

# Random number generators and application

Marcin Chrząszcz  
mchrzasz@cern.ch



Experimental Methods in Particle Physics,  
1 November, 2018

# Random and pseudorandom numbers

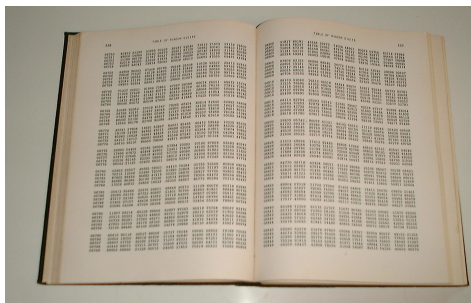
John von Neumann:

“Any one who considers arithmetical methods of producing random digits is, of course, in a state of sin. For, as has been pointed out several times, there is no such thing as a random number — there are only methods to produce random numbers, and a strict arithmetic procedure of course is not such a method.”

- ⇒ Random number: a given value that is taken by a random variable
  - by definition cannot be predicted.
- ⇒ Sources of truly random numbers:
  - Mechanical
  - Physical
- ⇒ Disadvantages of physical generators:
  - To slow for typical applications, especially the mechanical ones!
  - Not stable; small changes in boundary conditions might lead to completely different results!

# Random numbers - history remark

⇒ In the past there were books with random numbers:



⇒ It's obvious that they didn't become very popular ;)

⇒ This methods are coming back!

→ Storage device are getting more cheap and bigger (CD, DVD).

→ 1995: G. Marsaglia, 650MB of random numbers, "White and Black Noise".

# Pseudorandom numbers

- ⇒ Pseudorandom numbers are numbers that are generated accordingly to strict mathematical formula.
- ↪ Strictly speaking they are non random numbers, how ever they have all the statistical properties of random numbers.
- ↪ Discussing those properties is a wide topic so let's just say that without knowing the formula they are generated by one cannot say if those numbers are random or not.
- ⇒ Mathematical methods of producing pseudorandom numbers:
  - Good statistical properties of generated numbers.
  - Easy to use and fast!
  - Reproducible!
- ⇒ Since mathematical pseudorandom genrators are dominantly:  
pseudorandom  $\rightsquigarrow$  random.

# Middle square generator; von Neumann

⇒ The first mathematical generator (middle square) was proposed by von Neumann (1964).

↪ Formula:  $X_n = \lfloor X_{n-1}^2 \cdot 10^{-m} \rfloor - \lfloor X_{n-1}^2 \cdot 10^{-3m} \rfloor \cdot 10^{2m}$

↪ where  $X_0$  is a constant (seed),  $\lfloor \cdot \rfloor$  is the cut-off of a number to integer.

⇒ Example:

Let's put  $m = 2$  and  $X_0 = 2045$ :

$$\begin{array}{ccc} \text{↪ } X_0^2 = & \underbrace{04} & 1820 & \underbrace{25} \\ & \text{rej} & & \text{rej} \end{array} \quad \Rightarrow X_1 = 1820$$

$$\begin{array}{ccc} \text{↪ } X_1^2 = & \underbrace{03} & 3124 & \underbrace{00} \\ & \text{rej} & & \text{rej} \end{array} \quad \Rightarrow X_1 = 3124$$

↪ Simple generator but unfortunately quite bad generator. Firstly the sequences are very short and strongly dependent on the  $X_0$  number.

# Linear generators Lecture2/Linear\_gen1

⇒ This was a first generator written and it's a good example how to not write generators.

⇒ It's highly non stable!

```
mchrzasz-ThinkPad-W530% python gen.py 14714 4
21650.0
46872.0
219698.0
4826721.0
2329723538.0
5.42761170924e+14
2.94589685716e+25
8.67830820626e+46
7.53130325698e+89
Traceback (most recent call last):
  File "gen.py", line 29, in <module>
    sys.exit(main())
  File "gen.py", line 22, in main
    tmp=X0**2
OverflowError: (34, 'Numerical result out of range')
```

# Linear generators

⇒ General equation:

$$X_n = (a_1 X_{n-1} + a_2 X_{n-2} + \dots + a_k X_{n-k} + c) \bmod m,$$

↷ where  $a_i, c, m$  are parameters of a generator (integer numbers).

↷ Generator initialization ⇔ setting those parameters.

