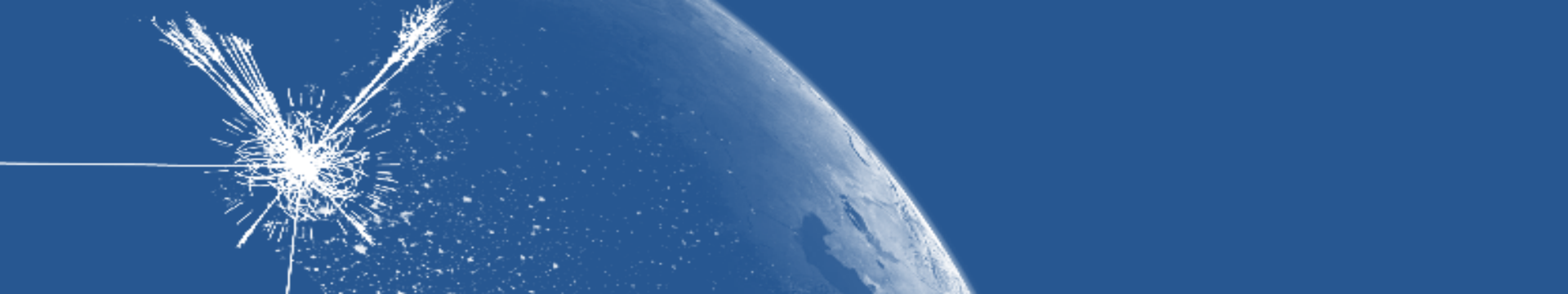# Introduction to probability and statistics (4)

Andreas Hoecker (CERN)

CERN Summer Student Lecture, 17–21 July 2017

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

# Outline (4 lectures)

1st lecture:
- Introduction
- Probability

2nd lecture:
- Probability axioms and hypothesis testing
- Parameter estimation
- Confidence levels (some catch up to do…)

3rd lecture:
- Maximum likelihood fits
- Monte Carlo methods
- Data unfolding

4th lecture:
- Multivariate techniques and machine learning
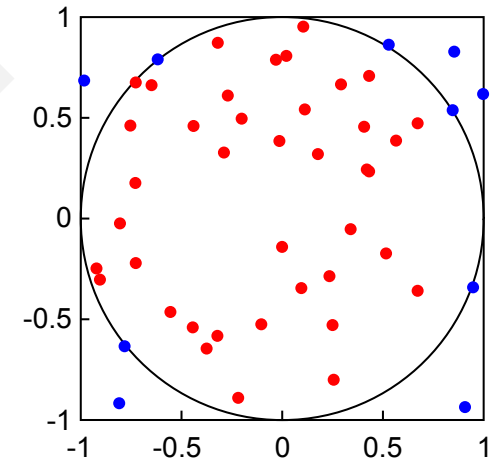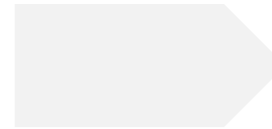
# Monte Carlo techniques

# Why "Monte Carlo" techniques ?

Monte Carlo (MC) techniques are computational algorithms that rely on repeated random sampling to obtain numerical results

They are used when analytical solutions are too complex or not even known
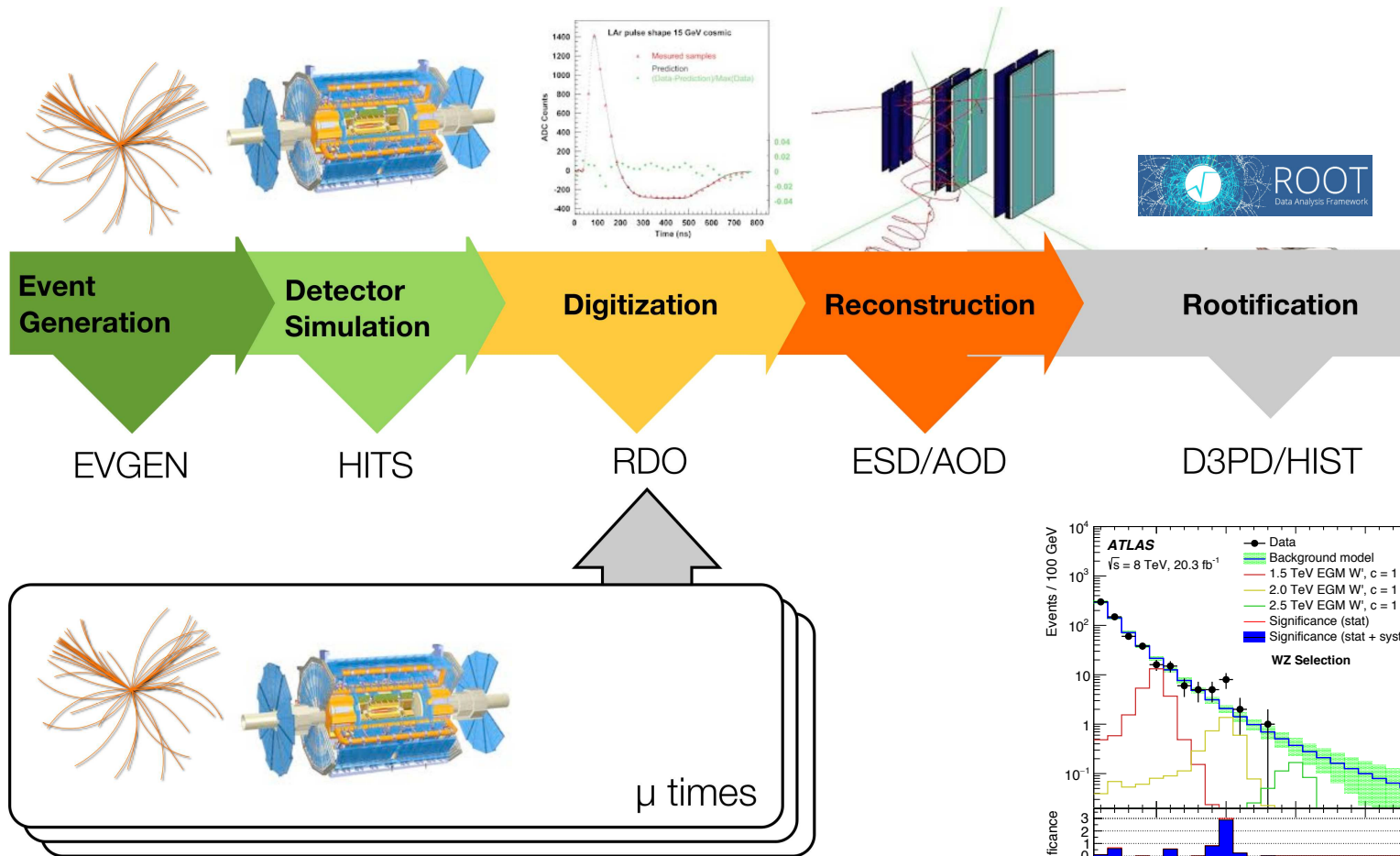
Examples:

