



Foundations of statistics (1)

Andreas Hoecker (CERN)

CERN Summer Student Lecture, 26–29 July 2018

If you have questions, please do not hesitate to contact me: andreas.hoecker@cern.ch



Outline (4 lectures)

Slides with title: *Digression — ...*
→ not discussed here, for later reading

1st lecture:

- Introduction
- Probability

2nd lecture:

- Probability axioms and hypothesis testing
- Parameter estimation
- Confidence levels

3rd lecture:

- Maximum likelihood fits
- Monte Carlo methods
- Data unfolding

4th lecture:

- Multivariate techniques and machine learning



Outline (4 lectures)

1st lecture:

- Introduction
- Probability



...a bit dry

2nd lecture:

- Probability axioms and hypothesis testing
- Parameter estimation
- Confidence levels

3rd lecture:

- Maximum likelihood fits
- Monte Carlo methods
- Data unfolding

4th lecture:

- Multivariate techniques and machine learning

Acknowledgements and some further reading

- I am grateful to Helge Voss, whose comprehensive statistics lectures serve as model for the present lectures. His latest lectures on multivariate analysis and neural networks at the School of Statistics, 2016: <https://indico.in2p3.fr/event/12667/other-view?view=standard>
- Glen Cowan's book "Statistical data analysis" represents a great introduction:
<http://www.pp.rhul.ac.uk/~cowan/sda>
- PDG reviews by G. Cowan on probability (<http://pdg.lbl.gov/2013/reviews/rpp2013-rev-probability.pdf>) and statistics (<http://pdg.lbl.gov/2013/reviews/rpp2013-rev-statistics.pdf>) :
- Luca Lista's lectures "Practical Statistics for Particle Physicists", at the CERN-JINR school 2016:
<https://indico.cern.ch/event/467465/other-view?daysPerRow=5&view=nicecompact>
- Kyle Cranmer's article "Practical Statistics for the LHC", arXiv:1503.07622
- Harrisson Prosters 2015 CERN Academic Training lecture "Practical Statistics for LHC Physicists":
<https://indico.cern.ch/event/358542/>
- Machine learning introduction: T. Hastie, R. Tibshirani, J. Friedman, "The Elements of Statistical Learning", Springer 2001; C.M. Bishop, "Pattern Recognition and Machine Learning", Springer 2006
- Very incomplete list. Some more references are given throughout the lecture.

Why we need probability in the particle world

Since Pierre-Simon Laplace's times (1749–1827) the universe's fate was deterministic and, in spite of technical difficulties, was considered predictable if the complete equation of state were known.

Challenged by Heisenberg's uncertainty principle (1927), Albert Einstein proclaimed "*Gott würfelt nicht*" ("God does not play dice"), but so-called *hidden variables* to bring back determinism through the back door into the quantum world were never found.

In quantum mechanics, particles are represented by wave functions. The size of the wave function gives the probability that the particle will be found in a given position. The rate, at which the wave function varies from point to point, gives the speed of the particle.

Quantum phenomena like particle reactions occur according to certain probabilities. Quantum field theory allows us to compute cross-sections of particle production in scattering processes, and decays of particles. It cannot, however, predict how a single event will come out. We use probabilistic "Monte Carlo" techniques to simulate event-by-event realisations of quantum probabilities.



Pierre-Simon de Laplace



Albert Einstein



Werner Heisenberg



Statistics of large systems

Statistical physics uses probability theory and statistics to make statements about the approximate physics of large populations of stochastic nature, neglecting individuals.

Heavy-ion collisions at the LHC are modelled using notions of hydrodynamics (the strongly interacting medium behaves like perfect fluid)

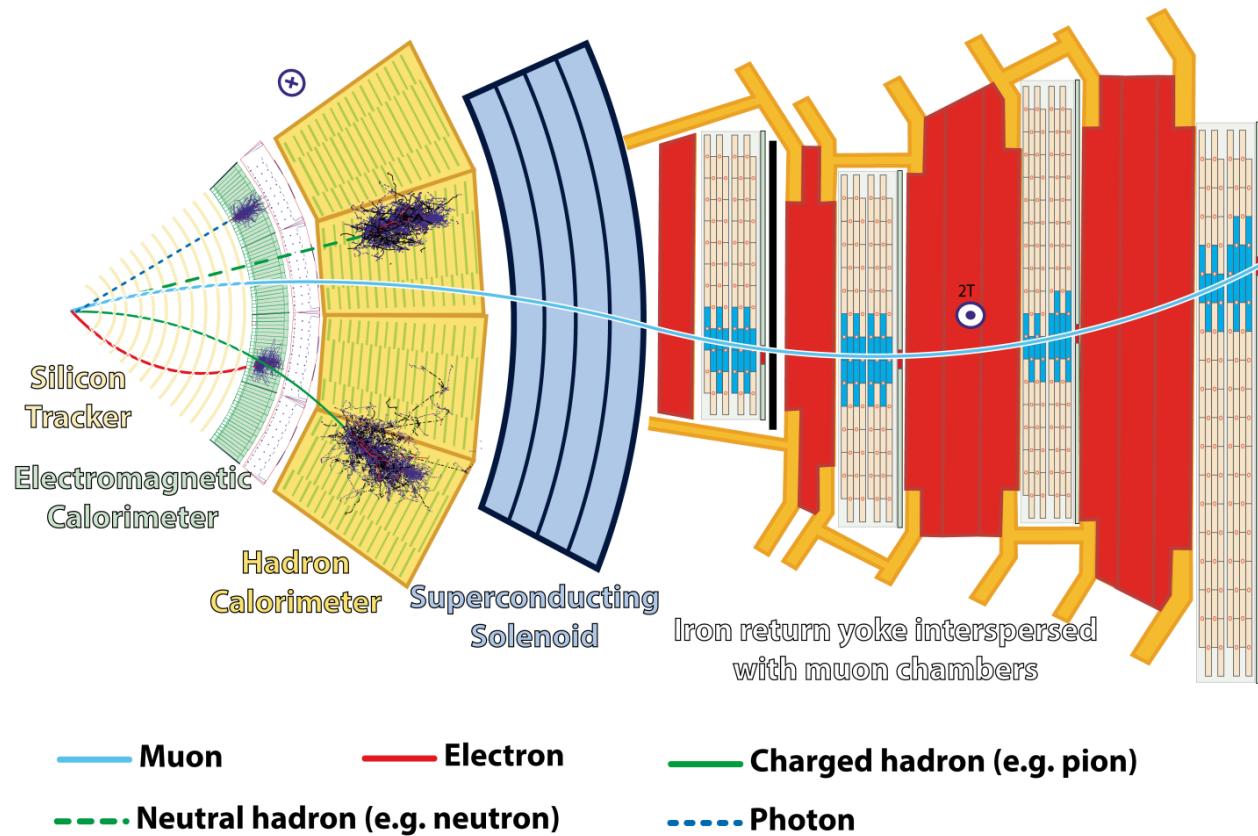
Statistical mechanics provides a framework for relating the microscopic properties of individual atoms and molecules to the macroscopic properties of materials that can be observed in everyday life. It thus explains thermodynamics as a natural result of statistics, classical mechanics, and quantum mechanics at the microscopic level.



Display of ATLAS Run-2 Heavy-Ion collision

Probability and statistics are fundamental ingredients & tools in all modern sciences

Statistics in measurement processes



In addition to the intrinsic probabilistic character of particle reactions, the measurement process through the interaction of particles with active detector materials contributes statistical degrees of freedom leading to measurement errors and to genuine systematic effects (eg, detector misalignment), that need to be considered in the statistical analysis

Statistics in measurement processes

Multiple scattering (diffusion) of light passing through a wetter and wetter windscreens (left to right).



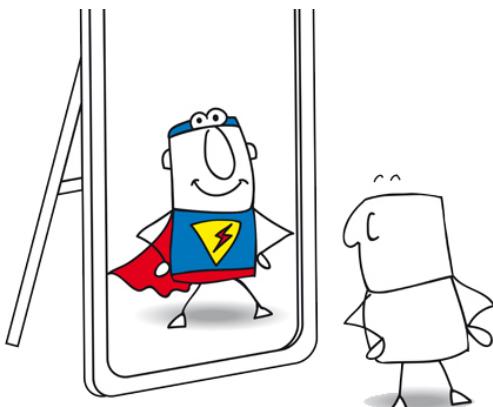
Track fitting in the LHC environment is challenging. It must deal with ambiguities, hit overlaps, multiple scattering, bremsstrahlung, multiple vertices, etc. Track fitters take Gaussian noise (e.g., Kalman filter) and non-Gaussian noise (e.g., Gaussian sum filter) into account.

Fitting is a statistical procedure that takes into account *a priori* known uncertainties.

Measurements and hypothesis testing

From a measured data sample, we want to determine parameters of a known model (eg, the top-quark mass in the Standard Model), we want to discover and measure missing pieces of the model (eg, the Higgs boson, neutrino masses), *and* we want to watch out for the unknown (test the data versus the predictions of a known model), or exclude parameters of suggested new physics models

Supersymmetry ?



QUARKS	mass →	$\approx 2.3 \text{ MeV}/c^2$	charge →	$2/3$	spin →	$1/2$	u	c	t	g	Higgs boson
							up	charm	top	gluon	boson
LEPTONS	mass →	$\approx 4.8 \text{ MeV}/c^2$	charge →	$-1/3$	spin →	$1/2$	d	s	b	γ	photon
							down	strange	bottom		
GAUGE BOSONS	mass →	$0.511 \text{ MeV}/c^2$	charge →	-1	spin →	$1/2$	e	μ	τ	Z	Z boson
							electron	muon	tau		
	mass →	$<2.2 \text{ eV}/c^2$	charge →	0	spin →	$1/2$	ν _e	ν _μ	ν _τ	W	W boson
							electron neutrino	muon neutrino	tau neutrino		

The Standard Model:

- The model needed to be developed
- Its particles needed to be observed
- Its parameters needed to be measured
- Is it all there is ?