⇒ Very old generators. (often used in Pascal, or first C versions):

$$k = 1 : X_n = (aX_{n-1} + c) \bmod m,$$

$$c = \begin{cases} = 0, & \text{multiplicative generator} \\ \neq 0, & \text{mix generator} \end{cases}$$

⇒ The period can be achieved by tuning the seed parameters:

$$P_{\max} = \begin{cases} 2^{L-2}; & \text{for } m = 2^L \\ m - 1; & \text{for } m = \text{prime number} \end{cases}$$

# Shift register generator

⇒ General equation:

$$b_n = (a_1 X_{n-1} + a_2 X_{n-2} + \dots + a_k X_{n-k} + c) \bmod 2,$$

where  $a_i \in \{0, 1\}$

⇒ Super fast and easy to implement due to:  $(a + b) \bmod 2 = a \text{ xor } b$

| a | b | a xor b |
|---|---|---------|
| 0 | 0 | 0       |
| 1 | 0 | 1       |
| 0 | 1 | 1       |
| 1 | 1 | 0       |

⇒ Maximal period is  $2^k - 1$ .

⇒ Example (Tausworths generator):

$a_p = a_q = 1$ , other  $a_i = 0$  and  $p > q$ . Then:  $b_n = b_{n-p} \text{ xor } b_{n-q}$

⇒ How to get numbers from bits (for example):

$$U_i = \sum_{j=1}^L 2^{-j} b_{is+j}, \quad s < L.$$



# Non linear generators

⇒ The natural solutions to problems of linear generators are the non-linear generators (second part of 1980s).

⇒ Eichenauera i Lehna (1986):

$$X_n = (aX_{n-1}^{-1} + b) \bmod m,$$

⇒ Eichenauera-Hermann (1993)

$$X_n = [a(n + n_0) + b]^{-1} \bmod m,$$

⇒ L. Blum, M. Blum, Shub (1986):

$$X_n = X_{n-1}^2 \bmod m,$$

→ Very popular in cryptography.

⇒ Pros and cons:

- They all pass all statistical tests.
- Much slower than linear generators.

# RANLUX generator

⇒ All described generators are based on some mathematical algorithms and recursion. The typical scheme is of constructing a MC generator:

- Think of a formula that takes some initial values.
- Generate large number of random numbers and put them through statistical tests.
- If the test are positive we accept the the generator.

⇒ Now let's think: why the hell numbers obtained that way are showing some random number properties?

# RANLUX generator

⇒ All described generators are based on some mathematical algorithms and recursion. The typical scheme is of constructing a MC generator:

- Think of a formula that takes some initial values.
- Generate large number of random numbers and put them through statistical tests.
- If the test are positive we accept the the generator.

⇒ Now let's think: why the hell numbers obtained that way are showing some random number properties? There is no science behind it, it's pure luck!

⇒ M.Luscher (1993) hep-lat/9309020

⇒ Generator RANLUX based on Kolomogorow entropy and Lyapunov exponent. **Effectively we are building inside the generator the chaos theory.**

⇒ RANLUX and Mersenne Twister (TRandom1, TRandom3) are the 2 most powerful generators in the world that passed every known statistical test.

# Chaos theory in a nut shell

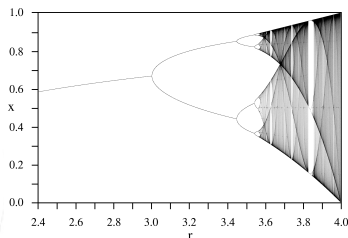
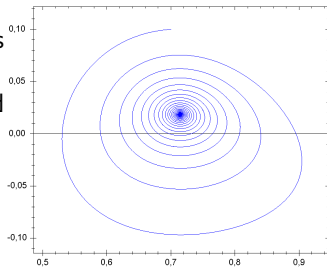
⇒ We know that the solution of classical systems is described by trajectory in phase spaces. Now the problem with this picture starts to be when around one point in this phase space we are getting more and more trajectories that are drifting a part later on.

⇒ The Lyapunov exponent tells us how a two solutions drift apart with time:

$$|\delta X(t)| \approx e^{\lambda t} |\delta X_0|$$

⇒ Kolomogorow entropy:

$$h_K = \int_P \lambda d\mu$$



## Reverting the c.d.f.

- ⇒ Let  $U$  be a random variable from  $\mathcal{U}(0, 1)$
- ⇒ Now let  $F$  be a non decreasing function such that:

$$F(-\infty) = 0 \quad F(\infty) = 1$$

then:

$$X = F^{-1}(U)$$

has a p.d.f. distribution with a c.d.f. function of  $F$ .