- Numerical integration of complex, multidimensional integrals (eg phase-space integration of matrix elements describing particle physics processes )

- Simulation of LHC particle collisions ("events") as measured by the particle detectors. This involves:

  – Matrix element generation of collision

  – Decay of produced particles and propagation of stable particles through detector material

  – Electronic response of active detector layers, and reconstructions of signals

  – Physics analysis



Numerical estimation of circle area by taking ratio of red to red+blue points times the square's area

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

# Typical Monte Carlo Production Chain



Event Generation → Detector Simulation → Digitization → Reconstruction → Rootification

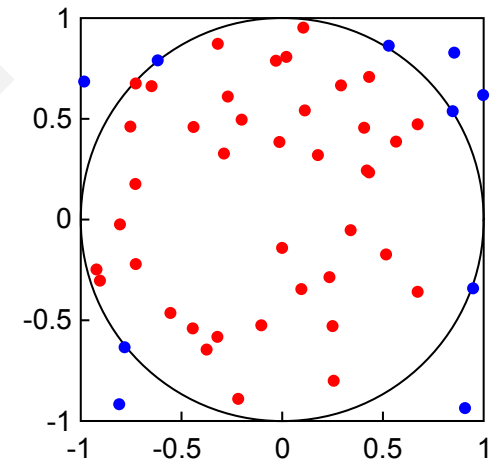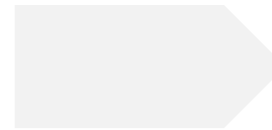EVGEN    HITS    RDO    ESD/AOD    D3PD/HIST

μ times

# Why "Monte Carlo" techniques ?

Monte Carlo (MC) techniques are computational algorithms that rely on repeated random sampling to obtain numerical results

They are used when analytical solutions are too complex or not even known

Examples:

- Numerical integration of complex, multidimensional integrals (eg phase-space integration of matrix elements describing particle physics processes )

- Simulation of LHC particle collisions ("events") as measured by the particle detectors. This involves:

  - Matrix element generation of collision

  - Decay of produced particles and propagation of stable particles through detector material

  - Electronic response of active detector layers, and reconstructions of signals

  - Physics analysis

- Simpler: error propagation and estimation of error on a measured quantity with unknown property ($\rightarrow$ next slide)



Numerical estimation of circle area by taking ratio of red to red+blue points times the square's area

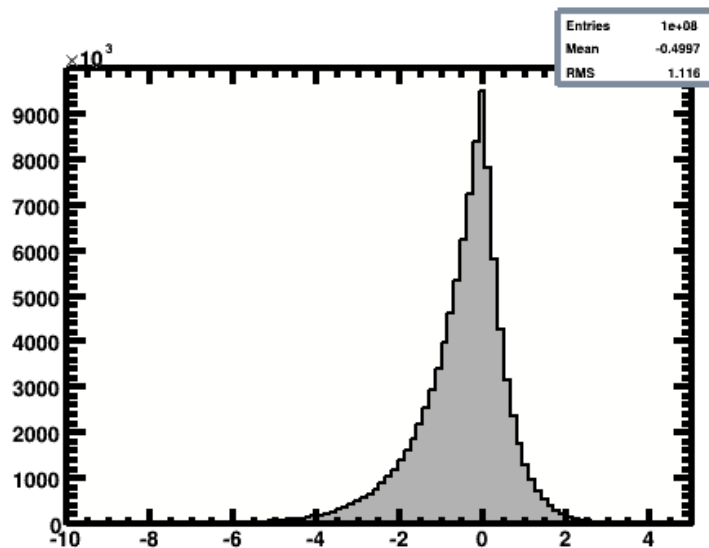Figure from: https://en.wikipedia.org/wiki/Monte_Carlo_integration

# Bootstrap method

Consider the following problem: a quantity $x$ was measured $N$ times: $x_i$ $(i = 1 \dots N)$

One wants to determined a derived quantity $y(x_1, \dots, x_N)$, and needs an error for it.

$\rightarrow$ Error propagation (remember: $\sigma_y = \frac{dy(x)}{dx}\Big|_{x=\bar{x}} \cdot \sigma_x$), but it requires to know the PDF of $x$

Assume the distribution of the measured $x_i$ looks like this:



| Entries | 1e+08 |
|---|---|
| Mean | -0.4997 |
| RMS | 1.116 |

- This **is** a PDF, and the best available information

- One can obtain a new set to "simulate" the measurements by applying *resampling with replacement*

- That is: one draws $M$ events from the ensemble allowing to re-draw the same event multiple times

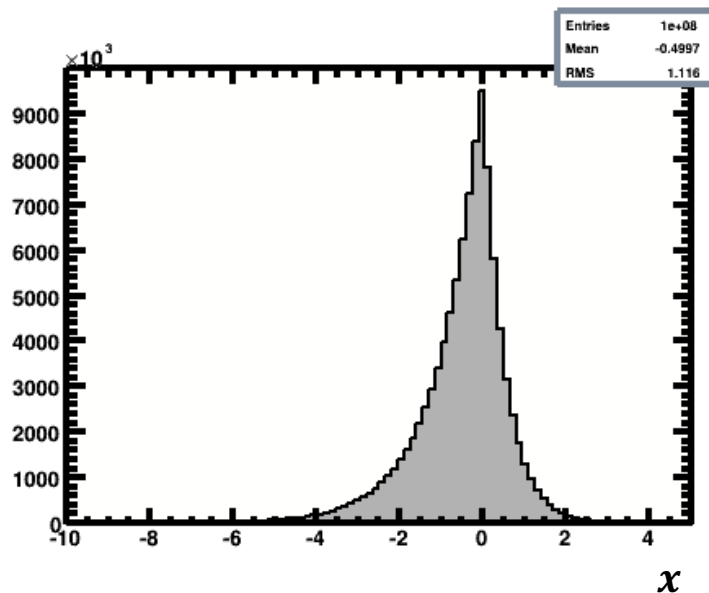- One does this many times

$\rightarrow$ **Bootstrapping**

# Bootstrap method

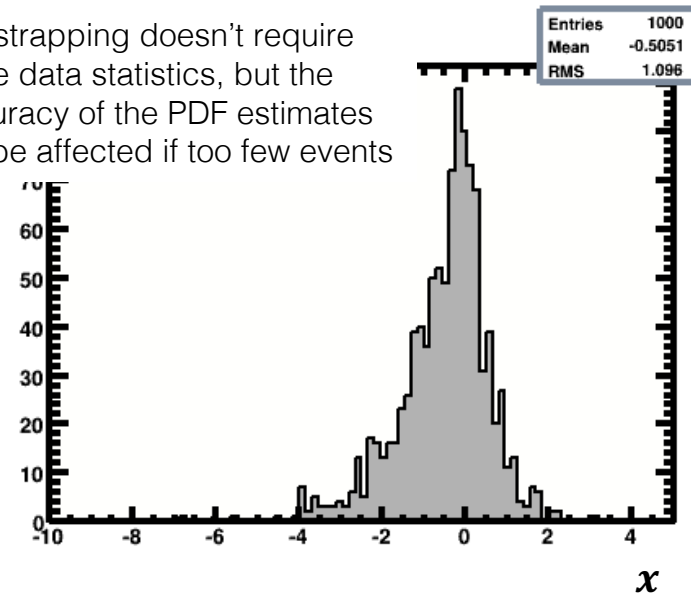Consider the following problem: a quantity $x$ was measured $N$ times: $x_i$ ($i = 1 \dots N$)

One wants to determined a derived quantity $y(x_1, \dots, x_N)$, and needs an error for it.

