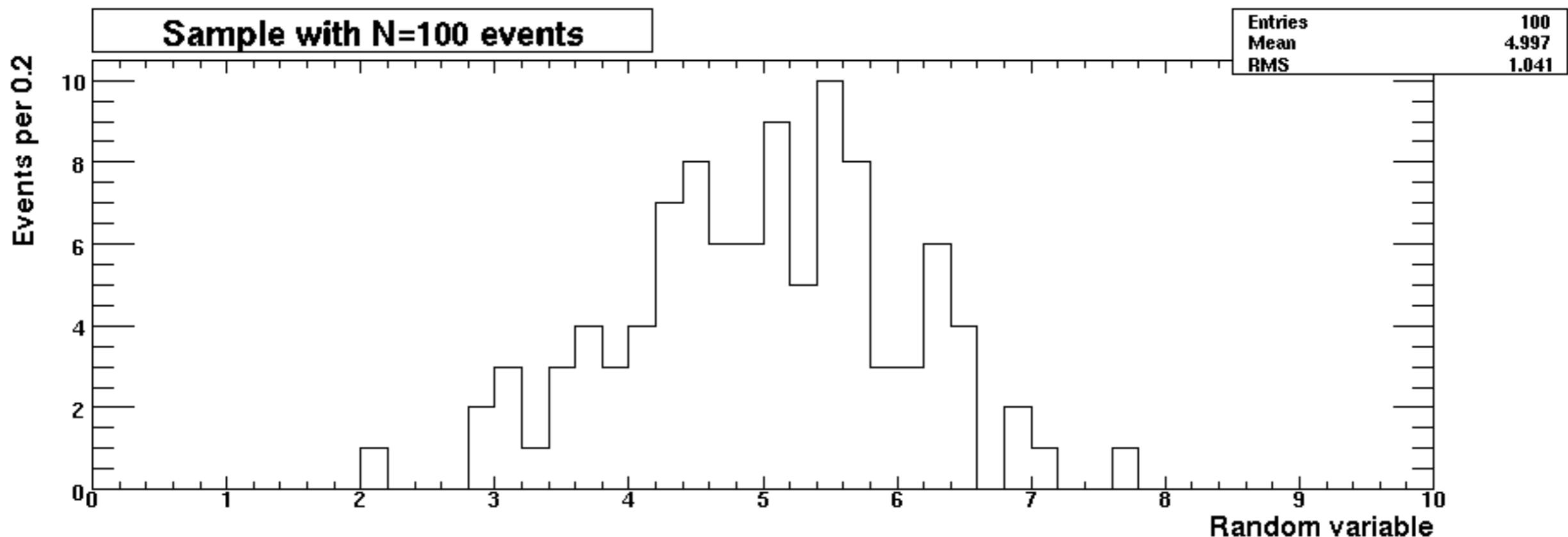
- σ is the standard deviation

If you only remember one thing today, remember this

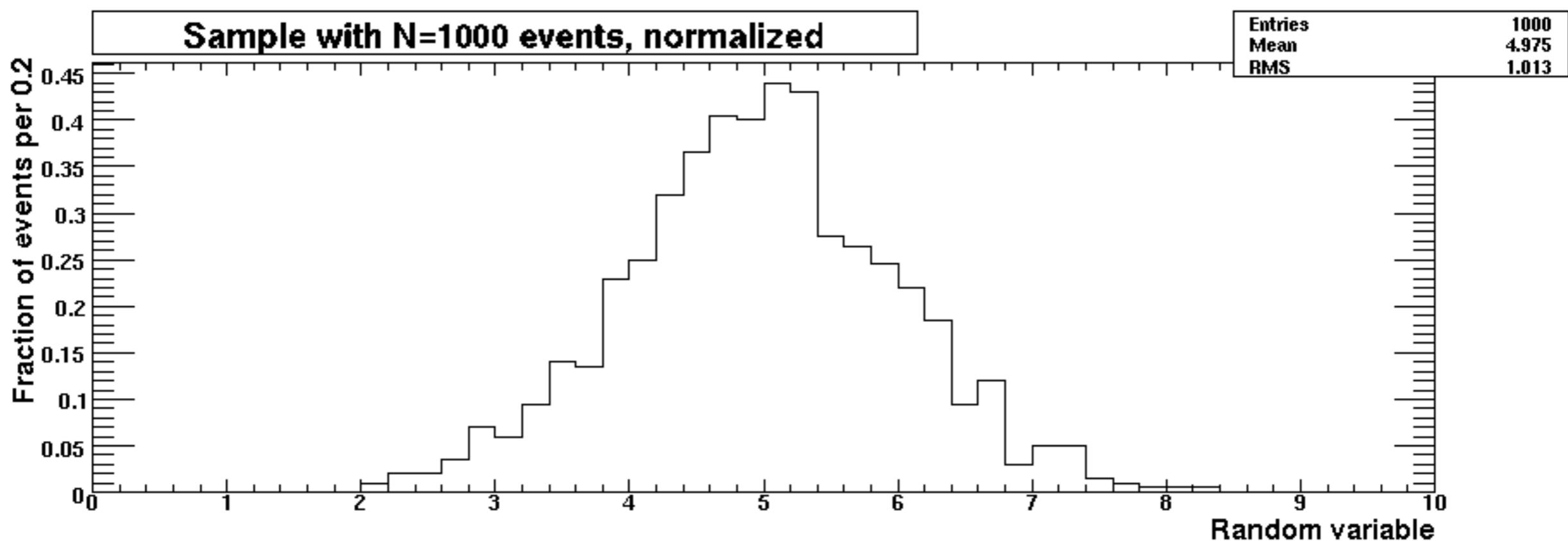
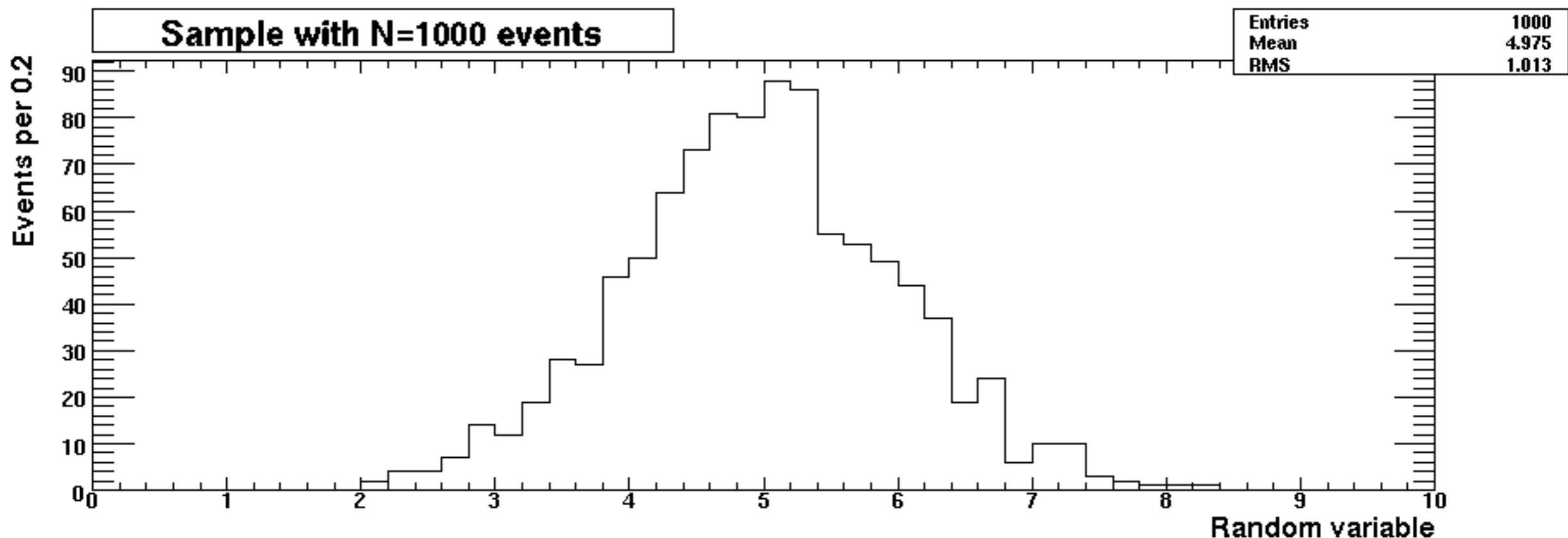
If we know the PDF, we know how to determine μ and σ

Example in jupyter notebook

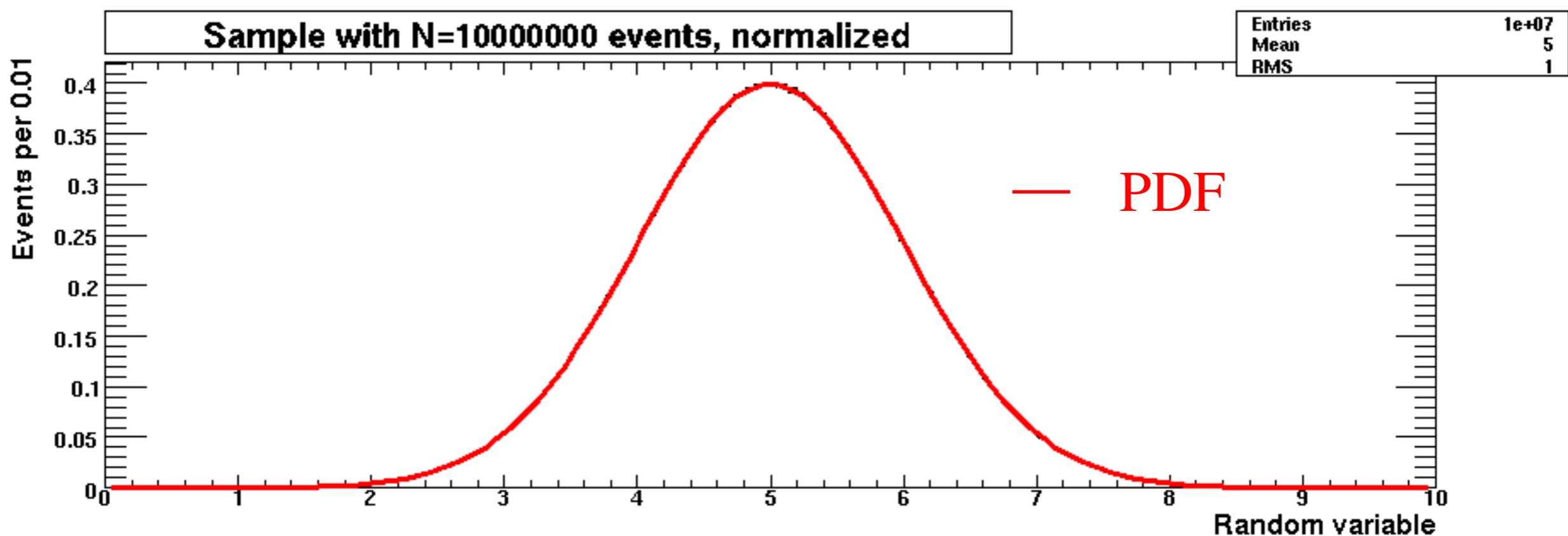
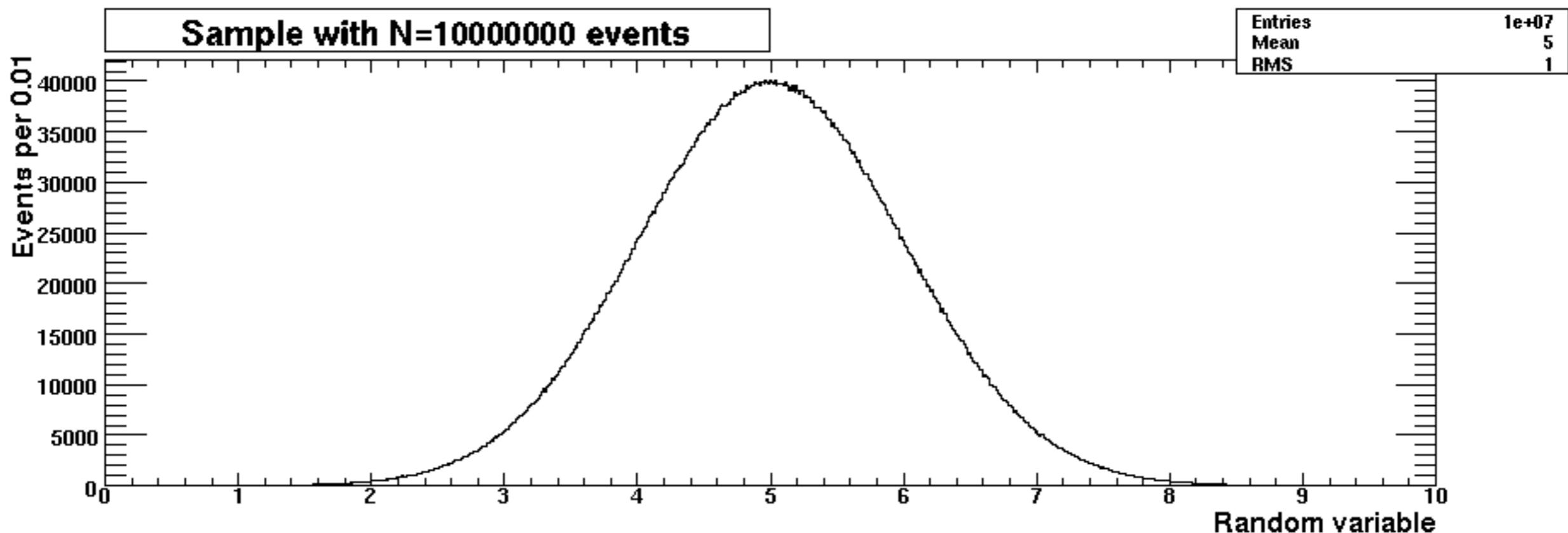
Sample from a Continuous PDF



