



DE LA RECHERCHE À L'INDUSTRIE

UNCERTAINTY QUANTIFICATION

A broad introduction on statistical aspects

HPC and Uncertainty Treatment with Open TURNS and Uranie, 10/05/2021

Jean-Baptiste Blanchard (jean-baptiste.blanchard@cea.fr), Fabrice Gaudier
(fabrice.gaudier@cea.fr)

CEA DES/ISAS/DM2S/STMF/LGLS SACLAY

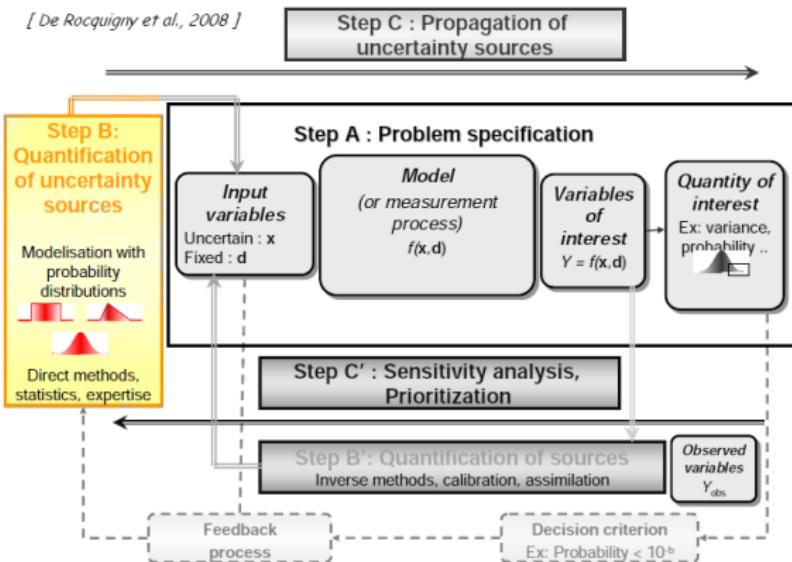
Brief reminders

Descriptive statistics

Data modelisation with PDF

Goodness-of-fit techniques

[De Rocquigny et al., 2008]

**Main steps:**

- A: problem definition**
 - Uncertain input variables
 - Variable/quantity of interest
 - Model construction
- B: uncertainty quantification**
 - Choice of pdfs
 - Choice of correlations
- B': quantification of sources**
 - Inverse methods using data to constrain input values and uncertainties
- C: uncertainty propagation**
 - Evolution of output variability w.r.t input uncertainty
- C': sensitivity analysis**
 - Uncertainty source sorting

These steps are usually model dependent, it might be useful to iterate to help converging to proper conclusions

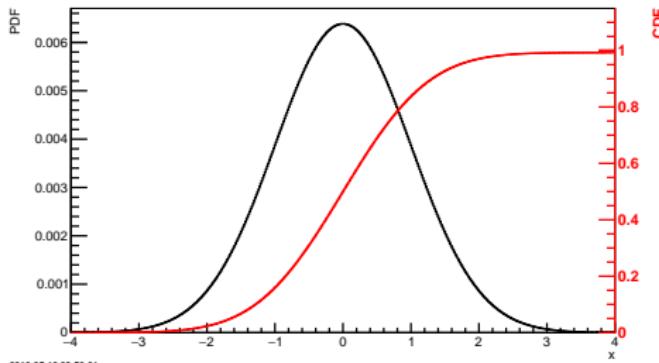
For every random variable $X : \Omega \rightarrow \mathbb{R}$

- ▶ **PDF** (Probability Density Function): if the random variable X has a density f_X , where f_X is a non-negative Lebesgue-integrable function, then

$$P \{a \leq X \leq b\} = \int_a^b f_X(s)ds$$

- ▶ **CDF** (Cumulative Distribution Function): the function $F_X : \mathbb{R} \rightarrow [0, 1]$, given by

$$F_X(a) = \int_{-\infty}^a f_X(s)ds, \quad a \in \mathbb{R}$$



2019-05-13 08:56:01

So far we focus on probabilistic approaches:

- ▶ What is a probability and in which space does it lives
- ▶ The axioms that govern them
- ▶ How to interpret them as well as describe probabilities
- ▶ What might link different realisations

What we're usually handing

In real life, we're handling a sample meaning a restricted set of information

- ▶ What can we learn from a sample
- ▶ How to describe it ? Represent it ?
- ▶ Can we test some hypothesis about what might have provided this sample ?

This branch dealing with these issues is the **statistics**

Brief reminders

Descriptive statistics

Univariate case

Bivariate case

CLT and statistical tests

Data modelisation with PDF

Goodness-of-fit techniques

The effect of the "location" parameter is to translate the graph relative to the standard distribution

- **Mean μ :**

$$\mu = \frac{1}{n_S} \sum_{i=1}^{n_S} x_i$$

- **Mode M:** Value where the probability is the greatest value
- **α -Quantile q_α with $\alpha \in [0, 1]$:** defined as

$$\mathbb{P}[X \leq q_\alpha] = \alpha$$

- **Median $q_{0.5}$:** it is the 0.5-quantile defined as

$$\mathbb{P}[X \leq q_{0.5}] = 0.5 = \mathbb{P}[X \geq q_{0.5}]$$

- **Quartiles:** $q_{0.25}, q_{0.5}, q_{0.75}$
- **Extreme values :** Min and Max

The effect of a "dispersion" parameter is to stretch/shrink the standard distribution

- ▶ **Variance** $\text{Var}(X)$: measure of spread in the data about the mean
 $\text{Var}(X) = \mathbb{E}[(X - \mathbb{E}(X))^2]$, and can be estimated by:

$$\text{Var}(X) = \frac{1}{n_S - 1} \sum_{i=1}^{n_S} (x_i - \mu)^2$$

- ▶ **Standard Deviation** σ : to have an information in the same unit as the variable

$$\sigma = \sqrt{\text{Var}(X)}$$

- ▶ **Coefficient of Variation** δ : σ does not indicate the degree (%) of dispersion around the mean value μ , a non-dimensional term can be introduced:

$$\delta = \frac{\sigma}{\mu}$$

- ▶ **Range** R :

$$R = \text{Max} - \text{Min}$$

- ▶ **Inter-quartile interval H**:

$$H = q_{0.75} - q_{0.25}$$

Any parameter of a PDF that affect the shape of a distribution rather than simply shifting it or stretching/shrinking it.