$\rightarrow$ Error propagation (remember: $\sigma_y = \frac{dy(x)}{dx}\Big|_{x=\bar{x}} \cdot \sigma_x$), but it requires to know the PDF of $x$

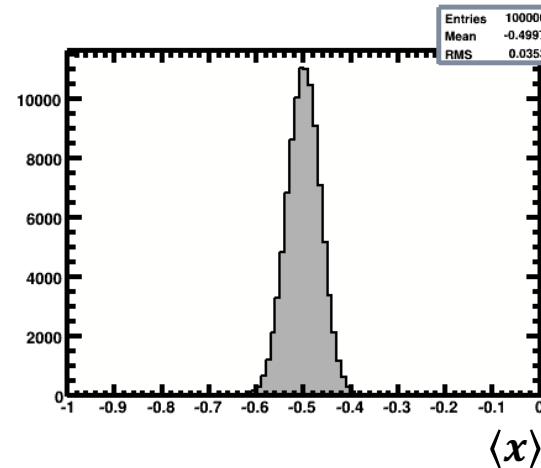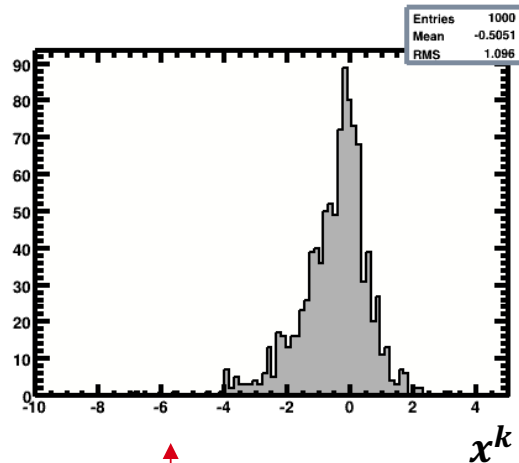Assume the distribution of the measured $x_i$ looks like this:



Boostrapping doesn't require large data statistics, but the accuracy of the PDF estimates will be affected if too few events

| Entries | 1e+08 |
| Mean | -0.4997 |
| RMS | 1.116 |

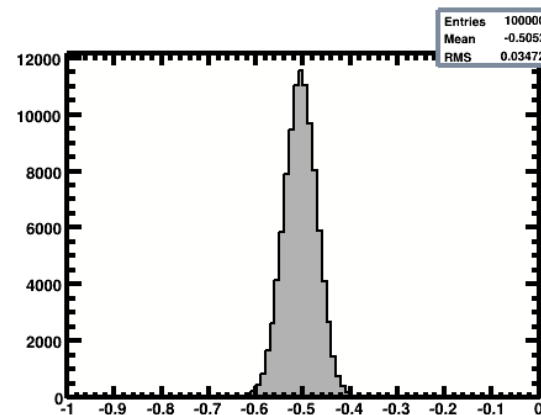| Entries | 1000 |
| Mean | -0.5051 |
| RMS | 1.096 |

# Bootstrap method — does it really work ?

Let's try with our toy example: simulate 100,000 experiments with 1,000 events each sampled from some analytic PDF (Nature's unknown truth) that we want to approximate

One example experiment $k$

Distribution of mean values among all experiments

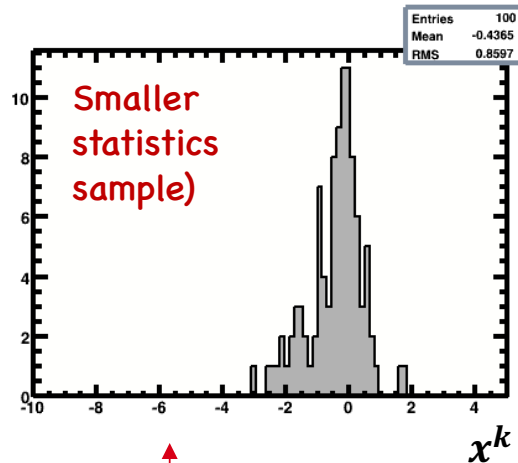Now, use **this experiment** to sample 100000 bootstrap experiments

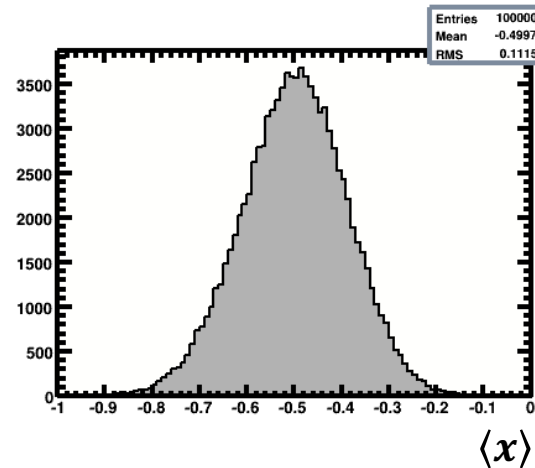Satisfying result: RMS reproduced within 1.7%

# Bootstrap method — does it really work ?

Let's try with our toy example: simulate 100,000 experiments with **100** events each sampled from some analytic PDF (Nature's unknown truth) that we want to approximate
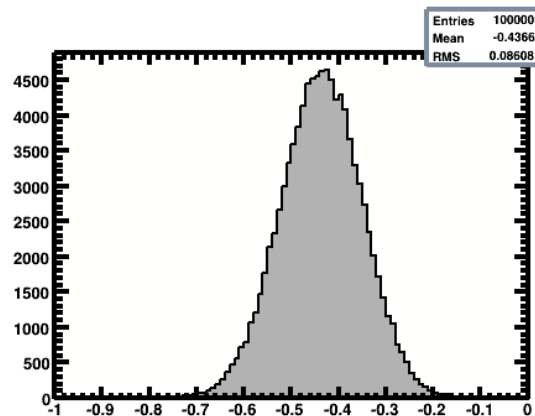
One example experiment $k$