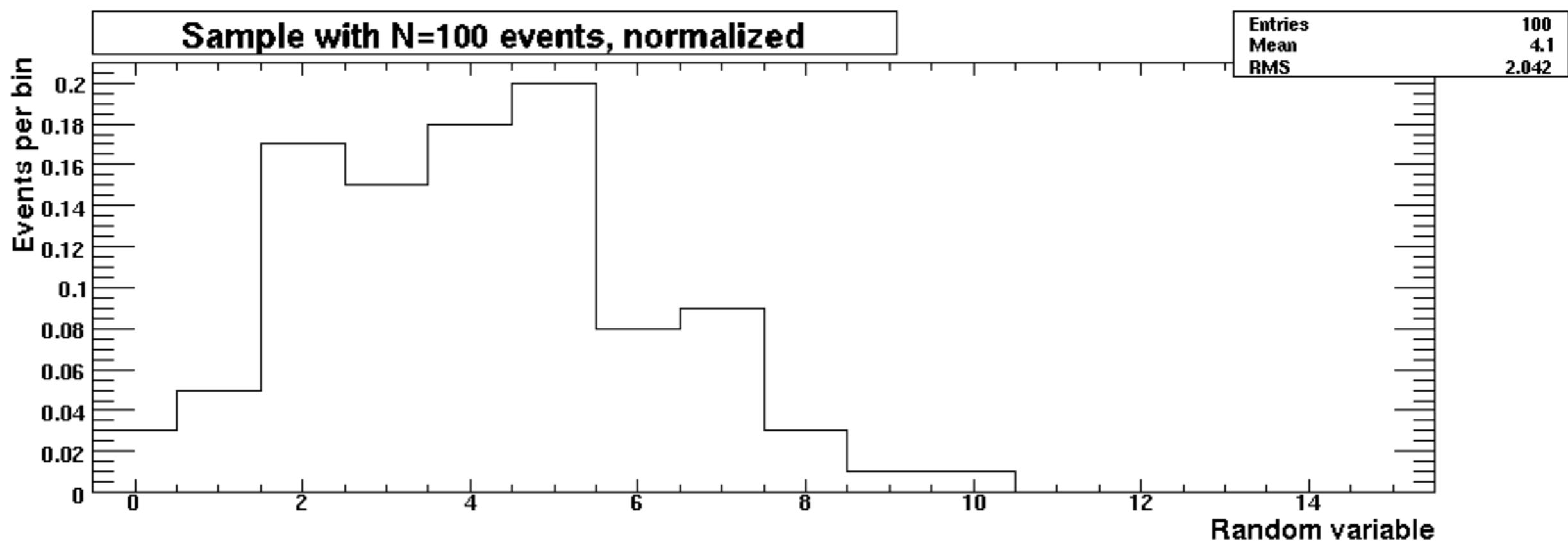
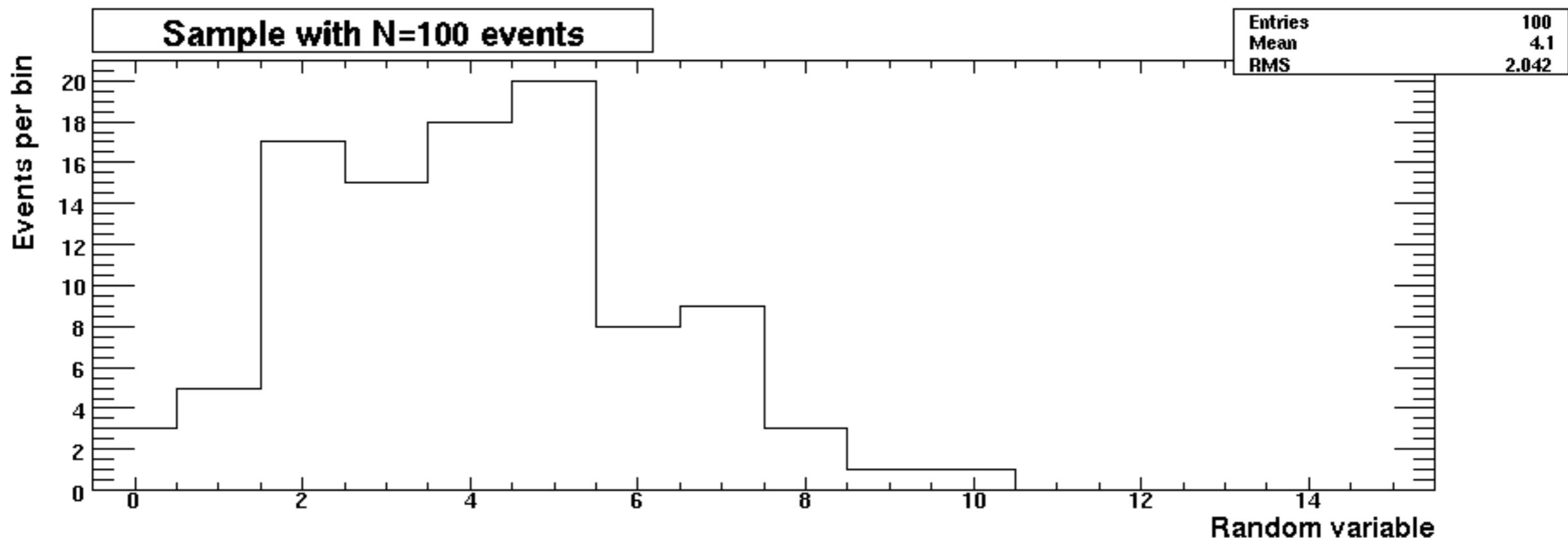
Sample From a Continuous PDF



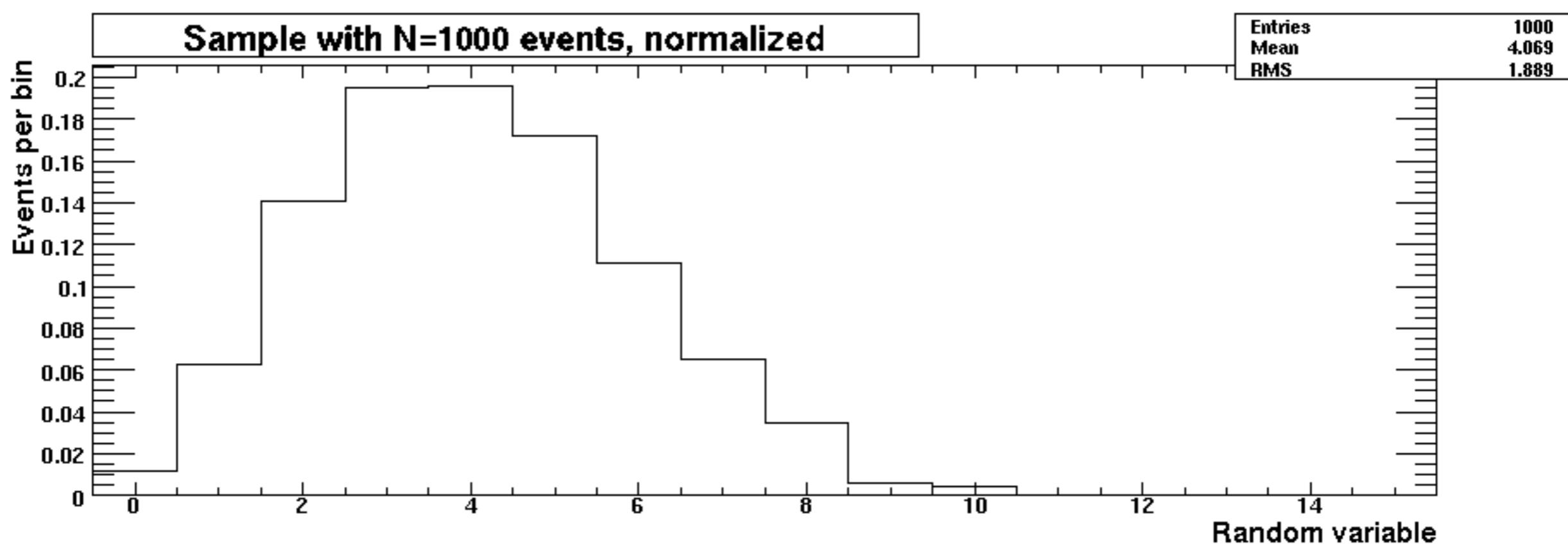
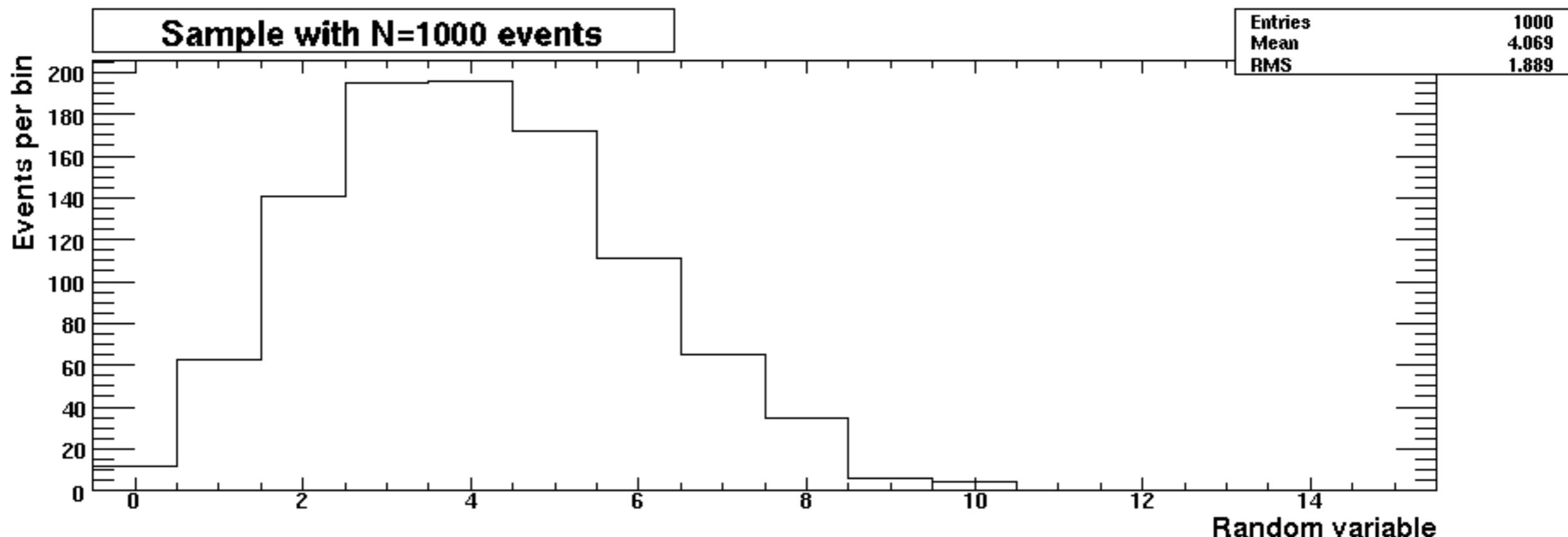
Sample From a Continuous PDF



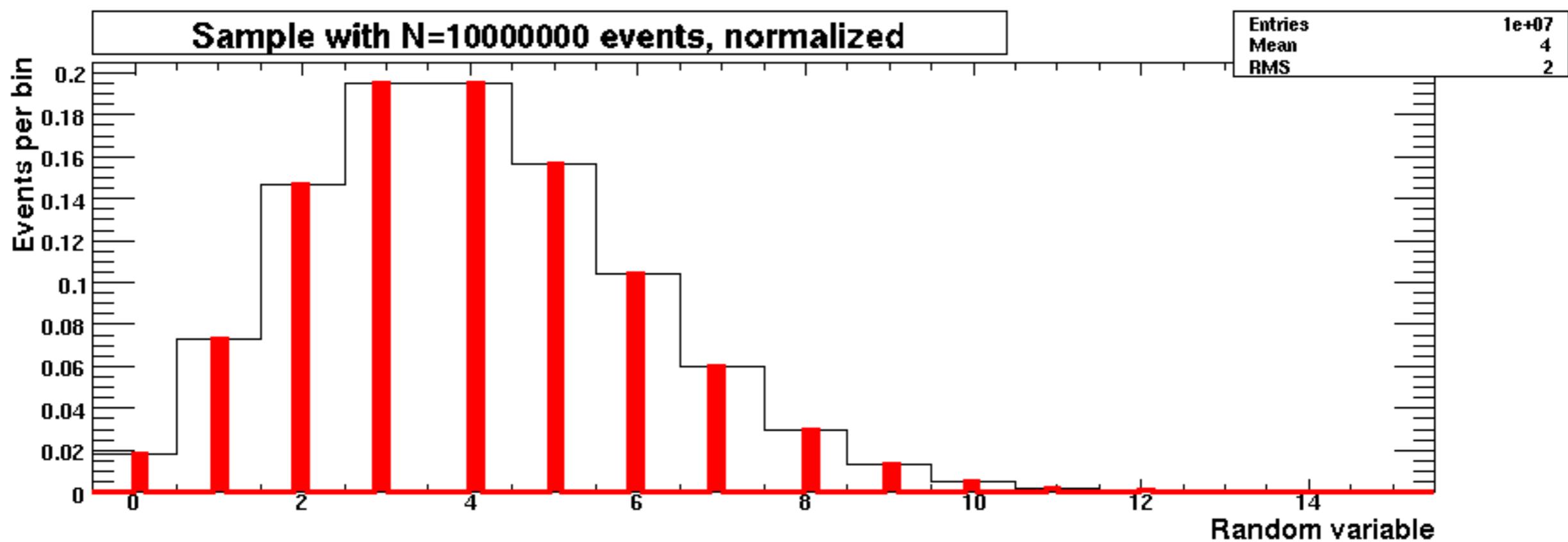
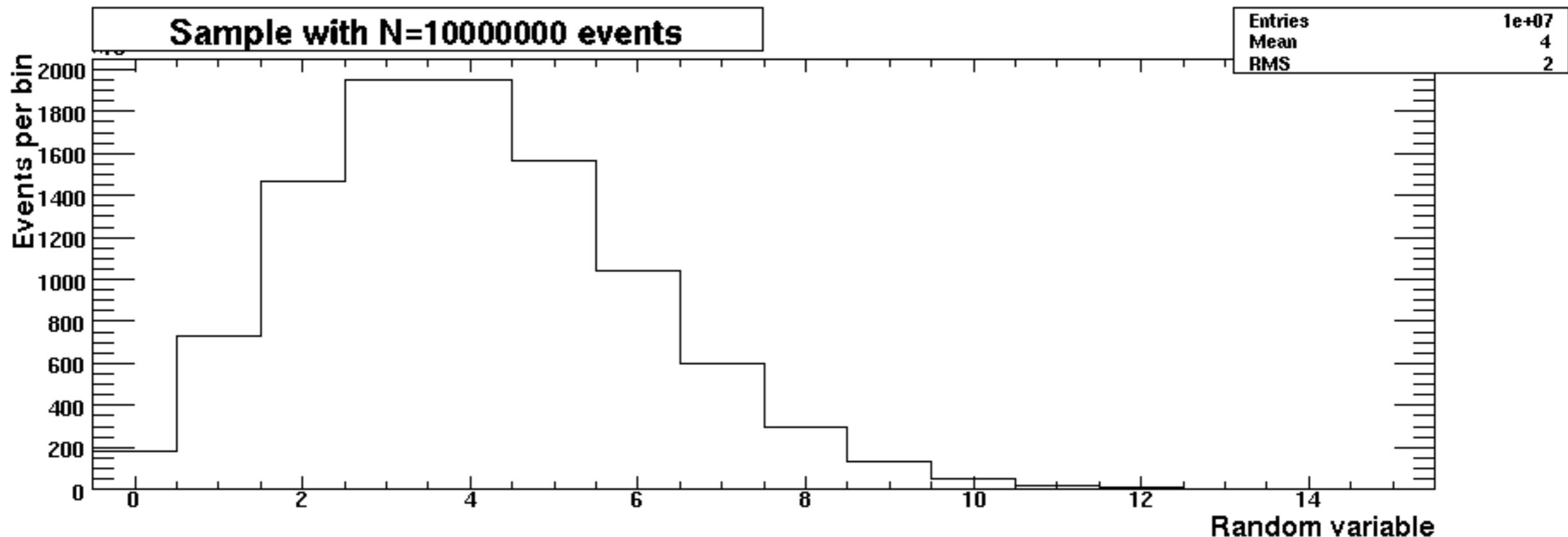
Sample From a Discrete PDF