- **Moment order p:** $\mu_p = \mathbb{E}[(X - \mathbb{E}(X))^p]$

$$\mu_p = \frac{1}{n_S} \sum_{i=1}^{n_S} (x_i - \mu)^p$$

- **Skewness:** γ_1 is a measure of the asymmetry of the PDF

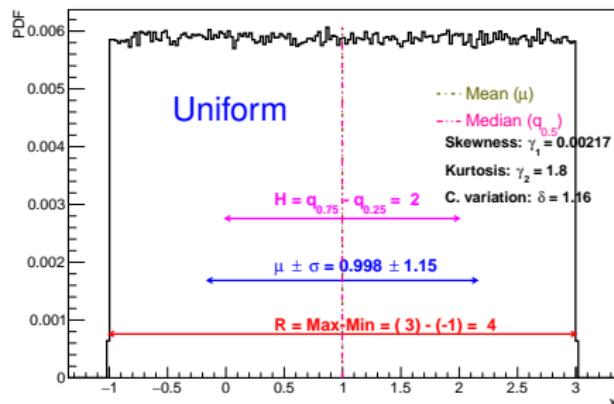
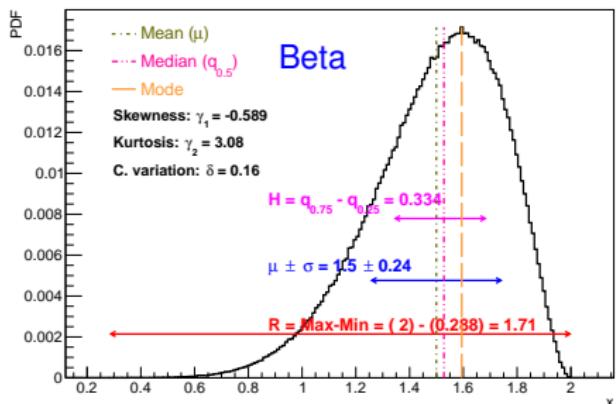
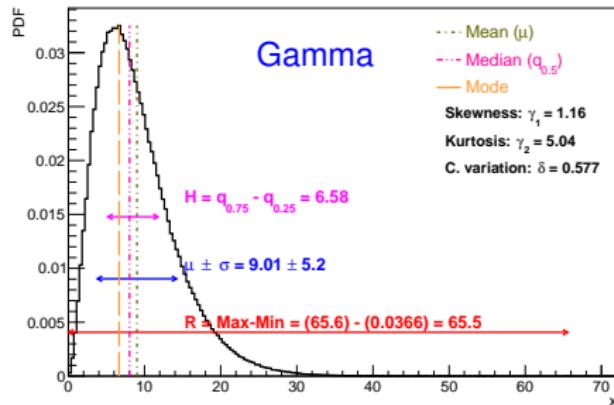
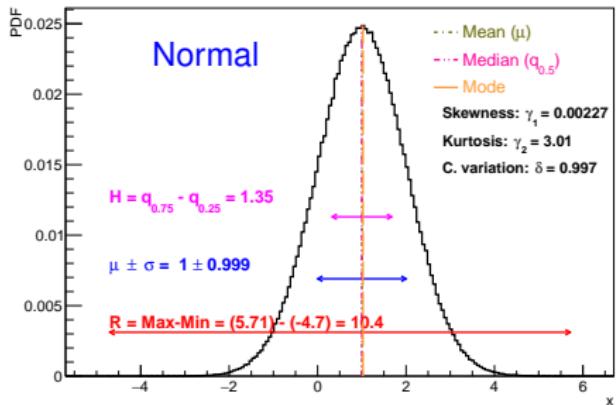
$$\gamma_1 = \mathbb{E} \left[\left(\frac{X - \mu}{\sigma} \right)^3 \right] = \frac{\mu_3}{\sigma^3} = \frac{\mathbb{E}(X^3) - 3\mu\sigma^2 - \mu^3}{\sigma^3}$$

- **Kurtosis:** γ_2 is a measure of the "peakedness" of the PDF

$$\gamma_2 = \frac{\mu_4}{\sigma^4};$$

- Normalised γ_2 : sometimes -3.0 is added to it as $\gamma_2 = 3.0$ for $\mathcal{N}(\mu, \sigma)$

Uni-variate case: illustration of some parameters



Usual convention

- ▶ X is a random variable (RV), whose realisation is noted x
- ▶ \mathbf{x} is a vector of realisation of size n_S , x_i being its i-th element.

Many possible ways to represent data, among which:

Histograms

$$\forall a, b \in \mathbb{R}^2 \text{ with } a < b, H_{[a,b]}(\mathbf{x}) = \sum_{i=1}^{n_S} \mathbb{1}_{[a,b]}(x_i)$$

Normalised w.r.t total number of events, weights...

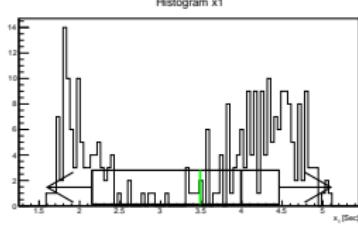
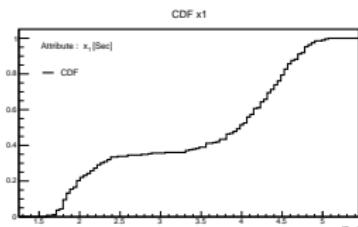
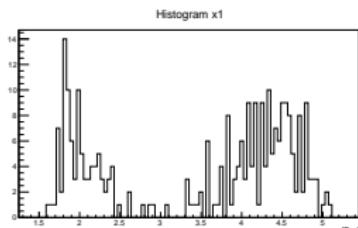
- Sturges: $N_{\text{bin}} = \log_2(n_S) + 1$
- Scott: $N_{\text{bin}} = (x_{\max} - x_{\min}) \times \sqrt[3]{n_S} / (3.5 \times \hat{\sigma}_x) \dots$

Empirical Cumulative Density Function (eCDF)

$$F_{n_S}(x) = \frac{1}{n_S} \sum_{i=1}^{n_S} \mathbb{1}(x_i \leq x)$$

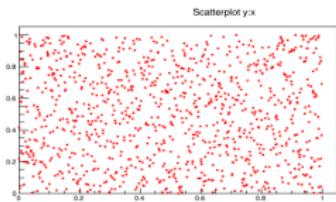
BoxPlot: Simple way to look at many information:

- Minimum and maximum (arrows)
- quartiles: 0.25, 0.5, 0.75 quantiles (black lines)
- Mean: green line

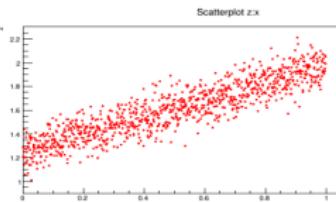


Detect and describe statistical dependencies between variables

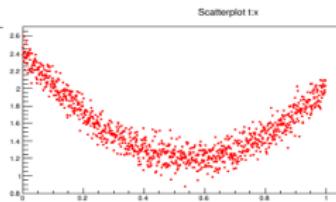
- ▶ independent variables \Rightarrow uncorrelated variables
- ▶ uncorrelated variables $\not\Rightarrow$ independent variables



uncorrelated



linear correlation



nonlinear correlation

The covariance is a measure of how much two RV change together:

$$\text{Cov}(X, Y) = \mathbb{E}[X - \mathbb{E}[X]] \times \mathbb{E}[Y - \mathbb{E}[Y]]$$

and the covariance estimated from a sample (x_i, y_i) is defined as

$$\widehat{\text{Cov}}(x, y) = \frac{1}{n_S} \sum_{i=1}^{n_S} (x_i - \bar{x})(y_i - \bar{y})$$

The sign of this coefficient is the tendency of the linear relationship between the variables, but the magnitude is not easy to interpret.