**Smaller statistics sample)**

| Entries | 100 |
| Mean | -0.4365 |
| RMS | 0.8597 |

$x^k$

| Entries | 100000 |
| Mean | -0.4997 |
| RMS | 0.1115 |

Distribution of mean values among all experiments

$\langle x \rangle$

Now, use **this experiment** to sample 100000 bootstrap experiments

| Entries | 100000 |
| Mean | -0.4366 |
| RMS | 0.08608 |

Mediocre result: RMS reproduced within 30%

# Jackknife resampling method (also called: leave-one-out cross validation)

*Old method (~1950), basically replaced by bootstrap. Nevertheless instructive to know*

Let's again consider: a quantity $\boldsymbol{x}$ was measured $\boldsymbol{N}$ times: $\boldsymbol{x_i}\ (\boldsymbol{i = 1 \dots N})$

One wants to determine a derived quantity $\boldsymbol{y = y(x_1, \dots, x_N)}$, and needs an error for it:

- Study how $\boldsymbol{y}$ changes when *leaving out one measurement* at the time
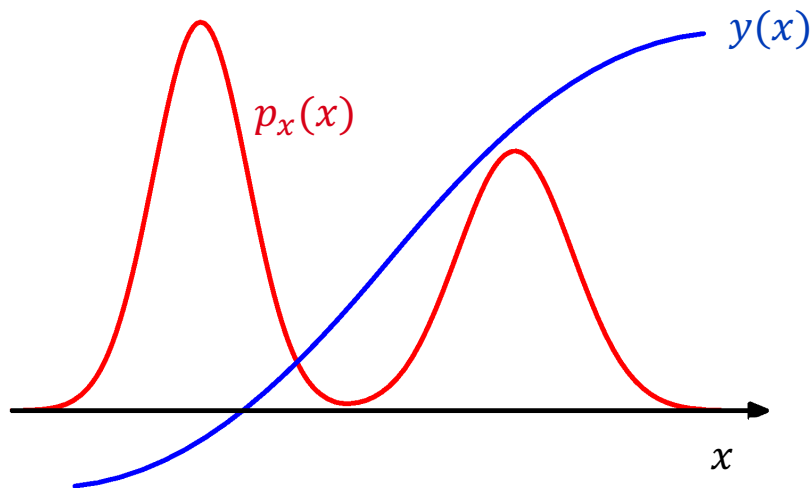
    let: $y_i = y_i(x_1, \dots x_{i-1}, x_{i+1}, \dots x_N)$,

    and compute the bias-corrected jackknife estimator of $y$: $y_i^{\text{Jack}} = Ny - (N-1)y_i$

- Plot $y_i^{\text{Jack}}$ for all $i = 1 \dots N$ and treat them as if they were independent samples of the measured quantity.

- Compute mean or variance from $y_i^{\text{Jack}}$ ensemble

# Monte Carlo (MC) integration

Want to *numerically* compute an expectation value: $E[y] = \int y(x)p_x(x)dx$



- Simplest solution: $n$-equidistant stepwise summation

- Works in 1, possibly *few* dimensions $D$

- Bad curse of dimensionality: exponential growth of $n$ with $D$

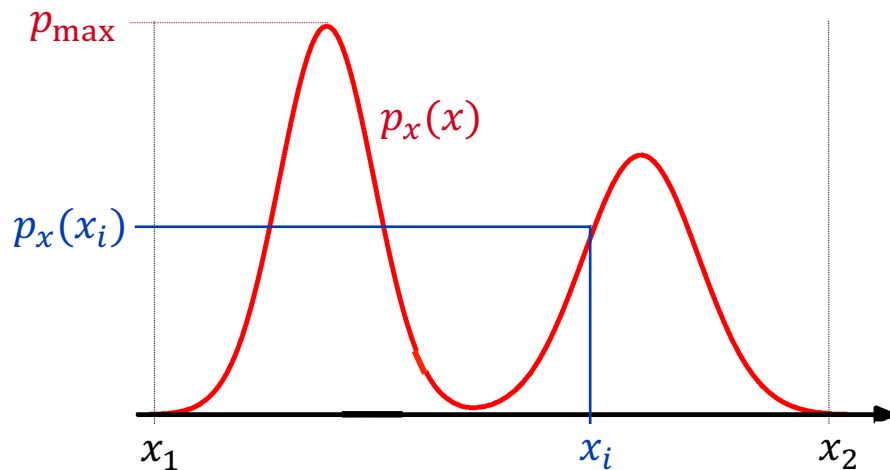- Random MC phase-space sampling converges faster for large $D$

MC integration to compute $E[y]$ requires MC sampling according to PDF $p_x(x)$

That given, one finds: $\int y(x)p_x(x)dx \approx \frac{1}{N_{\text{samples}}} \sum_{i=1}^{N_{\text{samples}}} y(x_i)$

# "Hit-or-miss" rejection sampling

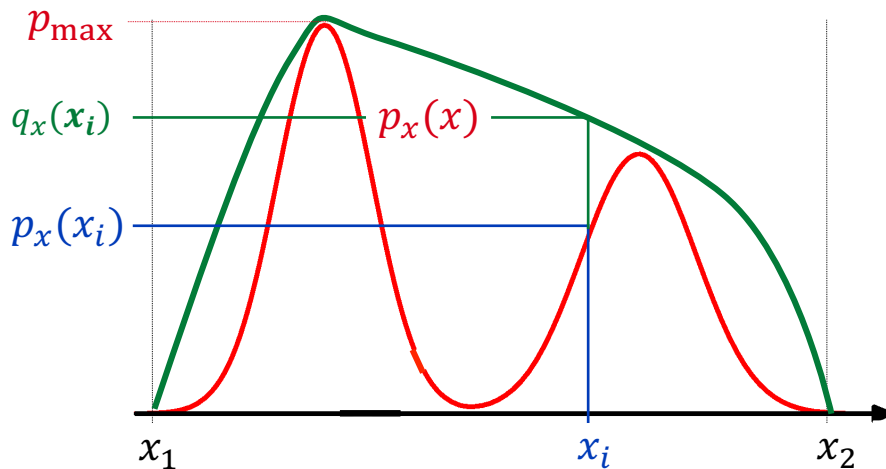Simplest way to generate random numbers (to "sample") according to PDF $p_x(x)$



1. Generate uniform random number in interval $[x_1, x_2] \to \boldsymbol{x_i}$ and $\boldsymbol{p_x(x_i)}$

2. Generate another uniform random number in interval $[0, p_{\max}] \to \boldsymbol{p_i}$

3. If $\boldsymbol{p_i < p_x(x_i)}$: **accept** $\boldsymbol{x_i}$; **else: reject**

$\boldsymbol{q_x(x)} = p_{\max}$, the (here uniform) PDF of generated $\boldsymbol{x}$ values defines *proposal distribution*

$\to$ One could be smarter to have a larger "accept" rate (efficiency)

# Rejection sampling

One can choose a (known) proposal distribution $q_x(x)$ closer to $p_x(x)$



1. Generate random number according to $q_x(x)$ in interval $[x_1, x_2] \to x_i$ and $p_x(x_i)$

2. Generate another uniform random number in interval $[0, q_x(x_i)] \to p_i$

3. If $p_i < p_x(x_i)$: **accept** $x_i$; **else: reject**

Fraction of accepted events now larger than before (there are techniques to adapt automatically the proposal distribution during the generation)

$\to$ Can be even more clever if only integration needed, no random event generation

# Markov chain Monte Carlo (MCMC) method

So far, the accuracy of the sampling depended on how closely $\boldsymbol{q_x(x)}$ follows $\boldsymbol{p_x(x)}$

This is a problem for sparsely known $\boldsymbol{p_x(x)}$ in case of complex multi-*D* structure. Every random point is chosen independently of every other one.

