

# Probabilidad y estadística para la física experimental

José Ocariz

*LA-CoNGA physics*

*ocariz@in2p3.fr*

11 de mayo de 2021



Latin American alliance for  
Capacity building in Advanced physics  
**LA-CoNGA physics**



Cofinanciado por el  
programa Erasmus+  
de la Unión Europea

UAN  
UNIVERSIDAD NACIONAL



UNIVERSIDAD  
NACIONAL DE  
INGENIERIA

UNMSM



UNIVERSIDAD SIMÓN BOLÍVAR

Université  
de Paris

TECHNISCHE  
UNIVERSITÄT  
DRESDEN

UNIVERSITÄT  
DUISBURG  
ESSEN



CLARA

DBACCESS

frontier x  
ANALYTICS





- ▶ El libro clásico de referencia (912 páginas) :
  - ▶ Stuart, K. Ord, S. Arnold, Kendall's Advanced theory of statistics Volume 2A : Classical Inference and the Linear Model, John Wiley & Sons, 2009
- ▶ Libros de estadísticas, escritos por físicos de partículas :
  - ▶ L. Lyons, Statistics for Nuclear and Particle Physics, Cambridge, 1986
  - ▶ G. Cowan, Statistical Data Analysis, Clarendon Press, Oxford 1998  
(ver también [http://www.p.rhu1.ac.uk/~cowan/stat\\_course.htm](http://www.p.rhu1.ac.uk/~cowan/stat_course.htm))
  - ▶ R.J. Barlow, A Guide to the Use of Statistical Methods in the Physical Sciences John Wiley & Sons, 1989
  - ▶ F. James, Statistical Methods in Experimental Physics, World Scientific, 2006
- ▶ El PDG también es una fuente conveniente para referencias rápidas :
  - ▶ 2020 Review of Particle Physics, P.A. Zyla et al. (Particle Data Group), Prog. Theor. Exp. Phys. 2020, 083C01 (2020) "Mathematical Tools" section  
(ver también <https://pdg.lbl.gov/>)
- ▶ Las grandes colaboraciones internacionales tienen foros y grupos de trabajo, con muchos enlaces y referencias.



- ▶ La probabilidad matemática es un concepto axiomático abstracto, desarrollado por Kolmogorov (1933) y otros
- ▶ La teoría de la probabilidad es el marco conceptual para el estudio de los procesos aleatorios
- ▶ Un proceso es llamado aleatorio si satisface dos condiciones :
  - ▶ su realización (un “evento”) no puede ser predicha con total certeza ;
  - ▶ si el proceso se repite bajo las mismas condiciones, cada nueva realización puede ser diferente
- ▶ Es usual clasificar las fuentes de incertidumbre según su origen :
  - ▶ *reducibles* : errores en la la medición, p.e. limitaciones prácticas que en principio pueden ser mejoradas (mejores instrumentos, mejor control de las condiciones experimentales) ;
  - ▶ *cuasi-irreducibles* : errores aleatorios en la medición, como efectos térmicos o de turbulencia ;
  - ▶ *fundamentales* : cuando el proceso físico es intrínsecamente incierto (mecánica cuántica).

En física subatómica experimental, figuran los tres tipos de incertidumbre. Notar en particular :

- ▶ los eventos resultantes de colisiones de partículas son independientes, y son un ejemplo perfecto de procesos aleatorios de origen cuántico
- ▶ las partículas inestables obedecen probabilidades de desintegración descritas por la mecánica cuántica

Ejercicio : dar ejemplos de procesos físicos para cada una de las fuentes de incertidumbre



Sea  $\Omega$  el universo total de posibles realizaciones de un proceso aleatorio, y sean  $X, Y \dots$  elementos de  $\Omega$ . Una función de probabilidad  $\mathcal{P}$  se define como un mapa en los números reales :

$$\begin{aligned}\mathcal{P} : \{\Omega\} &\rightarrow [0 : 1] , \\ X &\rightarrow \mathcal{P}(X) .\end{aligned}$$