The Pearson coefficient (ρ or ρ_P) is a normalised version of the covariance as it is divided by the standard deviations.

It is a measure of the linear correlation (dependence) between two RV X and Y ,

$$\rho_P(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y}.$$

Its estimation on a sample (x_i, y_i) can be written as \widehat{r}_P :

$$\widehat{r}_P = \frac{\sum_{i=1}^{n_S} (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^{n_S} (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^{n_S} (y_i - \bar{y})^2}}$$

Properties of this coefficient

- ▶ $\widehat{r}_P \in [-1, 1]$
- ▶ $\widehat{r}_P = \pm 1 \Leftrightarrow$ perfect linear description between X and Y , the data points lying exactly on a positive (negative) identity line.
- ▶ $\widehat{r}_P = 0$, X and Y are said to be (linearly) uncorrelated (but not necessarily independents !!)

Let's consider an i.i.d. sequence (*independent and identically distributed*) X_1, \dots, X_n from distribution function F with expectation μ and variance σ^2 .

Investigating the average

Let's call $\bar{X}_n = n^{-1} \sum_{i=1}^n X_i$. This is a RV, so

- ▶ $\mathbb{E}[\bar{X}_n] = n^{-1} \mathbb{E}\left[\sum_{i=1}^n X_i\right] = n^{-1} \left(\sum_{i=1}^n \mu\right) = \mu$
- ▶ $\text{Var}[\bar{X}_n] = n^{-2} \text{Var}\left[\sum_{i=1}^n X_i\right] = n^{-2} \left(\sum_{i=1}^n \sigma^2\right) = n^{-1} \sigma^2$
- With increasing number of samples, n , \bar{X}_n deviates less and less from μ .

Chebyshev's inequality

For an arbitrary RV Y and any $a \in \mathbb{R}^+$

$$\mathbb{P}(|Y - \mathbb{E}[Y]| \geq a) \leq \frac{1}{a^2} \text{Var}(Y)$$

Applying Chebyshev's inequality to our example, with $E[\bar{X}_n] = \mu$, $\text{Var}[\bar{X}_n] = \sigma^2/n$ and $\varepsilon > 0$, leads to

$$\mathbb{P}(|\bar{X}_n - \mu| \geq \varepsilon) \leq \frac{\sigma^2}{n\varepsilon^2}$$

The law of large numbers

If \bar{X}_n is the average of n independent random variables with expectation μ , and variance σ^2 , then for $\varepsilon \in \mathbb{R}^+$

$$\lim_{n \rightarrow \infty} \mathbb{P}(|\bar{X}_n - \mu| \geq \varepsilon) = 0 \quad (\text{weak})$$

Another formulation can be found :

$$\mathbb{P}\left(\lim_{n \rightarrow \infty} \bar{X}_n = \mu\right) = 1 \quad (\text{strong})$$

Many mathematicians have contributed to the CLT and its proof \Leftrightarrow many different statements of the theorem are accepted.

Formulation

Suppose X_1, X_2, \dots, X_n are i.i.d. with mean μ and a finite variance σ^2 . Then,

$$\frac{S_n - n\mu}{\sqrt{n\sigma^2}} \rightarrow \mathcal{N}(0, 1),$$

where $S_n = X_1 + X_2 + \dots + X_n$, $n \geq 1$ and \rightarrow represents convergence in distribution.

This can help providing

- ▶ Confidence interval
- ▶ test-of-hypothesis

Illustration of this (next few slides)

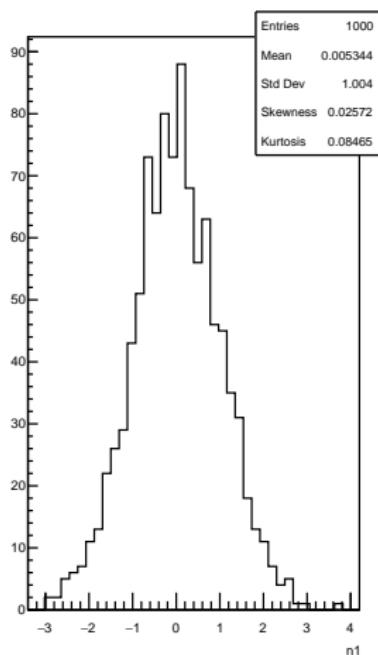
Sample of 1000 points for 200 random variables: exponential (e) / normal (n) / uniform (u) laws

→ Their properties are chosen so that $\mu = 0$ and $\sigma = 1$

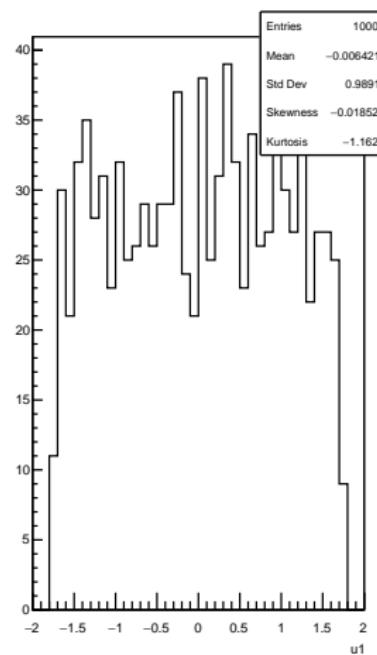
- ▶ Check the laws
- ▶ Check the distribution of $S_n = \sum_n X_n$ when increasing n
- ▶ Check the distribution of $S_n^t = \sum_n X_n^t$ for all $t = e/n/u$