Ingredients

Due to the intrinsic randomness of the data, probability theory is required to extract the information that addresses our questions. Statistics is used for the actual data analysis.

The interpretation of probability depends on the statement we want to make.

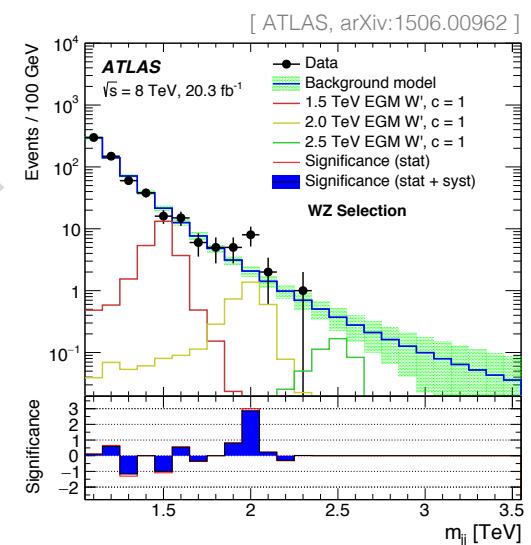
- For repeatable experiments, probability can be a measure of how frequently a statement is true
- In a more subjective approach, one could express a degree of belief of a statement

Repeatable experiments are for example:

- Playing a dice and finding 6
- Fluctuations in a background distribution and finding a peak of some size or more (note that the contrary: the probability that *the* peak is due to a background fluctuation, is *not* repeatable)

Non-repeatable statements are for example:

- The probability that dark matter is made of axions
- The probability that the new 125 GeV boson is the Higgs boson



Short excursion: Look-elsewhere effect (aka: trials factor)

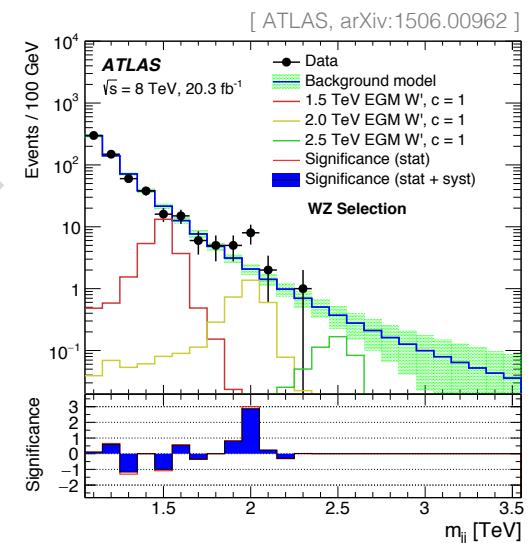
If we compute the probability (relative frequency) of an observed “bump” to happen under the background hypothesis, we need to take into account that the bump might occur anywhere in the spectrum as we do not know where the new physics would lie

→ Need to correct observed result by the *trials factor* (eg, using pseudo Monte Carlo event)

Ignoring the trials factor would be similar to ignoring that the probability of finding 6 when playing dice increases from 0.17 with a single dice to, eg, 0.60 with 5 dice

Repeatable experiments are for example:

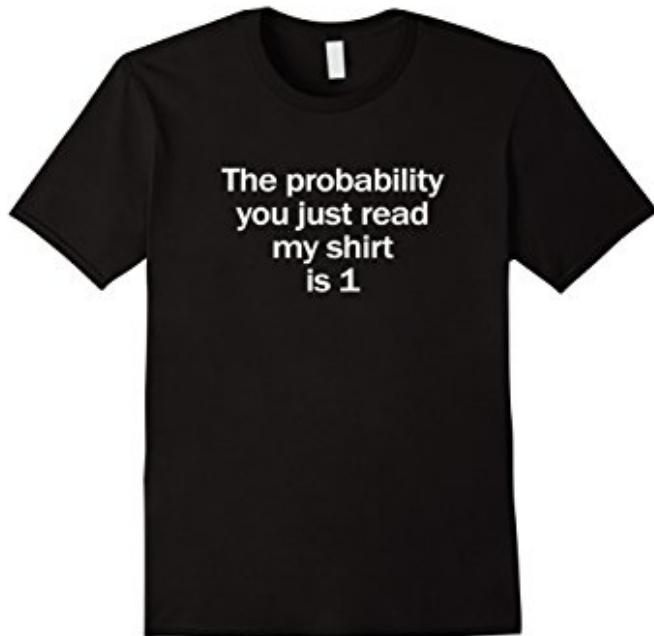
- Playing a dice and finding 6
- Fluctuations in a background distribution and finding a peak of some size or more (note that the contrary: the probability that *the* peak is due to a background fluctuation, is *not* repeatable)



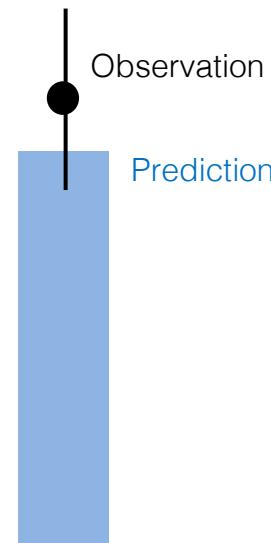
Random variables

In a statistical context, instead of “data” that follow a distribution, one often (typically) speaks of a “random variable”

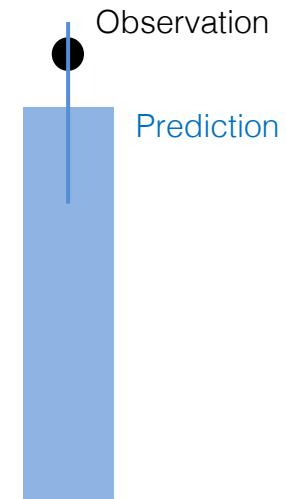
The data follow an (unknown) distribution of a random variable. Although we plot error bars around observed data points, a given observation does not have an uncertainty. It is fixed. It would be more correct to plot the error bar around the prediction. By plotting the error bar around the observation, we assume by convenience the observed value to be the truth.



Custom:



Actually:



Statistical distributions

Measurement results typically follow some “distribution”, ie, the data do not appear at fixed values, but are “spread out” in a characteristic way

Which type of distribution it follows depends on the particular case

- It is important to know the occurring distributions to be able to pick the correct one when interpreting the data (example: *Poisson* vs. *Compound Poisson*)
- ...and it is important to know their characteristics to extract the correct information

Probability distribution / density of a random variable

Random variable \mathbf{k} (discrete) or \mathbf{x} (continuous): quantity or point in sample space

Discrete variable

$$P(k_i) = p_i$$

Continuous variable

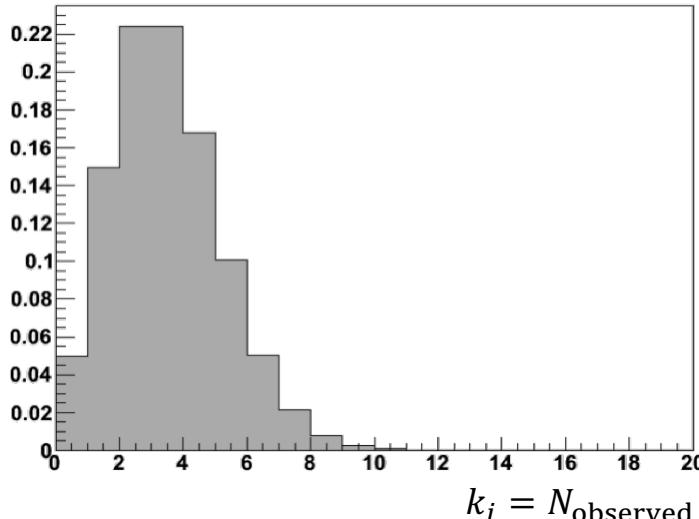
$$P(x \in [x, x + dx]) = p_x(x)dx$$

Normalisation (your parameter/event space covers all possibilities - *unitarity*)

