### PSI Zuoz Summer School 2022

### Statistics

Nicolas Berger (LAPP Annecy)



### PSI Zuoz Summer School 2022

## Statistics **Eor Physicists**

Nicolas Berger (LAPP Annecy)



### **Lecture Plan**

### Statistics basic concepts (Today)

[Basic ingredients (PDFs, etc.)]

Parameter estimation (maximum likelihood, least-squares, ...)

Model testing ( $\chi^2$  tests, hypothesis testing, p-values, ...)

### **Computing statistical results** (Today/Tomorrow)

Discovery Confidence intervals Upper limits Systematics and profiling [Bayesian techniques]

The class will be based on both lectures and hands-on tutorial

### Hands-on sessions

The hands-on session will be based on Jupyter notebooks built using the **numpy/scipy/pyplot** stack.

If you have a computer with you, **please install** anaconda before the start of the class. This provides a consistent installation of python, JupyterLab, etc.

 $\rightarrow$  Alternatively, you can also install JupyterLab as a standalone package.

 $\rightarrow$  Another solution is to run the notebooks on the public jupyter servers at mybinder.org. This will probably be slower but avoids a local install.

Warmup		notebook [solutions]	binder [solutions]
Lecture 1	Lecture Notes	notebook	binder
Lecture 2	Lecture notes	notebook	binder

The **warmup** item includes material that will not be covered in detail in the class, as well as an introduction to the notebooks. Please have a look before the beginning of the classes if you are unfamiliar with this.

### **Statistics are everywhere**

*"There are three types of lies - lies, damn lies, and statistics." – Benjamin Disraeli* 



### And Physics ?

"If your experiment needs statistics, you ought to have done a better experiment" – E. Rutherford

### Introduction

Statistical methods play a critical role in many areas of physics



GeV

Data

Background ZZ(\*)

Background Z+jets, tt

ATLAS

 $H \rightarrow ZZ^{(^*)} \rightarrow 4I$ 

### Introduction

Sometimes difficult to distinguish a bona fide discovery from a **background fluctuation**...



### Introduction

Sometimes difficult to distinguish a bona fide discovery from a **background fluctuation**...



### **Randomness in High-Energy Physics**

Experimental data is produced by incredibly complex processes





### **Randomness in High-Energy Physics**



**Randomness** involved in all stages

- $\rightarrow$  **Classical** randomness: detector reponse
- $\rightarrow$  Quantum effects in particle production, decay

Hard scattering

PDFs, Parton shower, Pileup

More details in Anna Sfyrla's lectures!

Decays

**Detector response** 

### Reconstruction



**Example**: measuring the energy of a photon in a calorimeter





**Example**: measuring the energy of a photon in a calorimeter





**Example**: measuring the energy of a photon in a calorimeter









Cannot predict the measured value for a given event

### ⇒ Random process ⇒ Need a probabilistic description

### Quantum Randomness: H→ZZ\*→4l



### Quantum Randomness: H→ZZ\*→4l



**Rare process**: Expect 1 signal event every ~6 days



# http://www.phdcomics.com/comics/archive.php?comicid=1489

### View online

### Quantum Randomness: H→ZZ\*→4l



"Will I get an event today ?"  $\rightarrow$  only **probabilistic** answer

### **Statistical Modeling**

Probabilistic treatment of possible outcomes ⇒ **Probability Distribution** 

**Example**: two-coin toss

 $\rightarrow$  Fractions of events in each bin i converge to a limit p<sub>i</sub>

### **Probability distribution** :

 $\{P_i\}$  for i = 0, 1, 2

### Properties

- P<sub>i</sub> > 0
- Σ P<sub>i</sub>=1



Probabilistic treatment of possible outcomes ⇒ **Probability Distribution** 

**Example**: two-coin toss

 $\rightarrow$  Fractions of events in each bin i converge to a limit p<sub>i</sub>

### **Probability distribution** :

{ P<sub>i</sub> } for i = 0, 1, 2

### Properties

- P<sub>i</sub> > 0
- Σ P<sub>i</sub>=1



Probabilistic treatment of possible outcomes ⇒ **Probability Distribution** 

**Example**: two-coin toss

 $\rightarrow$  Fractions of events in each bin i converge to a limit p<sub>i</sub>

### **Probability distribution** :

{ P<sub>i</sub> } for i = 0, 1, 2

### Properties

- P<sub>i</sub> > 0
- Σ P<sub>i</sub>=1



Probabilistic treatment of possible outcomes ⇒ Probability Distribution

**Example**: two-coin toss

→ Fractions of events in each bin i converge to a limit p<sub>i</sub>

### **Probability distribution** :

{ P<sub>i</sub> } for i = 0, 1, 2

### Properties

- P<sub>i</sub> > 0
- Σ P<sub>i</sub>=1



### **Continuous Variables: PDFs**

**Continuous variable**: can consider **per-bin** probabilities p<sub>i</sub>, i=1.. n<sub>bins</sub>



Generalizes to **multiple variables** :  $P(x,y) > 0, \int P(x,y) dx dy = 1$ 

Contours: P(x,y)

Bin size  $\rightarrow$  0 : **Probability distribution function P(x)** 

High PDF value

 $\Rightarrow$  High chance to get a measurement here

P(x) > 0,  $\int P(x) dx = 1$ 



### **Continuous Variables: PDFs**

**Continuous variable**: can consider **per-bin** probabilities p<sub>i</sub>, i=1.. n<sub>bins</sub>



Generalizes to **multiple variables** :  $P(x,y) > 0, \int P(x,y) dx dy = 1$ 

Contours: P(x,y)

Bin size  $\rightarrow$  0 : **Probability distribution function P(x)** 

High PDF value

 $\Rightarrow$  High chance to get a measurement here

P(x) > 0,  $\int P(x) dx = 1$ 



### **Continuous Variables: PDFs**

**Continuous variable**: can consider **per-bin** probabilities p<sub>i</sub>, i=1.. n<sub>bins</sub>



Generalizes to **multiple variables** :  $P(x,y) > 0, \int P(x,y) dx dy = 1$ 

Contours: P(x,y)

Bin size  $\rightarrow$  0 : **Probability distribution function P(x)** 

High PDF value

 $\Rightarrow$  High chance to get a measurement here

P(x) > 0,  $\int P(x) dx = 1$ 



### **Random Variables**

X, Y... are **Random Variables** (continuous or discrete), a.ka. **observables** :  $\rightarrow$  X can take any value x, with probability **P(X=x)**.

 $\rightarrow$  P(X=x) is the **PDF** of X, a.k.a. the **Statistical Model**.

→ The **Observed data** is **one value**  $x_{obs}$  of X, drawn from P(X=x).







### **PDF Properties: Mean**

**E(X) = <X> : Mean** of X – expected outcome on average over many measurements

$$\langle X \rangle = \sum_{i} x_{i} P_{i}$$
 or  
 $\langle X \rangle = \int x P(x) dx$ 

 $\rightarrow$  Property of the **PDF** 

For measurements  $x_1 \dots x_n$ , then can compute the **Sample mean**:

$$\bar{x} = \frac{1}{n} \sum_{i} x_{i}$$

- $\rightarrow$  Property of the sample
- $\rightarrow$  approximates the PDF mean.