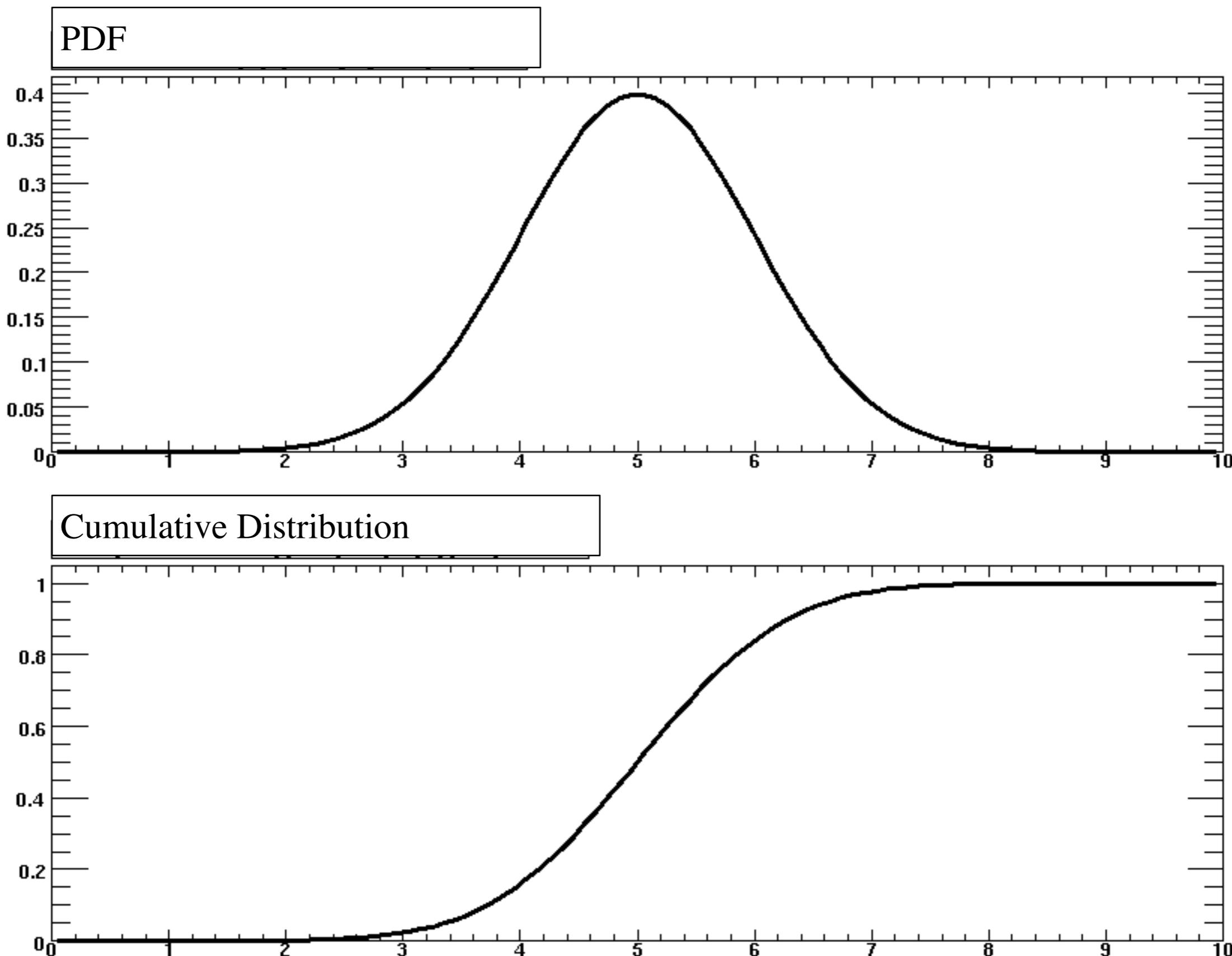
Sample From a Discrete PDF



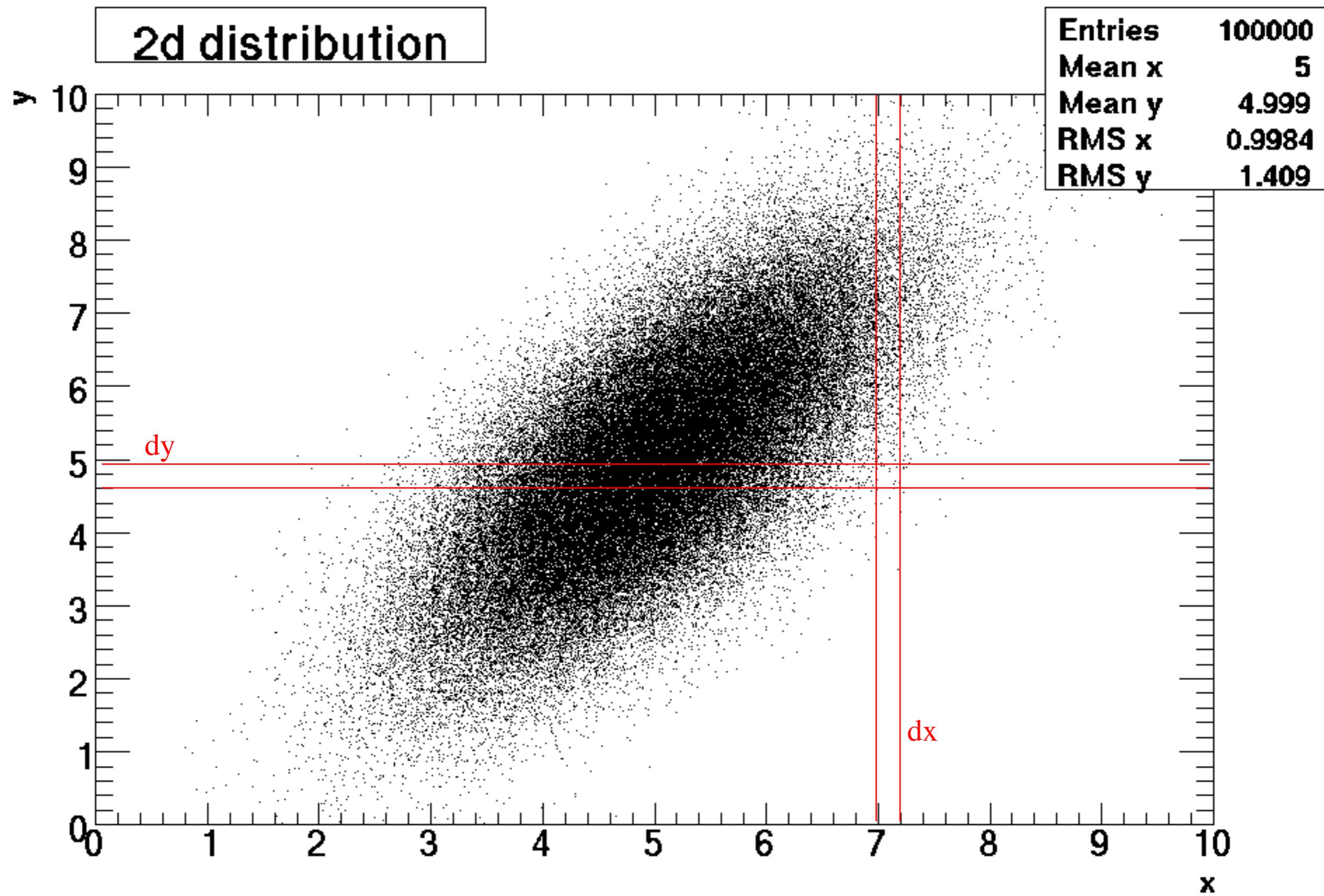
Sample From a Discrete PDF



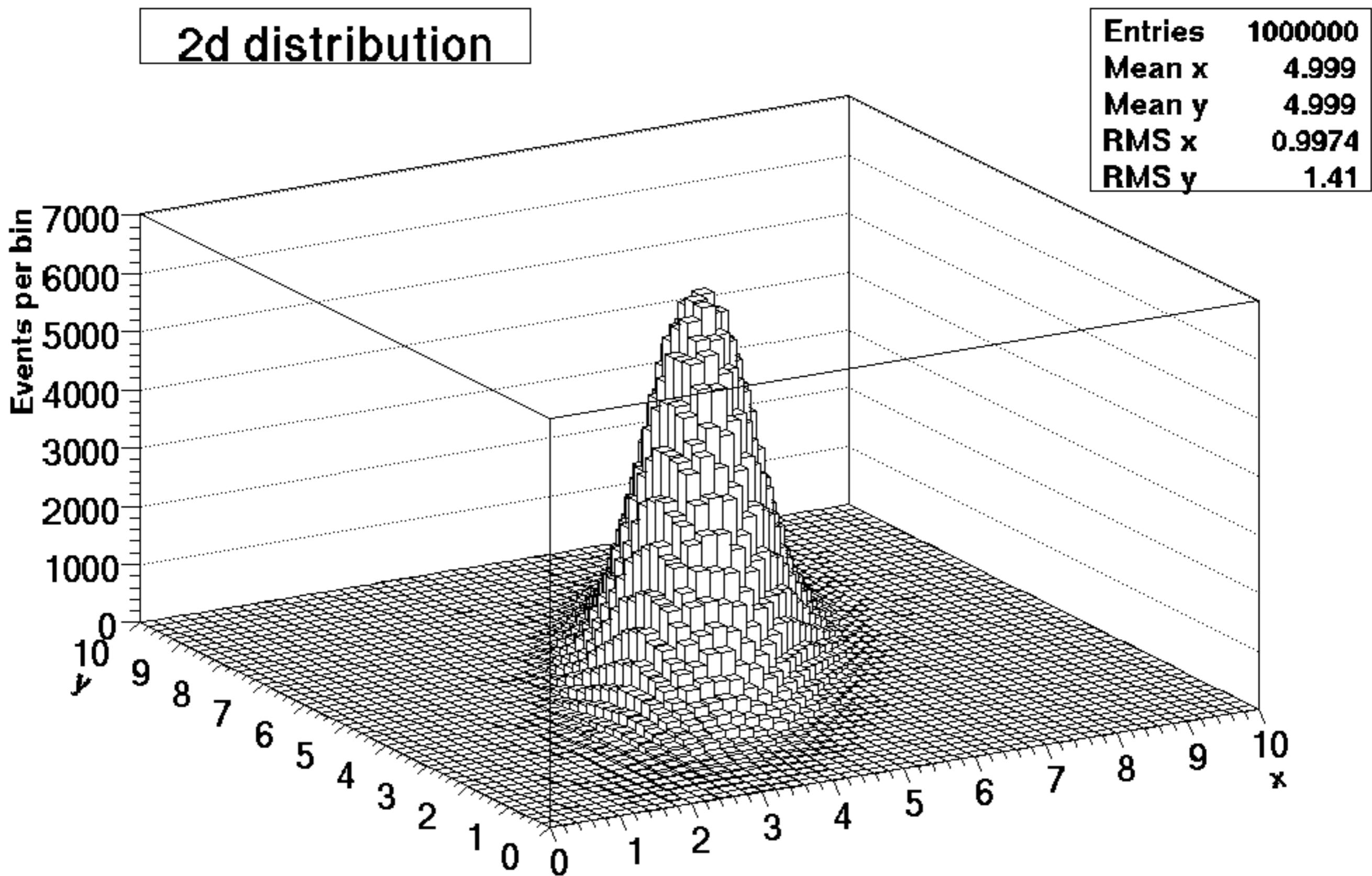
Cumulative Distribution



2d Distribution



2d Distribution



Most important distributions

- or distribution
- distribution
-
-
- distribution
-
- There are many others
 -
 -
 -
 -

Binomial Distribution

- Random process with two possible outcomes

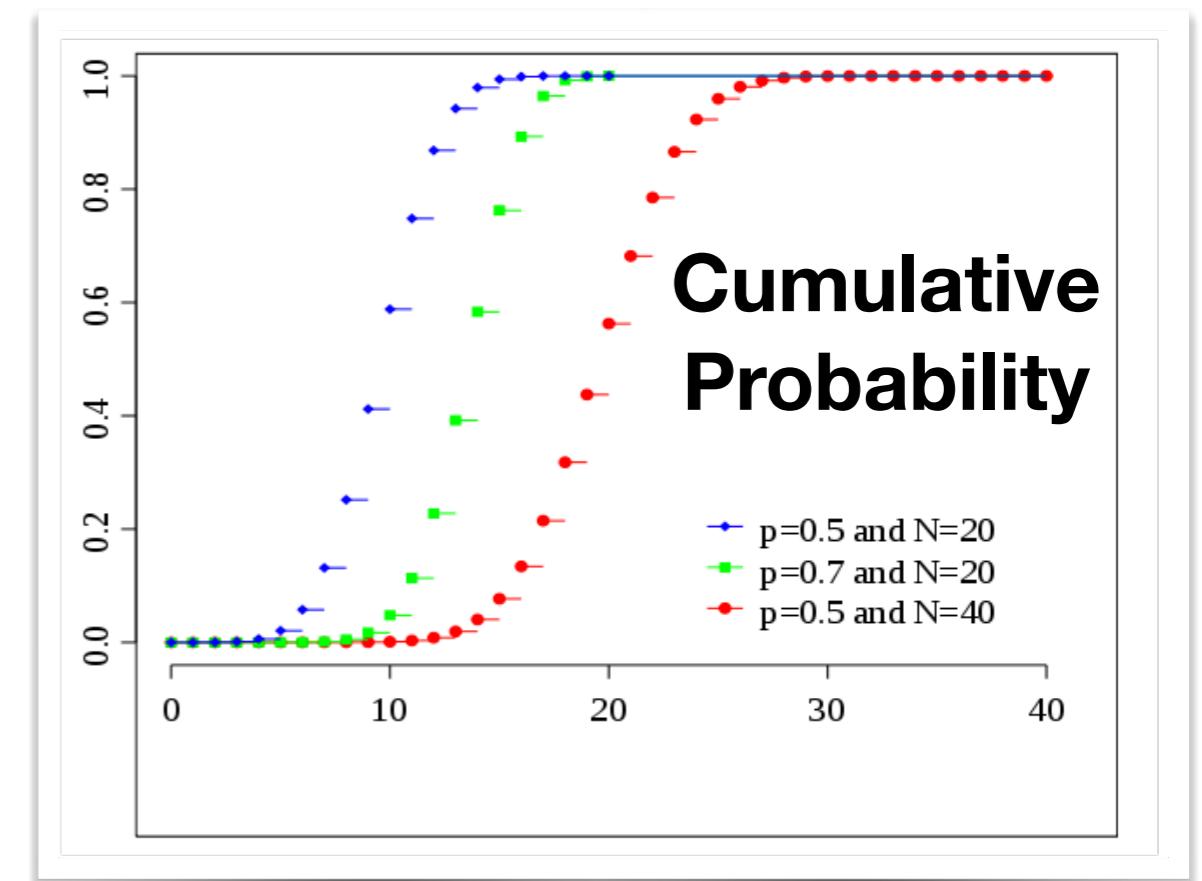
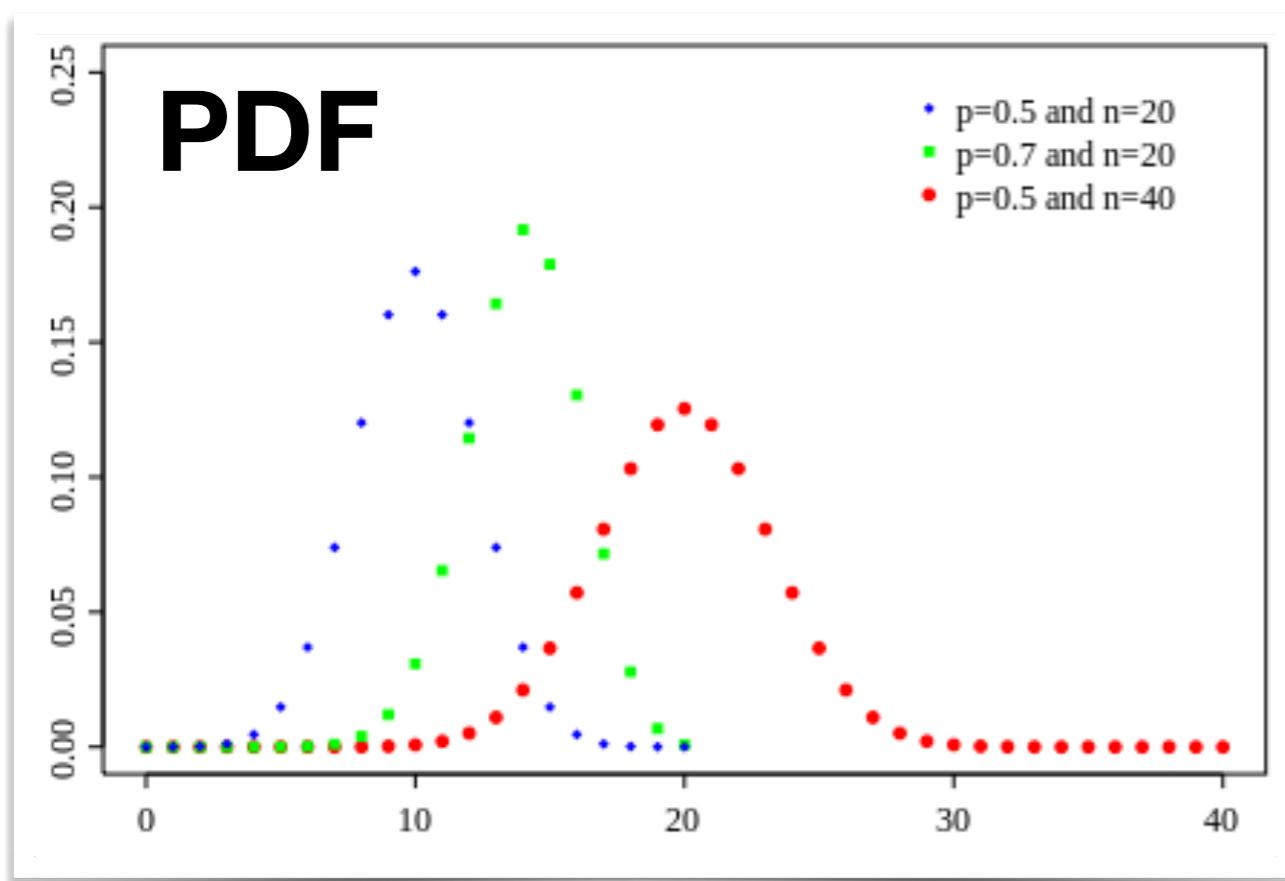
$$\bullet p = \quad q =$$