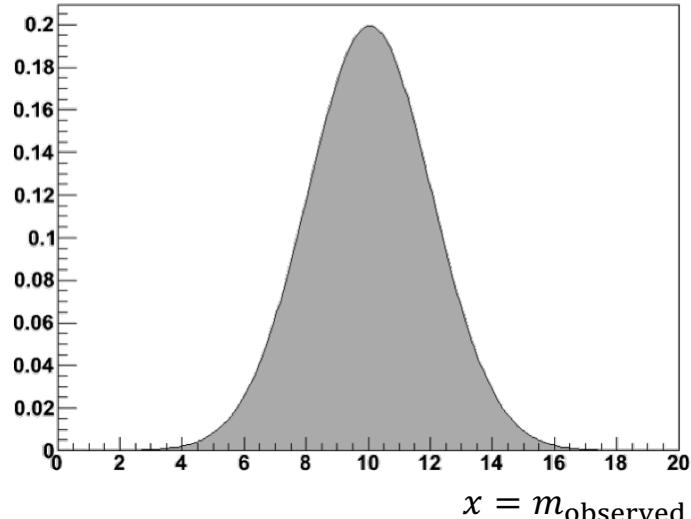
$$\sum_{i=0 \dots \infty} P(k_i) = 1$$

$$\int_{-\infty}^{\infty} p_x(x)dx = 1$$

Poisson distribution



Gaussian (or *Normal*) distribution



Cumulative distribution

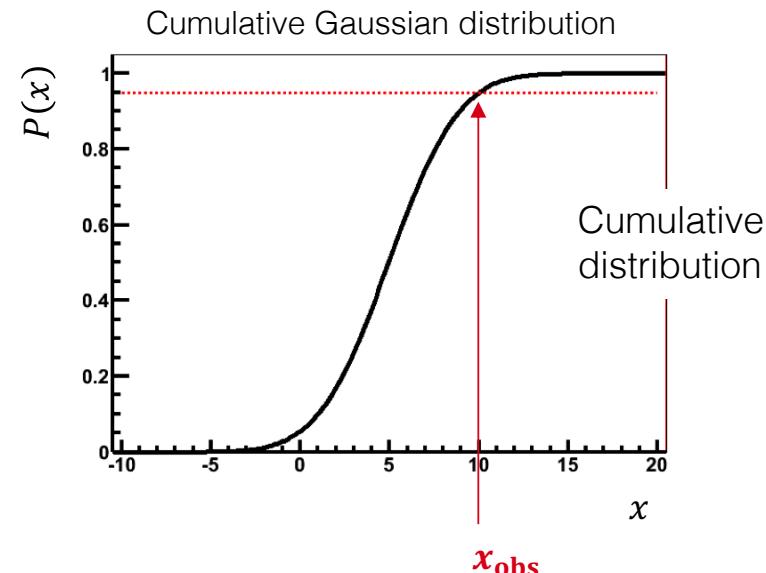
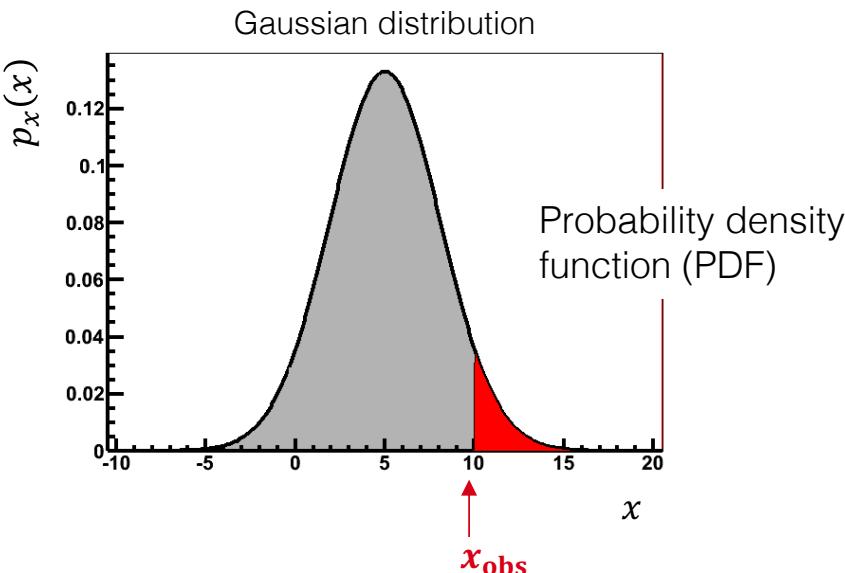
$p_x(x)$: probability density distribution for some “measurement” x under the assumption of some model and its parameters

The cumulative distribution $P(x)$ is *the probability to observe a random value x smaller than the one observed, x_{obs}*

→ Examples for cumulative distributions: χ^2 , p-values, confidence limits (will come back to this)

$$p_x(x) = dP(x)/dx$$

$$\int_{-\infty}^x p_x(x') dx' \equiv P(x)$$



Selected probability (density) distributions

Imagine a monkey discovered a huge bag of alphabet noodles. The monkey blindly draws noodles out of the bag and places them in a row. The text reads:

ROGER FEDERER

The probability for this to happen is about 10^{-17}



Infinite monkey theorem: provided enough time,
the monkey will type Shakespeare's Hamlet

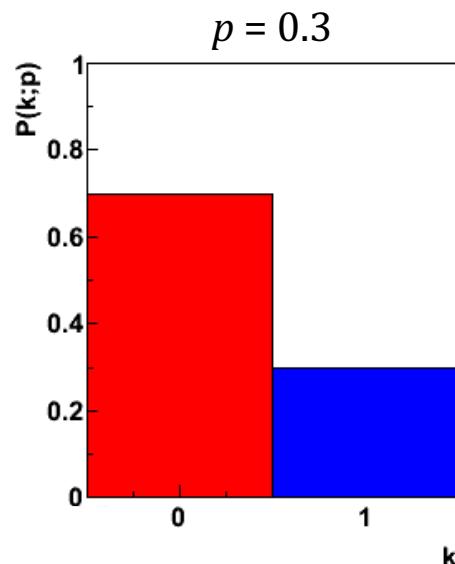
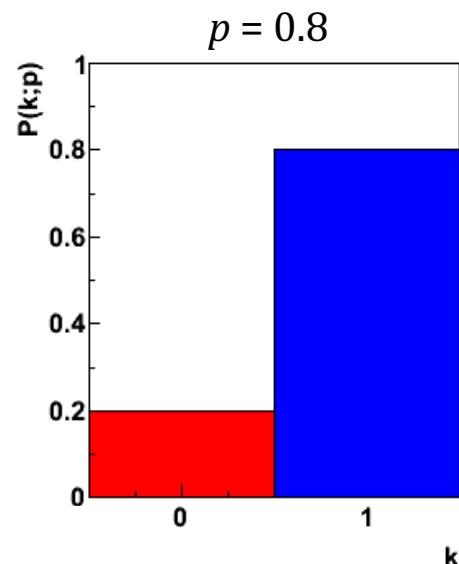
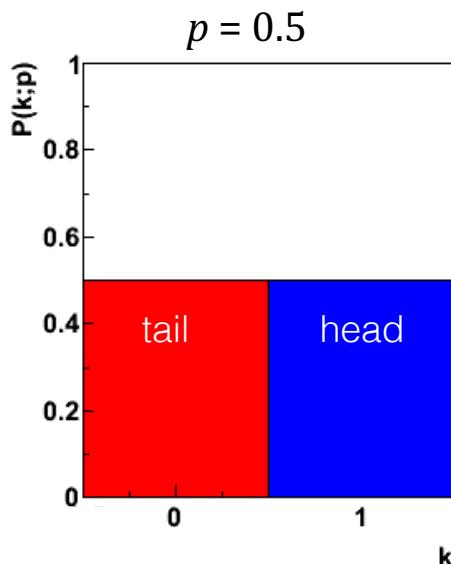
Bernoulli distribution

Experiment with two possible discrete outcomes: $k = 1$ / $k = 0$ (or yes / no or head / tail, etc)

What is the probability of one or the other ?

$$P(\text{head}) = p \text{ (where } 0 \leq p \leq 1\text{)}, \quad P(\text{tail}) = 1 - P(\text{head}) = 1 - p$$

$$\Rightarrow P(k; p) = p^k (1 - p)^{1-k} \text{ for } k \in \{0,1\}$$



Binomial distribution (very important!)

Now let's get more complex: throw N coins (or similar binary choices)

How often (likely) is $k \times \text{head}$ and $(N - k) \times \text{tail}$?

- Each coin: $P(\text{head}) = p, P(\text{tail}) = 1 - p$
- Pick k particular coins → the probability of all having **head** is:

$$P(k \times \text{head}) = P(\text{head}) \cdot P(\text{head}) \cdot \dots \cdot P(\text{head}) = P(\text{head})^k = p^k$$

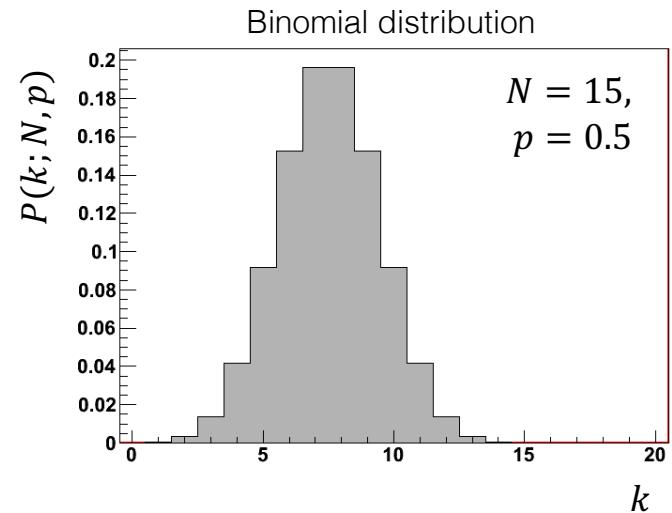
- Multiply this by the probability that all remaining $N-k$ coins land on **tail**:

$$P(\text{head})^k \cdot P(\text{tail})^{N-k} = p^k (1-p)^{N-k}$$

- This was for a particular choice of k coins
- Now include all $\binom{N}{k}$ permutations for *any* k coins