### **PDF Properties: (Co)variance**

Variance of X:

$$\operatorname{Var}(X) = \langle (X - \langle X \rangle)^2 \rangle$$

- → Average square of deviation from mean → RMS(X) =  $\sqrt{Var(X)} = \sigma_x$  standard deviation
- Can be approximated by **sample variance**:

$$\hat{\sigma}^2 = \frac{1}{n-1} \sum_i (x_i - \bar{x})^2$$

**Covariance of X and Y:** 

$$\operatorname{Cov}(X,Y) = \langle (X - \langle X \rangle) (Y - \langle Y \rangle) \rangle$$

 $\rightarrow$  Large if variations of X and Y are "synchronized"

**Correlation coefficient** 

$$\mathbf{b} = \frac{\operatorname{Cov}(X, Y)}{\sqrt{\operatorname{Var}(X)\operatorname{Var}(Y)}}$$







### **PDF Properties: (Co)variance**

Variance of X:

$$\operatorname{Var}(X) = \langle (X - \langle X \rangle)^2 \rangle$$

- → Average square of deviation from mean → RMS(X) =  $\sqrt{Var(X)} = \sigma_x$  standard deviation
- Can be approximated by **sample variance**:

$$\hat{\sigma}^2 = \frac{1}{n-1} \sum_i (x_i - \bar{x})^2$$

**Covariance of X and Y:** 

$$\operatorname{Cov}(X,Y) = \langle (X - \langle X \rangle) (Y - \langle Y \rangle) \rangle$$

 $\rightarrow$  Large if variations of X and Y are "synchronized"

**Correlation coefficient** 

$$o = \frac{\operatorname{Cov}(X, Y)}{\sqrt{\operatorname{Var}(X)\operatorname{Var}(Y)}}$$







### **PDF Properties: (Co)variance**

Variance of X:

$$\operatorname{Var}(X) = \langle (X - \langle X \rangle)^2 \rangle$$

- → Average square of deviation from mean → RMS(X) =  $\sqrt{Var(X)} = \sigma_x$  standard deviation
- Can be approximated by **sample variance**:

$$\hat{\sigma}^2 = \frac{1}{n-1} \sum_i (x_i - \bar{x})^2$$

**Covariance of X and Y:** 

$$\operatorname{Cov}(X,Y) = \langle (X - \langle X \rangle) (Y - \langle Y \rangle) \rangle$$

 $\rightarrow$  Large if variations of X and Y are "synchronized"

**Correlation coefficient** 

$$\mathbf{b} = \frac{\operatorname{Cov}(X, Y)}{\sqrt{\operatorname{Var}(X)\operatorname{Var}(Y)}}$$







### "Linear" vs. "non-linear" correlations

For non-Gaussian cases, the **Correlation coefficient p** is not the whole story:



Source: Wikipedia

In particular, variables can still be correlated even when  $\rho=0$ : "Non-linear" correlations.

### **Gaussian PDF**

**Gaussian distribution:** 

$$G(x; X_0, \sigma) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x-X_0)^2}{2\sigma^2}}$$

- → Mean :  $X_0$ → Variance :  $\sigma^2$  (⇒ RMS =  $\sigma$ )
- Generalize to N dimensions:  $\rightarrow$  Mean :  $X_0$
- → Covariance matrix :

$$C = \begin{bmatrix} \operatorname{Var}(X_1) & \operatorname{Cov}(X_1, X_2) \\ \operatorname{Cov}(X_2, X_1) & \operatorname{Var}(X_2) \end{bmatrix}$$
$$= \begin{bmatrix} \sigma_1^2 & \rho \sigma_1 \sigma_2 \\ \rho \sigma_1 \sigma_2 & \sigma_2^2 \end{bmatrix}$$





### **Gaussian Quantiles**

Consider

$$z = \left(\frac{x - x_0}{\sigma}\right) \quad \text{``pull'' of x}$$

 $G(x;x_0,\sigma)$  depends only on  $z \sim G(z; 0, 1)$ 

Probability  $P(|x - x_0| > Z\sigma)$  to be away from the mean:

Gaussian Cumulative Distribution Function (CDF) :

$$\Phi(z) = \int_{-\infty}^{z} G(u; 0, 1) \, du$$

Z	$P( x - x_0  > Z\sigma)$
1	0.317
2	0.045
3	0.003
4	<b>3 x 10</b> <sup>-5</sup>
5	<b>6 x 10</b> <sup>-7</sup>

P(|x - x<sub>0</sub>| < 1♂) = 68.3 %



### **Gaussian Quantiles**

Consider

$$z = \left(\frac{x - x_0}{\sigma}\right) \quad \text{``pull'' of x}$$

 $G(x;x_0,\sigma)$  depends only on  $z \sim G(z; 0, 1)$ 

Probability  $P(|x - x_0| > Z\sigma)$  to be away from the mean:



$$\Phi(z) = \int_{-\infty}^{z} G(u; 0, 1) \, du$$



P(|x - x₀| < 2♂) = 95.4 % 0.4 0.35 0.3 0.25 0.2 0.15 0.1 0.05 1 -3 -2 -1 0 2 3 42**3**
## **Gaussian Quantiles**

Consider

$$z = \left(\frac{x - x_0}{\sigma}\right) \quad \text{``pull'' of x}$$

 $G(x;x_0,\sigma)$  depends only on  $z \sim G(z; 0, 1)$ 

Probability  $P(|x - x_0| > Z\sigma)$  to be away from the mean:

Gaussian Cumulative Distribution Function (CDF) :

$$\Phi(z) = \int_{-\infty}^{z} G(u; 0, 1) \, du$$





# **Central Limit Theorem**

(\*) Assuming  $\sigma_x < \infty$ and other regularity conditions

For an observable X with **any**<sup>(\*)</sup> **distribution**, one has

$$\bar{x} = \frac{1}{n} \sum_{i=1}^{n} x_i \stackrel{n \to \infty}{\sim} G(\langle X \rangle, \frac{\sigma_X}{\sqrt{n}})$$

What this means:

- The average of many measurements is always Gaussian, whatever the distribution for a single measurement
- The mean of the Gaussian is the average of the single measurements
- The **RMS** of the Gaussian **decreases as** √**n** : smaller fluctuations when averaging over many measurements

Another version:

$$\sum_{i=1}^{n} x_{i} \stackrel{n \to \infty}{\sim} G(n \langle X \rangle, \sqrt{n} \sigma_{X})$$

Mean scales like n, but RMS only like  $\sqrt{n}$ 

Draw events from a parabolic distribution (e.g. decay  $\cos \theta^*$ )



**Distribution becomes Gaussian**, although very non-Gaussian originally **Distribution becomes narrower** as expected (as  $1/\sqrt{n}$ )

Draw events from a parabolic distribution (e.g. decay  $\cos \theta^*$ )



**Distribution becomes Gaussian**, although very non-Gaussian originally **Distribution becomes narrower** as expected (as  $1/\sqrt{n}$ )

Draw events from a parabolic distribution (e.g. decay  $\cos \theta^*$ )



**Distribution becomes Gaussian**, although very non-Gaussian originally **Distribution becomes narrower** as expected (as  $1/\sqrt{n}$ )

Draw events from a parabolic distribution (e.g. decay  $\cos \theta^*$ )



**Distribution becomes Gaussian**, although very non-Gaussian originally **Distribution becomes narrower** as expected (as  $1/\sqrt{n}$ )

Draw events from a parabolic distribution (e.g. decay  $\cos \theta^*$ )



**Distribution becomes Gaussian**, although very non-Gaussian originally **Distribution becomes narrower** as expected (as  $1/\sqrt{n}$ )





**Distribution becomes Gaussian**, although very non-Gaussian originally **Distribution becomes narrower** as expected (as  $1/\sqrt{n}$ )

Draw events from a parabolic distribution (e.g. decay  $\cos \theta^*$ )



**Distribution becomes Gaussian**, although very non-Gaussian originally **Distribution becomes narrower** as expected (as  $1/\sqrt{n}$ )

Draw events from a parabolic distribution (e.g. decay  $\cos \theta^*$ )



**Distribution becomes Gaussian**, although very non-Gaussian originally **Distribution becomes narrower** as expected (as  $1/\sqrt{n}$ )

# **Chi-squared**

Multiple Independent Gaussian variables x<sub>i</sub>: Define

$$\chi^2 = \sum_{i=1}^n \left( \frac{x_i - x_i^0}{\sigma_i} \right)^2$$

Measures global distance from reference point  $(x_1^{0} \dots x_n^{0})$ 

Distribution depends on n :

Rule of thumb:

 $\chi^2/n$  should be  $\preceq 1$ 



# **Chi-squared**

Multiple Independent Gaussian variables x<sub>i</sub>: Define

$$\chi^2 = \sum_{i=1}^n \left( \frac{x_i - x_i^0}{\sigma_i} \right)^2$$

Measures global distance from reference point  $(x_1^{0} \dots x_n^{0})$ 

Distribution depends on n :

Rule of thumb:

 $\chi^2/n$  should be  $\preceq 1$ 



## **Histogram Chi-squared**

Histogram  $\chi$ 2 with respect to a reference shape:

- Assume an independent Gaussian distribution in each bin
- Degrees of freedom = (number of bins) (number of fit parameters)