Ese mapeo debe satisfacer los siguientes axiomas :

$$\begin{aligned}\mathcal{P}(\Omega) &= 1 , \\ \text{si } X \cap Y &= \emptyset , \text{ entonces } \mathcal{P}(X \cup Y) = \mathcal{P}(X) + \mathcal{P}(Y) ,\end{aligned}$$

de los cuales se pueden derivar varias propiedades útiles, p.e. (donde  $\bar{X}$  es el complemento de  $X$ )

$$\begin{aligned}\mathcal{P}(\bar{X}) &= 1 - \mathcal{P}(X) , \\ \mathcal{P}(X \cup \bar{X}) &= 1 , \\ \mathcal{P}(\emptyset) &= 1 - \mathcal{P}(\Omega) = 0 , \\ \mathcal{P}(X \cup Y) &= \mathcal{P}(X) + \mathcal{P}(Y) - \mathcal{P}(X \cap Y) ,\end{aligned}$$



## Probabilidad condicional, teorema de Bayes

La probabilidad condicional  $\mathcal{P}(X | Y)$  se define como la probabilidad de  $X$ , dado  $Y$

- ▶ equivale a restringir el universo  $\Omega$  a la muestra  $Y$ .

El ejemplo más sencillo de probabilidad condicional es para realizaciones independientes :

- ▶ dos elementos  $X$  e  $Y$  son independientes (sus realizaciones no estén relacionadas en ninguna manera) si

$$\mathcal{P}(X \cap Y) = \mathcal{P}(X)\mathcal{P}(Y).$$

- ▶ por lo tanto, si  $X$  e  $Y$  son independientes, se satisface la condición

$$\mathcal{P}(X | Y) = \mathcal{P}(X)$$

El teorema de Bayes cubre el caso general : en vista de la relación  $\mathcal{P}(X \cap Y) = \mathcal{P}(Y \cap X)$ , se obtiene que

$$\mathcal{P}(X | Y) = \frac{\mathcal{P}(Y | X)\mathcal{P}(X)}{\mathcal{P}(Y)} .$$

Un corolario útil del teorema de Bayes : si  $\Omega$  puede dividirse en un número de submuestras disjuntas  $X_i$  (una "partición"), entonces

$$\mathcal{P}(X | Y) = \frac{\mathcal{P}(Y | X)\mathcal{P}(X)}{\sum_i \mathcal{P}(Y | X_i)\mathcal{P}(X_i)} .$$



## Variables aleatorias, funciones de densidad de probabilidad (I)

El escenario más relevante para nosotros es cuando la realización de un proceso aleatorio se presenta en forma numérica (p.e. corresponde a una medición) : a cada elemento  $X$  corresponde una variable  $x$  (real o entera). Para  $x$  continuo, su función de densidad de probabilidad (PDF)  $P(x)$  se define como :

$$\mathcal{P}(X \text{ en } [x, x + dx]) = P(x)dx ,$$

donde  $P(x)$  es definida-positiva para todo de  $x$ , y satisface la condición de normalización

$$\int_{-\infty}^{+\infty} dx' P(x') = 1 .$$

Para  $x_i$  discreto, la definición es similar :

$$\mathcal{P}(X \text{ en } x_i) = p_i ,$$

con  $\sum_j p_j = 1$  y  $p_k \geq 0 \forall k$  .

Probabilidades finitas se obtienen por integración sobre un rango no-infinitesimal,

$$\mathcal{P}(a \leq X \leq b) = \int_a^b dx' P(x') .$$



En ocasiones es conveniente referirse a la función de densidad acumulativa (CDF) :

$$C(x) = \int_{-\infty}^x dx' P(x') ,$$

de modo que las probabilidades finitas corresponden a evaluar la CDF en los bordes del rango de interés :

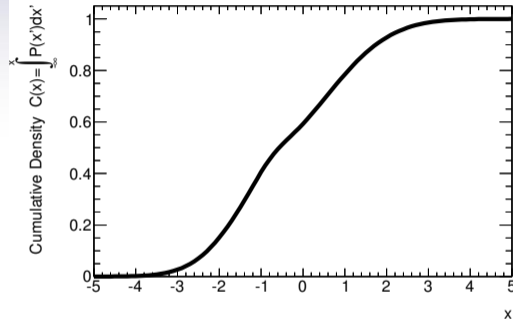
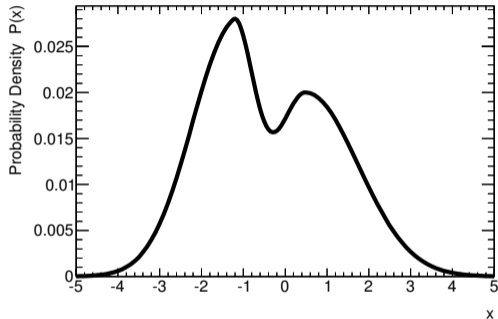
$$\mathcal{P}(a \leq X \leq b) = \int_a^b dx' P(x') = C(b) - C(a) .$$

Una PDF no puede ser completamente arbitraria :

- ▶ debe satisfacer la condición de normalización previamente indicada
- ▶ debe ser definida positiva
- ▶ debe ser de soporte acotado, con valores despreciables fuera de una región finita

Fuera de esas condiciones, una PDF puede ser arbitraria, p.e. exhibir uno o varios máximos locales, tener discontinuidades...