**Markov chain:** (eg, "random walk")    Andrey Markov (1856–1922)

- Consecutive random steps depend on previous ones in random variable space

- Allows to favor stepping into regions where $\boldsymbol{p_x(x)}$ is large

Several Markov chain algorithms: *Metropolis*, *Gibbs*, … $\rightarrow$ next pages
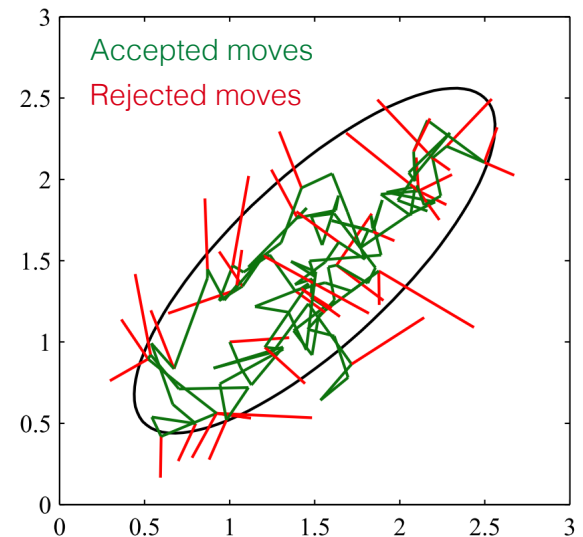
# Metropolis sampling algorithm (1953)

Autocorrelation sampling of PDF $p_x(x)$

1.  Start anywhere in (multidimensional) $x$ space and sample this point: $p_1 = p_x(x_1)$

2.  Provide proposal distribution $q_x(x_2|x_1)$ to move from $x_1 \rightarrow x_2$

    •   $q_x(x_2|x_1)$ could be Gaussian with appropriate metric in $x$ space to cover full space

    •   Accept $x_2$ if: $p_2 > p_1$ else: according to probability $p_2/p_1$

    •   If not accepted, set $x_2 = x_1$ (no walk)

3.  Iterate step 2. for $x_3$ vs. $x_2$, etc.

Sample points $x$ will wander closer and closer to the peak of the PDF, still jumping enough from time to time to sample the whole space

(Algorithm requires sufficient iterations. Test by checking stability of derived result, or by comparing several sampling ensembles obtained with different start values)
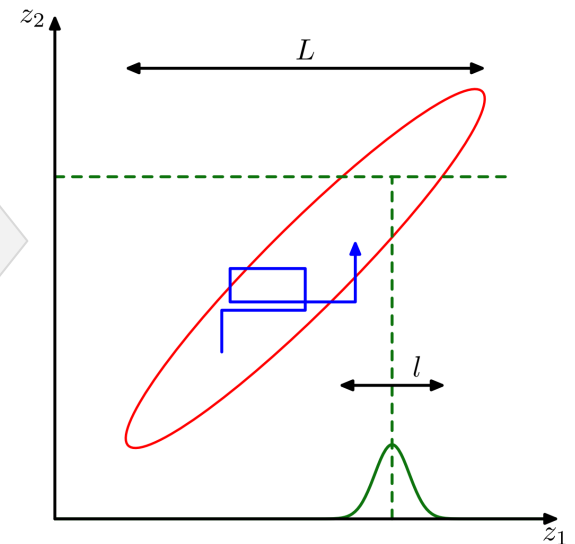


This subfigure from PRML, Bishop (2006)
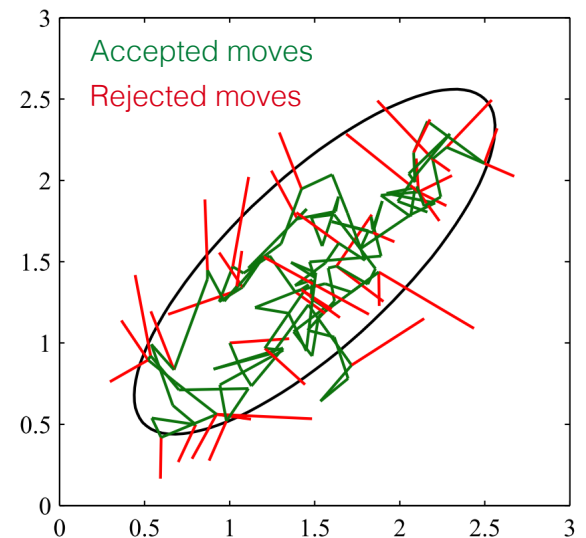
# Gibbs sampling algorithm <span>(1984)</span>

A method with no rejection:

1. Instead of moving along all dimensional components (and reject low-probability moves), the Gibbs sampler moves along 1 component according to the PDF conditioned on all other components.

2. Cycle through all components

Markov chain Monte Carlo usually converge fast and, if metric well chosen, cover the full space

However: care needs to be taken as the sample points $x$ are correlated (with either sampling method, although there exist tricks to reduce the correlations): it depends on the application whether or not this is an issue





This subfigure from PRML, Bishop (2006)
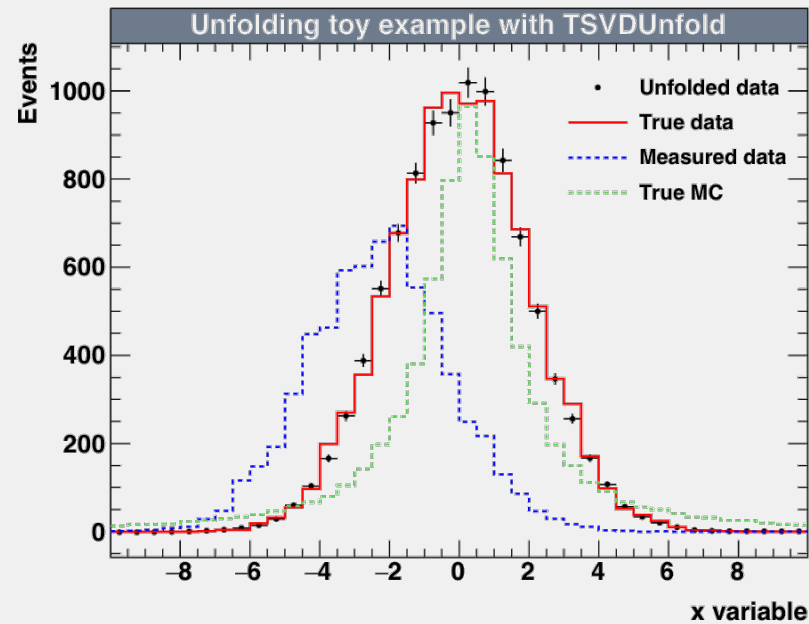
# Data unfolding