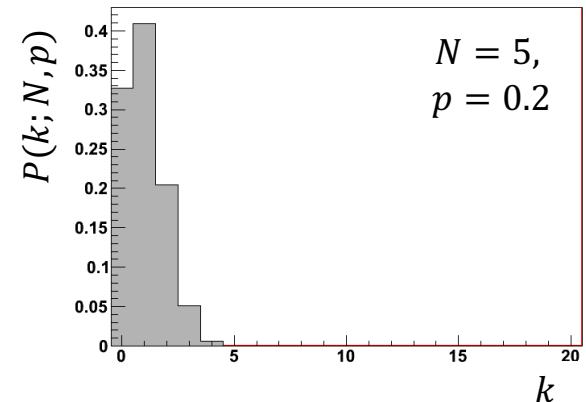
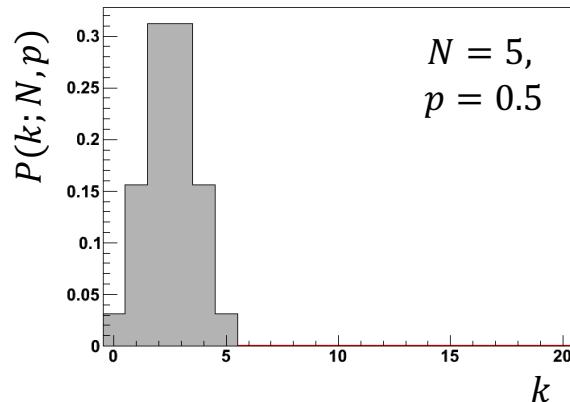
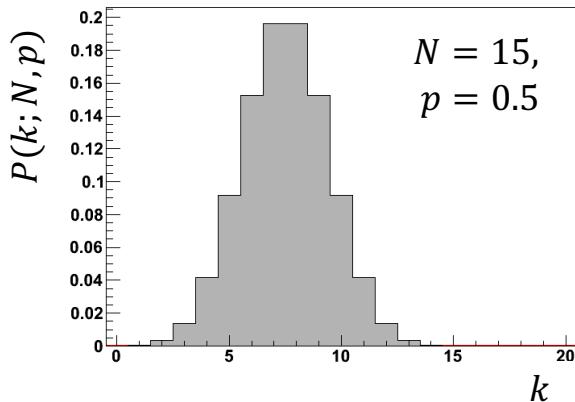
$$P(k; N, p) = p^k (1-p)^{N-k} \binom{N}{k}$$

where $\binom{N}{k} = \frac{N!}{k!(N-k)!}$ is the binomial coefficient



Binomial distribution (continued)

Example binomial distributions:



- *Expectation value*: sum over all possible outcomes and *average* (i.e.: weighted average)

$$E[k] = \sum_k kP(k; N, p) = Np$$

- *Variance*: (see next slide for definition)

$$V[k] = Np(1 - p)$$

Characteristic quantities of distributions

Quantity	Discrete variable	Continuous variable
Expectation (mean) value E	$E[k] = \langle k \rangle = \sum_k kP(k)$	$E[x] = \langle x \rangle = \int x \cdot p_x(x) dx$
Variance (spread) $V = \sigma^2$	$E[(k - \langle k \rangle)^2] = E[k^2] - (E[k])^2$	same with $k \rightarrow x$
Higher moments: skew	$E[(k - \langle k \rangle)^3]$	same with $k \rightarrow x$

Note that “expectation value” and “variance” are properties of the full data population. Unbiased estimates can be derived from N samples extracted from the population:

Sample variance (unbiased)	$\frac{1}{N-1} \sum_{i=1}^N (k_i - \langle k \rangle)^2$	same with $k \rightarrow x$
----------------------------	--	-----------------------------

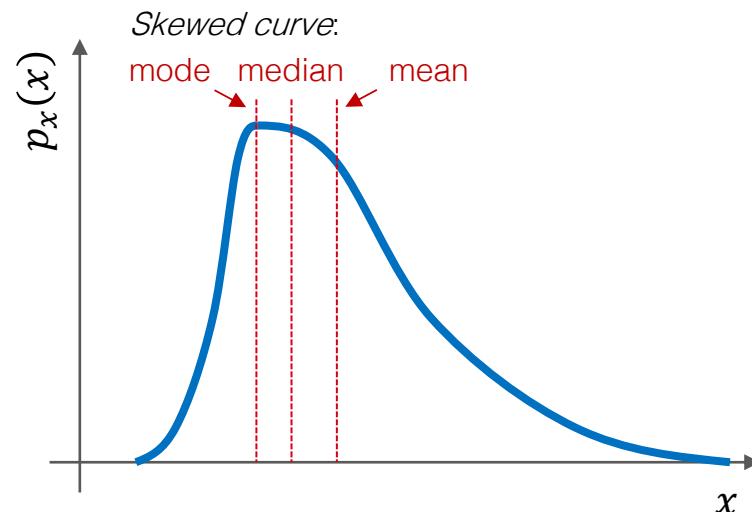
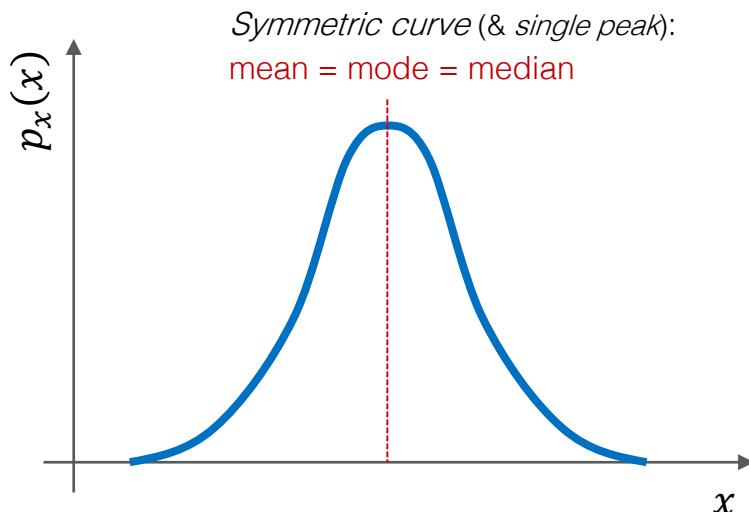
Characteristic quantities of distributions (continued)

Mean, Mode, Median:

- **Mean**: $\langle x \rangle$ — defined before
- **Mode**: most probable value x_{mode} : $p_x(x_{\text{mode}}) \geq p_x(x), \forall x$
- **Median**: *2-quantile*: 50% of x values are larger than x_{median} , 50% are smaller

Can generalise *k-quantile*: points at regular intervals of the cumulative distribution.

Boundaries of binning chosen such that each bin contains the $1/k$ -th of total integral of distribution.



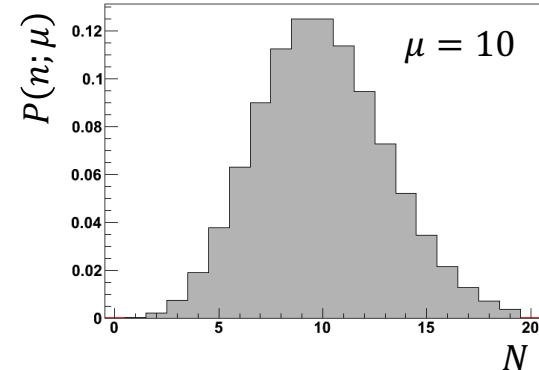
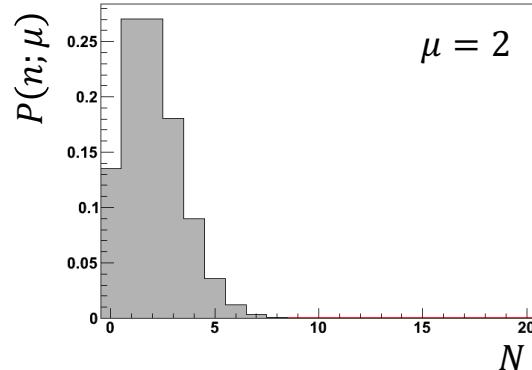
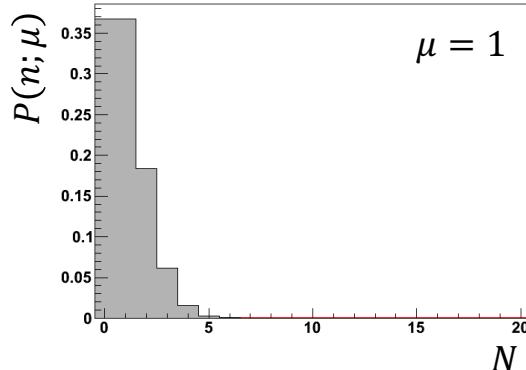
Poisson distribution

Recall: for individual events each with two possible outcomes → Binomial distribution

How about: number of counts in radioactive decay experiment during given time interval Δt ?

- Events happen “randomly” but there is no such 2nd outcome; Δt is continuous, no discrete num. of trials
- μ : average number of counts in Δt . What is the probability of observing N counts?
- Limit of Binomial distribution for large number of trials and small p : $N \rightarrow \infty$ & $p \rightarrow 0$ so that $Np \rightarrow \mu$.

→ Poisson distribution: $P(N; \mu) = \frac{\mu^N}{N!} e^{-\mu}$



Expectation value: $E[N] = \sum_N N \cdot P(N; \mu) = \mu$, Variance: $V[N] = \mu$

Poisson is good approximation for Binomial distribution for $N \gg \mu$ ($= Np$)

Gaussian (also: “Normal”) distribution

In limit of large μ a Poisson distribution approaches a symmetric Gaussian distribution

- This is the case not only for the Poisson distributions, but for almost any sufficiently large sum of samples with different sub-properties (mean & variance) → **Central Limit Theorem** (will discuss later)
- Gaussian distribution is of utter use, and luckily has simple properties

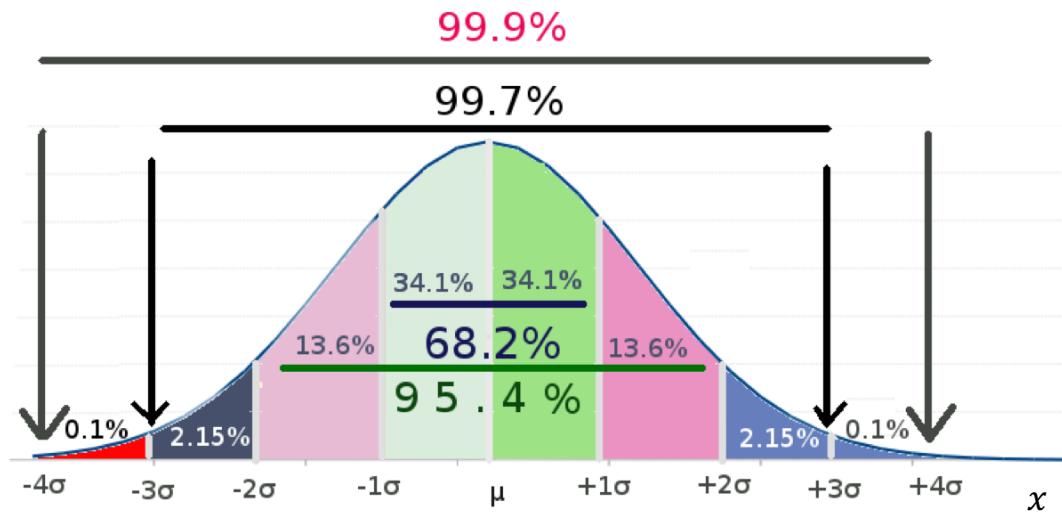
→ Gauss distribution: $P(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi} \cdot \sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$