- After n trials, prob of getting outcome #I exactly k times is

$$f(k, p) = \text{where } \binom{n}{k} =$$



$$\mu = \quad \sigma =$$

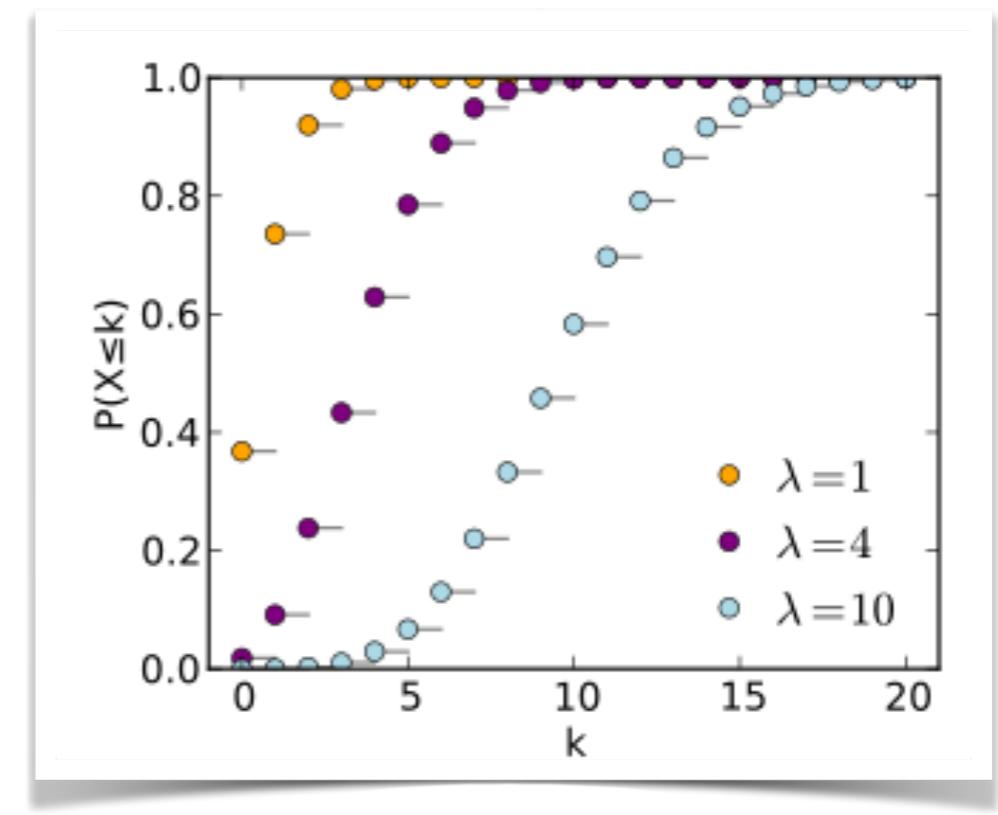
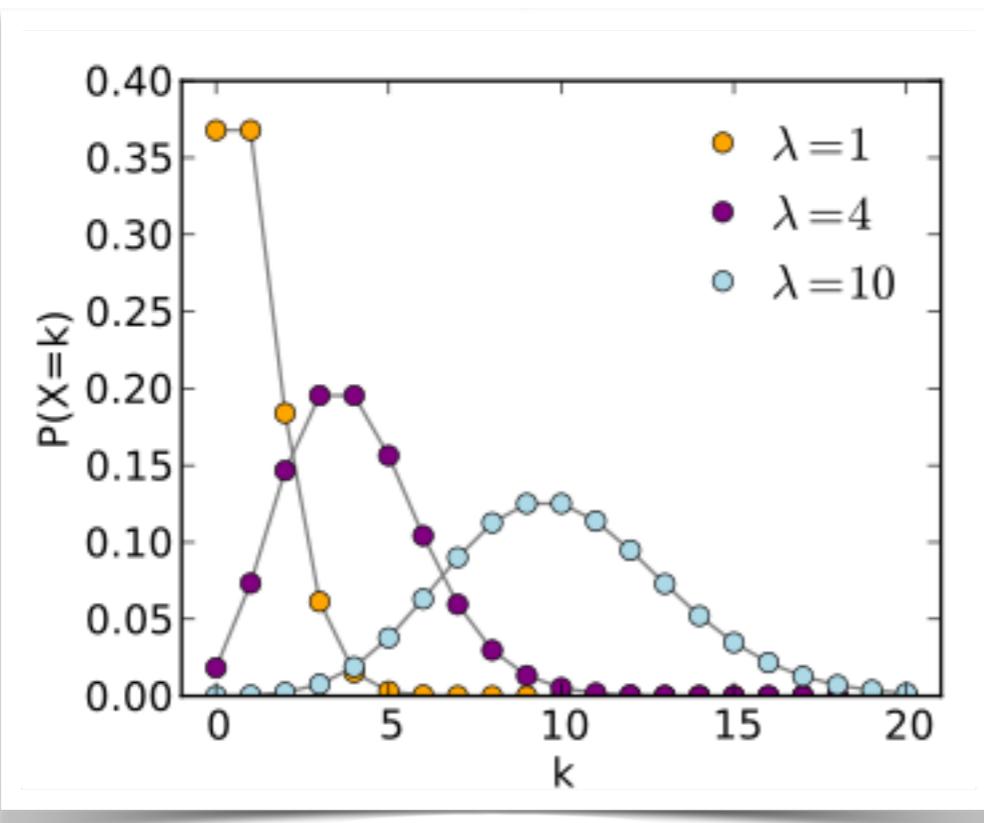


Poisson Distribution

- Prob of finding exactly k events in the interval if the events occur with an average rate λ
- $f(k; \lambda) =$
- For large λ approaches a Gaussian

$$\mu =$$

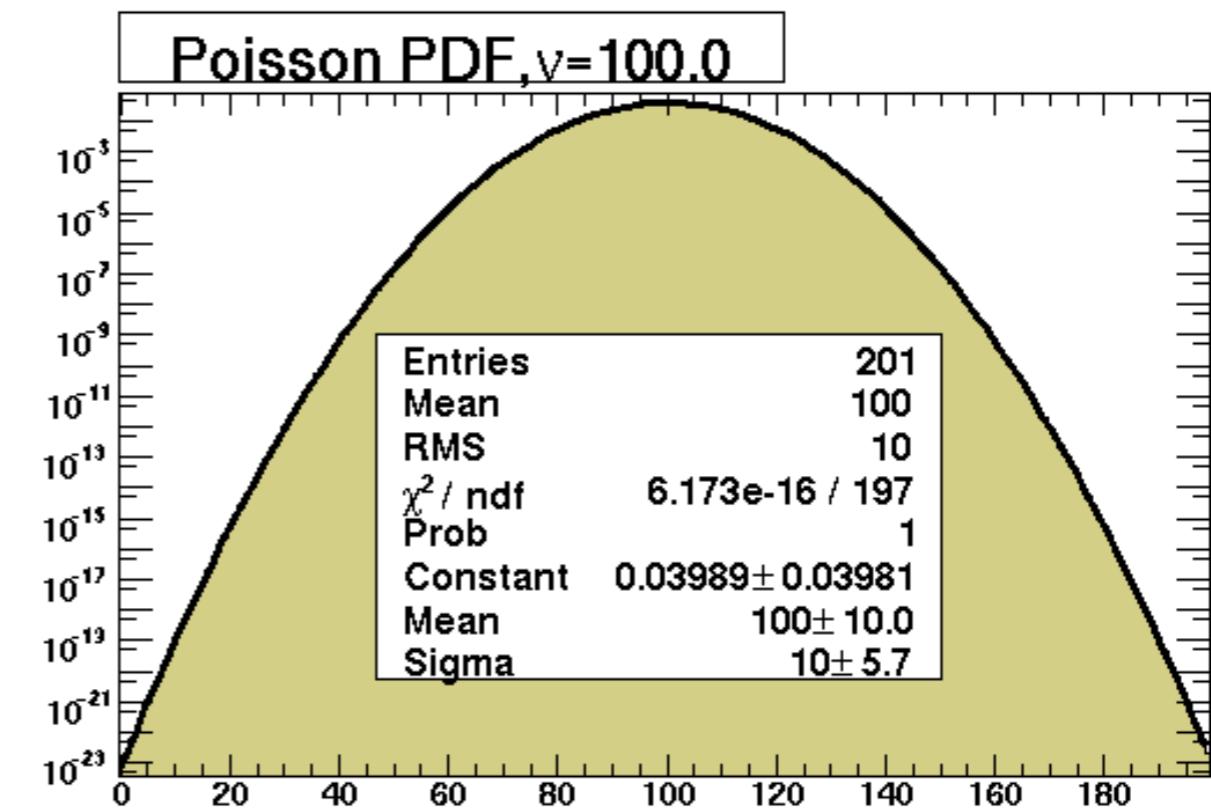
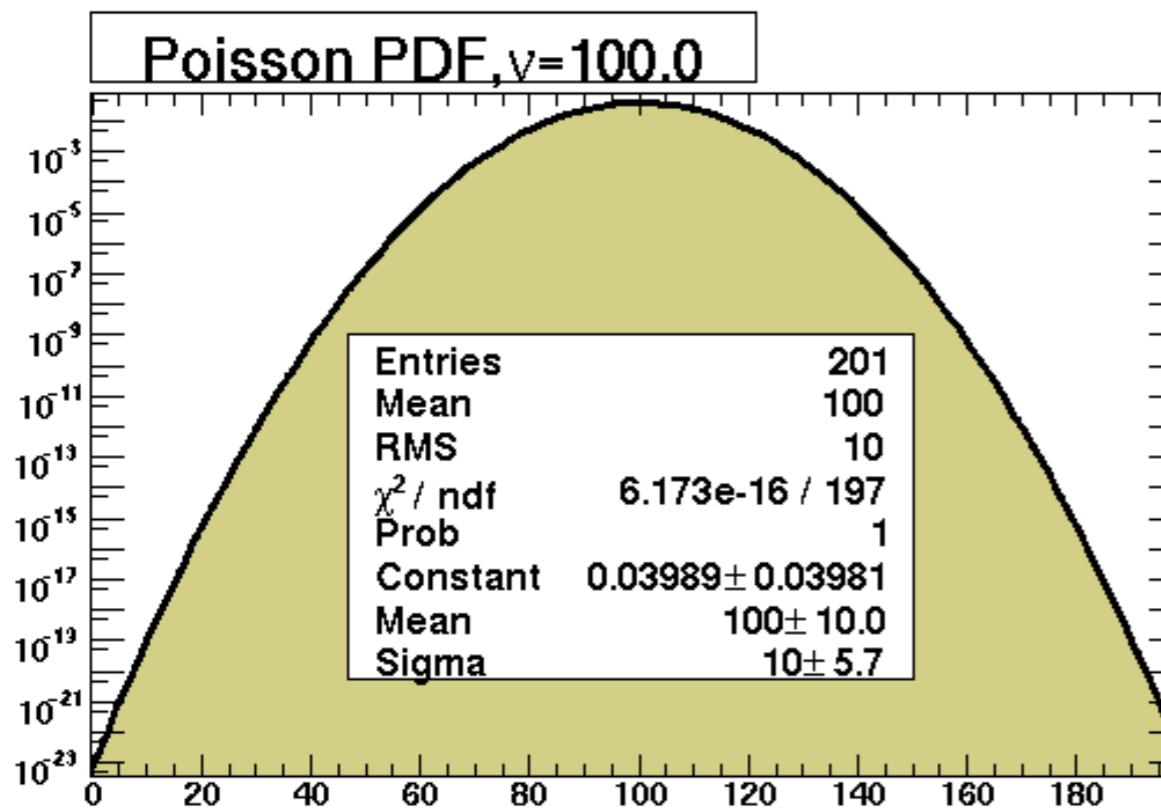
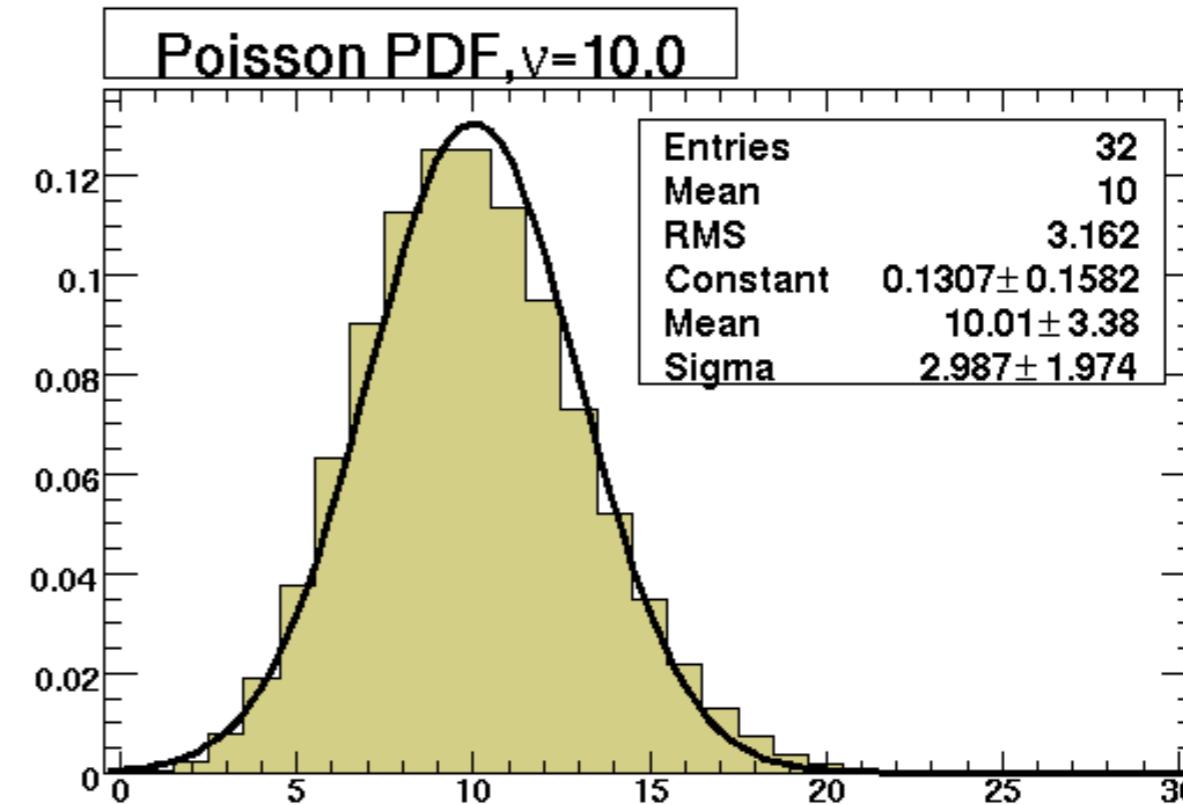
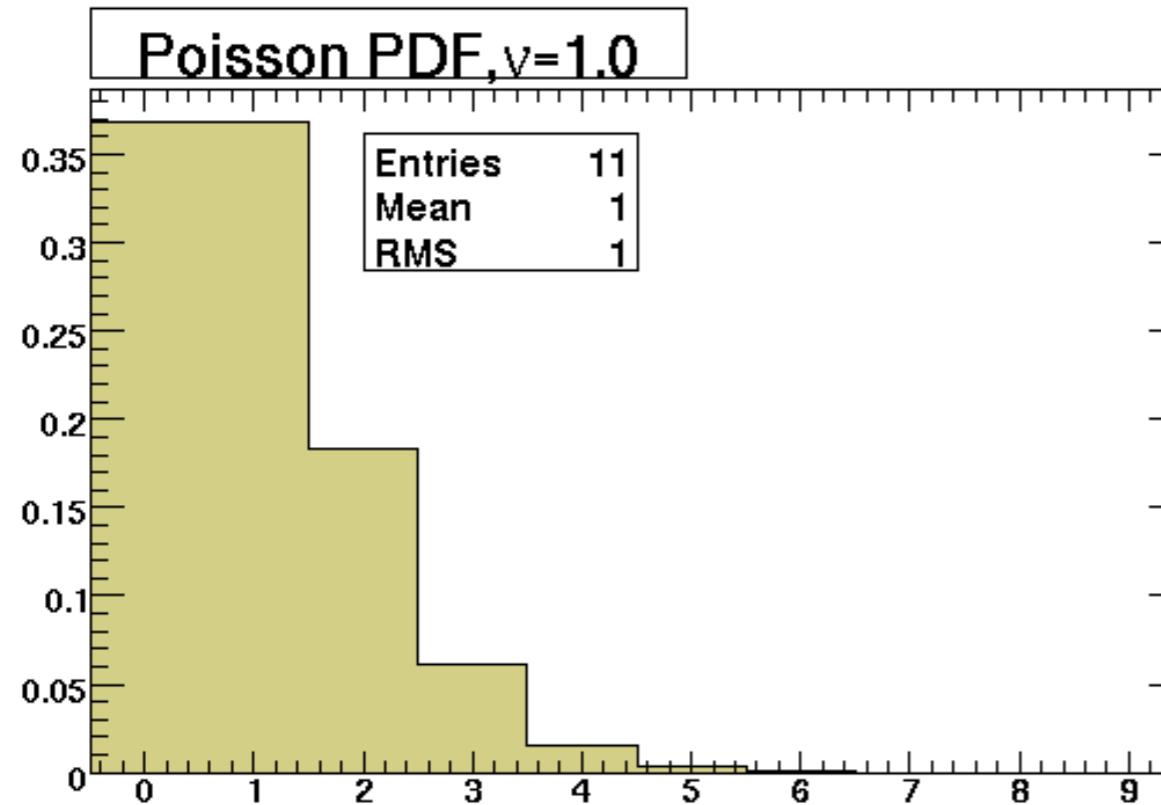
$$\sigma =$$