# Data unfolding — introduction

"Unfolding" means correcting measured data for any effects related to the measurement device. The unfolded data can be directly compared to theory or among experiments

Consider a measured histogram $\boldsymbol{y}^{\text{data}} = \{y_1^{\text{data}}, \ldots, y_m^{\text{data}}\}$, a corresponding Monte Carlo histogram $\boldsymbol{y}^{\text{MC}}$ of the same process as the data that underwent full detector simulation, its *truth* distribution (ie, before detector simulation) $\boldsymbol{x}^{\text{MC}} = \{x_1^{\text{MC}}, \ldots, x_n^{\text{MC}}\}$, and the $m \times n$ matrix $\boldsymbol{A}^{\text{MC}}$ obtained from MC that describes the "smearing" process due to the measurement:

$$\boldsymbol{A}^{\text{MC}} \cdot \boldsymbol{x}^{\text{MC}} = \boldsymbol{y}^{\text{MC}}$$

Note that in general $\boldsymbol{y}^{\text{data}} \neq \boldsymbol{y}^{\text{MC}}$ (the physics leading to $\boldsymbol{y}^{\text{data}}$ is what we want to measure), but we assume $\boldsymbol{A}^{\text{MC}} = \boldsymbol{A}^{\text{data}}$ (we know the detector response).
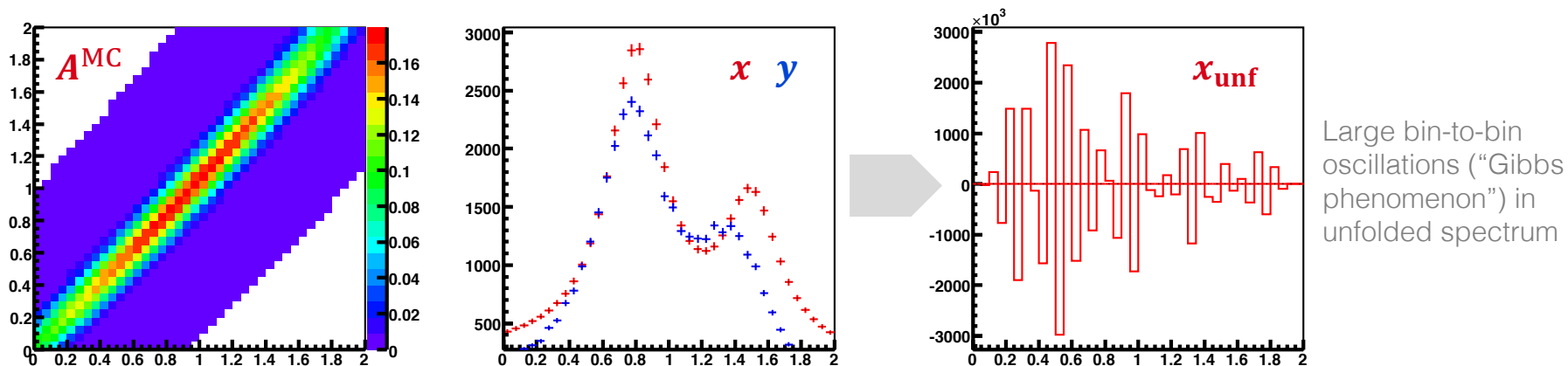
Hence, to obtain the truth information $\boldsymbol{x}^{\text{data}}$, one "just" needs to invert $\boldsymbol{A}$:

$$\boldsymbol{x}^{\text{data}} = \left(\boldsymbol{A}^{\text{MC}}\right)^{-1} \cdot \boldsymbol{y}^{\text{data}}$$

This is where the trouble begins …

# Data unfolding — introduction

The distribution $\boldsymbol{y}^{\text{data}}$ and the matrix $\boldsymbol{A}^{\text{MC}}$ have finite statistics. An attempt to solve the problem directly and "exactly" will end up looking like this:



Large bin-to-bin oscillations ("Gibbs phenomenon") in unfolded spectrum

The poor solution, bin-by-bin corrections, $x_i^{\text{data}} = \frac{x_i^{\text{MC}}}{y_i^{\text{MC}}} \cdot y_i^{\text{data}}$, only works if $\boldsymbol{A}^{\text{MC}}$ is square and ~diagonal so that the ratio $x_i^{\text{MC}}/y_i^{\text{MC}}$ corrects for mainly efficiency effects, or if $y_i^{\text{data}} \cong y_i^{\text{MC}}$.

A better solution is to regularise the matrix inversion problem …

# Data unfolding — regularisation

Regularisation damps the oscillations, by suppressing statistically insignificant bins in the data distribution and response matrix.

In simplified form, one can write the unfolding problem as a minimisation of

$$\chi^2(x^{\text{data}}) = \left(A^{\text{MC}} \cdot x^{\text{data}} - y^{\text{data}}\right)^T \left(A^{\text{MC}} \cdot x^{\text{data}} - y^{\text{data}}\right) + \tau \cdot \left(Cx^{\text{data}}\right)^T \left(Cx^{\text{data}}\right)$$

where $C$ is a matrix and $Cx^{\text{data}}$ is the sum of squares of the 2nd derivative of $x^{\text{data}}$

Minimising $\chi^2$ wrt. the first term only corresponds to the (bad) exact inversion solution. The second term regularises the inversion by damping the oscillations.