CLT illustration: RV one-by-one

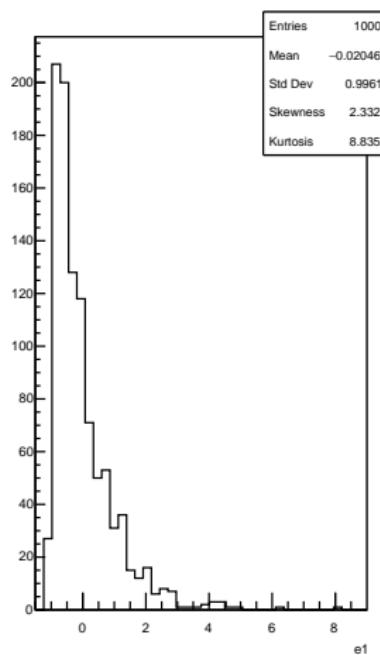
Histogram n1



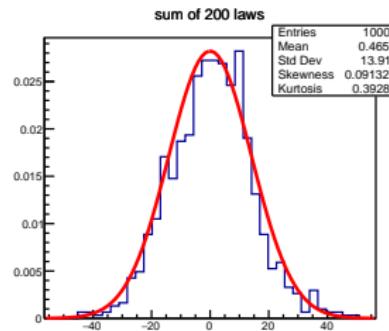
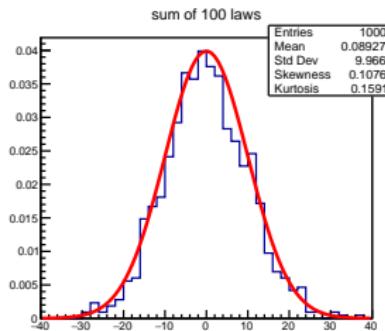
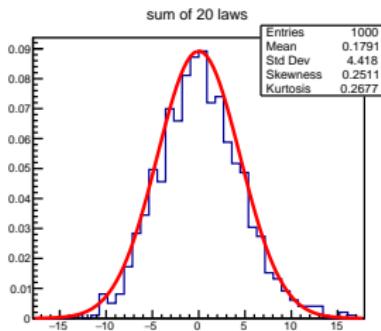
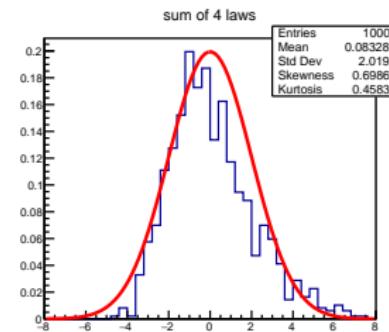
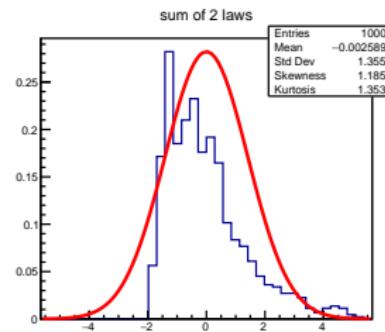
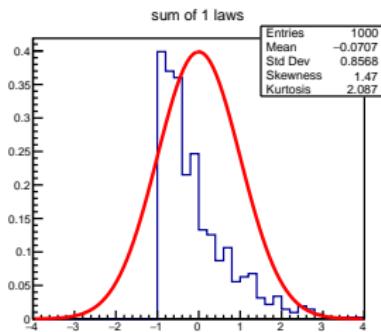
Histogram u1



Histogram e1

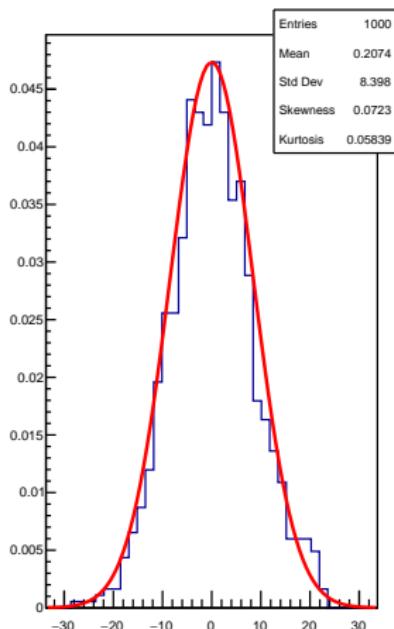


CLT illustration: summing laws of all types

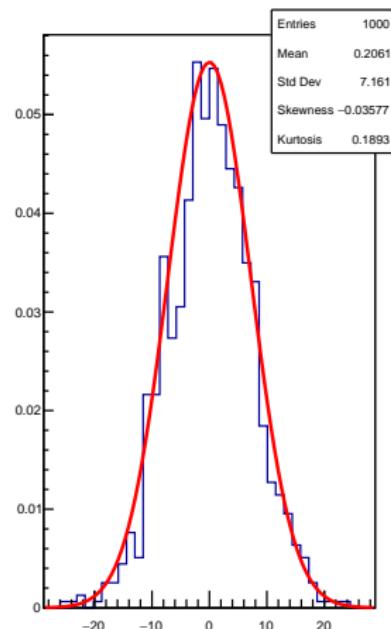


CLT illustration: sums split by types

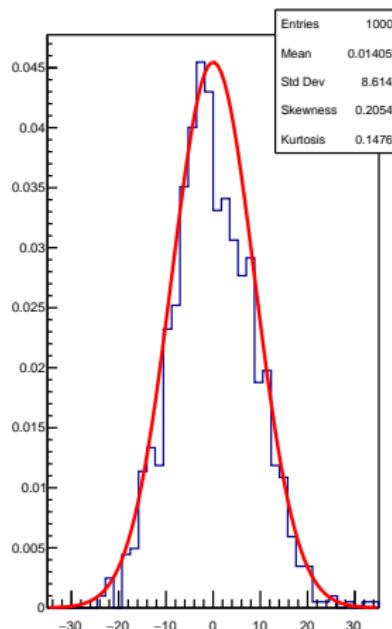
sum of 71 n laws



sum of 52 u laws



sum of 77 e laws



The aim of a test-of-hypothesis is to check the validity of a given hypothesis, providing a certain chosen confidence level ($1 - \alpha$).

Principle in few key steps

- Purpose: significance, goodness-of-fit, independence, conformity...

A factory build tubes whose lifetime $\sim N(1200, 300)$. 100 tubes are produced with a new process

$\bar{x} = 1265$. Is this significant ? Is the new μ_N greater than 1200 ?

The aim of a test-of-hypothesis is to check the validity of a given hypothesis, providing a certain chosen confidence level ($1 - \alpha$).

Principle in few key steps

- Purpose: significance, goodness-of-fit, independence, conformity...

A factory build tubes whose lifetime $\sim N(1200, 300)$. 100 tubes are produced with a new process

$\bar{x} = 1265$. Is this significant? Is the new μ_N greater than 1200?

- Hypothesis

1 H_0 is the null-hypothesis (to be tested): $\mu_0 = 1200$

2 H_1 is the alternative hypothesis: $\mu_1 > 1200$

- Confidence level: choose probability α

A usual choice is to set $\alpha = 0.05$, resulting in a 95% CL

	H_0 accepted	H_0 rejected
H_0 true	Correct ($1 - \alpha$)	Type-I error (α)
H_0 false	Type-II error (β)	Correct ($1 - \beta$)

The aim of a test-of-hypothesis is to check the validity of a given hypothesis, providing a certain chosen confidence level ($1 - \alpha$).

Principle in few key steps

- Purpose: significance, goodness-of-fit, independence, conformity...

A factory build tubes whose lifetime $\sim N(1200, 300)$. 100 tubes are produced with a new process

$\bar{x} = 1265$. Is this significant? Is the new μ_N greater than 1200?