Symmetric distribution:

- Expectation value: $E[x] = \mu$
- Variance: $V[x] = \sigma^2$
- Probability content:

$$\int_{-\sigma}^{+\sigma} P(x; \mu, \sigma) dx = 68.2\%$$

$$\int_{-2\sigma}^{+2\sigma} P(x; \mu, \sigma) dx = 95.4\%$$



Gaussian (also: “Normal”) distribution

In limit of large μ a Poisson distribution approaches a symmetric Gaussian distribution

- This is the case not only for the Poisson distributions, but for almost any sufficiently large sum of samples with different sub-properties (mean & variance) → **Central Limit Theorem** (will discuss later)
- Gaussian distribution is of utter use, and luckily has simple properties

→ Gauss distribution: $P(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi} \cdot \sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$

Symmetric distribution:

- Expectation value: $E[x] = \mu$
- Variance: $V[x] = \sigma^2$

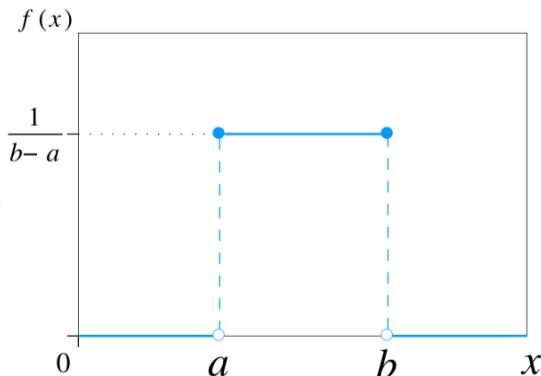
Poisson distribution:

- Expectation value: $E[x] = \mu$
- Variance: $V[x] = \mu$

→ For large μ , the standard deviation (σ) of the expected event counts is $\sqrt{\mu}$!

Some other distributions

Uniform (“flat”) distribution



Exponential distribution

- Particle decay density versus time
(in the particle’s rest frame!)

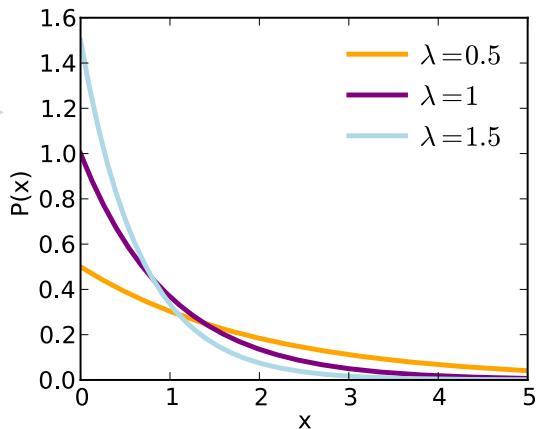
Relativistic Breit-Wigner distribution

- Distribution of resonance of unstable particle as function of centre-of-mass energy in which the resonance is produced
(originates from the propagator of an unstable particle)

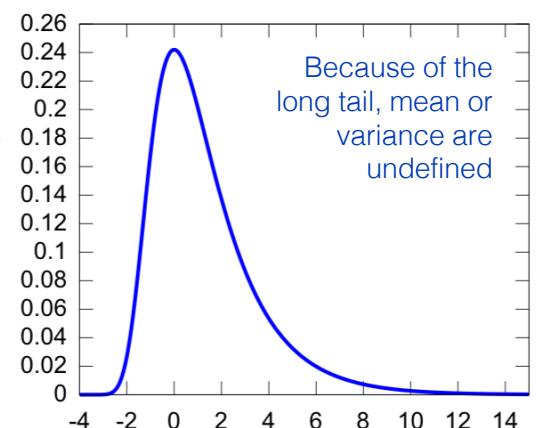
Chi-squared (χ^2) distribution

- Sum of squares of Gaussian distributed variables;
used to derive goodness of a fit to describe data

Landau distribution



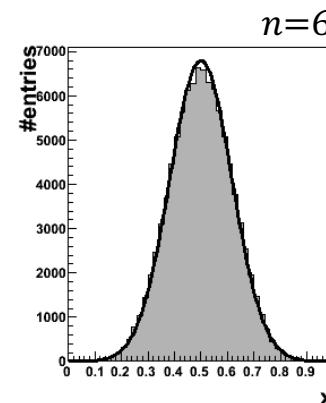
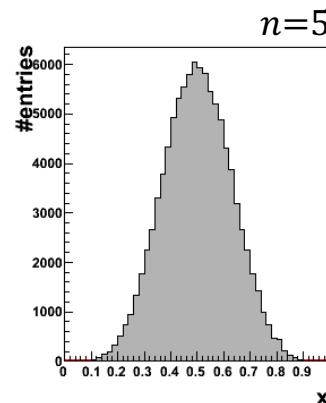
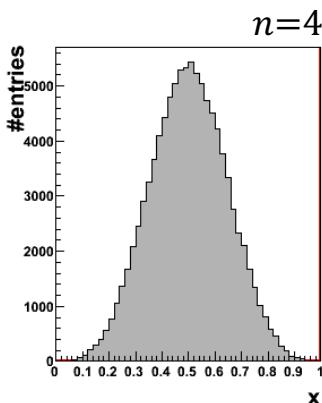
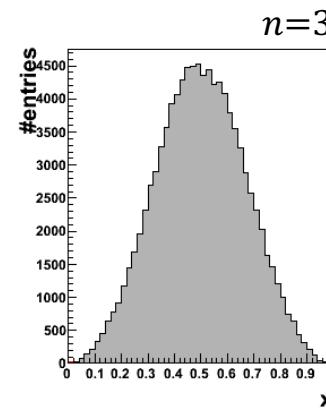
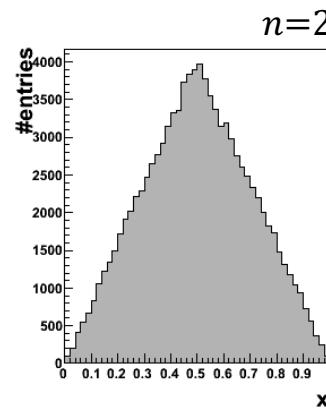
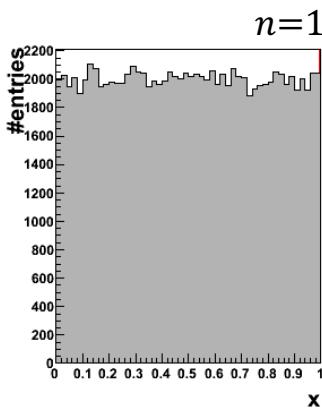
Many more, see <http://pdg.lbl.gov/2015/reviews/rpp2015-rev-probability.pdf>
for definitions and properties.



Central limit theorem (CLT)

CLT: the sum of n independent samples x_i ($i = 1, \dots, n$) drawn from any PDF $D(x_i)$ with well defined expectation value and variance is Gaussian distributed in the limit $n \rightarrow \infty$

$$D: E_D[x_i] = \mu; V_D[x_i] = \sigma_D^2, \text{ and: } y = \frac{1}{n} \sum_{i=1}^n x_i \Rightarrow E_{\text{Gauss}}[y] = \mu; V_{\text{Gauss}}[y] = \frac{\sigma_D^2}{n}$$



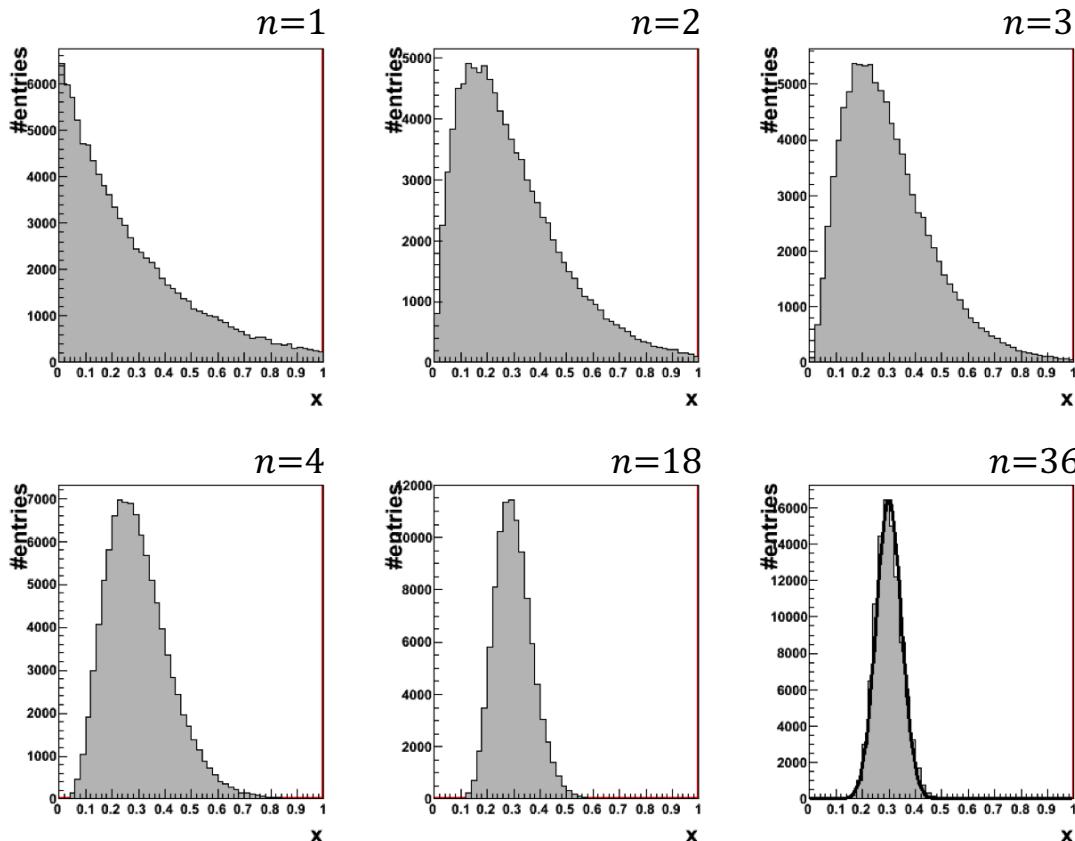
Averaging reduces
the variance

Example: summing up uniformly distributed ensembles within $[0, 1]$

Central limit theorem (CLT)