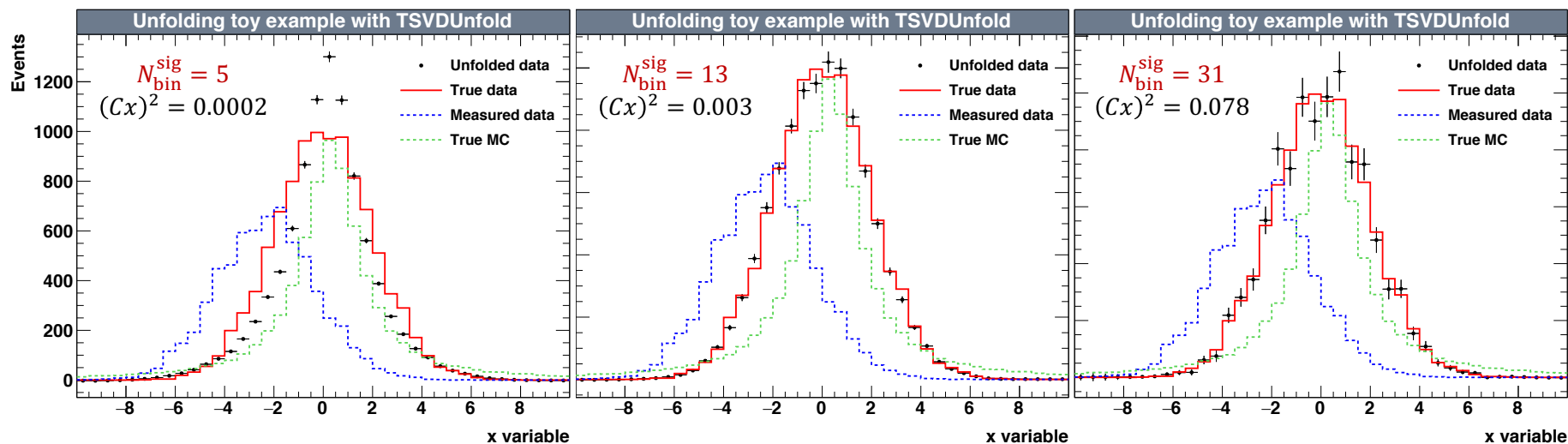The parameter $\tau$ regulates the strength of the damping:

- If $\tau$ too small $\rightarrow$ oscillations

- If $\tau$ too large $\rightarrow$ information in $x^{\text{data}}$ is suppressed
  ($x^{\text{data}}$ becomes too "smooth" and will be biased towards $x^{\text{MC}}$)

- The right choice captures the significant information and discards the rest

# Data unfolding — example

The parameter $\tau$ regulates the strength of the damping:

- If $\tau$ too small $\rightarrow$ oscillations

- If $\tau$ too large $\rightarrow$ information in $\boldsymbol{x}^{\mathrm{data}}$ is suppressed
  ($\boldsymbol{x}^{\mathrm{data}}$ becomes too "smooth" and will be biased towards $\boldsymbol{x}^{\mathrm{MC}}$)

- The right choice captures the significant information and discards the rest

# Folding versus unfolding

Unfolding is an ill-defined problem which necessarily leads to some obstruction of information in the data and transfer of statistical uncertainty to a systematic one after regularisation (this is similar to a non-parametric fit to data)

Technically simpler and mathematically well defined is the folding of a theoretical prediction $x^{\text{theo}}(\boldsymbol{\theta})$, depending on a set of parameters $\boldsymbol{\theta}$, through the detector response and direct comparison with the measured data. It allows the statistical test:

$$\chi^2(x^{\text{theo}}(\theta)) = \left(A^{\text{MC}} \cdot x^{\text{theo}}(\theta) - y^{\text{data}}\right)^T \left(A^{\text{MC}} \cdot x^{\text{theo}}(\theta) - y^{\text{data}}\right)$$

Folding requires that the experiments either perform the test, or publish $A^{\text{MC}}$ and $y^{\text{data}}$

Folding does not allow a model-independent combination or comparison among experiments. In most case, unfolding is the only viable solution for easy and long-term use of the experimental results.
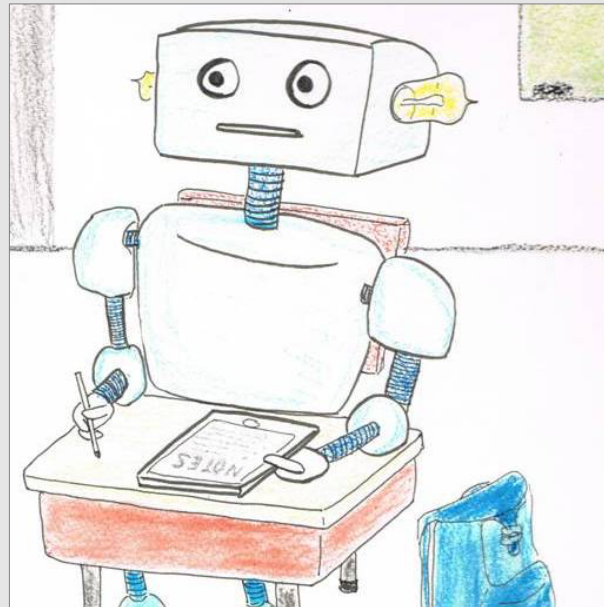
Bootstrapping methods allow to straightforwardly re-sample measured data for the purpose of error propagation

Brief introduction to Monte Carlo integration and the sampling of random data according to any arbitrary PDF

Markov-Chain Monte Carlo integration is a very effective method that "automatically" samples the important regions (where the PDF is large) more often than tails. Try yourself!

Unfolding is a delicate mathematical operation that requires careful regularisation. Folding can help in some cases.

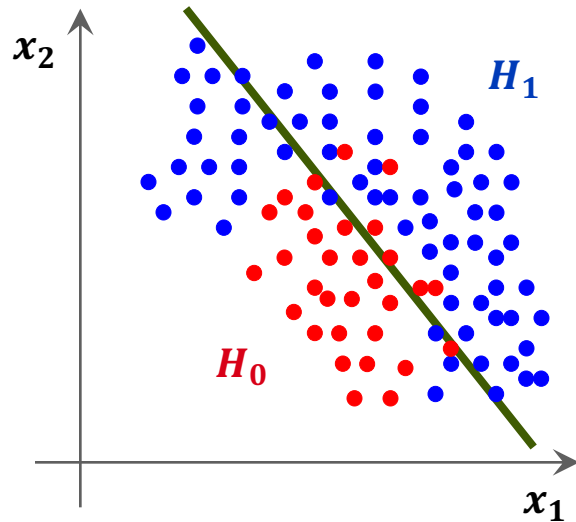# Multivariate techniques and machine learning

# Event classification

Suppose data sample with two types of events: $H_0$, $H_1$

- We have found discriminating input variables $x_1$, $x_2$, …

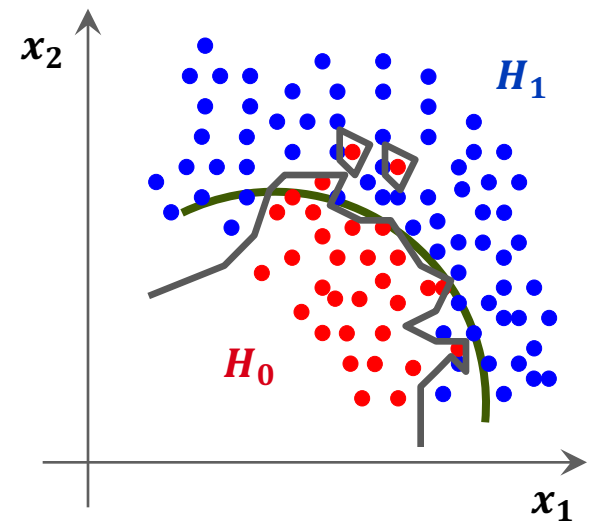- What decision boundary should we use to select events of type $H_1$?