- Hypothesis

1 H_0 is the null-hypothesis (to be tested): $\mu_0 = 1200$

2 H_1 is the alternative hypothesis: $\mu_1 > 1200$

- Confidence level: choose probability α

A usual choice is to set $\alpha = 0.05$, resulting in a 95% CL

- Statistical test to be computed

$$\text{Use classical test } \hat{Z} = \frac{\bar{x} - \mu_0}{\sigma / \sqrt{N}} = \frac{1265 - 1200}{30} = 2.17$$

	H_0 accepted	H_0 rejected
H_0 true	Correct ($1 - \alpha$)	Type-I error (α)
H_0 false	Type-II error (β)	Correct ($1 - \beta$)

The aim of a test-of-hypothesis is to check the validity of a given hypothesis, providing a certain chosen confidence level ($1 - \alpha$).

Principle in few key steps

- Purpose: significance, goodness-of-fit, independence, conformity...

A factory build tubes whose lifetime $\sim N(1200, 300)$. 100 tubes are produced with a new process

$\bar{x} = 1265$. Is this significant? Is the new μ_N greater than 1200?

- Hypothesis

- 1 H_0 is the null-hypothesis (to be tested): $\mu_0 = 1200$
- 2 H_1 is the alternative hypothesis: $\mu_1 > 1200$

- Confidence level: choose probability α

A usual choice is to set $\alpha = 0.05$, resulting in a 95% CL

- Statistical test to be computed

$$\text{Use classical test } \hat{Z} = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}} = \frac{1265 - 1200}{30} = 2.17$$

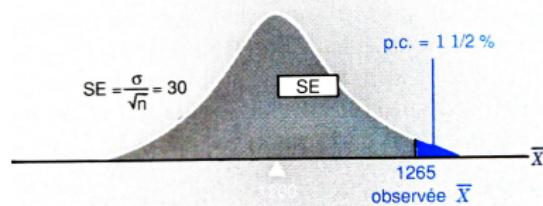
- Result interpretation

Two possibilities:

- 1 Look at table and see that $\hat{Z} > Z_{0.05}$ ($= 1.64$) $\Rightarrow H_0$ rejected!
- 2 Look at table and see that $\hat{Z} = 2.17 \Leftrightarrow P_c = 0.015 \Rightarrow H_0$ rejected!
- 3 Look at table and see that to get $\hat{Z}_C = 1.64 \Leftrightarrow \bar{x}_c = 1249 \Rightarrow H_0$ rejected!

Example of tables: <https://www.statisticshowto.com/tables/z-table/>

	H_0 accepted	H_0 rejected
H_0 true	Correct ($1 - \alpha$)	Type-I error (α)
H_0 false	Type-II error (β)	Correct ($1 - \beta$)



Assuming that $X, Y \sim \mathcal{N}(\mu_1, \mu_2, \sigma_1, \sigma_2, \rho)$ a bi-variate normal distribution

1 H_0 : test whether X and Y are independent, meaning $\rho = 0$

2 H_1 : it exists a relation between X and Y , meaning $\rho \neq 0$

→ Statistical test to be computed *Use t-statistic test* $\hat{t} = \frac{\rho\sqrt{n_S - 2}}{\sqrt{1 - \rho^2}}$

Using a 15-sample database showing weight and height for 2-year old children.

X: Height (cm)	82.9	83.4	82.4	82.1	84.8	86.7	84.	89.	85.	85.4	87.7	87.7	86.4	86.4
Y: Weight (kg)	8.7	9.2	9.5	10.1	10.4	10.5	10.8	11.	11.5	11.6	12.4	13.6	13.8	13.9

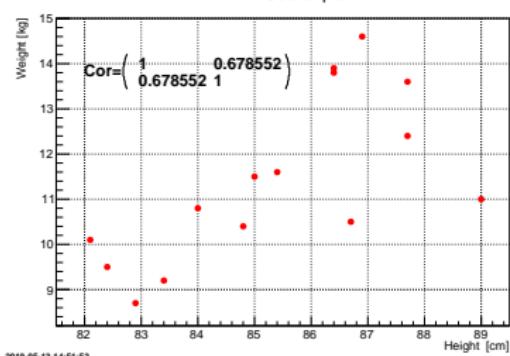
Scatterplot Y:X

Setting the test

$$\hat{t} = \frac{0.6786 \times \sqrt{15 - 2}}{\sqrt{1 - 0.6786^2}} = 3.33067$$

Interpret these results

- For a chosen $\alpha = 0.05$, $t_{5\%}(13) = 2.16$
 $\Rightarrow \hat{t} > t_{5\%}(13) \Leftrightarrow H_0$ rejected !
- It exists a relation between X and Y at 5% significance level
- For a chosen $\alpha = 0.01$, $t_{1\%}(13) = 3.012$
 $\Rightarrow \hat{t} > t_{1\%}(13) \Leftrightarrow H_0$ rejected !
- It exists a relation between X and Y at 1% significance level
- Looking at table, $3.01 < \hat{t} < 3.37$
 \Rightarrow Critical probability $0.005 < P_c < 0.01$



122-points sample.

Setting the test

1 $n_S = 122 \Rightarrow \text{degree of freedom} = 120 (n_S - 2)$

2 $\hat{\rho} = 0.0668$

$$\hat{t} = \frac{0.0668 \times \sqrt{120}}{\sqrt{1 - 0.0668^2}} = 0.733$$

<https://archimede.mat.ulaval.ca/stt1920/STT-1920-Loi-de-Student.pdf>

Interpret these results

122-points sample.

Setting the test

1 $n_S = 122 \Rightarrow \text{degree of freedom} = 120 (n_S - 2)$

2 $\hat{\rho} = 0.0668$

$$\hat{t} = \frac{0.0668 \times \sqrt{120}}{\sqrt{1 - 0.0668^2}} = 0.733$$

<https://archimede.mat.ulaval.ca/stt1920/STT-1920-Loi-de-Student.pdf>

Interpret these results

- ▶ For $\alpha = 0.05$, $t_{5\%}(120) = 1.98$ and for $\alpha = 0.01$, $t_{1\%}(120) = 2.62$
 $\Rightarrow H_0$ accepted !
- ▶ Looking at table assuming student \rightarrow normal.
 \Rightarrow Critical probability is $0.46 < P_c^P < 0.47$

\Leftrightarrow You know nothing Jon Snow !

122-points sample.