En contraste, la CDF es una función monótonicamente creciente de  $x$ .  
(ver ejemplo en la lámina siguiente)



Un ejemplo arbitrario de PDF con un máximo global y un segundo máximo local, y su CDF correspondiente. Fuera del intervalo en el gráfico, el valor de la PDF es totalmente despreciable.





## PDFs multidimensionales (I)

Para un evento descrito por un conjunto  $n$ -dimensional de elementos  $\vec{X} = \{X_1, X_2, \dots, X_n\}$  y su correspondiente conjunto de variables aleatorias  $\vec{x} = \{x_1, x_2, \dots, x_n\}$ , tenemos su PDF multidimensional :

$$P(\vec{x})d\vec{x} = P(x_1, x_2, \dots, x_n)dx_1dx_2 \dots dx_n .$$

PDFs de menor dimensionalidad pueden derivarse por integración de ciertas variables. Por ejemplo, para una variable específica  $x = x_j$  su densidad de probabilidad marginal unidimensional  $P_X(x)$  es :

$$P_X(x)dx = dx \int_{-\infty}^{+\infty} dx_1 \dots \int_{-\infty}^{+\infty} dx_{j-1} \int_{-\infty}^{+\infty} dx_{j+1} \dots \int_{-\infty}^{+\infty} dx_{n-1} .$$

Caso bidimensional, con elementos  $X, Y$  y variables aleatorias  $\vec{X} = \{x, y\}$ . La probabilidad finita en un rango bidimensional rectangular es

$$P(a \leq X \leq b ; c \leq Y \leq d) = \int_a^b dx \int_c^d dy P(x, y) .$$

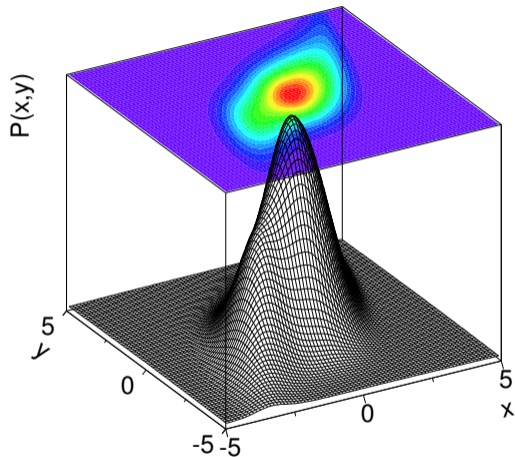
Para un valor fijo de  $Y$ , la función de densidad condicional de  $X$  es

$$P(x | y) = \frac{P(x, y)}{\int dy P(x, y)} = \frac{P(x, y)}{P_Y(y)} .$$

De nuevo, la relación  $P(x, y) = P_X(x) \cdot P_Y(y)$  solamente es válida para  $X, Y$  independientes.



Ejemplo de una función de densidad  
bidimensional con variables  
no-independientes,  
 $P(x, y) \neq P_X(x) \cdot P_Y(y)$ .





Modelo : descripción de un proceso aleatorio con funciones de densidad

Modelo paramétrico : sus PDFs pueden describirse completamente usando un número finito de parámetros

- ▶ este requisito no es obligatorio; las PDFs pueden también ser no-paramétricas (equivalente a suponer que se necesita un número infinito de parámetros), o pueden ser mixtas

Una implementación sencilla de una PDF paramétrica es cuando sus parámetros son argumentos analíticos de la función de densidad ; la notación

$$P(x, y, \dots ; \theta_1, \theta_2, \dots)$$

indica la dependencia funcional o *forma* de la PDF en términos de variables  $x_1, y_2, \dots$  y parámetros  $\theta_1, \theta_2, \dots$

Consideremos una variable aleatoria  $X$  con PDF  $P(x)$ . Para una función genérica  $f(x)$ , su *valor de expectación*  $E[f]$  es su promedio ponderado sobre el rango cubierto por  $x$  :

$$E[f] = \int dx P(x) f(x) .$$

Como describiremos más adelante, los parámetros de una PDF pueden ser estimados a partir de ciertos valores de expectación.