**BLUE histogram vs. flat reference**  $\chi^2 = 12.9$ ,  $p(\chi^2=12.9, n=10) = 23\%$ 

# **Histogram Chi-squared**

Histogram  $\chi 2$  with respect to a reference shape:

- Assume an independent Gaussian distribution in each bin
- Degrees of freedom = (number of bins) (number of fit parameters)



BLUE histogram vs. flat reference  $\chi^2 = 12.9$ ,  $p(\chi^2=12.9, n=10) = 23\%$ RED histogram vs. flat reference  $\chi^2 = 38.8$ ,  $p(\chi^2=38.8, n=10) = 0.003\%$ 

27

# **Histogram Chi-squared**

Histogram  $\chi 2$  with respect to a reference shape:

- Assume an independent Gaussian distribution in each bin
- Degrees of freedom = (number of bins) (number of fit parameters)



BLUE histogram vs. flat reference  $\chi^2 = 12.9$ ,  $p(\chi^2=12.9, n=10) = 23\%$ RED histogram vs. flat reference  $\chi^2 = 38.8$ ,  $p(\chi^2=38.8, n=10) = 0.003\%$ RED histogram vs. correct reference  $\chi^2 = 9.5$ ,  $p(\chi^2=9.5, n=10) = 49\%$ 

### **Error Bars**

Strictly speaking, *the uncertainty is given by the model* :

- $\rightarrow$  **Bin central value** ~ mean of the bin PDF
- $\rightarrow$  **Bin uncertainty** ~ RMS of the bin PDF

The data is just what it is, a simple observed point.

⇒ One should in principle **show the error bar on the prediction**.

 $\rightarrow$  In practice, the usual convention is to have error bars on the data points.



### **Error Bars**

Strictly speaking, *the uncertainty is given by the model* :

- $\rightarrow$  **Bin central value** ~ mean of the bin PDF
- $\rightarrow$  **Bin uncertainty** ~ RMS of the bin PDF

The data is just what it is, a simple observed point.

⇒ One should in principle **show the error bar on the prediction**.

 $\rightarrow$  In practice, the usual convention is to have error bars on the data points.



# **Statistical Modeling**

# **Example 1: Z counting**

Measure the cross-section (event rate) of the  $Z \rightarrow$  ee process





#### $\sigma^{\text{fid}} = 0.781 \pm 0.004 \text{ (stat)} \pm 0.018 \text{ (syst) nb}$

Fluctuations in the data counts

Other uncertainties (assumptions, parameter values)

"Single bin counting" : only data input is n<sub>data</sub>.

### Example 2: ttH→bb

#### arXiv:2111.06712



### Event counting in different regions: *Multiple-bin counting*

#### Lots of information available

- $\rightarrow$  Potentially higher sensitivity
- $\rightarrow$  How to make optimal use of it ?

# **Example 3: unbinned modeling**



All modeling done using continuous distributions:

$$P_{\text{total}}(m_{\gamma\gamma}) = \frac{S}{S+B} P_{\text{signal}}(m_{\gamma\gamma}; m_H) + \frac{B}{S+B} P_{\text{bkg}}(m_{\gamma\gamma})$$

- $\rightarrow$  In principle, binomial process
- $\rightarrow$  In practice, **P**  $\ll$  **1**, **N**  $\gg$  **1**,  $\Rightarrow$  Poisson approximation.
- $\rightarrow$  *i.e.* **very rare** process, but **very many trials** so still expect to see good events



- $\rightarrow$  In principle, binomial process
- $\rightarrow$  In practice, **P**  $\ll$  **1**, **N**  $\gg$  **1**,  $\Rightarrow$  Poisson approximation.
- $\rightarrow$  *i.e.* **very rare** process, but **very many trials** so still expect to see good events



- $\rightarrow$  In principle, binomial process
- $\rightarrow$  In practice, **P**  $\ll$  **1**, **N**  $\gg$  **1**,  $\Rightarrow$  Poisson approximation.
- $\rightarrow$  *i.e.* **very rare** process, but **very many trials** so still expect to see good events



- $\rightarrow$  In principle, binomial process
- $\rightarrow$  In practice, **P**  $\ll$  **1**, **N**  $\gg$  **1**,  $\Rightarrow$  Poisson approximation.
- $\rightarrow$  *i.e.* **very rare** process, but **very many trials** so still expect to see good events



- $\rightarrow$  In principle, binomial process
- $\rightarrow$  In practice, **P**  $\ll$  **1**, **N**  $\gg$  **1**,  $\Rightarrow$  Poisson approximation.
- $\rightarrow$  *i.e.* **very rare** process, but **very many trials** so still expect to see good events



- $\rightarrow$  In principle, binomial process
- $\rightarrow$  In practice, **P**  $\ll$  **1**, **N**  $\gg$  **1**,  $\Rightarrow$  Poisson approximation.
- $\rightarrow$  *i.e.* **very rare** process, but **very many trials** so still expect to see good events



# **Statistical Model for Counting**

#### **Observable: number of events n**

Typically both Signal and Background present:

$$P(n; S, B) = e^{-(S+B)} \frac{(S+B)^n}{n!}$$



Model has **parameters S** and **B**.

B can be known a priori or not (S usually not...)

 $\rightarrow$  Example: **assume B is known**, use **measured n** to find out about **S**.



# **Multiple counting bins**



**Shapes f** typically obtained from simulated events (*Monte Carlo*)

 $\rightarrow$  HEP: typically excellent modeling from simulation, although some uncertainties need to be accounted for.

However not always possible to generate sufficiently large MC samples MC stat fluctuations can create artefacts, especially for  $S \ll B$ .

## **Model Parameters**

Model typically includes:

• Parameters of interest (POIs) : what we want to measure

 $\rightarrow$  S, m<sub>w</sub>, ...

• Nuisance parameters (NPs) : other parameters needed to define the model

 $\rightarrow$  Background levels (B)

 $\rightarrow$  For binned data,  $f^{sig}$ ,  $f^{bkg}$ 

NPs must be either:

- → Known a priori (within uncertainties) or
- $\rightarrow$  Constrained by the data



# **Takeaways**

Random data must be described using a statistical model:



Description	Observable	Likelihood
Counting	n	Poisson $P(\mathbf{n}; \mathbf{S}, \mathbf{B}) = e^{-(\mathbf{S} + \mathbf{B})} \frac{(\mathbf{S} + \mathbf{B})^{\mathbf{n}}}{\mathbf{n}!}$
Binned shape analysis	n <sub>i</sub> , i = 1 N <sub>bins</sub>	Poisson product $P(\mathbf{n}_{i}; \mathbf{S}, \mathbf{B}) = \prod_{i=1}^{n_{\text{bins}}} e^{-(\mathbf{S} f_{i}^{\text{sig}} + \mathbf{B} f_{i}^{\text{bkg}})} \frac{(\mathbf{S} f_{i}^{\text{sig}} + \mathbf{B} f_{i}^{\text{bkg}})^{\mathbf{n}_{i}}}{\mathbf{n}_{i}!}$
Unbinned shape analysis	m <sub>i</sub> , i = 1 n <sub>evts</sub>	Extended Unbinned Likelihood $P(\boldsymbol{m_i}; \boldsymbol{S}, \boldsymbol{B}) = \frac{e^{-(\boldsymbol{S} + \boldsymbol{B})}}{\boldsymbol{n_{\text{evts}}}!} \prod_{i=1}^{\boldsymbol{n_{\text{evts}}}} \boldsymbol{S} P_{\text{sig}}(\boldsymbol{m_i}) + \boldsymbol{B} P_{\text{bkg}}(\boldsymbol{m_i})$

Model can include multiple **categories**, each with a separate description Includes **parameters of interest** (POIs) but also **nuisance parameters** (NPs) **Next step**: use the model to obtain information on the POIs

# Hypothesis Testing and discovery



# **Discovery Testing**

We see an unexpected feature in our data, is it a signal for new physics or a fluctuation ?

e.g. Higgs discovery : **"We have 5σ" !** 



GeV

Events/5 ( 02 02

15

10

Data

///// Syst.Unc.