Example: Measure Efficiency

- Generate a sample of
- Apply selection; suppose events passed
- Estimate
 - $\hat{\epsilon} =$
 - $\sigma(\hat{\epsilon}) =$
 - $\sigma(\hat{\epsilon}) \neq$

Poisson Distribution



Gaussian Distribution

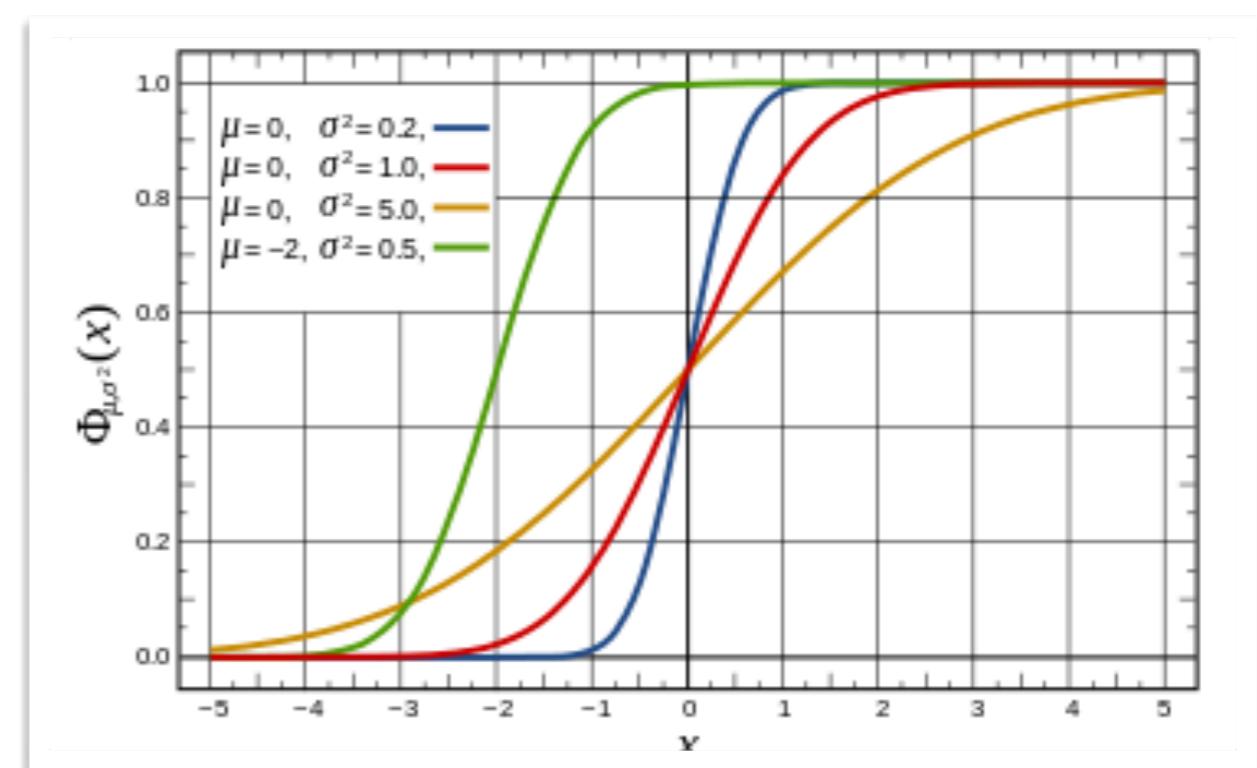
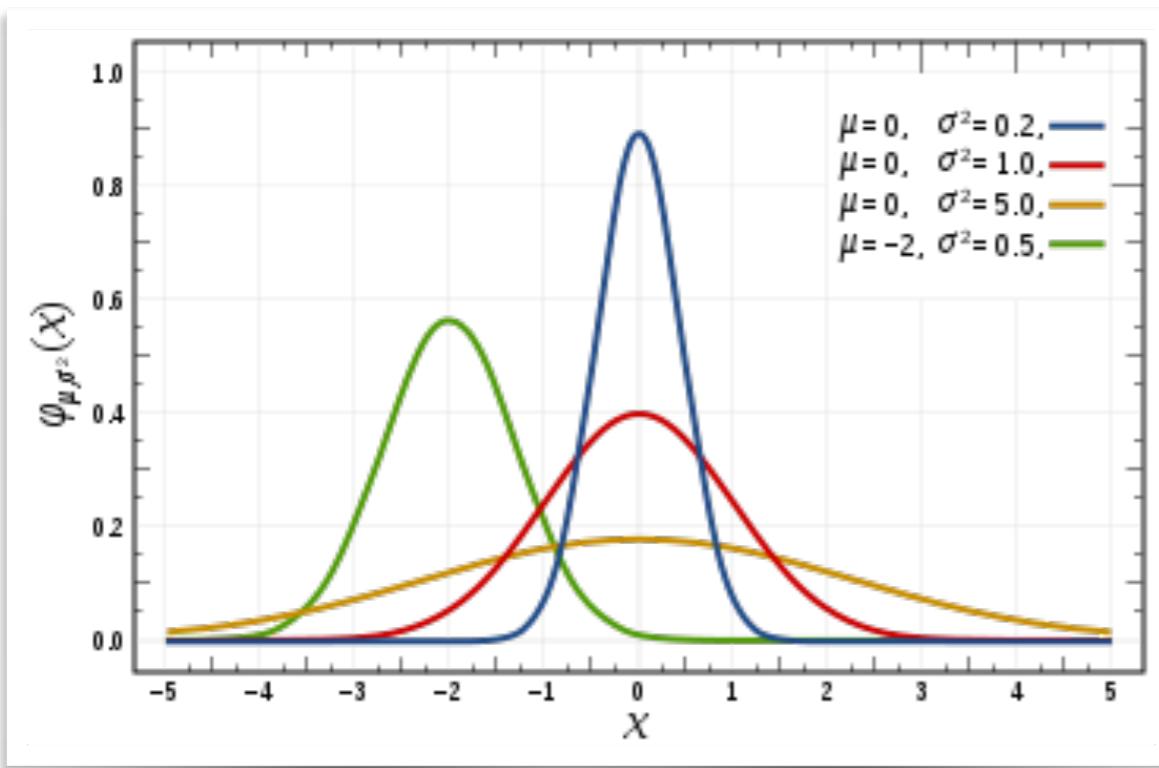
- **Central Limit Theorem**