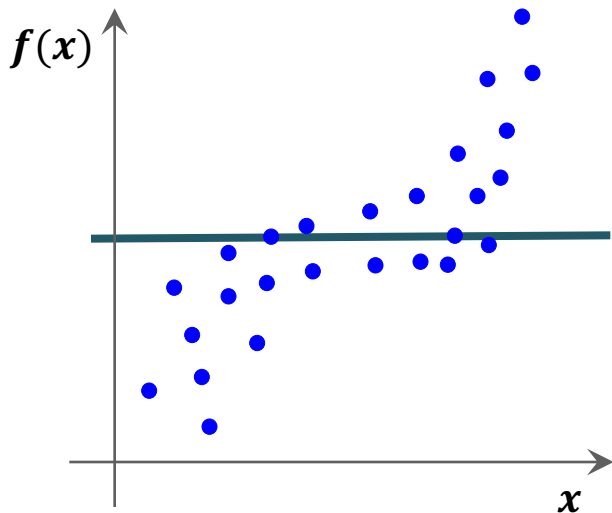Low variance (stable), high bias methods          High variance, small bias methods
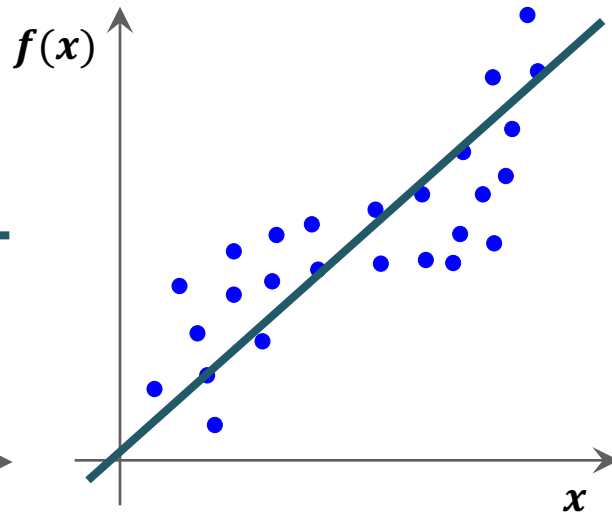
# Parameter regression

How to estimate a *functional behaviour* from a set of measurements? HEP examples:

- Energy deposit in a the calorimeter, distance between overlapping photons, …

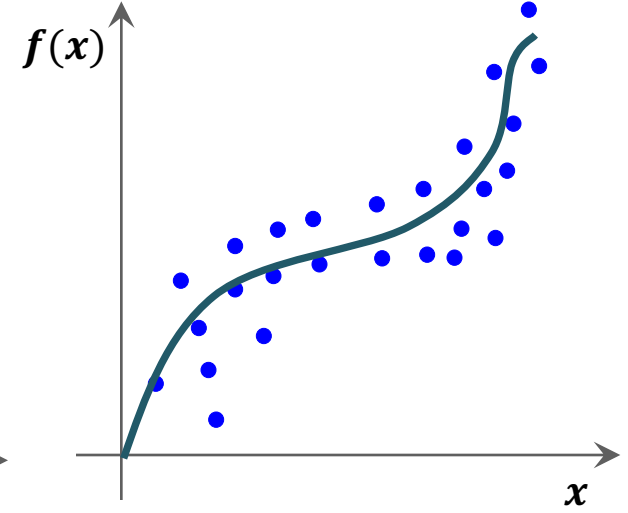- Entry location of a particle in the calorimeter or on a silicon pad, …



| Constant ? | Linear function ? | A non-linear one ? |

Looks trivial? What if we have many input variables?

Note: the goal is *not to fit* given data, but to learn $f(x)$ vs. $x$ to *predict* **target** $f(x)$ for new measurements $x$

These are the simplest applications of statistical machine learning (ML).
Most particle physics utilisations so far fall into this category

However, there is no limit of use cases for complex machine learning…



Image recognition



Table tennis (KUKA advertisement)



Autonomous driving (Google)

+ speech recognition, language understanding, syntax parsing, translation, face recognition, road hazard detection, …
and: **PHYSICS !**

…as long as the ML algorithms are smart and efficient enough, there is sufficient computing power, and a complete set of training data

Also in particle physics we can apply ML to more complex problems. Among these: track reconstruction, calibration, tuning

*Machine learning is giving computers the ability to learn without explicitly programming them* (Arthur Samuel, 1959)

It is fundamentally different from applying a set of fixed rules (a "program") to solve a problem
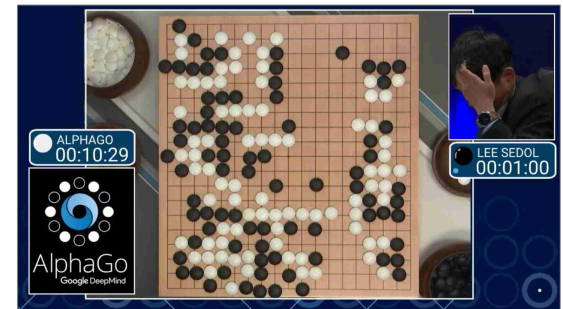
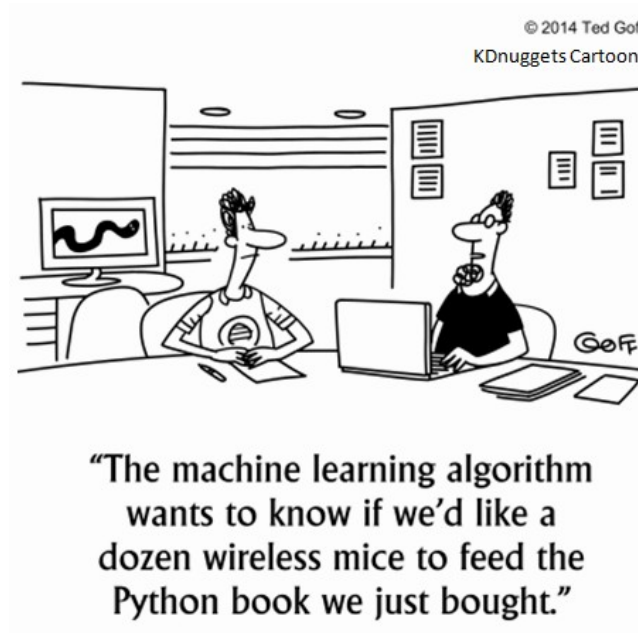Chess computers in the 1980's and 90's

Alpha Go in 2016



Rule-based programming

Machine learning

(Deep reinforcement learning)

Not so long ago, real-life artificial intelligence used to be like this:



Such things still happen, but the improvements have nevertheless been astounding

# More tomorrow…