Setting the test

1 $n_S = 122 \Rightarrow \text{degree of freedom} = 120 (n_S - 2)$

2 $\hat{\rho} = 0.0668$

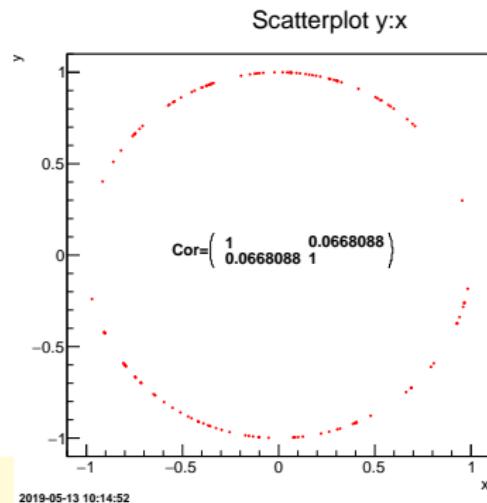
$$\hat{t} = \frac{0.0668 \times \sqrt{120}}{\sqrt{1 - 0.0668^2}} = 0.733$$

<https://archimede.mat.ulaval.ca/stt1920/STT-1920-Loi-de-Student.pdf>

Interpret these results

- ▶ For $\alpha = 0.05$, $t_{5\%}(120) = 1.98$ and for $\alpha = 0.01$, $t_{1\%}(120) = 2.62$
 $\Rightarrow H_0$ accepted !
- ▶ Looking at table assuming student \rightarrow normal.
 \Rightarrow Critical probability is $0.46 < P_c^P < 0.47$

↔ You know nothing Jon Snow !



Brief reminders

Descriptive statistics

Data modelisation with PDF

Parametric PDFs

Non-parametric PDFs

Goodness-of-fit techniques

Data modelisation with PDF

Problem

- ▶ Let (x_1, \dots, x_{n_S}) an i.i.d sample of a PDF $f(x, \theta)$ where $\theta \in \Theta$ is a vector of parameters for this family
- ▶ The true value of the parameters θ^* , from which the data come from, is unknown
- ▶ Build an estimator $\hat{\theta}$ which would be as close to the true value θ^* as possible.

Two usual methods are:

1 Maximum Likelihood (MLE)

The method of maximum likelihood selects the set of values of the model parameters that maximizes the likelihood function. This function measures the "agreement" of the selected model with the observed data.

2 Moments Method (MM)

- One starts with deriving equations that relate the population moments to the parameters θ
- The moments are estimated from the given sample
- The equations are then solved for the parameters θ , using the sample moments in place of the (unknown) population moments

Build an estimator $\hat{\theta}$ for the model's parameters of the $f(x, \theta)$ from the data $(x_i)_{1 \leq i \leq n_S}$

We use the **Likelihood** function $\mathcal{L}(\theta; x_1, \dots, x_{n_S})$:

$$\mathcal{L}(\theta; x_1, \dots, x_{n_S}) = f(x_1, \dots, x_{n_S} | \theta) = \prod_{i=1}^{n_S} f(x_i | \theta)$$

In practice it is often more convenient to work with the logarithm of the likelihood function, called the **log-likelihood**:

$$\ln(\mathcal{L}(\theta; x_1, \dots, x_{n_S})) = \sum_{i=1}^{n_S} \ln(f(x_i | \theta))$$

or the **average log-likelihood**:

$$\hat{l}(\theta; x_1, \dots, x_{n_S}) = \frac{1}{n_S} \ln(\mathcal{L}(\theta; x_1, \dots, x_{n_S}))$$

MLE estimates $\hat{\theta}_{MLE}$ by finding the value of θ that maximizes the \hat{l} function

$$\hat{\theta}_{MLE} = \arg \max_{\theta \in \Theta} \hat{l}(\theta; x_1, \dots, x_{n_S})$$

... if any maximum exists

Let (x_1, \dots, x_{n_S}) be an i.i.d sample from a normal law $\mathcal{N}(\mu, \sigma)$

If one defines $\theta = (\mu, \sigma)$ the unknown parameters, the density can be written:

$$f(x|\theta) = \frac{1}{\sigma\sqrt{2\pi}} \exp^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$\text{The Likelihood is : } \mathcal{L}(\theta; x_1, \dots, x_{n_S}) = \prod_{i=1}^{n_S} f(x_i|\theta) = \left(\frac{1}{2\pi\sigma^2} \right)^{\frac{n_S}{2}} \exp^{-\frac{\sum_{i=1}^{n_S} (x_i - \mu)^2}{2\sigma^2}}$$

The average log-likelihood $\hat{l}(\theta; x_1, \dots, x_{n_S})$ can be written as:

$$\hat{l}(\theta; x_1, \dots, x_{n_S}) = -\frac{1}{2} \ln 2\pi - \ln \sigma - \frac{1}{2n_S\sigma^2} \sum_{i=1}^{n_S} (x_i - \mu)^2$$

► MLE for the mean parameter :

$$\frac{\partial \hat{l}}{\partial \mu} = \frac{1}{n_S\sigma^2} \sum_{i=1}^{n_S} (x_i - \mu) = 0 \Leftrightarrow \hat{\mu}_{MLE} = \bar{x} = \frac{1}{n_S} \sum_{i=1}^{n_S} x_i$$

► MLE for the variance parameter :

$$\frac{\partial \hat{l}}{\partial \sigma} = -\frac{1}{\sigma} + \frac{1}{n_S\sigma^3} \sum_{i=1}^{n_S} (x_i - \mu)^2 = 0 \Leftrightarrow \hat{\sigma}_{MLE}^2 = \frac{1}{n_S} \sum_{i=1}^{n_S} (x_i - \hat{\mu}_{MLE})^2$$

Build an estimator $\hat{\theta}$ for the model's parameters of $f(x, \theta)$ using data $(x_i)_{1 \leq i \leq n_S}$

Suppose the first k moments of the true PDF can be expressed as functions of θ :

$$\mu_1 = \mathbb{E}[X] = g_1(\theta_1, \theta_2, \dots, \theta_k)$$

$$\mu_2 = \mathbb{E}[X^2] = g_2(\theta_1, \theta_2, \dots, \theta_k)$$

...

$$\mu_k = \mathbb{E}[X^k] = g_k(\theta_1, \theta_2, \dots, \theta_k)$$

We compute the same first k moments from the sample $(x_i)_{1 \leq i \leq n_S}$

$$\hat{\mu}_j = \frac{1}{n_S} \sum_{i=1}^{n_S} x_i^j$$

The moments method estimator for (θ_j) denoted by $\hat{\theta}_{MM}$ is defined as the solution (if there is one) to the system of equations:

$$\hat{\mu}_1 = g_1(\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_k)$$

$$\hat{\mu}_2 = g_2(\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_k)$$

...

$$\hat{\mu}_k = g_k(\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_k)$$

Comments

- ▶ The moments method is fairly simple and yields consistent estimators (under very weak assumptions), though these estimators are often biased
- ▶ Estimates by the moments method may be used as the first approximation to the solutions of the likelihood equations, and successive improved approximations may then be found by the Newton Raphson method. In this way the moments method and the method of maximum likelihood are symbiotic
- ▶ In some cases, as in the example of the gamma distribution, the likelihood equations may be intractable without computers, whereas the moments method estimators can be quickly and easily calculated by hand