- 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

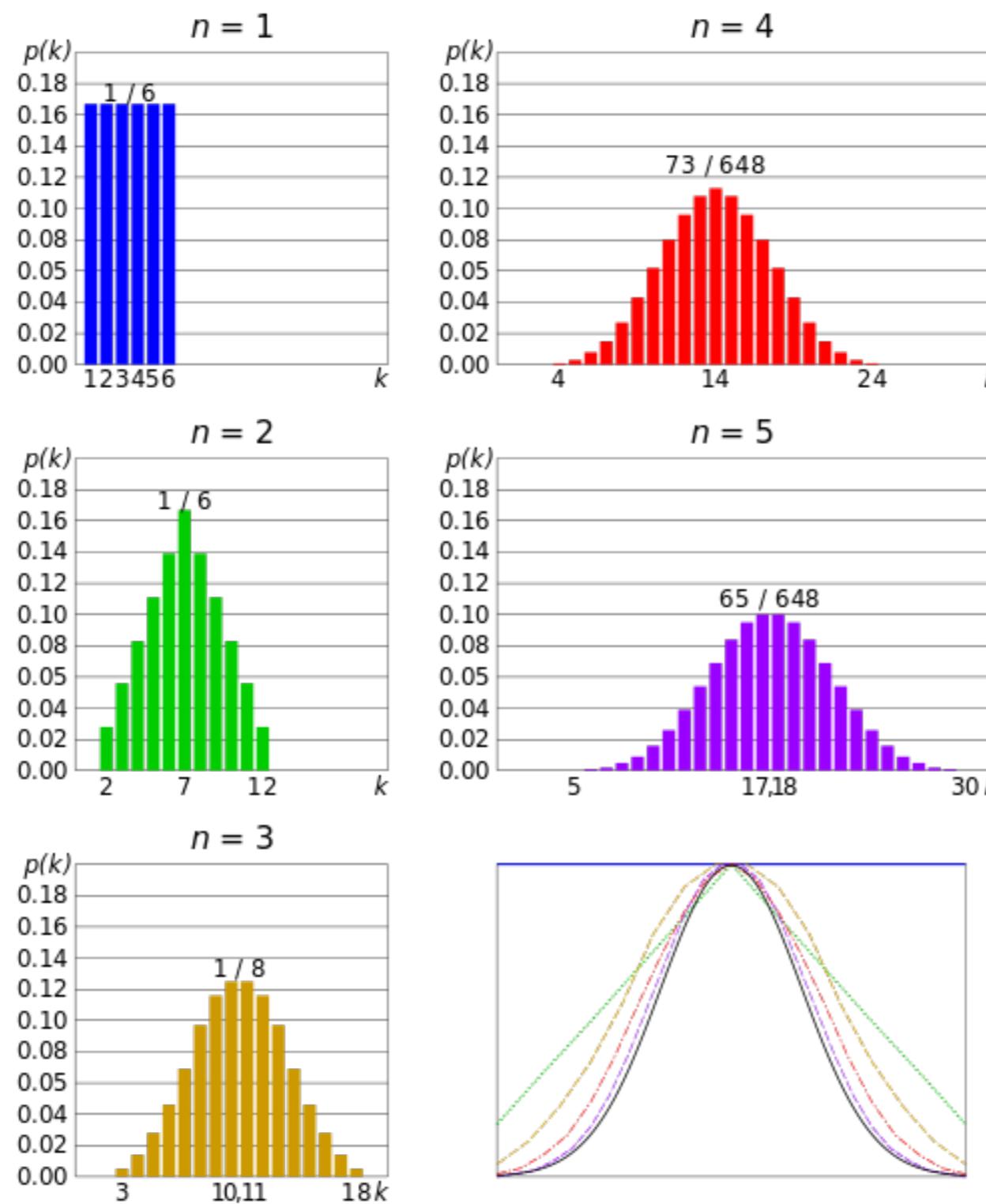
- $S/n =$

- the distribution of S/n approaches the distribution as

- $f(x; \mu, \sigma) =$



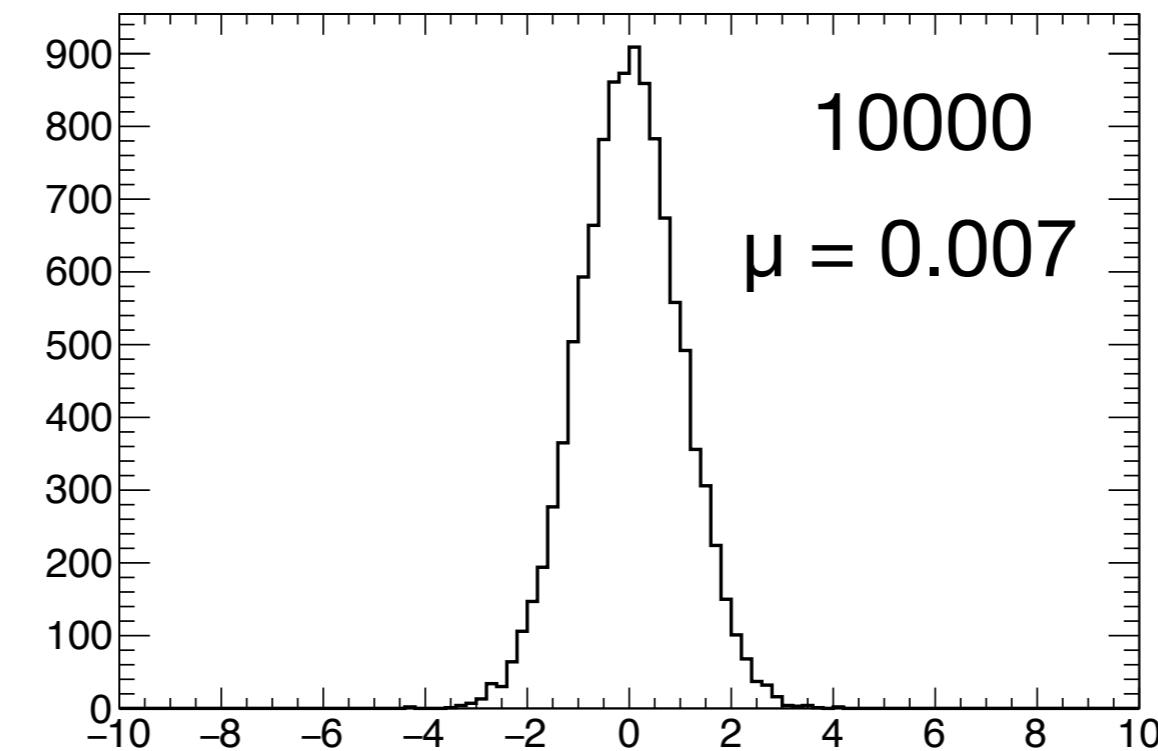
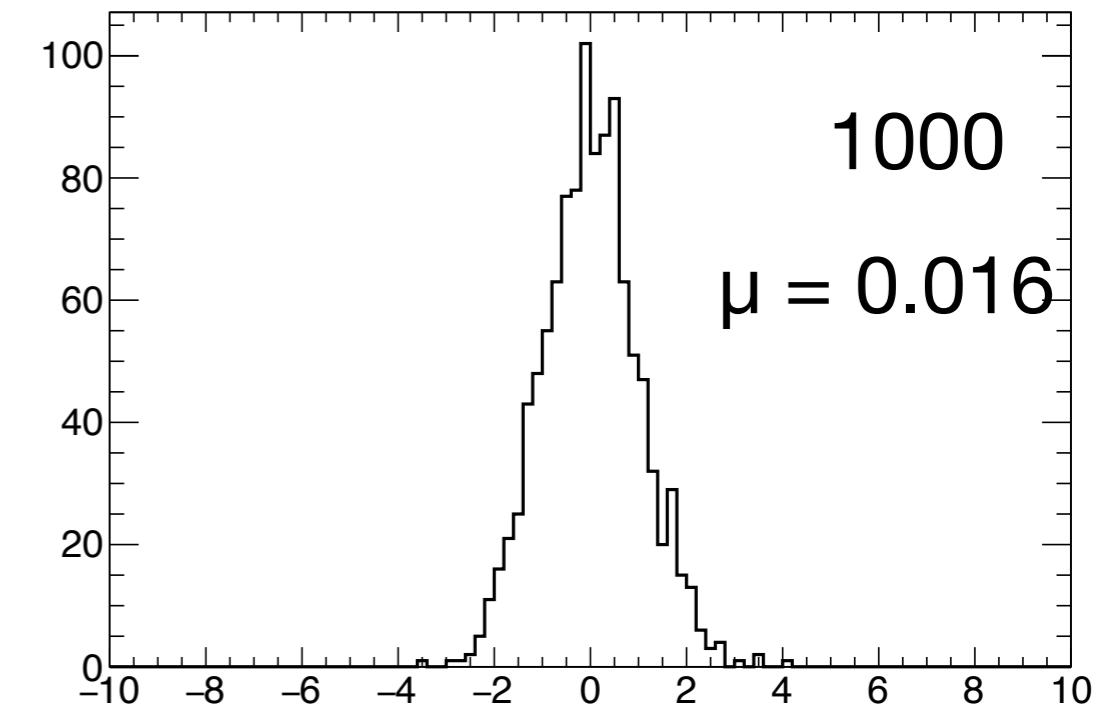
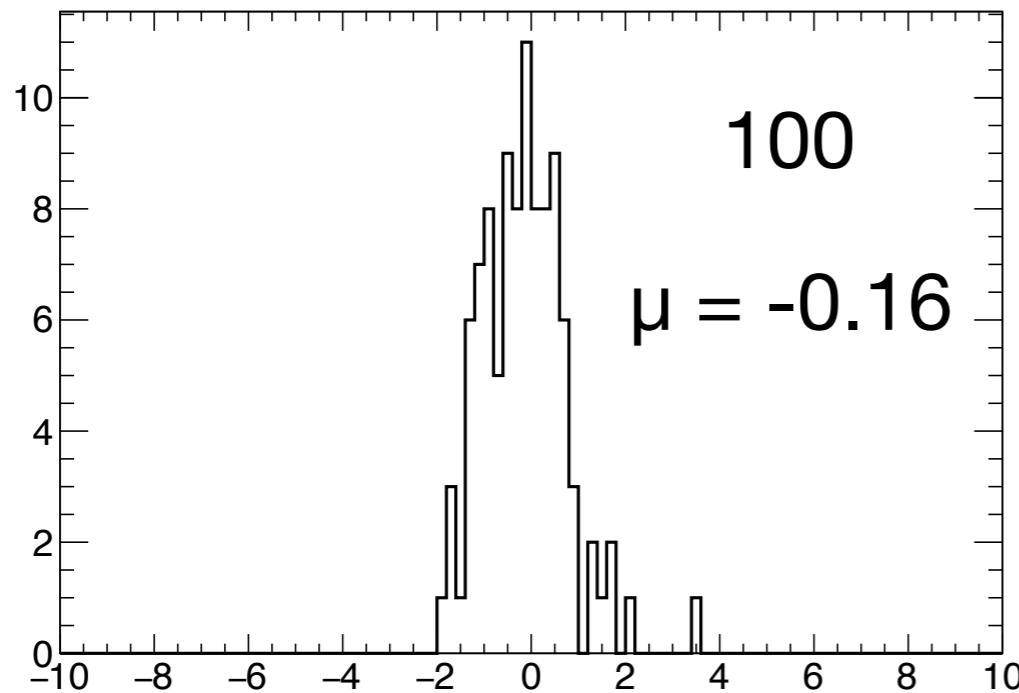
Central Limit Theorem



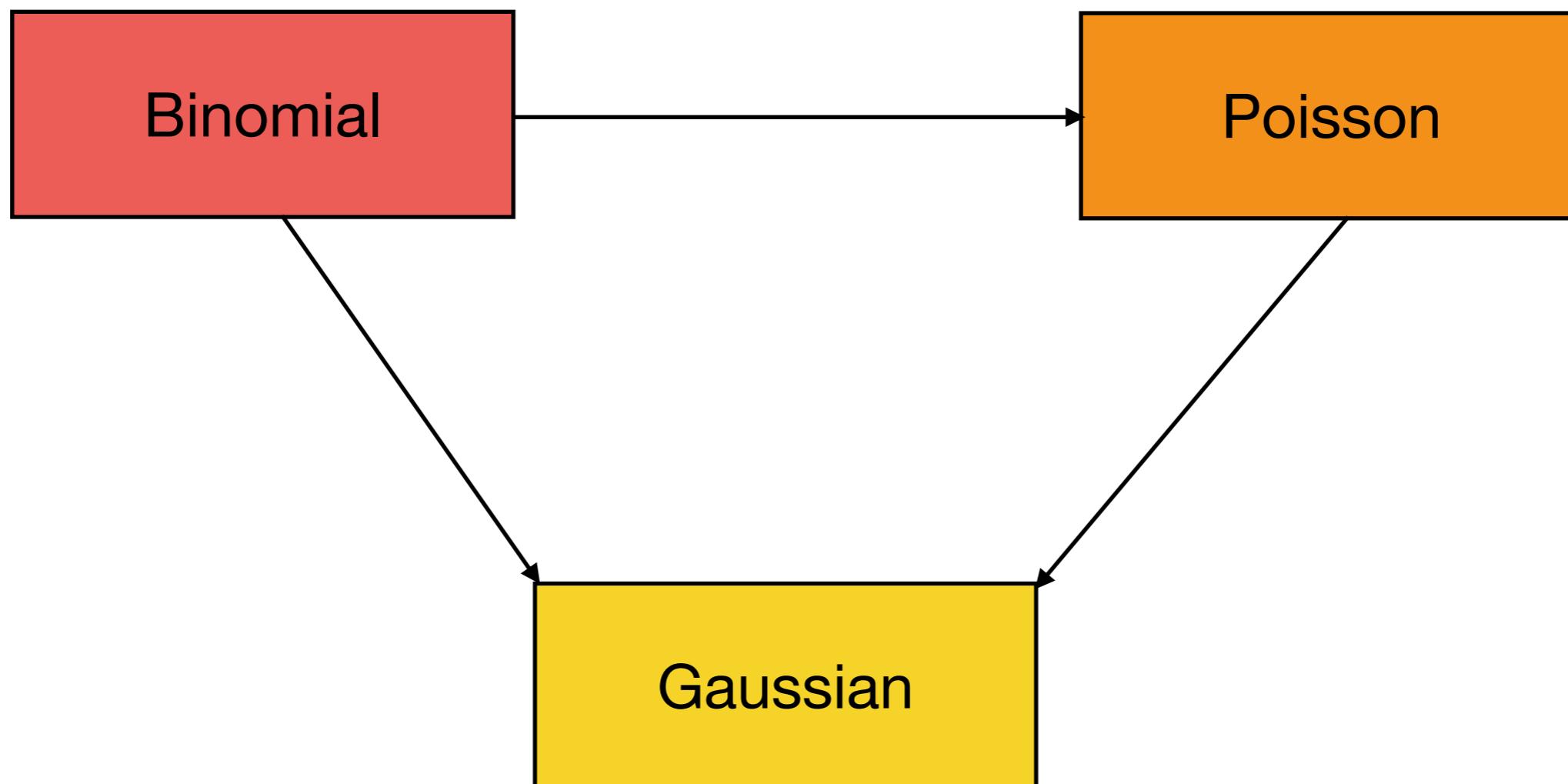
http://en.wikipedia.org/wiki/File:Dice_sum_central_limit_theorem.svg

Example

Example in jupyter notebook



Relationships between pdfs



Point Estimation

- Standard problem: set of described by
 - $f(x) \equiv$
- Point estimation:
 - $\hat{\theta} =$
 - Estimator of

Estimators

- Typical goal: estimate from experimental data and understand the uncertainty on that measurement
- **Characteristics** of an estimator
 -
 -
 -
 -
 -
- **Uncertainty**: how far the might be from our estimate due to statistical fluctuations in the

Basic Estimators

- Estimators for μ and σ^2
- Shape of the distribution
 - Most distributions are normal, but may be skewed or bimodal
 - Distribution of sample mean is not readily available
- Properties of estimators
 - Unbiased estimator for μ is the sample mean, \bar{x} , for data x_1, x_2, \dots, x_n
 - Convenient for estimating linear functions of parameters
 - Automatic measure of spread for data
 - Be careful of outliers
 - (e.g. when $x_i = 1000$ becomes $\bar{x} = 1000$)

Mean and Variance from a Sample

- Estimators (equally data)

- $\hat{\mu} =$

- $\hat{\sigma}^2 =$

- Variances of these

- $V[\hat{\mu}] =$

- $V[\hat{\sigma}^2] =$

- $\sigma[\hat{\sigma}] =$

Likelihood Function

- Likelihood $\mathcal{L}(x; \theta)$:
 - yield a
 - that a
 - for a
- Need to both the and the value for in that theory
- With an ensemble of measurements, overall likelihood is obtained from the of the measurements
- Here θ represents one or more parameters

Log Likelihood

- To estimate the parameter(s), maximise the likelihood
 - Set derivative to zero
 - Typically easier to maximise the
- $\frac{\partial \mathcal{L}}{\partial \theta} =$
- If there are several we can minimise with respect to each of them

Likelihood Example: Poisson

- independent trials with results
- Likelihood function for observing if true mean is
 - $\mathcal{L}(n_i; \mu) =$
- Product over N measurements
 - $\mathcal{L}(\text{data}, \mu) =$
 - $\Rightarrow \ln \mathcal{L} =$

Best estimator is
the mean value

Likelihood Example: Gaussian

- $G(x | \mu, \sigma) =$
- Take the derivative of the log likelihood

$$\bullet \frac{\partial}{\partial \mu} (\ln \mathcal{L})|_{\hat{\mu}=\mu} =$$

- The unbiased estimator for σ is
- $\hat{\sigma} =$

Binned vs unbinned likelihood functions

- Likelihood formalism works for any
 - Product of the $f(t)$ is a
 - **Example measurement:** Measure the τ of a particle of a given species for an N of such particles produced at t_0 such that the $f(t)$ at time t :
$$\bullet f(t) = \frac{1}{\tau} e^{-t/\tau}$$
 - Consider two ways to construct the likelihood
 - For decay τ measure t_i and take the $\ln(t_i/\tau)$ of all times t_i ()
 - Make a histogram of the $\ln(t_i/\tau)$ in bins of width $\Delta \ln(t_i/\tau)$
 - Measurement is the $\ln(t_i/\tau)$ in each bin ()

The Likelihood and χ^2

- If the data is Gaussian, we have

- $\ln \mathcal{L} =$

- Compare to

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

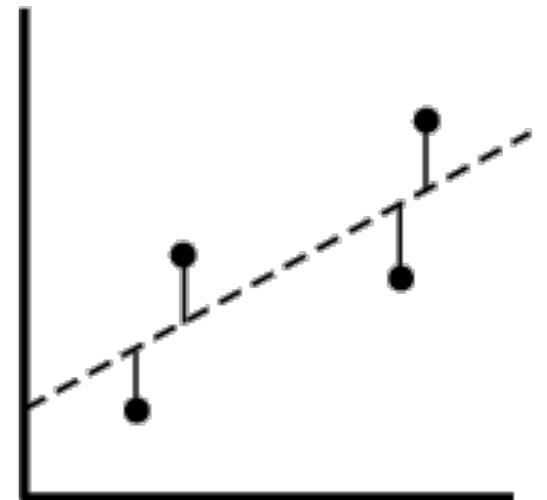
- By inspection

- $\chi^2 =$

- Note: the likelihood formulation works for not just

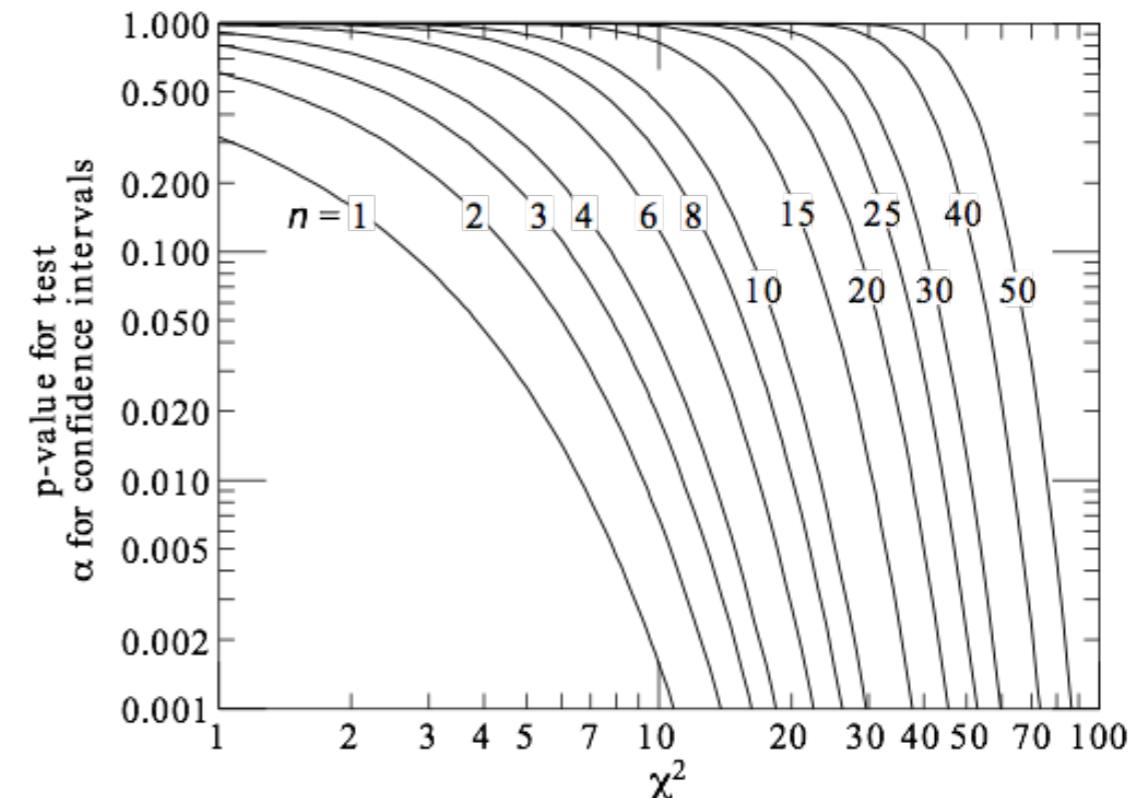
Method of Least Squares

- Assume we have enough statistics for our measurements such that we can assume we are in the Gaussian regime
- Goal: Find the parameters of a model that describes the data
- Minimize the difference of the observed data from the predicted values
 - Account for the scatter
 - Scatter defined by χ^2
 - $\chi^2 = \sum (y_i - \hat{y}_i)^2$
 - Can write the χ^2 in terms of our observables
 - $\chi^2 = \sum (y_i - \hat{y}_i)^2$
- Minimise χ^2 with respect to θ
- Useful when minimising $\ln \mathcal{L}$ is slow (high statistics samples)