⇒ Prove:

$$F(x) = \mathcal{P}(U \leq F(x)) = \mathcal{P}(F^{-1}(U) \leq x) = \mathcal{P}(X \leq x) \quad \square$$

⇒ So it looks very simple if  $x_1, X_2, \dots, X_n$  are random variables from  $\mathcal{U}(0, 1)$  then:  $\{X_i = F^{-1}(x_i)\}, i = 1, \dots, n$  is the sequence that has a c.d.f. distribution of  $F$ .

# Reverting the c.d.f., examples

- ⇒ The exponential distribution:  $E(0, 1)$ .
- ⇒ The p.d.f.:  $\rho(X) = e^{-X}$ ,  $X \geq 0$ .
- ⇒ The c.d.f.:  $F(x) = \int_0^x e^{-X} dX = 1 - e^{-x}$ .
- ⇒ Now let  $R \in \mathcal{U}(0, 1)$ :  $R = F(X) = 1 - e^{-X} \longrightarrow X = -\ln(1 - R)$
- ⇒ Now we can play a trick: if  $R \in \mathcal{U}(0, 1)$  then  $1 - R$  also in  $\mathcal{U}(0, 1)$ .
- ⇒ In the end we get:  $X = -\ln(R)$

⇒ Use the reverting to generate the following distributions:

- E 6.1  $\rho(X) = \frac{c}{X}$  on the interval  $[a, b]$ , where  $0 < a < b < \infty$ .
- E6.2 The Breit-Wigner function:

$$\rho_{\theta, \lambda}(X) = \frac{\lambda}{\pi} \frac{1}{(X - \theta)^2 + \lambda^2}$$

Hit: First do  $C(0, 1)$  then transform the variables.

## Reverting the c.d.f., general case



Lets assume:  $F$  - non decreasing function such that:  $F(-\infty) = 0$  and  $F(\infty) = 1$ .  
Then a random variable  $X$ :

$$X = \inf\{x : U \leq F(x)\}$$

has a distribution with a c.d.f. of  $F$

```
int gen_discrete(double rn, double *p){
// generating a discrete distribution
// P{X = k } = p[k], k=0,1,...
// rn random number from U(0,1)
int k=0;
double sum=p[0];
while(suma < rn) suma +=p[++k]
return k;
}
```

⇒ E6.3 Using the above example please generate the p.d.f. accordingly to:

$$\mathcal{P}(X = k) = p_k = A \sin\left(\frac{\pi}{10} \left[k + \frac{1}{2}\right]\right)$$

where  $A$  is the normalization constant (calculate it!). From the generated numbers make a histogram and compare to the exact function.

# Reverting the c.d.f., pros and cons

⇒ Pros:

- Very accurate method.
- Fast and easy.
- To generate one random number from c.d.f. you need one random number from  $\mathcal{U}(0, 1)$ .

⇒ Cons:

- Usually we require that the c.d.f. is revertible analytically( small number of functions).
- If we use the numerical method of reverting the function, then it's much slower and less accurate.
- Super hard fro multidimensional distributions.

⇒ Mathematical digression: Numerical reverting the c.d.f.:

- We look for the zero of the function:  $F_u(X) \equiv F(X) - U$ , where  $U \in \mathcal{U}(0, 1)$ .
- $X_0 = F^{-1}(U)$



# HEP simulation

⇒ There is some ambiguity what particle physicist call MC. Normally those are mathematical theories but when we say MC we usually mean MC simulation of a physics process. ⇒ There are plenty of things that need to be simulated:

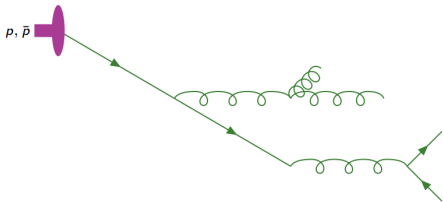
$p, \bar{p}$  

$t = -\infty$ , incoming protons

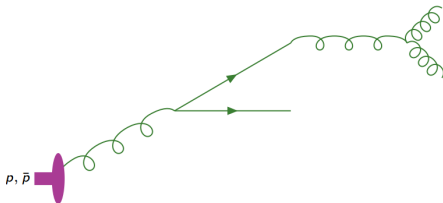
$p, \bar{p}$  

# HEP simulation

⇒ There is some ambiguity what particle physicist call MC. Normally those are mathematical theories but when we say MC we usually mean MC simulation of a physics process. ⇒ There are plenty of things that need to be simulated:

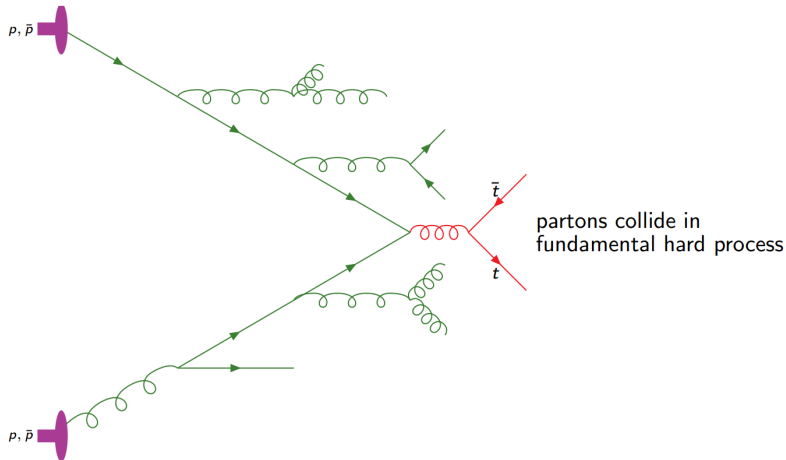


partons from the protons radiate



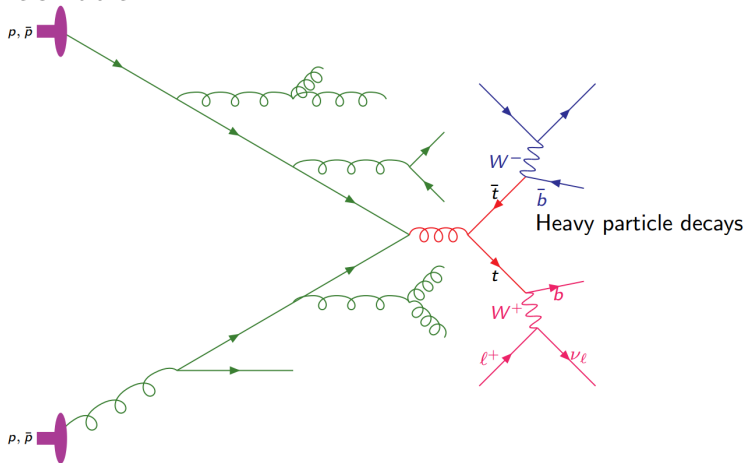
# HEP simulation

⇒ There is some ambiguity what particle physicist call MC. Normally those are mathematical theories but when we say MC we usually mean MC simulation of a physics process. ⇒ There are plenty of things that need to be simulated:



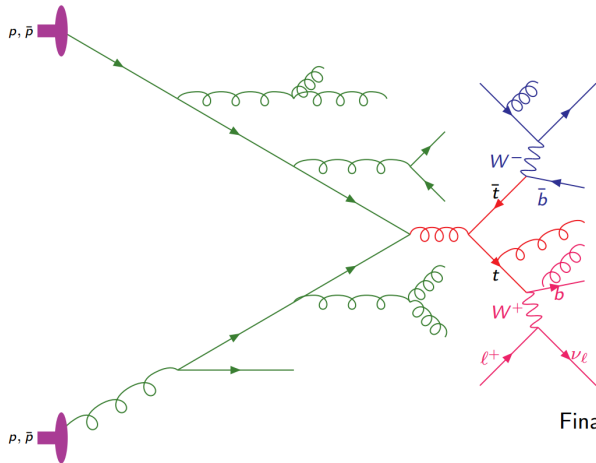
# HEP simulation

⇒ There is some ambiguity what particle physicist call MC. Normally those are mathematical theories but when we say MC we usually mean MC simulation of a physics process. ⇒ There are plenty of things that need to be simulated:



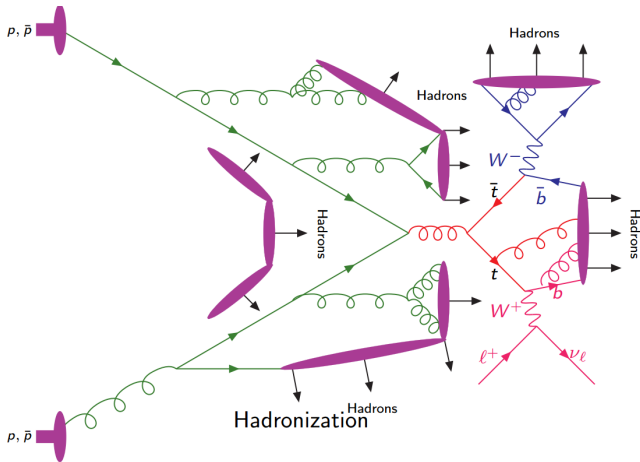
# HEP simulation

⇒ There is some ambiguity what particle physicist call MC. Normally those are mathematical theories but when we say MC we usually mean MC simulation of a physics process. ⇒ There are plenty of things that need to be simulated:



# HEP simulation

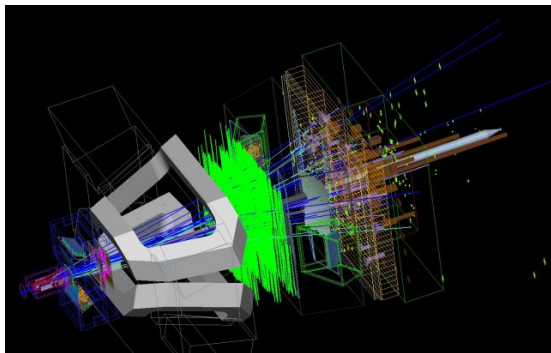
⇒ There is some ambiguity what particle physicist call MC. Normally those are mathematical theories but when we say MC we usually mean MC simulation of a physics process. ⇒ There are plenty of things that need to be simulated:



# Detector simulation

- ⇒ Things do not get simpler on the detector side simulation.
- ⇒ Lots of effects need to be taken into account:

- Bremsstrahlung
- Interactions with different detector materials
- Particle identification
- Showers



- ⇒ Example of generators:
- FLUKA
- Geant

# Method of Moments

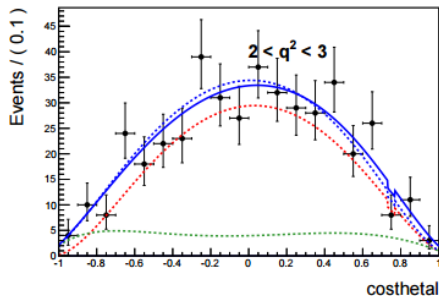
⇒ Now real cool things!

⇒ Let's consider we want to study a rare decay:  $B^\pm \rightarrow K^\pm \mu \mu$ . The decay is described by the following PDF:

$$\frac{1}{\Gamma} \frac{d^2\Gamma}{dq^2 d\cos\theta_l} = \frac{3}{4}(1 - F_H)(1 - \cos^2\theta_l) + F_H/2 + A_{FB} \cos\theta_l$$

⇒ PDF by construction is normalized:  $\int_{-1}^1 \frac{1}{\Gamma} \frac{d^2\Gamma}{dq^2 d\cos\theta_l} = 1$

- Normally we do a likelihood fit and we are done.
- There is a second way!





# Method of Moments

⇒ Let's calculate the integrals:

$$\int_{-1}^1 \frac{1}{\Gamma} \frac{d^2\Gamma}{dq^2 d \cos \theta_l} \cdot \cos \theta_l = \frac{2}{3} A_{FB}$$

$$\int_{-1}^1 \frac{1}{\Gamma} \frac{d^2\Gamma}{dq^2 d \cos \theta_l} \cdot \cos^2 \theta_l = \frac{1}{5} + \frac{2F_H}{15}$$

⇒ So we can get our parameters that we searched for by doing a integration. So now what?