CLT: the sum of n independent samples x_i ($i = 1, \dots, n$) drawn from any PDF $D(x_i)$ with well defined expectation value and variance is Gaussian distributed in the limit $n \rightarrow \infty$

$$D: E_D[x_i] = \mu; V_D[x_i] = \sigma_D^2, \text{ and: } y = \frac{1}{n} \sum_{i=1}^n x_i \Rightarrow E_{\text{Gauss}}[y] = \mu; V_{\text{Gauss}}[y] = \frac{\sigma_D^2}{n}$$



Example: summing up exponential distributions

Central Gaussian limit works even if D doesn't look Gaussian at all

Central limit theorem (CLT)

CLT is key concept in probability theory: it implies that probabilistic and statistical methods that work for Gaussian distributions can be applicable to many problems involving other types of distributions

Crucial for experimental particle physics: can propagate systematic uncertainties in a Gaussian manner via quadratic addition (see later)

Multidimensional random variables

What if a measurement consists of two variables?

Let:

A = measurement x in $[x, x + dx]$

B = measurement y in $[y, y + dy]$

Joint probability: $P(A \cap B) = p_{xy}(x, y)dx dy$

(where $p_{xy}(x, y)$ is joint probability density function — PDF)

If the two variables are independent:

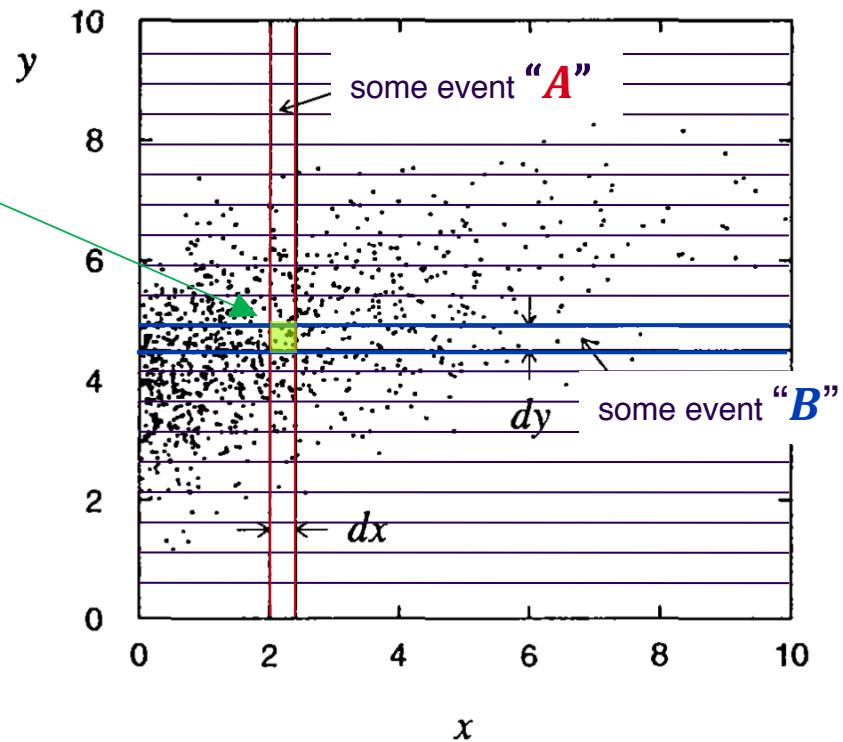
$$P(A \cap B) = P(A) \cdot P(B)$$

$$p_{xy}(x, y) = p_x(x) \cdot p_y(y)$$

Marginal PDF: if one is not interested in dependence on y (or cannot measure it),

- integrate out (“marginalise”) y , ie, project onto x
- resulting one-dimensional PDF: $p_x(x) = \int p_{xy}(x, y)dy$

From: Glen Cowan,
Statistical data analysis



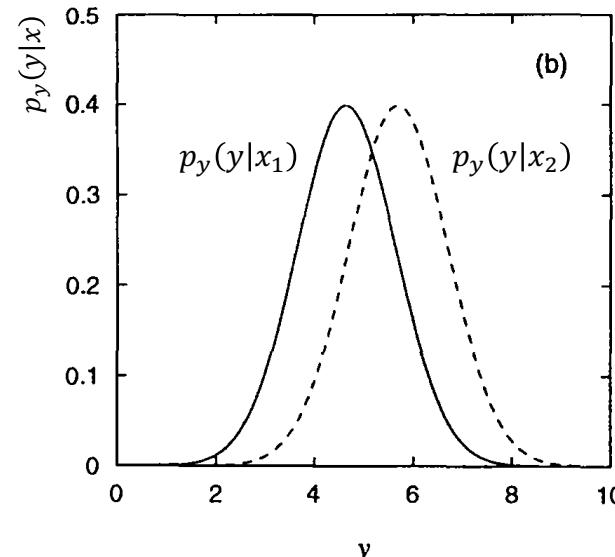
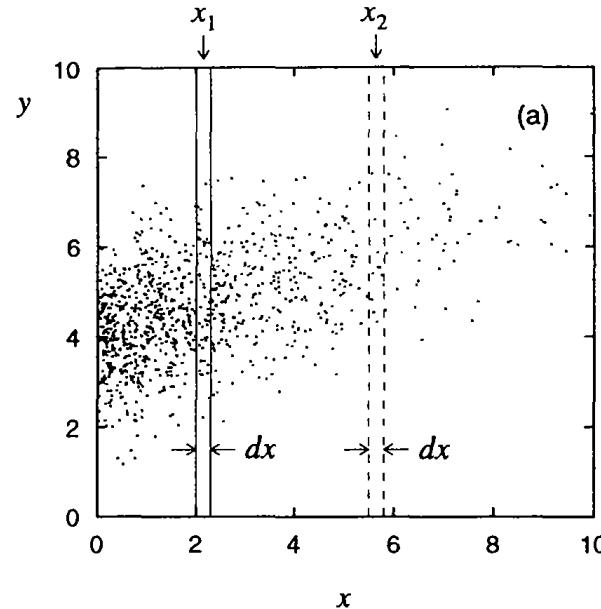
Conditioning versus marginalisation

Conditional probability $\mathbf{P}(\mathbf{A}|\mathbf{B})$: [read: $P(A|B)$ = “probability of A given B ”]

$$P(A \cap B) = P(A|B) \cdot P(B) \quad \Leftrightarrow \quad \mathbf{P}(\mathbf{A}|\mathbf{B}) = \frac{P(A \cap B)}{P(B)} = \frac{p_{xy}(x, y)dx dy}{p_y(y)dy}$$

Rather than integrating over the whole y region (marginalisation),
look at one-dimensional (1D) slices of the two-dimensional (2D) PDF $p_{xy}(x, y)$:

$$p_y(y|x_1) = p_{xy}(x = \text{const} = x_1, y)$$



From: Glen Cowan,
Statistical data analysis

Covariance and correlation

Recall, for 1D PDF $\mathbf{p}_x(x)$ we had: $E[x] = \mu_x$; $V[x] = \sigma_x^2$

For a 2D PDF $\mathbf{p}_{xy}(x, y)$, one correspondingly has: $\mu_x, \mu_y, \sigma_x, \sigma_y$

How do x and y co-vary? $\rightarrow C_{xy} = \text{covariance}_{xy} = E[(x - \mu_x)(y - \mu_y)] = E[xy] - \mu_x\mu_y$

From this define the scale / dimension invariant *correlation coefficient*:

$$\rho_{xy} = \frac{C_{xy}}{\sigma_x \sigma_y}, \text{ where } \rho_{xy} \in [-1, +1]$$

- If x, y are independent: $\rho_{xy} = 0$, ie, they are *uncorrelated* (or they *factorise*)