Background ZZ(\*)

\_√s = 7 TeV:∫Ldt = 4.8 fb<sup>-1</sup>

√s = 8 TeV: ∫Ldt = 5.8 fb<sup>-1</sup>

Background Z+jets, tt Signal (m\_=125 GeV) ATLAS

 $H \rightarrow ZZ^{(*)} \rightarrow 4I$ 

# **Discovery Testing**

Say we have a Gaussian measurement with a background **B=100**, and we measure **n=120** 

Did we just discover something ? *Maybe :-)* (but not very likely)

The measured signal is S = 20.  $S = n_{obs} - B$ 

Uncertainty on B is  $\sqrt{B} = 10$  $\Rightarrow$  Significance Z = 2  $\Rightarrow$  we are  $\sim 2\sigma$  away from S=0.

#### Gaussian quantiles :

Z = 2 happens  $p_0 \sim 2.3\%$  of the time if S=0

 $p_0 = 1 - \Phi(Z)$ 

 $\Rightarrow$  Rare, but not exceptional

$$Z = \frac{S}{\sqrt{B}}$$

$$= \frac{S}{\sqrt{B}}$$

$$= 0$$

$$B = 100$$

$$B = 100$$

$$h$$

$$\Phi(Z) = \int_{-\infty}^{Z} G(u; 0, 1) du$$



Obs: 120

# **Discovery Testing**

n



Straightf	р <sub>о</sub>	Z	S	n <sub>obs</sub>
Need to	31%	0.5σ	5	105
more co	16%	1σ	10	110
Fvidence	2.3%	2σ	20	120
Discovery	0.1%	3σ	30	130

Straightforward in this Gaussian case Need to be able to do the same in more complex cases: • Determine S

• Compute Z and p<sub>0</sub> 41 /

# **Maximum Likelihood Estimation**
# What a PDF is for

Model describes the distribution of the observable: P(data; parameters) ⇒ Possible outcomes of the experiment, for given parameter values Can draw random events according to PDF : generate pseudo-data



# What a PDF is also for: Likelihood

Model describes the distribution of the observable: P(data; parameters) ⇒ Possible outcomes of the experiment, for given parameter values We want the other direction: use data to get information on parameters



**Likelihood:** L(parameters) = P(data; parameters)

 $\rightarrow$  same as the PDF, but seen as function of the parameters

# **Maximum Likelihood Estimation**

To estimate a parameter  $\mu$ , find the value  $\hat{\mu}$  that maximizes L( $\mu$ )

Maximum Likelihood Estimator (MLE) **û**:

$$\hat{\mathbf{L}} = arg max L(\boldsymbol{\mu})$$



**MLE**: the value of μ for which **this data** was **most likely to occur The MLE is a function of the data** – itself an **observable** *No guarantee* it is the true value (data may be "unlikely") but sensible estimate

45 /

#### **Gaussian case**



**Best-fit** of Gaussian PDF mean to observed data

#### **Gaussian case**



**Best-fit** of Gaussian PDF mean to observed data

46 /

#### **Gaussian case**



**Best-fit** of Gaussian PDF mean to observed data

46 /



-2 log Likelihood:

$$\lambda(\mu) = -2 \log L(\mu) = \sum_{i=1}^{N_{\text{bins}}} \left(\frac{n_i - \mu_i}{\sigma_i}\right)^2$$

However typically need to perform non-linear minimization in other cases.

- MINUIT (C++ library within ROOT, numerical gradient descent)
- **scipy.minimize** using NumPy/TensorFlow/PyTorch/... backends
  - $\rightarrow$  Many algorithms gradient-based, etc.



-2 log Likelihood:

$$\lambda(\mu) = -2 \log L(\mu) = \sum_{i=1}^{N_{\text{bins}}} \left( \frac{n_i - \mu_i}{\sigma_i} \right)^2$$

However typically need to perform non-linear minimization in other cases.

- MINUIT (C++ library within ROOT, numerical gradient descent)
- **scipy.minimize** using NumPy/TensorFlow/PyTorch/... backends
  - $\rightarrow$  Many algorithms gradient-based, etc.



-2 log Likelihood:

$$\lambda(\mu) = -2 \log L(\mu) = \sum_{i=1}^{N_{\text{bins}}} \left( \frac{n_i - \mu_i}{\sigma_i} \right)^2$$

However typically need to perform non-linear minimization in other cases.

- MINUIT (C++ library within ROOT, numerical gradient descent)
- **scipy.minimize** using NumPy/TensorFlow/PyTorch/... backends
  - $\rightarrow$  Many algorithms gradient-based, etc.



-2 log Likelihood:

$$\lambda(\mu) = -2 \log L(\mu) = \sum_{i=1}^{N_{\text{bins}}} \left( \frac{n_i - \mu_i}{\sigma_i} \right)^2$$

However typically need to perform non-linear minimization in other cases.

- MINUIT (C++ library within ROOT, numerical gradient descent)
- **scipy.minimize** using NumPy/TensorFlow/PyTorch/... backends
  - $\rightarrow$  Many algorithms gradient-based, etc.

# **Hypothesis Testing**

Null Hypothesis: assumption on POIs, say value of S (e.g. H<sub>0</sub> : S=0)

 $\rightarrow$  Goal : decide if H<sub>0</sub> is favored or disfavored using a test based on the data

Possible outcomes:	Data disfavors H <sub>o</sub> (Discovery claim)			Data favors H <sub>o</sub> (Nothing found)	
H <sub>o</sub> is false (New physics!)	Discovery!			Missed discovery	
H <sub>o</sub> is true (Nothing new)	False discovery			No new physics, None found	Image: second

"... the null hypothesis is never proved or established, but is possibly disproved, in the course of experimentation. Every experiment may be said to exist only to give the facts a chance of disproving the null hypothesis." – R. A. Fisher

# **Hypothesis Testing**

Hypothesis: assumption on model parameters, say value of S (e.g. H<sub>o</sub>: S=0)

	Data disfavor (Discovery cla	s H <sub>o</sub> aim)	Data favors H <sub>o</sub> (Nothing found)	
H <sub>0</sub> is false (New physics!)	Discovery!		Type-II error (Missed discovery)	
H <sub>₀</sub> is true (Nothing new)	Type-I error (False discovery)		No new physics, none found	
	Ĺ	— n-valu	e significance	

**Lower Type-I errors** ⇔ **Higher Type-II errors** and vice versa: cannot have everything!

 $\rightarrow$  Goal: test that minimizes Type-II errors for a given level of Type-I error.



#### **ROC Curves**

#### more powerful Better discriminators "Receiver operating characteristic" (ROC) Curve: ිස $\rightarrow$ Shows Type-I vs Type-II rates for different selections Better ε<sub>Type-I</sub> (= $\rightarrow$ All curves monotonically decrease from (0,1) to (1,0) $\rightarrow$ Better discriminators more bent • towards (1,1) 0 $1 - \varepsilon_{\text{Type-II}} (= \varepsilon_{\text{S}})$ 0.4 S = 0**BSM** $\rightarrow$ **Goal**: test that minimizes Type-II 0.35 0.3 0.25 0.2 Type-I error 0.15 Type-I<mark>/</mark> Error p-value 0.1 0.05 0<u></u>5 5 50 2 -3 -2 0 3 4 \_4 -1 Discriminant observable

Increasingly

errors for given level of Type-I error.