## Valores de expectación (II)

Por ser de uso frecuente, algunos valores de expectación tienen nombre propio.

Para PDFs unidimensionales, la *media* y la *varianza* se definen así :

$$\text{Media} \quad : \quad \mu = \quad E[x] = \int dx P(x)x ,$$

$$\text{Varianza} \quad : \quad \sigma^2 = \quad V[x] = E[x^2] - \mu^2 = E[(x - \mu)^2] ;$$

y la *desviación estándar*  $\sigma$  es la raíz cuadrada de la varianza.

Para PDFs multidimensionales, la matriz de *covarianza*  $C_{ij} = C(x_i, x_j)$  y la matriz adimensional de correlación lineal  $\rho_{ij}$  se definen así :

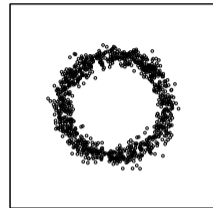
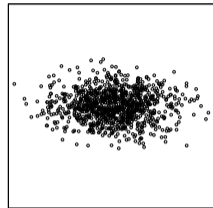
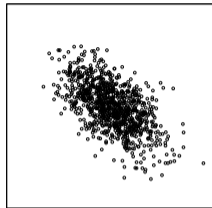
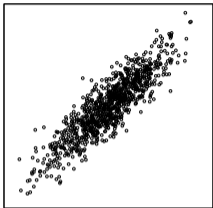
$$C_{ij} = E[x_i x_j] - \mu_i \mu_j = E[(x_i - \mu_i)(x_j - \mu_j)] , \quad \rho_{ij} = \frac{C_{ij}}{\sigma_i \sigma_j} .$$

Los coeficientes de correlación lineal tienen valores en el rango  $-1 \leq \rho_{ij} \leq 1$ , e indican la tendencia dominante de densidad en el patrón  $(x_i; x_j)$  pattern: se habla de correlaciones positivas y negativas (o anti-correlaciones).

Para variables aleatorias  $X_i, X_j$  independientes, es decir con  $P(x_i, x_j) = P_{X_i}(x_i)P_{X_j}(x_j)$ , se tiene

$$E[x_i x_j] = \int \int dx_i dx_j P(x_i, x_j) x_i x_j = \mu_i \mu_j , \quad \longrightarrow \quad \rho_{ij} = 0 .$$

(pero la converso no es necesariamente cierta, p.e. ejemplo en la lámina siguiente)



De izquierda a derecha :

- ▶  $\rho = +0,9$  ,
- ▶  $\rho = -0,5$  ;
- ▶  $\rho = 0$ , para variables independientes ,
- ▶ variables fuertemente correlacionadas con un patrón no lineal de correlación que “conspira” para arrojar una correlación lineal nula,  $\rho = 0$ .



- ▶ En la práctica, la verdadera dependencia funcional de una PDF es a menudo desconocida
- ▶ La información sobre su forma solamente puede extraerse a partir de una muestra de talla finita (digamos que contiene  $N$  eventos), es decir suponemos que la muestra disponible es una realización aleatoria a partir de una PDF desconocida.
- ▶ Si consideramos que esa PDF subyacente es de tipo paramétrico, la *caracterización de su forma* es un procedimiento para estimar los valores numéricos de sus parámetros, partiendo de una hipótesis “razonable” sobre la dependencia funcional sobre sus variables.
- ▶ Ahora, solamente un número finito de valores de expectación independientes pueden extraerse de una muestra de talla finita.
- ▶ No existe una receta única para la elección de los parámetros a ser estimados, con lo que el proceso es intrínsecamente incompleto.
- ▶ Se puede sin embargo confirmar en la práctica que el proceso de caracterización de forma es bastante poderoso, si los parámetros seleccionados proveen información útil y complementaria.

Ilustramos estas consideraciones con un ejemplo unidimensional para una variable aleatoria única  $x$ .