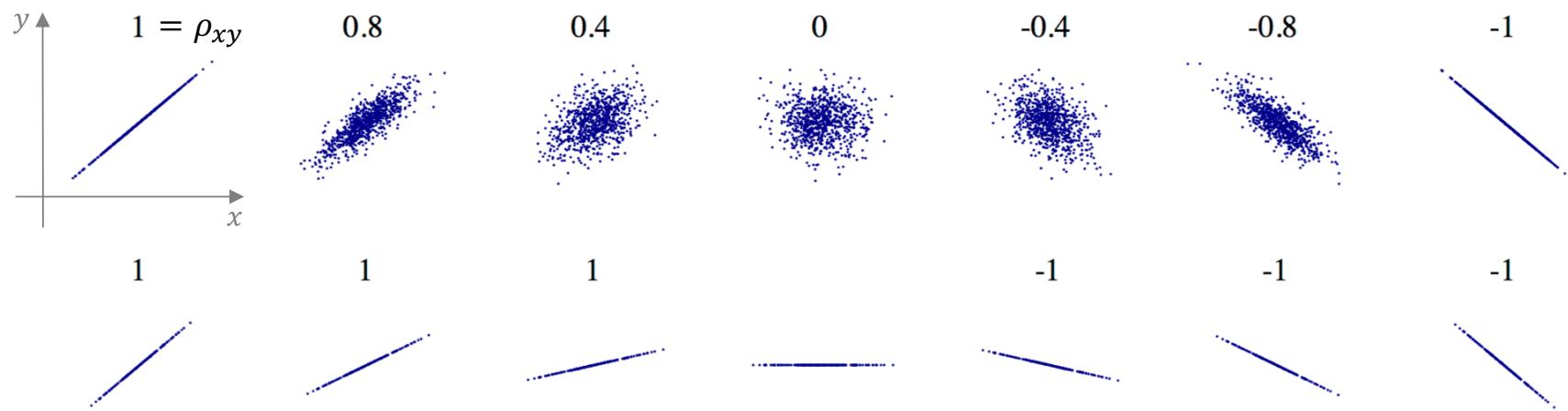
Proof: $E[xy] = \iint xy \cdot p_{xy}(x, y) dx dy = \iint xy \cdot p_x(x)p_y(y) dx dy = \int x \cdot p_x(x) dx \cdot \int y \cdot p_y(y) dy = \mu_x \mu_y$

- Note that the contrary is not always true: non-linear correlations can lead to $\rho_{xy} = 0$,
 \rightarrow see next page

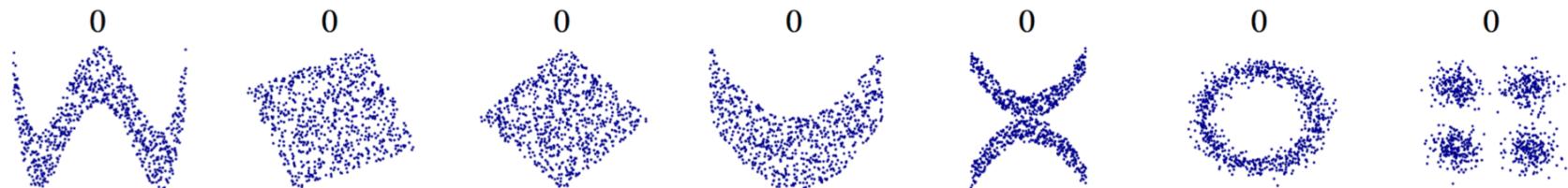
Correlations

Figure from: https://en.wikipedia.org/wiki/Correlation_and_dependence

The correlation coefficient measures the noisiness and direction of a linear relationship:



...it does not measure the slope ρ_{xy} (see above figures)



...and non-linear correlation patterns are not or only approximately captured by ρ_{xy} (see above figures)

Digression — Mutual information

Non-linear correlation can be captured by the “*mutual information*” quantity I_{xy} :

$$I_{xy} = \iint p_{xy}(x, y) \cdot \ln\left(\frac{p_{xy}(x, y)}{p_x(x)p_y(y)}\right) dx dy$$

Measure of mutual dependence between two variables:
“How much information is shared among them”

where $I_{xy} = 0$ only if x, y are fully statistically independent

Proof: if independent, then $p_{xy}(x, y) = p_x(x)p_y(y) \Rightarrow \ln(\dots) = 0$

NB: $I_{xy} = H_x - H_x(y) = H_y - H_y(x)$,

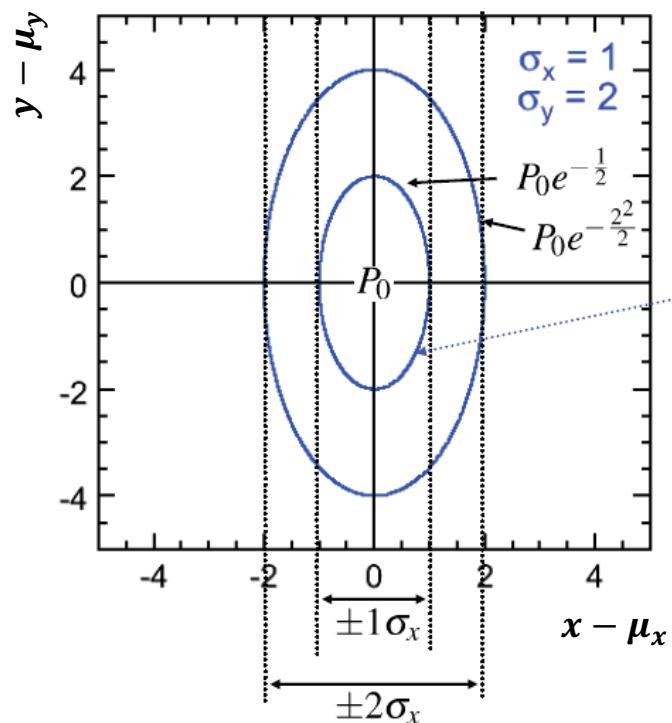
where $H_x = -\int p_x(x) \cdot \ln(p_x(x)) dx$ is *entropy*, $H_x(y)$ is *conditional entropy*



2D Gaussian (uncorrelated)

Two variable \mathbf{x}, \mathbf{y} are independent: [$p_{xy}(x, y) = p_x(x) \cdot p_y(y)$]

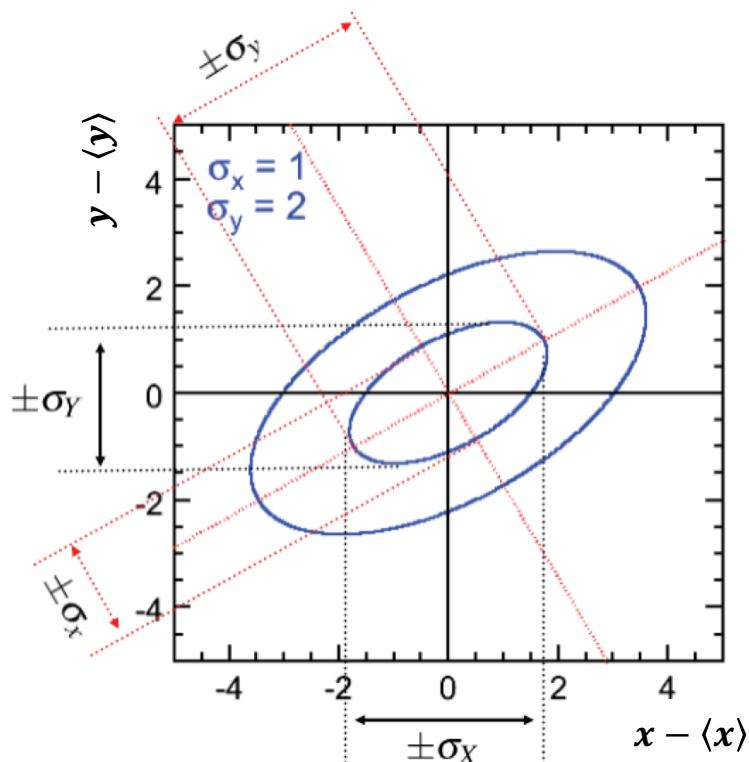
$$p_{xy}(x, y) = \frac{1}{\sqrt{2\pi}\sigma_x} e^{-\frac{(x-\mu_x)^2}{2\sigma_x^2}} \cdot \frac{1}{\sqrt{2\pi}\sigma_y} e^{-\frac{(y-\mu_y)^2}{2\sigma_y^2}}$$



2D Gaussian (correlated)

Two variable \mathbf{x}, \mathbf{y} are *not* independent: [$p_{xy}(x, y) \neq p_x(x) \cdot p_y(y)$]

$$p_{\vec{x}}(\vec{x}) = \frac{1}{2\pi\sqrt{\det(C)}} \cdot \exp\left(-\frac{1}{2}(\vec{x} - \vec{\mu})^T C^{-1} (\vec{x} - \vec{\mu})\right)$$



where:

$$C = \begin{pmatrix} \langle x^2 \rangle - \langle x \rangle^2 & \langle xy \rangle - \langle x \rangle \langle y \rangle \\ \langle xy \rangle - \langle x \rangle \langle y \rangle & \langle y^2 \rangle - \langle y \rangle^2 \end{pmatrix}$$

is the (symmetric) *covariance matrix*

Corresponding correlation matrix elements:

$$\rho_{ij} = \rho_{ji} = \frac{C_{ij}}{\sqrt{C_{ii} \cdot C_{jj}}}$$

Digression — SQRT decorrelation

Find variable transformation that diagonalises a covariance matrix C

Determine “square-root” C' of C (such that: $C = C' \cdot C'$) by first diagonalising C