 $\rightarrow$  Usually set predefined level of acceptable Type-I error (e.g. "5 $\sigma$ ")

#### **ROC Curves**

#### more powerful Better discriminators "Receiver operating characteristic" (ROC) Curve: ිස $\rightarrow$ Shows Type-I vs Type-II rates for different selections Better ε<sub>Type-I</sub> (= $\rightarrow$ All curves monotonically decrease from (0,1) to (1,0) $\rightarrow$ Better discriminators more bent • towards (1,1) 0 $1 - \varepsilon_{\text{Type-II}} (= \varepsilon_{\text{S}})$ 0.4 S = 0**BSM** $\rightarrow$ **Goal**: test that minimizes Type-II 0.35 0.3 errors for given level of Type-I error. 0.25 0.2 Type-I error 0.15 Type-I/ Error p-value 0 $\rightarrow$ Usually set predefined level of 0.05 acceptable Type-I error (e.g. "5 $\sigma$ ") 50 -3 -2 0 2 3 -1 4 5 Discriminant observable

Increasingly

#### **ROC Curves**

#### more powerful Better discriminators "Receiver operating characteristic" (ROC) Curve: ිස $\rightarrow$ Shows Type-I vs Type-II rates for different selections Better ε<sub>Type-I</sub> (= $\rightarrow$ All curves monotonically decrease from (0,1) to (1,0) $\rightarrow$ Better discriminators more bent \_ towards (1,1) 0 $1 - \varepsilon_{\text{Type-II}} (= \varepsilon_{\text{S}})$ 0.4 S = 0**BSM** $\rightarrow$ **Goal**: test that minimizes Type-II 0.35 0.3 errors for given level of Type-I error. 0.25 0.2 Type-I error 0.15 Type-I/ Error p-value 0 $\rightarrow$ Usually set predefined level of 0.05 acceptable Type-I error (e.g. "5 $\sigma$ ") 50 -3 -2 0 2 3 4 5

/

Discriminant observable

Increasingly

# **Discovery Testing**

n



St	р <sub>о</sub>	Z	S	n <sub>obs</sub>
Ne	31%	0.5σ	5	105
m	16%	1σ	10	110
) Fvi	2.3%	2σ	20	120
Dis	0.1%	3σ	30	130

Straightforward in this Gaussian case Need to be able to do the same in more complex cases: • Determine S

- Compute Z and P<sub>0</sub>51

# **Testing for Evidence in Gaussian counting**



52 /

## **Testing for Evidence in Gaussian counting**



52 /

#### **Neyman-Pearson Lemma**

When comparing two hypotheses  $H_0$  and  $H_1$ , the

optimal discriminator is the Likelihood ratio (LR)

$$\frac{L(S = 5; data)}{L(S = 0; data)}$$

e.g.

**Caveat**: Strictly true only for *simple hypotheses* (no free parameters)

As for MLE, choose the hypothesis that is more likely given the data we have.

- $\rightarrow$  Always need an **alternate hypothesis** to test against the **null**.
- $\rightarrow$  **Minimizes Type-II uncertainties** for given level of Type-I uncertainties

 $\rightarrow$  In the following: all tests based on LR, will focus on p-values (Type-I errors), trusting that Type-II errors are anyway as small as they can be...

$$\frac{L(\mathbf{H}_{1}; data)}{L(\mathbf{H}_{0}; data)}$$

53

# **Discovery: Test Statistic**

Cowan, Cranmer, Gross & Vitells, Eur.Phys.J.C71:1554,2011

**Discovery**:

- H<sub>0</sub>: background only (S = 0) against
- H<sub>1</sub>: presence of a signal (S > 0)



 $\rightarrow$  For H<sub>1</sub>, any S > 0 is possible, which to use ? The one preferred by the data,  $\hat{S}$ .

⇒ Use Likelihood ratio:

$$\frac{L(S=0)}{L(\hat{S})}$$

 $\rightarrow$  In fact use the **test statistic**  $q_0 = -2\log \frac{L(S=0)}{L(\hat{S})}$ 

**Note**: for  $\hat{S} < 0$ , set  $q_0 = 0$  to reject negative signals ("one-sided test statistic")  $\frac{54}{7}$ 

# **Discovery p-value**

Large values of 
$$-2 \log \frac{L(S=0)}{L(\hat{S})}$$
 if:

 $\Rightarrow$  observed  $\hat{S}$  is far from 0

$$\Rightarrow$$
 H<sub>0</sub>(S=0) disfavored compared to H<sub>1</sub>(S≠0).

 $\Rightarrow$  Large  $\hat{S}$  !

Compute *p-value* in the tail of the distribution

to exclude H<sub>o</sub> (... and claim a discovery!)



Need to know  $f(q_0 | S=0)$ , the distribution of the test statistic...

# **Asymptotic distribution of q**<sub>0</sub>

**Gaussian regime for \hat{S}** (e.g. large  $n_{evts}$ , Central-limit theorem) :

Wilk's Theorem:  $q_0$  distributed as  $\chi^2$  ( $n_{par}$ ) for S = 0

$$\Rightarrow$$
 n<sub>par</sub> = 1 :  $\sqrt{q_0}$  is distributed as a Gaussian

⇒ Can compute p-values from Gaussian quantiles

 $p_0 = 1 - \Phi(\sqrt{q_0})$ 

 $\Rightarrow$  Even more simply, the significance is:

 $Z=\sqrt{q_0}$ 

Typically works well already for for event counts of O(5) and above  $\Rightarrow$  Widely applicable

(\*) 1-line "proof": asymptotically L and S are Gaussian, so  $L(S) = \exp\left[-\frac{1}{2}\left(\frac{S-\hat{S}}{\sigma}\right)^2\right] \Rightarrow q_0 = \left(\frac{\hat{S}}{\sigma}\right)^2 \Rightarrow \sqrt{q_0} = \frac{\hat{S}}{\sigma} \sim G(0,1) \Rightarrow q_0 \sim \chi^2(n_{dof}=1)$ 



# **Homework 1: Gaussian Counting**

#### Count number of events n in data

- $\rightarrow$  Assume n large enough so process is Gaussian
- $\rightarrow$  Assume B is known, and we measure S

Likelihood: 
$$L(S; n_{obs}) = e^{-\frac{1}{2} \left( \frac{n_{obs} - (S+B)}{\sqrt{S+B}} \right)^2}$$

- → Find the best-fit value (MLE)  $\hat{S}$  for the signal (can use  $\lambda$  = -2 log L instead of L for simplicity)
- $\rightarrow$  Find the expression of  $q_0$  for  $\hat{S} > 0$ .
- $\rightarrow$  Find the expression for the significance





# **Homework 2: Poisson Counting**

Same problem but now *not* assuming Gaussian behavior:

$$L(S;n) = e^{-(S+B)}(S+B)^n$$

 $\rightarrow$  As before, compute  $\hat{S}$ , and  $q_0$ 

(Can remove the n! constant since we're only dealing with L ratios)

 $\rightarrow$  Compute Z =  $\sqrt{q_0}$ , assuming asymptotic behavior

Solution:  

$$Z = \sqrt{2 \left[ (\hat{S} + B) \log \left| 1 + \frac{\hat{S}}{B} \right| - \hat{S} \right]}$$

Exact result can be obtained using

pseudo-experiments  $\rightarrow$  close to  $\sqrt{q_0}$  result

Asymptotic formulas justified by Gaussian regime, but remain valid even for small values of S+B (down to 5 events!)



# **Discovery Thresholds**

Evidence :  $3\sigma \Leftrightarrow p_0 = 0.3\% \Leftrightarrow 1$  chance in 300

**Discovery:**  $5\sigma \Leftrightarrow p_0 = 3 \ 10^{-7} \Leftrightarrow 1 \ \text{chance in } 3.5\text{M}$ 

Why so high thresholds ? (from Louis Lyons):

 Look-elsewhere effect: searches typically cover multiple independent regions ⇒ Higher chance to have a fluctuation "somewhere"

 $N_{trials} \sim 1000 : local 5\sigma \Leftrightarrow O(10^{-4})$  more reasonable

 Mismodeled systematics: factor 2 error in syst-dominated analysis ⇒ factor 2 error on Z...



• **History**: 3σ and 4σ excesses do occur regularly, for the reasons above

**Extraordinary claims require extraordinary evidence!** 

# **Extra Slides**

# **Rare Processes ?**

**HEP** : almost always use Poisson

distributions. Why?

#### ATLAS :

• Event rate ~ 1 GHz

(L~10<sup>34</sup> cm<sup>-2</sup>s<sup>-1</sup>~10 nb<sup>-1</sup>/s,  $\sigma_{tot}$ ~10<sup>8</sup> nb, )

Trigger rate ~ 1 kHz

(Higgs rate ~ 0.1 Hz)

⇒ p ~ 10<sup>-6</sup> ≪ 1 (p<sub>H→γγ</sub> ~ 10<sup>-13</sup>)

A day of data: N ~  $10^{14} \gg 1$ 

⇒ Poisson regime! Similarly true in many other physics situations.

(Large N = design requirement, to get not-too-small  $\lambda$ =Np...)