Gamma distribution

Given an i.i.d sample $(x_i)_{1 \leq i \leq n_S}$ from a Gamma law, where $\theta = (\alpha, \beta)$ unknown

$$f(x) = \frac{x^{\alpha-1} e^{-\beta x} \beta^\alpha}{\Gamma(\alpha)}$$

1 first moment:

$$\mu = \frac{\alpha}{\beta}$$

2 second moment:

$$\sigma^2 = \text{Var}[X] = \mathbb{E}[X^2] - \mathbb{E}[X]^2 = \frac{\alpha}{\beta^2}$$

Looking for $(\hat{\alpha}, \hat{\beta})$, for which $\hat{\mu}_1 = \sum_{i=1}^{n_S} x_i = \frac{\hat{\alpha}}{\hat{\beta}}$ and $\hat{\mu}_2 = \sum_{i=1}^{n_S} x_i^2 = \frac{\hat{\alpha}}{\hat{\beta}^2}$

From 1 $\Rightarrow \hat{\alpha} = \hat{\mu}_1 \hat{\beta}$

Injecting this into 2 $\Rightarrow \hat{\beta} = \frac{\hat{\mu}_1}{\hat{\mu}_2 - \hat{\mu}_1^2}$

$$\hat{\alpha} = \frac{\hat{\mu}_1^2}{\hat{\mu}_2 - \hat{\mu}_1^2} \quad \text{and} \quad \hat{\beta} = \frac{\hat{\mu}_1}{\hat{\mu}_2 - \hat{\mu}_1^2}$$

From the point of view of the histogram,

$$f(x) = F'(x) \simeq \frac{F(x+h) - F(x-h)}{2 \times h} \quad \forall h > 0, h \text{ "small"}$$

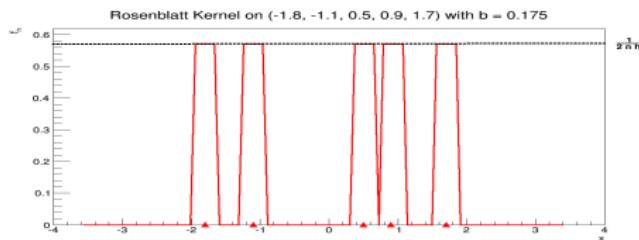
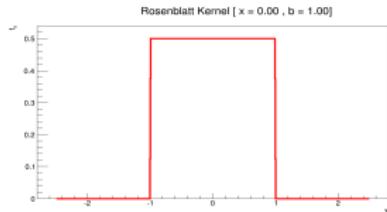
Then **Rosenblatt** (1956) suggests the estimator :

$$\hat{f}_{nS,h}(x) = \frac{\hat{F}_{nS}(x+h) - \hat{F}_{nS}(x-h)}{2 \times h}$$

which has another representation **Parzen** (1962)

$$\hat{f}_{nS,h}(x) = \frac{1}{nS} \sum_{i=1}^{nS} \frac{1}{h} K\left(\frac{x-x_i}{h}\right)$$

$$\text{with } K(u) = \frac{1}{2} \times \mathbb{I}_{[-1,1]}(u)$$



Kernel estimators - definitions

- A function $K : \mathbb{R} \rightarrow \mathbb{R}$ is said a **Kernel** if

$$\int K(u) \, du = 1.$$

- Often, but not necessarily,

- K is symmetric around the origin:
$$K(-u) = K(u) \quad \forall u$$
- K is positive:
$$K(u) > 0 \quad \forall u$$

- $\forall h > 0$,

$$\hat{f}_{nS,h}(x) = \frac{1}{nS} \sum_{i=1}^{nS} \frac{1}{h} K\left(\frac{X_i - x}{h}\right)$$

is a **kernel estimator** of the density f
$$(\int \hat{f}_{nS,h}(x) \, dx = 1)$$

- Kernel approach is a histogram which, for estimating the density of $f(x)$, has been shifted so that x , say, lies at the center of a mesh interval. And For evaluating the density at another point, say y , the mesh is shifted again, so that y is at the center of a mesh interval.
- The parameter h is a *smoothing* parameter called **bandwidth**; More greater h is, more the estimation $\hat{f}_{nS,h}$ is smooth.

- Rectangular (**Rosenblatt**) (black)

$$K(u) = \frac{1}{2} \times \mathbb{I}_{[-1,1]}(u)$$

- Triangular (**red**)

$$K(u) = (1 - |u|) \times \mathbb{I}_{[-1,1]}(u)$$

- Epanechnikov** (**blue**)

$$K(u) = \frac{3}{4}(1 - x^2) \times \mathbb{I}_{[-1,1]}(u)$$

- Biweight (**green**)

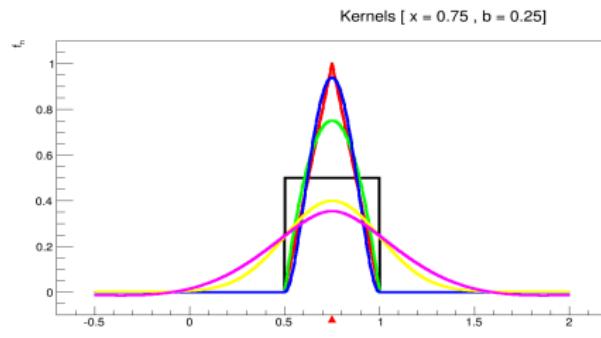
$$K(u) = \frac{15}{16}(1 - x^2)^2 \times \mathbb{I}_{[-1,1]}(u)$$

- Gaussian (**yellow**)

$$K(u) = \frac{\exp^{-x^2/2}}{\sqrt{2\pi}}$$

- Silverman** (**magenta**)

$$K(u) = \frac{1}{2} \exp^{-|u|/\sqrt{2}} \sin(|u|/\sqrt{2} + \pi/4)$$



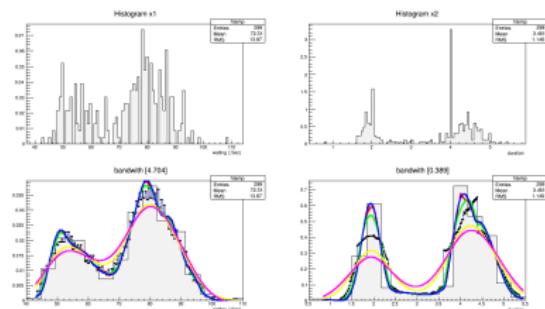
- Optimal bandwidth with the Silverman Rule (1996)

$$h_{nS} = 1.364 \times \alpha_K \times \text{MIN}\{\hat{\sigma}, \frac{\text{IQR}}{1.349}\} \times nS^{-1/5}$$

with

- $\hat{\sigma}$ is the sample standard deviation
- IQR is the "InterQuartile Range" ($\text{IQR} = q_{0.75} - q_{0.25}$)
- α_K is a constant that only depends on the used kernel

Kernel	$k(x)$	σ_K
Rectangular	$1/2$, $ x < 1$	1.3510
Triangular	$1 - x $, $ x < 1$	1.8882
Epanechnikov	$\frac{3}{4}(1 - x^2)$, $ x < 1$	1.7188
Biweight	$\frac{15}{16}(1 - x^2)^2$, $ x < 1$	2.0362
Gaussian	$\frac{\exp^{-x^2/2}}{\sqrt{2\pi}}$	0.7764