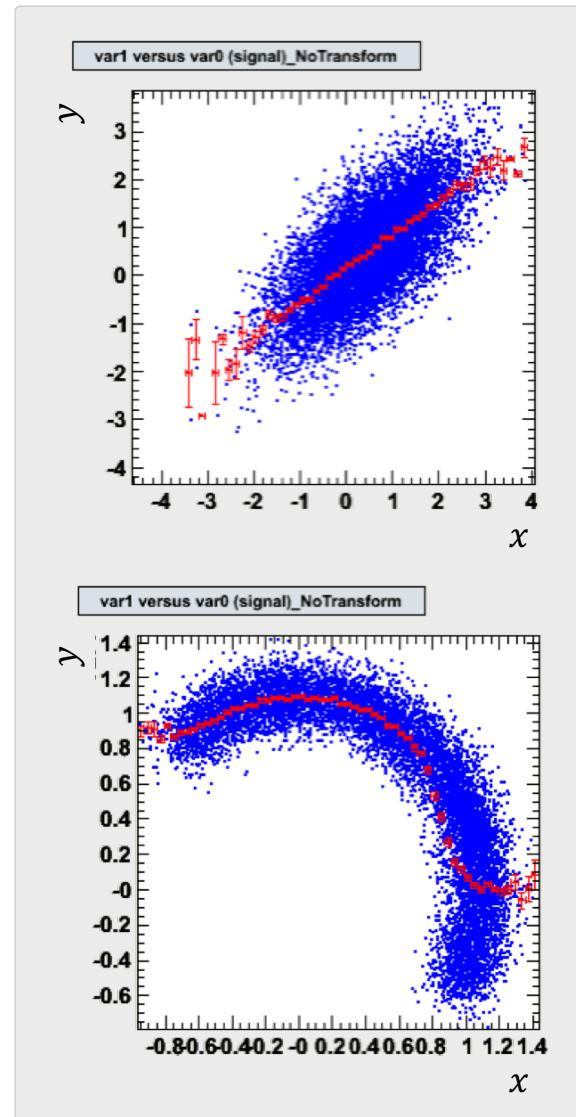
$$D = S^T \cdot C \cdot S \Leftrightarrow C' = S \cdot \sqrt{D} \cdot S^T$$

where D is diagonal, $\sqrt{D} = \{\sqrt{d_{11}}, \dots, \sqrt{d_{nn}}\}$,
and S an orthogonal matrix

Linear decorrelation of correlated vector x then obtained by

$$x' = (C')^{-1} \cdot x$$

Principle component analysis (PCA) is another convenient method to achieve linear decorrelation
(PCA is linear transformation that rotates a vector such that the maximum variability is visible. It identifies most important gradients)



Example:
original
correlations

Digression — SQRT decorrelation

Find variable transformation that diagonalises a covariance matrix C

Determine “square-root” C' of C (such that: $C = C' \cdot C'$) by first diagonalising C

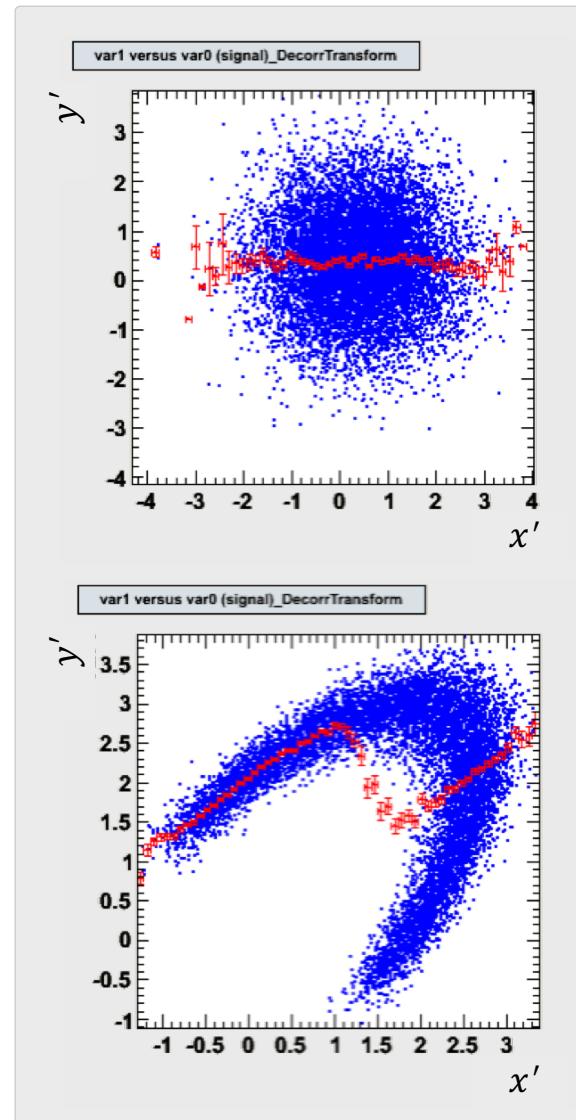
$$D = S^T \cdot C \cdot S \Leftrightarrow C' = S \cdot \sqrt{D} \cdot S^T$$

where D is diagonal, $\sqrt{D} = \{\sqrt{d_{11}}, \dots, \sqrt{d_{nn}}\}$,
and S an orthogonal matrix

Linear decorrelation of correlated vector x
then obtained by

$$x' = (C')^{-1} \cdot x$$

SQRT decorrelation works only for linear correlations!



Example:
after SQRT
decorrelation

Functions of random variables

Any function of a random variable is itself a random variable

E.g., x with PDF $p_x(x)$ becomes: $y = f(x)$

y could be a parameter extracted from a measurement

What is the PDF $p_y(y)$?

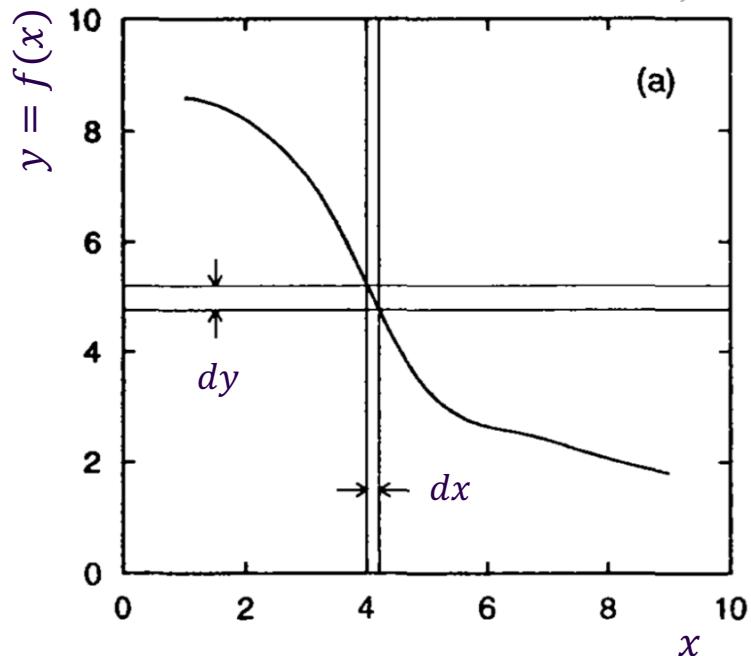
- Probability conservation: $p_y(y)|dy| = p_x(x)|dx|$
- For a 1D function $f(x)$ with existing inverse:

$$dy = \frac{df(x)}{dx} dx \Leftrightarrow dx = \frac{df^{-1}(y)}{dy} dy$$

- Hence: $p_y(y) = p_x(f^{-1}(y)) \left| \frac{dx}{dy} \right|$

Note: this is **not** the standard error propagation but the full PDF !

Glen Cowan: Statistical data analysis



Error propagation

Let's assume a measurement \mathbf{x} with *unknown* PDF $\mathbf{p}_x(\mathbf{x})$, and a transformation $\mathbf{y} = \mathbf{f}(\mathbf{x})$

- \bar{x} and \hat{V} are estimates of μ and variance σ^2 of $p_x(x)$

What are $E[y]$ and, in particular, σ_y^2 ? → Taylor-expand $f(x)$ around \bar{x} :

- $f(x) = f(\bar{x}) + \frac{df}{dx}\Big|_{x=\bar{x}} (x - \bar{x}) + \dots \Rightarrow E[f(x)] \simeq f(\bar{x}) \quad (\text{because: } E[x - \bar{x}] = 0 !)$

Now define $\bar{y} = f(\bar{x})$, and from the above follows:

$$\Leftrightarrow y - \bar{y} \simeq \frac{df}{dx}\Big|_{x=\bar{x}} (x - \bar{x})$$

$$\Leftrightarrow E[(y - \bar{y})^2] = \left(\frac{df}{dx}\Big|_{x=\bar{x}}\right)^2 E[(x - \bar{x})^2]$$

$$\Leftrightarrow \hat{V}_y = \left(\frac{df}{dx}\Big|_{x=\bar{x}}\right)^2 \hat{V}_x$$

$$\Leftrightarrow \sigma_y = \frac{df}{dx}\Big|_{x=\bar{x}} \cdot \sigma_x \quad \rightarrow \quad (\text{approximate}) \text{ error propagation}$$

Error propagation (continued)

In case of several variables, compute covariance matrix and partial derivatives

- Let $f = f(x_1, \dots, x_n)$ be a function of n randomly distributed variables
- $\left(\frac{df}{dx}\Big|_{x=\bar{x}}\right)^2 \hat{V}_x$ becomes: $\sum_{i,j=1}^n \frac{\partial f}{\partial x_i} \frac{\partial f}{\partial x_j} \Big|_{\bar{x}} \cdot \hat{V}_{i,j}$ (where: $\bar{x} = (\bar{x}_1, \dots, \bar{x}_n)$)
- with the covariance matrix:

$$\hat{V}_{i,j} = \begin{bmatrix} \sigma_{x_1}^2 & \cdots & \sigma_{x_1 x_n} \\ \vdots & \ddots & \vdots \\ \sigma_{x_n x_1} & \cdots & \sigma_{x_n}^2 \end{bmatrix}$$

- The resulting “error” (uncertainty) depends on the correlation of the input variables
- Typically (not always!) positive correlations lead to an increase of the total error,
 - and negative correlations decrease the total error

For complicated functional dependence $f = f(x_1, \dots, x_n)$, it is often more practical to use Monte Carlo techniques (“pseudo MC generation”) to propagate uncertainties



Summary for today

Probability and statistics are everywhere in science, and in particular profoundly contained in particle physics and in the physics of large ensembles

Overview of some important probability density distributions given, and also:

- Joint / marginal / conditional probabilities
- Covariance and correlations
- Error propagation

Tomorrow: how to use these concepts for hypothesis testing