# **Unbinned Shape Analysis**

**Observable**: set of values  $m_1 \dots m_n$ , one per event

- $\rightarrow$  Describe shape of the **distribution of m**
- $\rightarrow$  Deduce the **probability to observe m**<sub>1</sub>... m<sub>n</sub>



Vormalized events per GeV

0.25

0.2

0.15

0.1

0.05

m

110 120

130

140

150

160

Signal

Assume Poisson distribution with B = 0:  $P(n; S) = e^{-S} \frac{S^n}{n!}$ Say we observe n=5, want to infer information on the parameter S

- $\rightarrow$  Try different values of S for a fixed data value n=5
- $\rightarrow$  Varying parameter, fixed data: **likelihood**

$$L(S; n=5) = e^{-S} \frac{S^5}{5!}$$



Say we **observe n=5**, want to infer information on the parameter  $s^n = e^{-s} \frac{S^n}{n!}$   $\rightarrow$  Try different values of S for a fixed data use

- $\rightarrow$  Varying parameter, fixed data: **likelihood**





Say we **observe n=5**, want to infer information on the parameter  $s^n = e^{-s} \frac{S^n}{n!}$   $\rightarrow$  Try different values of S for a fixed data we

- $\rightarrow$  Varying parameter, fixed data: **likelihood**





Say we **observe n=5**, want to infer information on the parameter  $s^n = e^{-s} \frac{S^n}{n!}$   $\rightarrow$  Try different values of S for a fixed data we

- $\rightarrow$  Varying parameter, fixed data: **likelihood**





63

Say we **observe n=5**, want to infer information on the parameter  $s^n = e^{-s} \frac{S^n}{n!}$   $\rightarrow$  Try different values of S for a fixed data we

- $\rightarrow$  Varying parameter, fixed data: **likelihood**





# **MLEs in Shape Analyses**

Binned shape analysis:

$$L(\mathbf{S};\mathbf{n}_i) = P(\mathbf{n}_i;\mathbf{S}) = \prod_{i=1}^{N} \operatorname{Pois}(\mathbf{n}_i;\mathbf{S}f_i + B_i)$$

λT

Maximize global L(S) (each bin may prefer a different **S**) In practice easier to minimize

$$\lambda_{\text{Pois}}(\mathbf{S}) = -2\log L(\mathbf{S}) = -2\sum_{i=1}^{N} \log \text{Pois}(\mathbf{n}_i; \mathbf{S}f_i + B_i) \qquad \text{Needs a computer}$$

In the Gaussian limit

$$\lambda_{\text{Gaus}}(\mathbf{S}) = \sum_{i=1}^{N} -2\log G(\mathbf{n}_i; \mathbf{S}f_i + B_i, \sigma_i) = \sum_{i=1}^{N} \left| \frac{\mathbf{n}_i - (\mathbf{S}f_i + B_i)}{\sigma_i} \right|^2 \quad \chi^2 \text{ formula}$$

→ Gaussian MLE (min  $\chi^2$  or min  $\lambda_{Gaus}$ ) : Best fit value in a  $\chi^2$  (Least-squares) fit → Poisson MLE (min  $\lambda_{Pois}$ ) : Best fit value in a likelihood fit (in ROOT, fit option "L") In RooFit,  $\lambda_{Pois}$  ⇒ RooAbsPdf::fitTo(),  $\lambda_{Gaus}$  ⇒ RooAbsPdf::chi2FitTo().

#### In both cases, MLE ⇔ Best Fit



Classification BDT output

Н→үү



65

# **MLE Properties**

• Asymptotically Gaussian  $P(\hat{\mu}) \propto \exp\left(-\frac{(\hat{\mu}-\mu^*)^2}{2\sigma_{\hat{\mu}}^2}\right)$  for  $n \rightarrow \infty$ and unbiased  $\langle \hat{\mu} \rangle = \mu^*$  for  $n \rightarrow \infty$ Standard deviation of the distribution of  $\hat{\mu}$ 

for large enough datasets

- Asymptotically Efficient :  $\sigma_{\mu}$  is the lowest possible value (in the limit  $n \rightarrow \infty$ ) among consistent estimators.
  - $\rightarrow$  MLE captures all the available information in the data
- Also **consistent**:  $\hat{\mu}$  converges to the true value for large n,
- Log-likelihood : Can also minimize  $\lambda = -2 \log L$ 
  - $\rightarrow$  Usually more efficient numerically
  - $\rightarrow$  For Gaussian L,  $\lambda$  is parabolic:
- Can drop multiplicative constants in L (additive constants in  $\lambda$ )

 $\hat{\mathbf{u}} \xrightarrow{n \to \infty} \mathbf{u}^*$
### **Extra: Fisher Information**

### **Fisher Information:**

$$I(\mu) = \left| \left( \frac{\partial}{\partial \mu} \log L(\mu) \right)^2 \right| = - \left| \frac{\partial^2}{\partial \mu^2} \log L(\mu) \right|^2$$

Measures the **amount of information** available in the measurement of  $\mu$ .



For any estimator  $\tilde{\mu}$ .

- $\rightarrow$  cannot be more precise than allowed by information in the measurement.
- **Efficient** estimators reach the bound : **e.g. MLE in the large dataset limit.**

### **Some Examples**

#### Higgs Discovery: Phys. Lett. B 716 (2012) 1-29



### High-mass $X \rightarrow \gamma \gamma$ Search: JHEP 09 (2016) 1



# **Upper Limit Pathologies**

Upper limit: 
$$S_{up} \sim \hat{S} + 1.64 \sigma_{s}$$

**Problem**: for negative Ŝ, get **very** good observed limit.

 $\rightarrow$  For  $\hat{S}$  sufficiently negative, even  $S_{up} < 0$  !

How can this be ?

### → Background modeling issue ?... Or:

→ This is a 95% limit  $\Rightarrow$  5% of the time, the limit wrongly excludes the true value, e.g. S\*=0.

### **Options**

 $\rightarrow$  live with it: sometimes report limit < 0

 $\rightarrow$  Special procedure to avoid these cases, since if we assume S must be >0, we know a priori this is just a fluctuation.





The usual p-value under Usual solution in HEP : **CL**.  $\boldsymbol{p}_{S_0}$ H(S=S<sub>0</sub>) (=5%)  $p_{CL_s}$  –  $\rightarrow$  Compute modified p-value The p-value computed  $\Rightarrow$  **Rescale** exclusion at S<sub>0</sub> by exclusion at S=0. under H(S=0)  $\rightarrow$  Somewhat ad-hoc, but good properties... ц ц 95% limit, CL<sub>s+b</sub> **Ŝ compatible with 0** :  $p_{B} \sim O(1)$ 95% limit, CL  $p_{CLs} \sim p_{so} \sim 5\%$ , no change. **Far-negative**  $\hat{\mathbf{S}}$  : 1 -  $p_{R} \ll 1$  $p_{Cls} \sim p_{S0} / (1-p_B) \gg 5\%$  $\rightarrow$  lower exclusion  $\Rightarrow$  higher limit, σ<sub>s</sub> = 1 usually >0 as desired

#### Drawback: overcoverage

 $\rightarrow$  limit is claimed to be 95% CL, but actually >95% CL for small 1-p<sub>B</sub>.

# **CL**<sub>s</sub> : Gaussian Bands

Usual Gaussian counting example with known B: 95%  $CL_s$  upper limit on S:

$$S_{up} = \hat{S} + \left[ \Phi^{-1} \left( 1 - 0.05 \Phi(\hat{S}/\sigma_s) \right) \right] \sigma_s$$

Compute expected bands for S=0:

→ Asimov dataset  $\Leftrightarrow \hat{S} = 0$ →  $\pm n\sigma$  bands:

$$S_{up,exp}^{0} = 1.96 \sigma_{s}$$
  

$$S_{up,exp}^{\pm n} = \left(\pm n + \left[1 - \Phi^{-1}(0.05 \Phi(\mp n))\right]\right) \sigma_{s}$$



CLs :

- Positive bands somewhat reduced,
- Negative ones more so

Band width from  $\sigma_{s,A}^2 = \frac{S^2}{q_s(\text{Asimov})}$ depends on S, for non-Gaussian cases, different values for each band...