# Method of Moments

⇒ Let's calculate the integrals:

$$\int_{-1}^1 \frac{1}{\Gamma} \frac{d^2\Gamma}{dq^2 d \cos \theta_l} \cdot \cos \theta_l = \frac{2}{3} A_{FB}$$

$$\int_{-1}^1 \frac{1}{\Gamma} \frac{d^2\Gamma}{dq^2 d \cos \theta_l} \cdot \cos^2 \theta_l = \frac{1}{5} + \frac{2F_H}{15}$$

⇒ So we can get our parameters that we searched for by doing a integration. So now what?

⇒ Well nature is the best random number generator so let's take the data and treat and calculate the integral estimates:

$$\int_{-1}^1 \frac{1}{\Gamma} \frac{d^2\Gamma}{dq^2 d \cos \theta_l} \cdot \cos \theta_l = \frac{2}{3} A_{FB} = \frac{1}{N} \sum_{i=1}^N \cos \theta_{l,i}$$

$$\int_{-1}^1 \frac{1}{\Gamma} \frac{d^2\Gamma}{dq^2 d \cos \theta_l} \cdot \cos^2 \theta_l = \frac{1}{5} + \frac{2F_H}{15} = \frac{1}{N} \sum_{i=1}^N \cos^2 \theta_{l,i}$$

# Method of Moments

⇒ So what did we do?

- We have just estimated a parameters of interests without using any fit!!

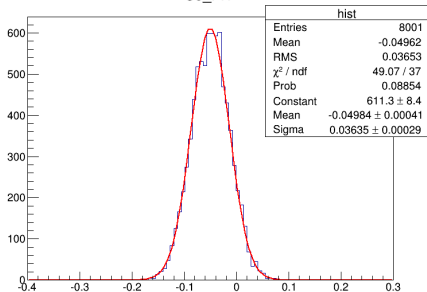
⇒ Pros and cones of method of moments:

- Are very immune to bias.
- Do not suffer from boundary problems.
- Require less statistic to work then likelihood fit.
- They always have a Gaussian error.
- Estimator has a larger uncertainty.

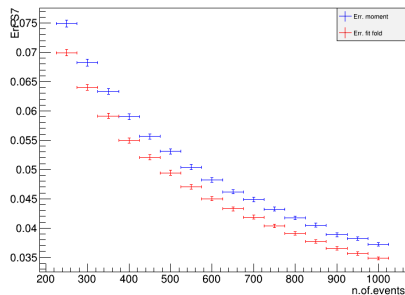
# Method of Moments, uncertainty estimator

⇒ It can be proven that Method of Moments estimator converges slower than the maximum likelihood fit.

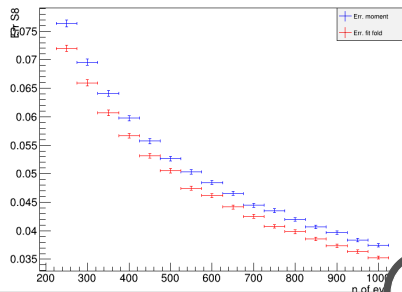
S8\_FIT



S7 Error



S8 Error



# Other application of MC - testing your analysis

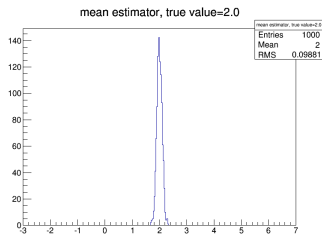
⇒ Probably the biggest application of MC methods in HEP are validations of your experimental methodology. The procedure is as follows:

- Define your analysis methodology: selection, efficiency corrections, parameters you want to measure.
- Simulate an assembly of simulation events for different values of parameters you want to measure.
- Do the analysis on this pseudo data.
- See if you are getting back what you have simulated.

# Testing your analysis, Lecture2/Test\_met

⇒ Probably the biggest application of MC methods in HEP are validations of your experimental methodology. The procedure is as follows:

- Define your analysis methodology: selection, efficiency corrections, parameters you want to measure.
- Simulate an assembly of simulation events for different values of parameters you want to measure.
- Do the analysis on this pseudo data.
- See if you are getting back what you have simulated.



# Wrap up

⇒ Things to remember:

- Computer cannot produce random numbers, only pseudorandom numbers.
- We use pseudorandom numbers as random numbers if they are statistically acting the same as random numbers.
- Linear generators are not commonly used nowadays.
- State of the art generators are the ones based on Kolomogorows theorem.
- MC methods used to simulate physics process, detector response and validating the estimators.

# Backup