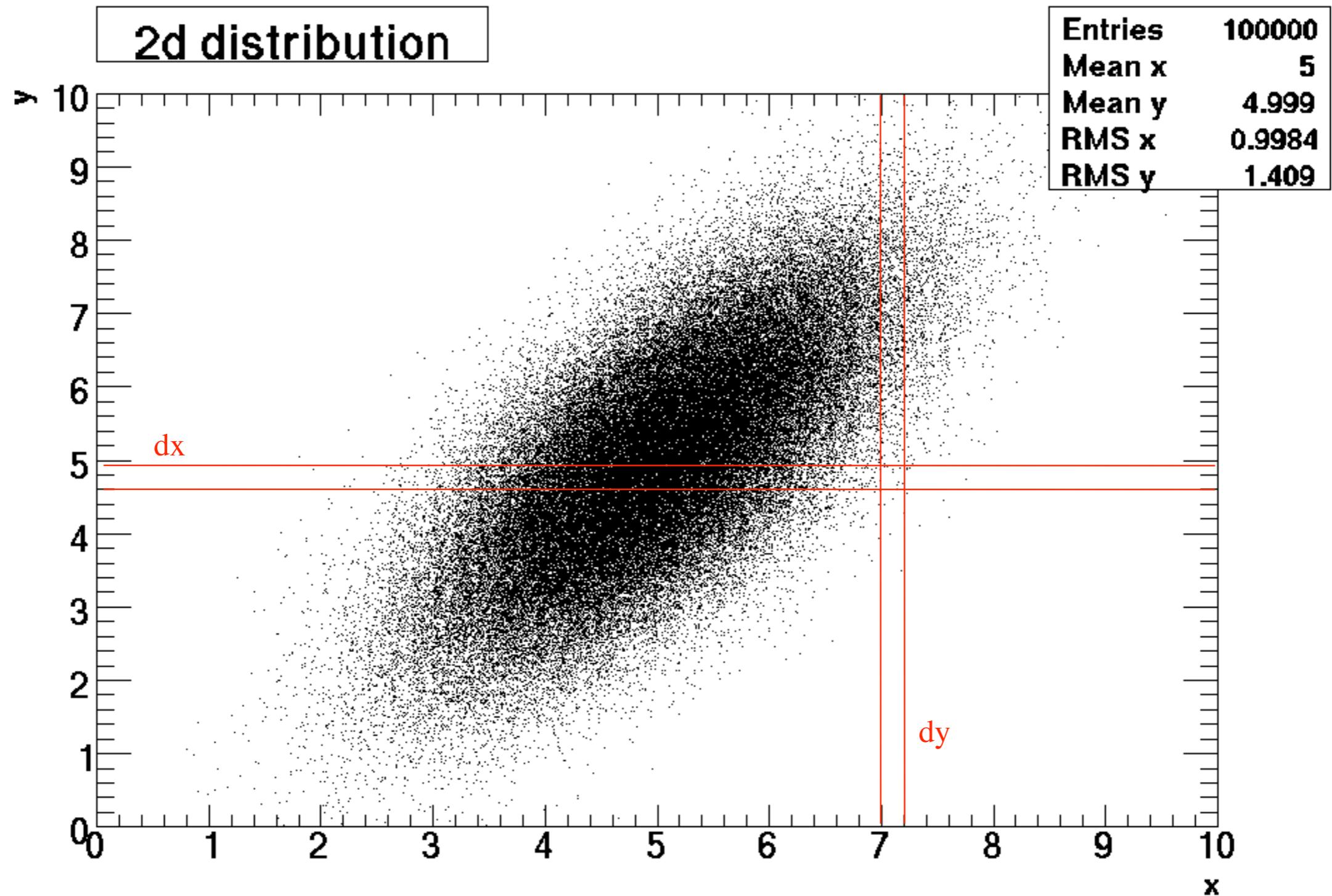
Example: chi-squared p-values

- One advantage of a is that the value of the can be
interpreted as a and
- iff on each data point are
the around their
are)
- In the plot below
- $n =$
- For a , expect
to be close to

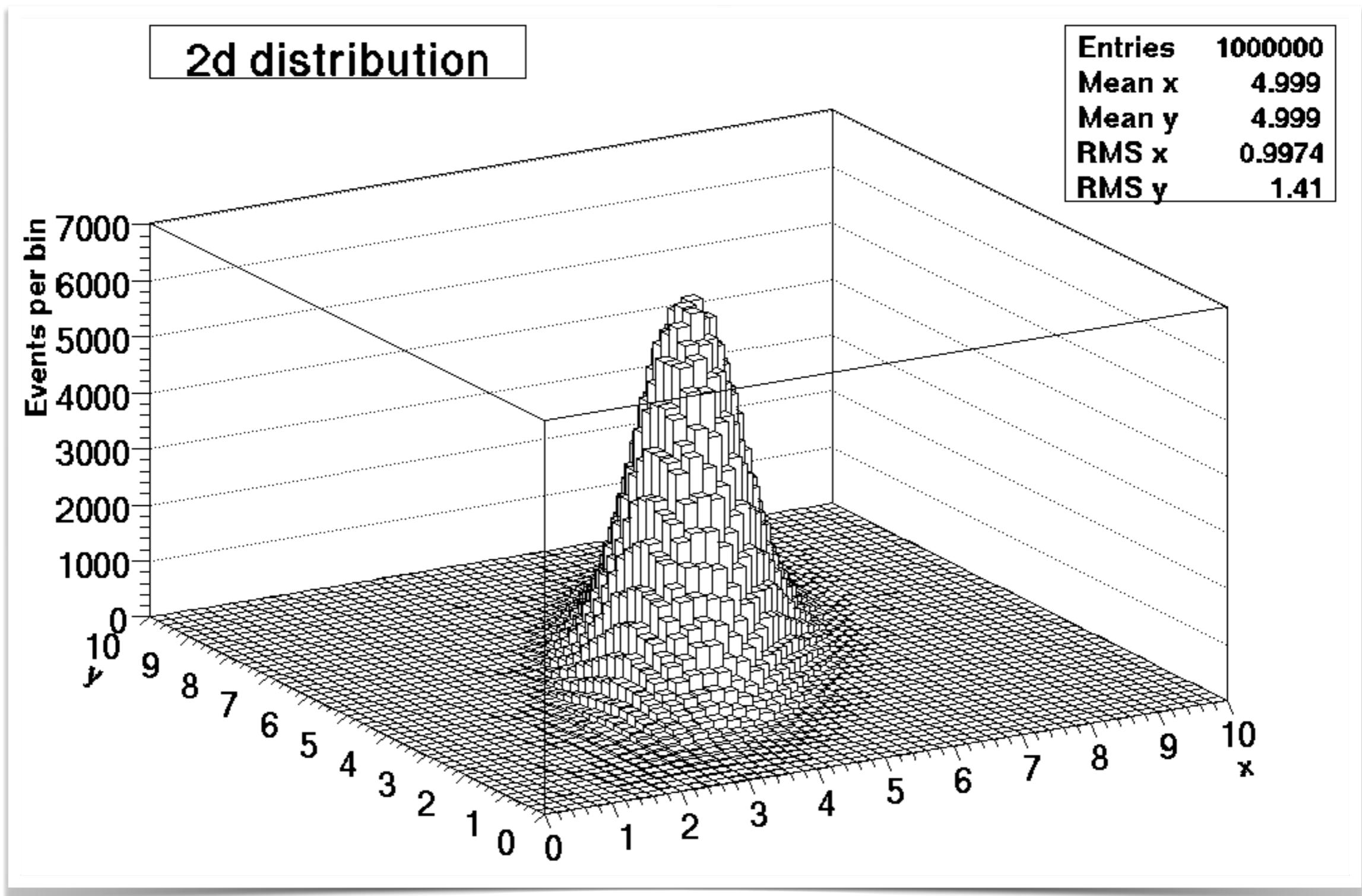


Example in jupyter notebook

2D distribution



2d distribution



Covariance and Correlation

Covariance Matrix

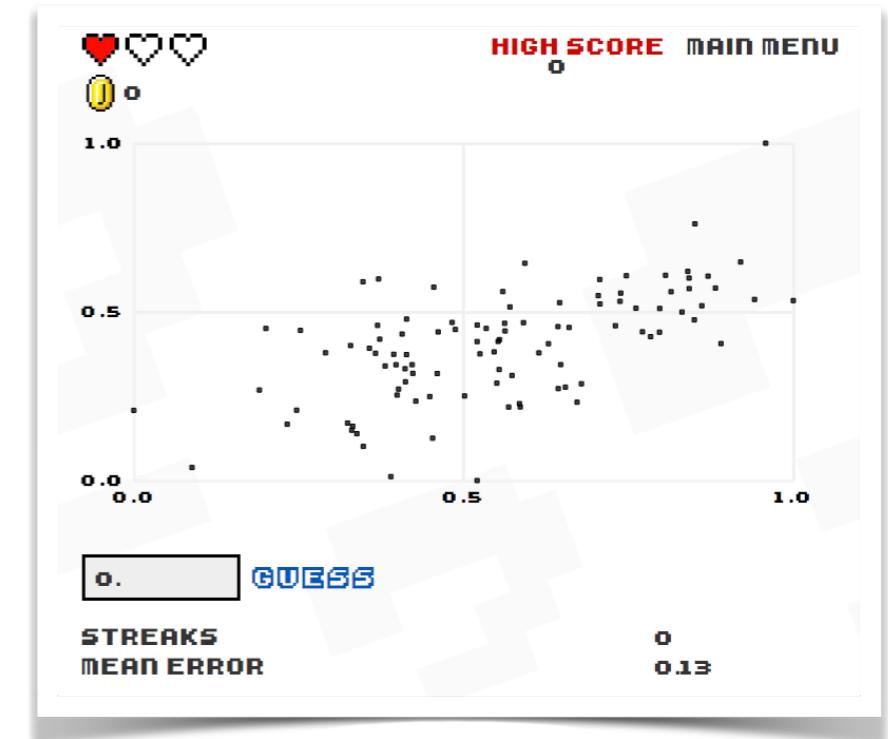
$$\text{cov}[x, y] =$$

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

Correlation (normalised covariance)