### **Comparison with LEP/TeVatron definitions**

Likelihood ratios are not a new idea:

- **LEP**: Simple LR with NPs from MC
  - Compare  $\mu$ =0 and  $\mu$ =1
- **Tevatron**: PLR with profiled NPs

Both compare to  $\mu=1$  instead of best-fit  $\hat{\mu}$ 



 $\rightarrow$  Asymptotically:

- **LEP/Tevaton**: q linear in  $\mu \Rightarrow$  **~Gaussian**
- LHC: q quadratic in  $\mu \Rightarrow -\chi 2$

 $\rightarrow$  Still use TeVatron-style for discrete cases

$$q_{LEP} = -2\log\frac{L(\mu=0,\widetilde{\theta})}{L(\mu=1,\widetilde{\theta})}$$
$$q_{Tevatron} = -2\log\frac{L(\mu=0,\widehat{\theta}_0)}{L(\mu=1,\widehat{\theta}_1)}$$



Probabilistic treatment of possible outcomes ⇒ **Probability Distribution** 

**Example**: two-coin toss

 $\rightarrow$  Fractions of events in each bin i converge to a limit p<sub>i</sub>

### **Probability distribution** :

{ P<sub>i</sub> } for i = 0, 1, 2

### Properties

- P<sub>i</sub> > 0
- Σ P<sub>i</sub>=1



Probabilistic treatment of possible outcomes ⇒ Probability Distribution

**Example**: two-coin toss

 $\rightarrow$  Fractions of events in each bin i converge to a limit p<sub>i</sub>

### **Probability distribution** :

 $\{P_i\}$  for i = 0, 1, 2

#### Properties

- P<sub>i</sub> > 0
- Σ P<sub>i</sub>=1



Probabilistic treatment of possible outcomes ⇒ Probability Distribution

**Example**: two-coin toss

 $\rightarrow$  Fractions of events in each bin i converge to a limit p<sub>i</sub>

### **Probability distribution :**

 $\{P_i\}$  for i = 0, 1, 2

### **Properties**

- P<sub>i</sub> > 0
- Σ P<sub>i</sub>=1



#### 100 trials

Probabilistic treatment of possible outcomes ⇒ Probability Distribution

Example: two-coin toss

 $\rightarrow$  Fractions of events in each bin i converge to a limit p<sub>i</sub>

### **Probability distribution** :

{ P<sub>i</sub> } for i = 0, 1, 2

### Properties

- P<sub>i</sub> > 0
- Σ P<sub>i</sub>=1



73

### **Continuous Variables: PDFs**

**Continuous variable**: can consider **per-bin** probabilities p<sub>i</sub>, i=1.. n<sub>bins</sub>



Bin size  $\rightarrow$  0 : **Probability distribution function P(x)**   $\rightarrow$  High values  $\Leftrightarrow$  high chance to get a measurement here

 $P(x) > 0, \int P(x) dx = 1$ 



Contours: P(x,y)

### **Continuous Variables: PDFs**

**Continuous variable**: can consider **per-bin** probabilities p<sub>i</sub>, i=1.. n<sub>bins</sub>



Bin size  $\rightarrow$  0 : **Probability distribution function P(x)**   $\rightarrow$  High values  $\Leftrightarrow$  high chance to get a measurement here

 $P(x) > 0, \int P(x) dx = 1$ 



Contours: P(x,y)

### **Continuous Variables: PDFs**

**Continuous variable**: can consider **per-bin** probabilities p<sub>i</sub>, i=1.. n<sub>bins</sub>



Bin size  $\rightarrow$  0 : **Probability distribution function P(x)**   $\rightarrow$  High values  $\Leftrightarrow$  high chance to get a measurement here

 $P(x) > 0, \int P(x) dx = 1$ 



Contours: P(x,y)

### **Random Variables**

X, Y... are **Random Variables** (continuous or discrete), a.ka. **observables** :  $\rightarrow$  X can take any value x, with probability **P(X=x)**.

 $\rightarrow$  P(X) is the **PDF** of X, a.k.a. the **Statistical Model**.

→ The **Observed data** is **one value**  $x_{obs}$  of X, drawn from P(X).







75

### **PDF Properties: Mean**

**E(X) = <X> : Mean** of X – expected outcome on average over many measurements

$$\langle X \rangle = \sum_{i} x_{i} P_{i} \qquad \text{Or}$$
  
 
$$\Rightarrow \text{ Property of the } \dot{P} DF$$
  
 
$$\langle X \rangle = \int x P(x) dx$$

For measurements  $x_1 \dots x_n$ , then can compute the **Sample mean**:



### PDF Mean



# **PDF Properties: (Co)variance**

Variance of X:

$$\operatorname{Var}(X) = \langle (X - \langle X \rangle)^2 \rangle$$

→ Average square of deviation from mean → RMS(X) =  $\sqrt{Var(X)} = \sigma_x$  standard deviation

Can be approximated by **sample variance**:

Covariance of X and Y:  

$$\hat{\sigma}^2 = \frac{1}{n-1} \sum_i (x_i - \bar{x})^2$$

 $\rightarrow$  Large if variations of X and Y are "synchronized"

$$\operatorname{Cov}(X,Y) = \langle (X - \langle X \rangle) (Y - \langle Y \rangle) \rangle$$





<u><</u> 1

$$\rho = \frac{\operatorname{Cov}(X, Y)}{\sqrt{\operatorname{Var}(X)\operatorname{Var}(Y)}} \quad -1 \le \rho$$

# **PDF Properties: (Co)variance**

Variance of X:

С

$$\operatorname{Var}(X) = \langle (X - \langle X \rangle)^2 \rangle$$

→ Average square of deviation from mean → RMS(X) =  $\sqrt{Var(X)} = \sigma_x$  standard deviation

Can be approximated by **sample variance**:

Covariance of X and Y:  

$$\hat{\sigma}^2 = \frac{1}{n-1} \sum_i (x_i - \bar{x})^2$$

 $\rightarrow$  Large if variations of X and Y are "synchronized"

$$\operatorname{Cov}(X,Y) = \langle (X - \langle X \rangle) (Y - \langle Y \rangle) \rangle$$





77

orrelation coefficient 
$$\rho = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X) \text{Var}(Y)}}$$
  $-1 \le \rho \le 1$ 

# **PDF Properties: (Co)variance**

Variance of X:

С

$$\operatorname{Var}(X) = \langle (X - \langle X \rangle)^2 \rangle$$

→ Average square of deviation from mean → RMS(X) =  $\sqrt{Var(X)} = \sigma_x$  standard deviation

Can be approximated by **sample variance**:

Covariance of X and Y:  

$$\hat{\sigma}^2 = \frac{1}{n-1} \sum_i (x_i - \bar{x})^2$$

 $\rightarrow$  Large if variations of X and Y are "synchronized"

$$\operatorname{Cov}(X,Y) = \langle (X - \langle X \rangle) (Y - \langle Y \rangle) \rangle$$





orrelation coefficient 
$$\rho = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X) \text{Var}(Y)}} \quad -1 \le \rho \le 1$$

# "Linear" vs. "non-linear" correlations

For non-Gaussian cases, the **Correlation coefficient**  $\rho$  is not the whole story:



Source: Wikipedia

In particular, variables can still be correlated even when  $\rho=0$ : "*Non-linear"* correlations.

### **Gaussian PDF**

Gaussian distribution:

$$\rightarrow Mea(\mathbf{n}^{X}; \mathbf{x}_{0}^{X}, \sigma) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x - X_{0})^{2}}{2\sigma^{2}}}$$
  

$$\rightarrow \text{Variance} : \sigma^{2} (\Rightarrow \text{RMS} = \sigma)$$



Generalize to N dimensions:

- $\rightarrow$  Mean : X<sub>o</sub>
- → Covariance matrix :

$$G(\mathbf{x}; \mathbf{X}_{0}, \mathbf{C}) = \frac{1}{[(2\pi)^{N} |\mathbf{C}|]^{1/2}} e^{-\frac{1}{2}(x - X_{0})^{T} \mathbf{C}^{-1}(x - X_{0})} \frac{1}{[(2\pi)^{N} |\mathbf{C}|]^{1/2}} e^{-\frac{1}{2}(x - X_{0})^{T} \mathbf{C}^{-1}(x - X_{0})} \frac{1}{\sigma_{1}^{2} - \sigma_{2}^{2}} \frac{1}{\sigma_{1}^{$$

$$C = \begin{bmatrix} \operatorname{Var}(X_1) & \operatorname{Cov}(X_1, X_2) \\ \operatorname{Cov}(X_2, X_1) & \operatorname{Var}(X_2) \end{bmatrix}$$
$$= \begin{bmatrix} \sigma_1^2 & \rho \sigma_1 \sigma_2 \\ \rho \sigma_1 \sigma_2 & \sigma_2^2 \end{bmatrix}$$