Geyser database for Gaussian Kernel (left) waiting b = 4.70, (right) duration b = 0.39

Brief reminders

Descriptive statistics

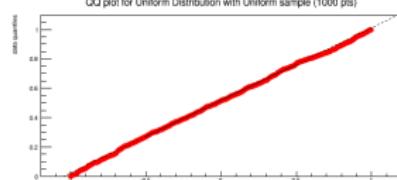
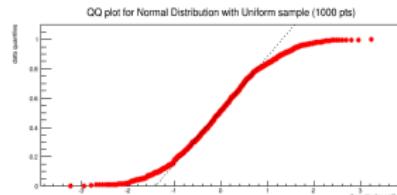
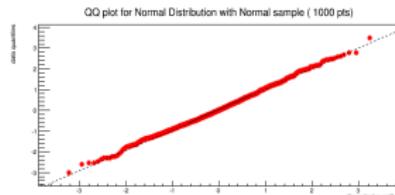
Data modelisation with PDF

Goodness-of-fit techniques

- Graphical methods
- Statistical tests

- Graphical methods
 - QQPlot
- Statistical Tests
 - Chi-Squared
 - Tests based on EDF Statistics
 - ★ Kolmogorov-Smirnov
 - ★ Cramer-von Misses
 - ★ Anderson-Darling

- a **QQ-plot** ("Q" stands for quantile) is a probability plot to compare two probability distributions by plotting their quantiles against each other
- A point (x, y) on the plot corresponds to one of the quantiles of the second distribution (y-coordinate) plotted against the same quantile of the first distribution (x-coordinate).
- If the two distributions being compared are similar, the points in the QQ-plot will approximately lie on the line $y = x$
- If the distributions are linearly related, the points in the QQ-plot will approximately lie on a line, but not necessarily on the line $y = x$.
- Select one axe for the theoretical distribution for Goodness-of-Fit test



In Goodness-of-Fit work, the commonly used statistical tests are:

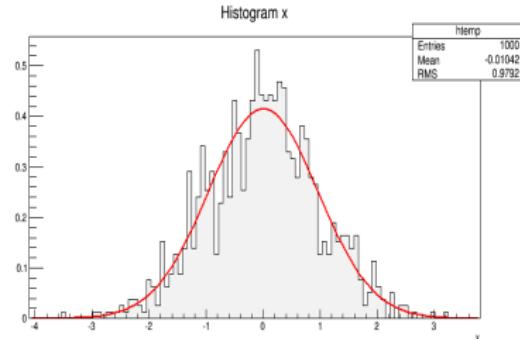
- Chi-Squared (χ^2)
- Tests based on EDF Statistics
 - Kolmogorov-Smirnov (**D**)
 - Cramer-von Mises (W^2)
 - Anderson-Darling (A^2)

The chi-squared test (χ^2)

- The χ^2 test is used to test if a sample (x_i) came from a specific distribution
- Useful when data are discrete, and applied to continuous distribution with a large number of observations
- The basic idea is to partitioned the range of the sample into k cells, and compare the observed frequency O_i with the expected frequency E_i in each cell i
- The statistic test is:

$$\chi^2 = \sum_{i=1}^k \frac{(O_i - E_i)^2}{E_i}$$

which follows a χ^2 distribution with $(k - 1 - t)$ degrees of freedom, where t is the number of parameters of the distribution to estimate



- The ratio nS/k must verify $nS/k \geq 5$
- The value of the χ^2 test statistic are dependent on how the data is binned
- χ^2 test is generally less powerful than *EDF* tests

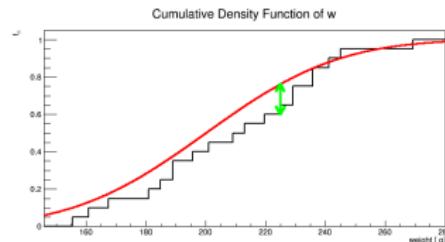
- Graphical methods have a wide appeal in deciding if a random sample appears to come from a given PDF
- We consider now tests of fit based on the *Empirical Distribution Function ("EDF")*
- *EDF* statistics are measures of the discrepancy between the empirical CDF and the theoretical CDF of the PDF
- They are based on the vertical differences between $F_{nS}(x)$ and $F(x)$, and divided into two classes :
 1. **the supremum statistics** : select the largest vertical difference between the two CDF; it is the **Kolmogorov-Smirnov test D**

$$D = \sup_x |F_{nS}(x) - F(x)|$$

2. **the quadratic statistics** : measure of discrepancy given by the Cramer-von Mises family

$$Q = nS \int_{-\infty}^{+\infty} (F_{nS}(x) - F(x))^2 \psi(x) dx$$

where ψ is a weight function



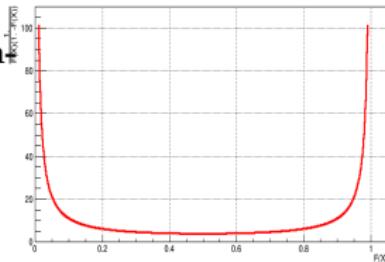
Tests based on EDF statistics (2/2)

- For $\psi(x) = 1$ we obtain the **Cramer-von Mises** Tests, denoted as W^2 :

$$W^2 = nS \int_{-\infty}^{+\infty} (F_{nS}(x) - F(x))^2 dx$$

- For $\psi(x) = \frac{1}{F(x)(1.0 - F(x))}$ we obtain the **Anderson-Darling** test, denoted A^2 :

$$A^2 = nS \int_{-\infty}^{+\infty} \frac{(F_{nS}(x) - F(x))^2}{F(x)(1.0 - F(x))} dx$$



- To compute these statistics, we use the *Probability Integral Transformation ("PIT")*
 - Let $X \sim F$ with F is the true CDF
 - If $Z = F(X)$, then $Z \sim \mathcal{U}[0., 1.]$
 - For The sample $(x_1, x_2, \dots, x_{nS})$, compute $z_i = F(x_i)$ and compare the empirical CDF of the z_i with the CDF of the uniform distribution

$$F^*(z) = z \quad , \quad 0 \leq z \leq 1$$

- EDF statistics computed from the EDF of the z_i compared with the uniform distribution will take the same values as if they were computed from the EDF of the x_i compared with F

- The χ^2 statistic is the lower powerfull for continuous PDF
- EDF statistics are usually much more powerfull than the χ^2 statistic (where data must be grouped, then loss of informations)
- the D statistic is the most well-known of the EDF statistics, but it is often much less powerfull than the quadratic statistics W^2 and A^2
- A^2 and W^2 give often similarly values, but A^2 is on the whole more powerfull when the distribution F departs from the true distribution in the tails (weight function)
- In Goodness-of-Fit work, departure in the tails is often important to detect, so A^2 is the recommended statistic



Thanks! Any questions?