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

$$\rho_{xy} =$$



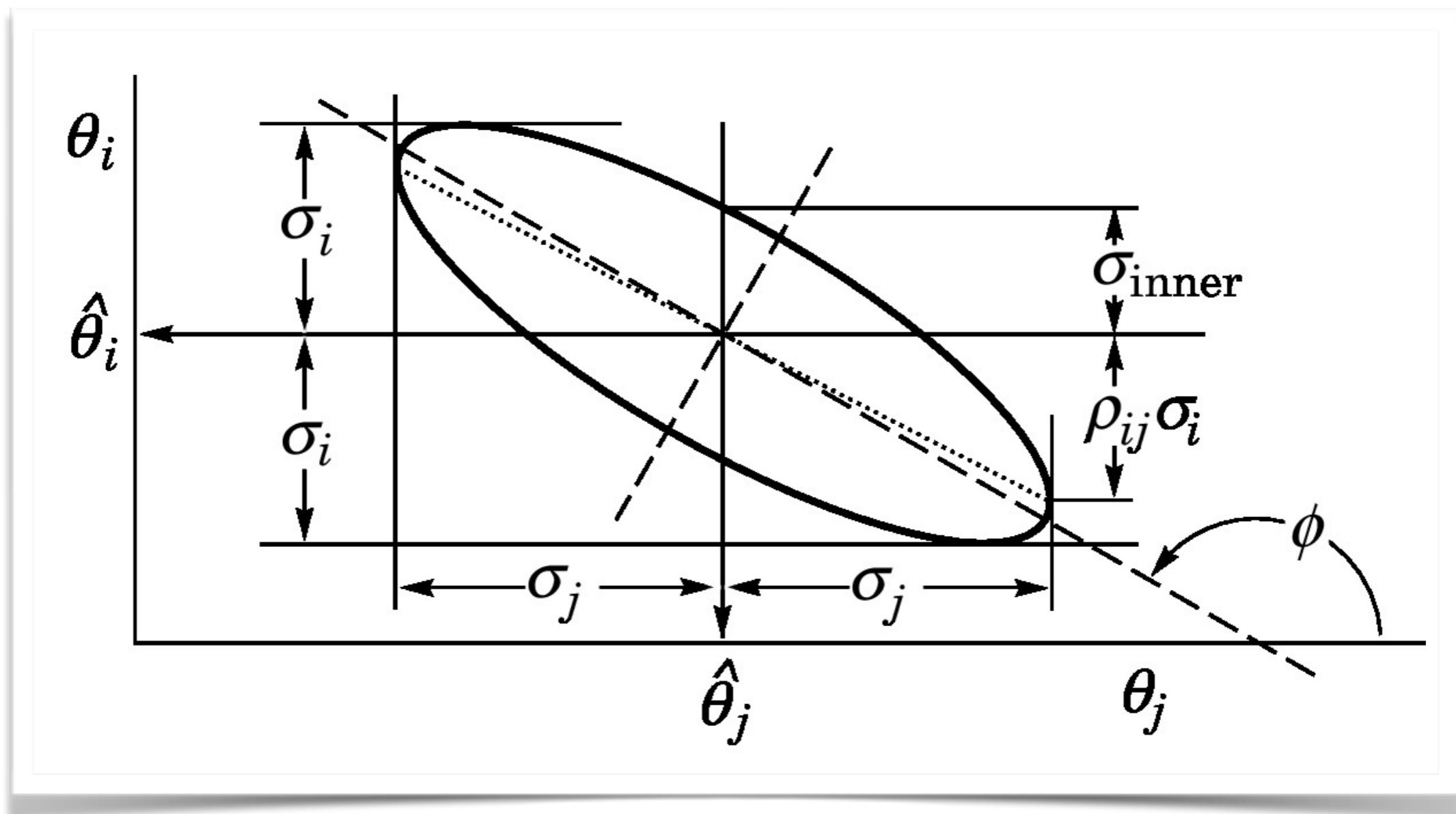
<http://guessthecorrelation.com/>

Covariance Matrix for a Gaussian

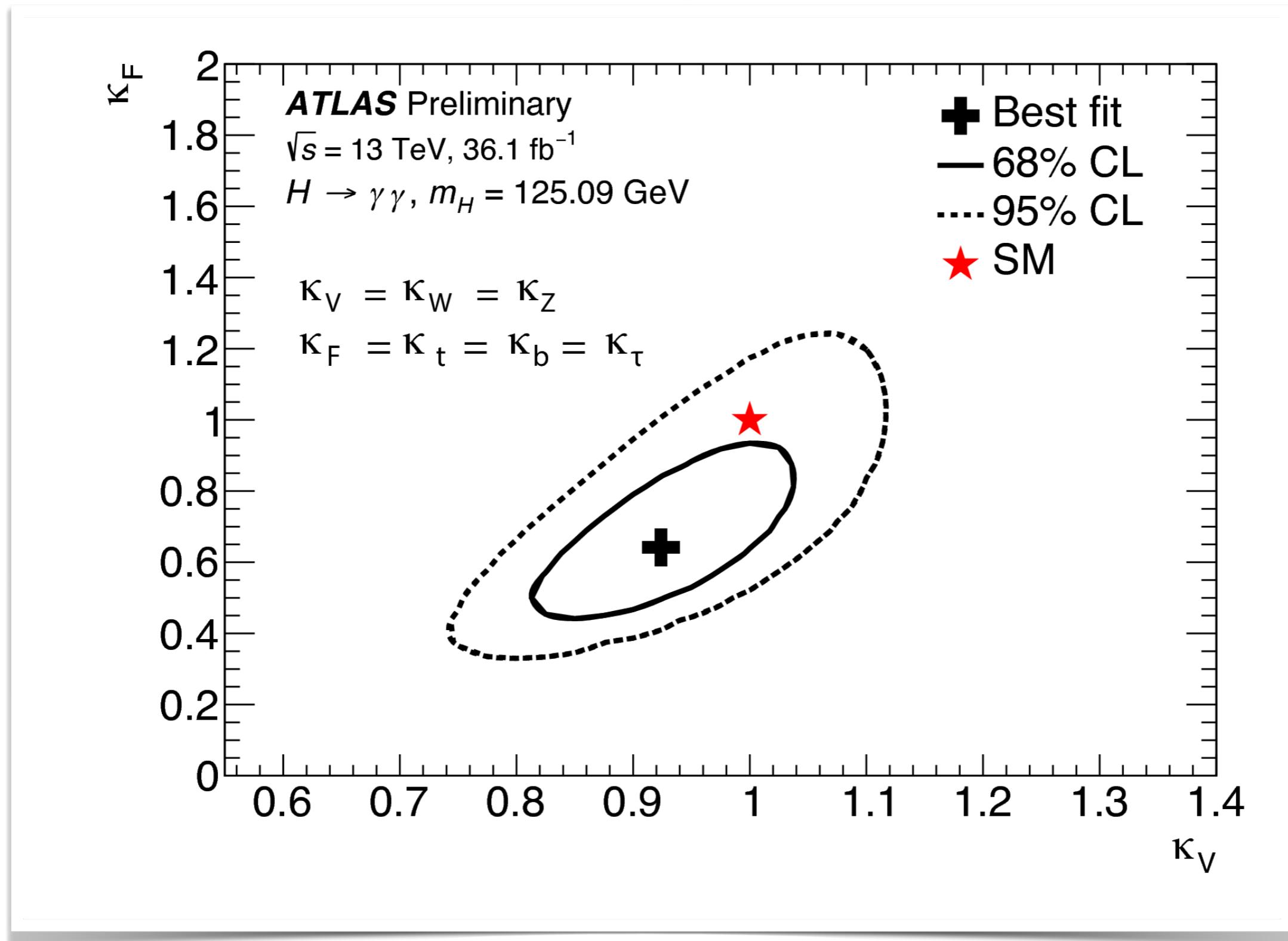
- If x and y are independent variables
 - $G(x, y | \mu_x, \sigma_x, \mu_y, \sigma_y) =$
 -
- Now, assume that x and y are correlated
 - Covariance matrix is defined by
 - $\langle \hat{V}^{-1} \rangle_{ij} =$
 - For a binned likelihood, where N is large and the likelihood can be reduced to a χ^2
 - $\langle \hat{V}^{-1} \rangle_{ij} =$

Correlated Uncertainties

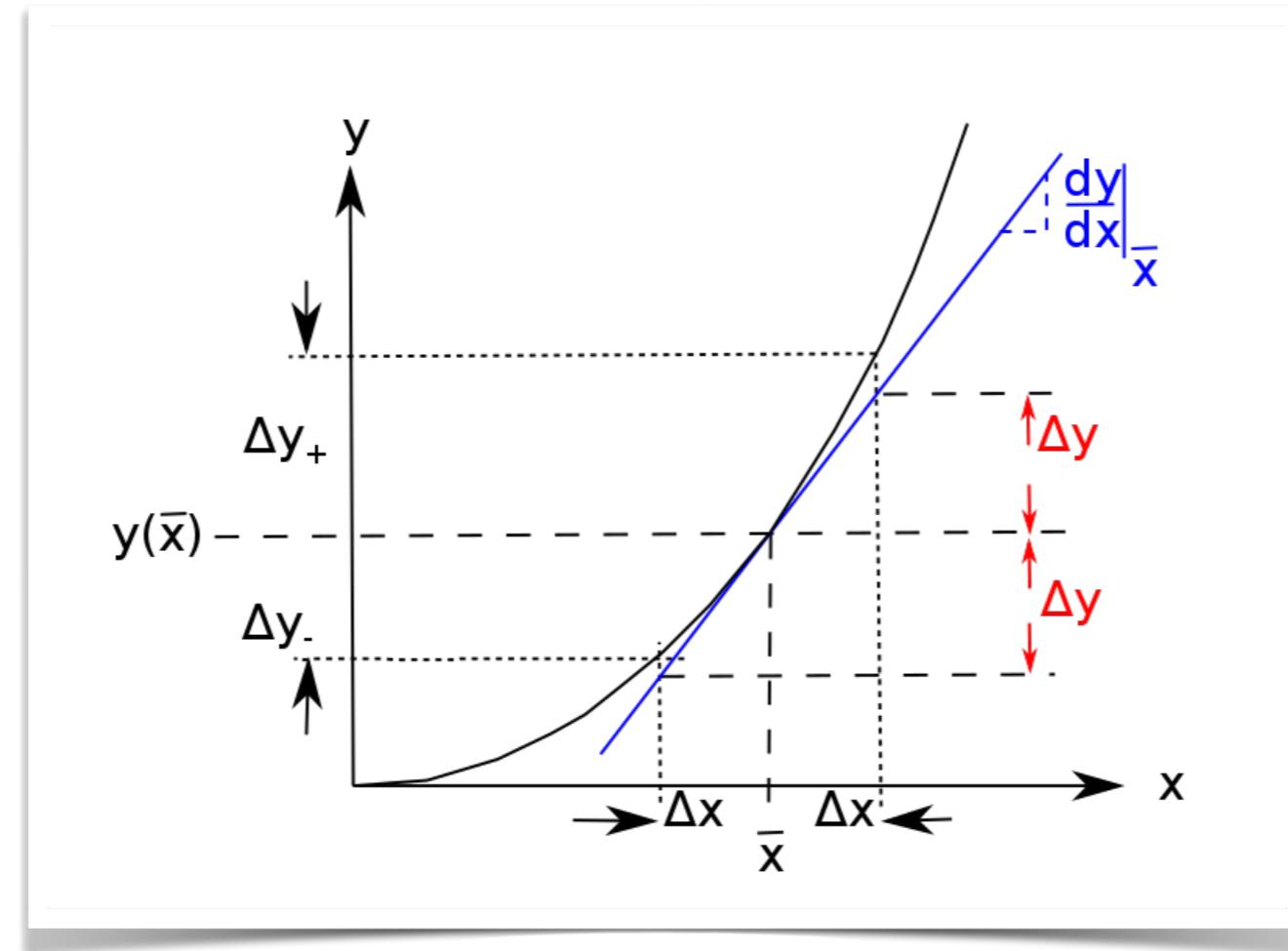
- Standard for two parameters with a correlation
- Slope related to correlation coefficient
- Correlation matrix typically determined from data numerically during fitting procedure



Correlation Example from Higgs



Propagation of Errors

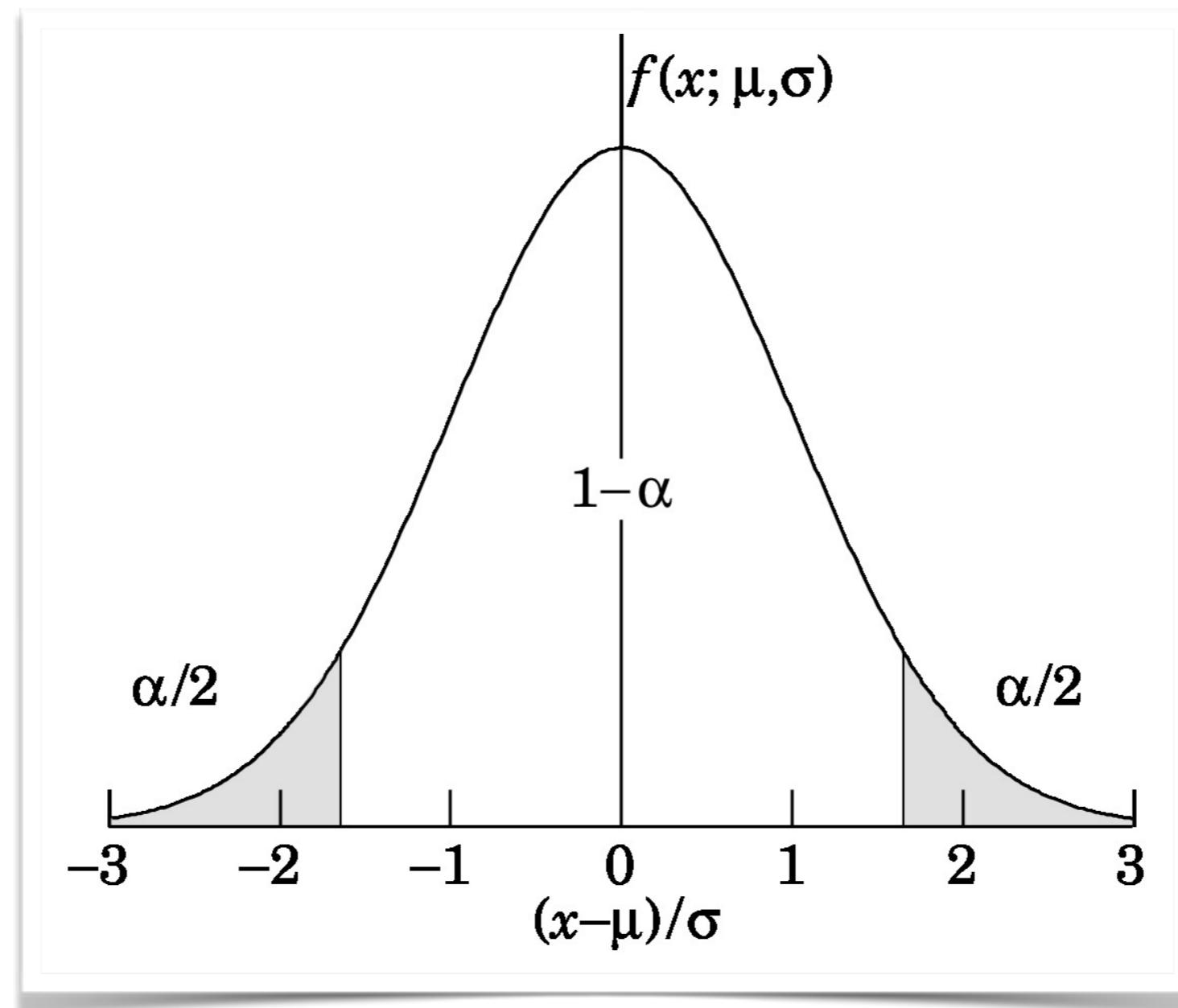


- Determine error on final measurement from known errors on input measurements
 - $\sigma_f^2 =$
- More dimensions are usually expressed as a matrix
- Useful reference: https://en.wikipedia.org/wiki/Propagation_of_uncertainty

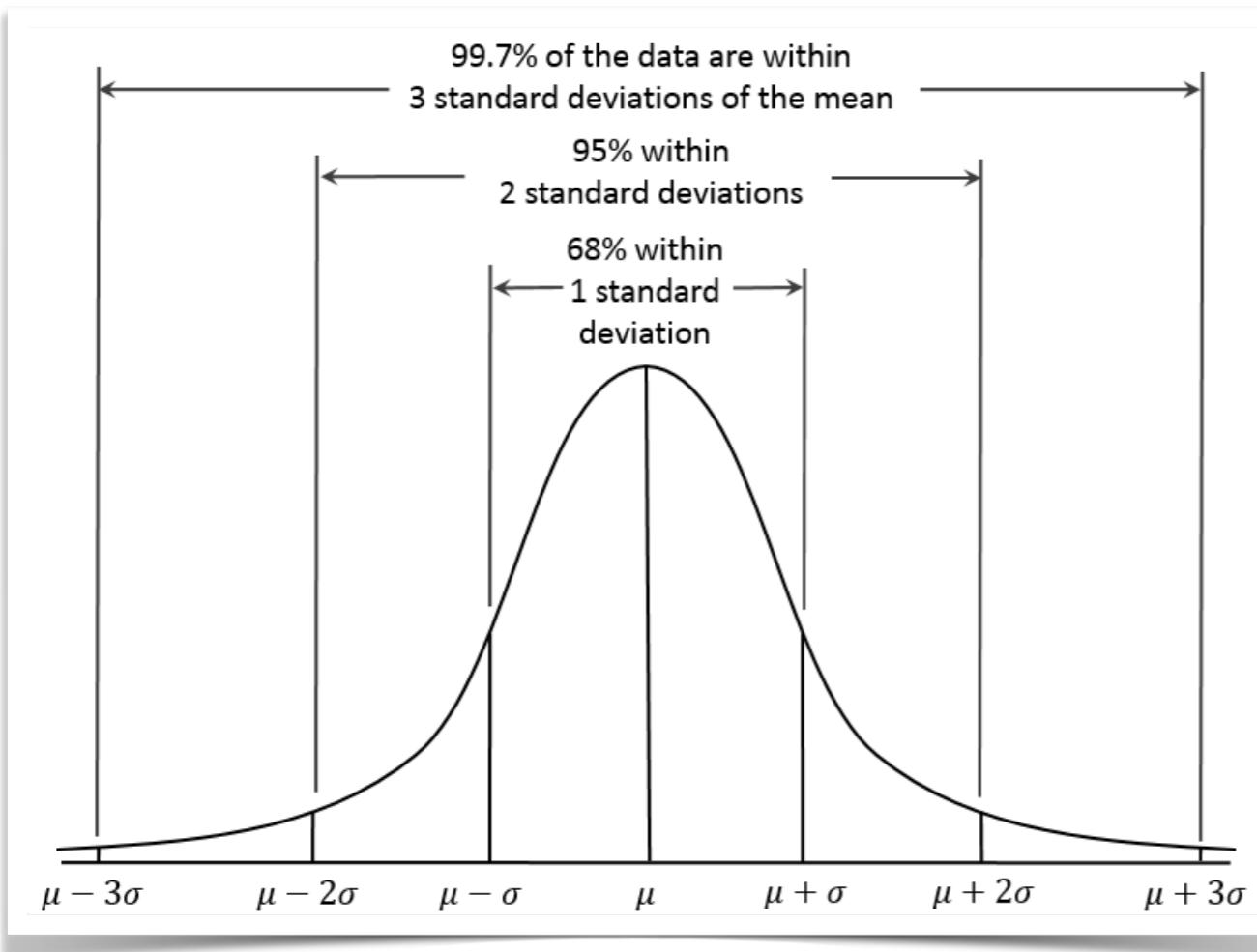
Confidence Interval

- Fraction of the result not between $\mu - z_{\alpha/2}$ and $\mu + z_{\alpha/2}$ is

- $1 - \alpha = \int_{-\infty}^{\infty} f(x; \mu, \sigma) dx$



Confidence Levels for a Gaussian



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

Confidence Levels: Higgs