### **Gaussian Quantiles**

Consider  $z = \left(\frac{x - x_0}{\sigma}\right)$  "pull" of x

 $G(x;x_0,\sigma)$  depends only on  $z \sim G(z; 0,1)$ 

Z $P(|x - x_0| > Z\sigma)$ 10.31720.04530.00343 x 10^{-5}56 x 10^{-7}

P(|x - x<sub>0</sub>| < 1♂) = 68.3 % Probability  $P(|x - x_0| > Z\sigma)$  to be away 0.4 from the mean: 0.35 0.3 Gaussian Cumulative Distribution Function (CDF) : 0.25 0.2 0.15  $\Phi(z) = \int_{-\infty}^{z} G(u; 0, 1) \, du$ 0.1 0.05 -3 -2 -1 0 2 3 1 4 8**ð** 

### **Gaussian Quantiles**

Consider  $z = \left(\frac{x - x_0}{\sigma}\right)$  "pull" of x

 $G(x;x_0,\sigma)$  depends only on  $z \sim G(z; 0,1)$ 

 $Z \qquad P(|x-x_0| > Z\sigma)$ 

P(|x - x<sub>0</sub>| < 2♂) = 95.4 % Probability  $P(|x - x_0| > Z\sigma)$  to be away 0.4 from the mean: 0.35 0.3 Gaussian Cumulative Distribution Function (CDF) : 0.25 0.2 0.15  $\Phi(z) = \int_{-\infty}^{z} G(u; 0, 1) \, du$ 0.1 0.05 -3 -2 -1 0 2 3 486 1

### **Gaussian Quantiles**

Consider  $z = \left(\frac{x - x_0}{\sigma}\right)$  "pull" of x

 $G(x;x_0,\sigma)$  depends only on  $z \sim G(z; 0,1)$ 

 $Z \qquad P(|x-x_0| > Z\sigma)$ 

P(|x - x<sub>0</sub>| < 3♂) = 99.7 % Probability  $P(|x - x_0| > Z\sigma)$  to be away 0.4 from the mean: 0.35 0.3 Gaussian Cumulative Distribution Function (CDF) : 0.25 0.2 0.15  $\Phi(z) = \int_{-\infty}^{z} G(u; 0, 1) \, du$ 0.1 0.05 -3 -2 2 3 -1 0 1 480

### **Central Limit Theorem**

(\*) Assuming  $\sigma_x < \infty$ and other regularity conditions

For an observable X with any distribution, one has(\*)

### What this means:

- hat this means:  $\overline{X} = \frac{1}{n} \sum_{i=1}^{n} x_i \stackrel{n \to \infty}{\sim} G(\langle X \rangle, \frac{\sigma_X}{\sqrt{n}})$ The average of many measurements is always Gaussian, whatever the distribution for a single measurement
- The mean of the Gaussian is the average of the single measurements
- The **RMS** of the Gaussian decreases as  $\sqrt{n}$  : smaller fluctuations when averaging over many measurements

Another version:

Mean scales like n, but RMS only like  $\sqrt{n}$ 

$$\sum_{i=1}^{n} x_{i} \stackrel{n \to \infty}{\sim} G(n \langle X \rangle, \sqrt{n} \sigma_{X})$$

Draw events from a parabolic distribution (e.g. decay  $\cos \theta^*$ )



**Distribution becomes Gaussian**, although very non-Gaussian originally **Distribution becomes narrower** as expected (as  $1/\sqrt{n}$ )

82

Draw events from a parabolic distribution (e.g. decay  $\cos \theta^*$ )



**Distribution becomes Gaussian**, although very non-Gaussian originally **Distribution becomes narrower** as expected (as  $1/\sqrt{n}$ )

Draw events from a parabolic distribution (e.g. decay  $\cos \theta^*$ )



**Distribution becomes Gaussian**, although very non-Gaussian originally **Distribution becomes narrower** as expected (as  $1/\sqrt{n}$ )

Draw events from a parabolic distribution (e.g. decay  $\cos \theta^*$ )



**Distribution becomes Gaussian**, although very non-Gaussian originally **Distribution becomes narrower** as expected (as  $1/\sqrt{n}$ )

Draw events from a parabolic distribution (e.g. decay  $\cos \theta^*$ )



**Distribution becomes Gaussian**, although very non-Gaussian originally **Distribution becomes narrower** as expected (as  $1/\sqrt{n}$ )

82

Draw events from a parabolic distribution (e.g. decay  $\cos \theta^*$ )



**Distribution becomes Gaussian**, although very non-Gaussian originally **Distribution becomes narrower** as expected (as  $1/\sqrt{n}$ )

82

Draw events from a parabolic distribution (e.g. decay  $\cos \theta^*$ )



**Distribution becomes Gaussian**, although very non-Gaussian originally **Distribution becomes narrower** as expected (as  $1/\sqrt{n}$ )

Draw events from a parabolic distribution (e.g. decay  $\cos \theta^*$ )



**Distribution becomes Gaussian**, although very non-Gaussian originally **Distribution becomes narrower** as expected (as  $1/\sqrt{n}$ )

# **Chi-squared**

Multiple Independent Gaussian variables x<sub>i</sub>: Define

$$\chi^2 = \sum_{i=1}^n \left( \frac{x_i - x_i^0}{\sigma_i} \right)^2$$

Measures global distance from reference point  $(x_1^0 \dots x_n^0)$ 

Distribution depends on n :

Rule of thumb:

 $\chi^2/n$  should be  $\preceq 1$ 



# **Chi-squared**

Multiple Independent Gaussian variables x<sub>i</sub>: Define

$$\chi^2 = \sum_{i=1}^n \left( \frac{x_i - x_i^0}{\sigma_i} \right)^2$$

Measures global distance from reference point  $(x_1^0 \dots x_n^0)$ 

Distribution depends on n :

Rule of thumb:

 $\chi^2/n$  should be  $\preceq 1$ 



83

### **Histogram Chi-squared**

Histogram  $\chi^2$  with respect to a reference shape:

- Assume an independent Gaussian distribution in each bin
- Degrees of freedom = (number of bins) (number of fit parameters)



# **BLUE histogram vs. flat reference** $\chi^2 = 12.9$ , $p(\chi^2=12.9, n=10) = 23\%$

### **Histogram Chi-squared**

Histogram  $\chi 2$  with respect to a reference shape:

- Assume an independent Gaussian distribution in each bin
- Degrees of freedom = (number of bins) (number of fit parameters)



BLUE histogram vs. flat reference  $\chi^2 = 12.9$ ,  $p(\chi^2=12.9, n=10) = 23\%$ RED histogram vs. flat reference

 $\chi^2 = 38.8$ ,  $p(\chi^2 = 38.8$ , n = 10) = 0.003%
## **Histogram Chi-squared**

Histogram  $\chi$ 2 with respect to a reference shape:

- Assume an independent Gaussian distribution in each bin
- Degrees of freedom = (number of bins) (number of fit parameters)



**BLUE histogram vs. flat reference**  $\chi^2 = 12.9$ ,  $p(\chi^2=12.9, n=10) = 23\%$ **RED histogram vs. flat reference**  $\chi^2 = 38.8$ ,  $p(\chi^2=38.8, n=10) = 0.003\%$ **RED histogram vs. correct reference**  $\chi^2 = 9.5$ ,  $p(\chi^2=9.5, n=10) = 49\%$ 

## **Error Bars**

Strictly speaking, the uncertainty is given by the model :

- $\rightarrow$  **Bin central value** ~ mean of the bin PDF
- $\rightarrow$  **Bin uncertainty** ~ RMS of the bin PDF

The data is just what it is, a simple observed point.

- $\Rightarrow$  One should in principle show the error bar on the prediction.
- $\rightarrow$  In practice, the usual convention is to have error bars on the data points.



## **Error Bars**

Strictly speaking, the uncertainty is given by the model :

- $\rightarrow$  **Bin central value** ~ mean of the bin PDF
- → **Bin uncertainty** ~ RMS of the bin PDF

The data is just what it is, a simple observed point.

- $\Rightarrow$  One should in principle **show the error bar on the prediction**.
- $\rightarrow$  In practice, the usual convention is to have error bars on the data points.