Ilustramos las consideraciones anteriores con un ejemplo unidimensional para una variable aleatoria única  $x$ . Consideremos el promedio empírico  $\bar{x}$ , definido como

$$\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i .$$

Mostraremos más adelante que  $\bar{x}$  es un buen *estimator* de la media  $\mu$  de la PDF subyacente  $P(x)$ . De modo análogo, la media cuadrática RMS (del inglés “root-mean-squared”), definida como

$$\text{RMS}^2 = \overline{x^2} - (\bar{x})^2 ,$$

es un estimador razonable de la varianza  $\sigma^2$  (veremos más adelante una mejor definición).

En términos intuitivos, el promedio y el RMS reflejan información útil y complementaria sobre la “localización” y la “dispersión” de la región de  $x$  mayor densidad de eventos, y esta región debe corresponder de manera aproximada a los intervalos de  $x$  donde la PDF tiene valores más grandes.

Obviamente, en un caso general esos dos parámetros son insuficientes para caracterizar una PDF más genérica, que requiere un procedimiento más sistemático.



## Caracterización de la forma de una PDF (III)

Para un procedimiento más sistemático, transformamos la  $x$ -dependencia de la PDF  $P(x)$  en una  $k$ -dependencia de la *función característica*  $C[k]$ , definida como

$$C[k] = E \left[ e^{ik \frac{x-\mu}{\sigma}} \right] = \sum_j \frac{(ik)^j}{j!} \mu_j .$$

Como se puede notar, la función característica es la transformación de Fourier de la PDF. Los coeficientes  $\mu_j$  de la expansión *are* se llaman *momentos reducidos*; por construcción, los primeros momentos son  $\mu_1 = 0$  y  $\mu_2 = 1$ ; en términos de la variable reescalada  $x' = (x - \mu)/\sigma$ , la PDF fue desplazada para tener promedio nulo, y escalada para tener varianza unitaria.

En principio, mientras mayor el número de momentos  $\mu_j$  sean estimados, más detallada será la caracterización de la forma de la PDF (pero una muestra finita solamente permite medir un número finito de momentos).

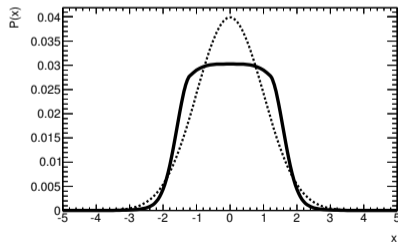
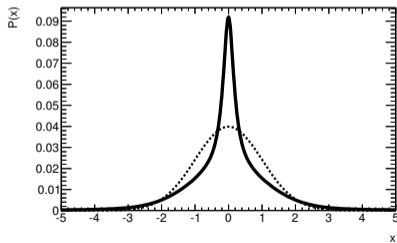
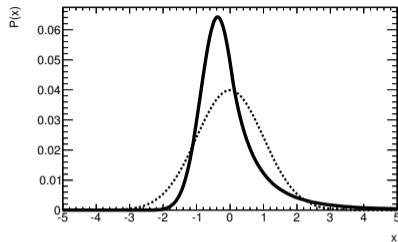
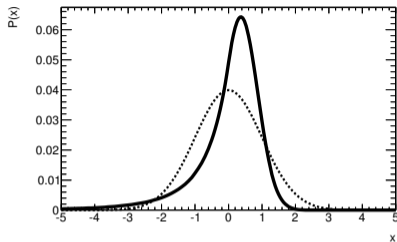
Los momentos 3 y 4 tienen nombres específicos, y sus valores son sencillos de interpretar en términos de la forma:

- ▶ el tercer momento es llamado oblicuidad (o *skewness* en inglés)
  - ▶ una distribución simétrica tiene skewness nula,
  - ▶ un valor negativo (positivo) indica una “anchura” mayor a la izquierda (derecha) de su media.
- ▶ el cuarto momento es llamado *kurtosis*
  - ▶ cantidad definida positiva, relacionada con cuán “picante” es la distribución
  - ▶ un valor grande indica un pico estrecho y “colas” de largo alcance: es una distribución leptokúrtica
  - ▶ un valor pequeño indica un pico central ancho y colas poco prominentes: es una distribución platykúrtica





# Ejemplos: skewness y kurtosis





La caracterización de la forma de una PDF a través de una estimación secuencial de parámetros de forma nos permitió introducir de manera cualitativa al concepto de estimación de parámetros (también llamado en inglés "point estimation"). Una receta más general sería la siguiente :

Consideremos una PDF  $n$ -dimensional,  $k$ -paramétrica,

$$P( x_1, x_2, \dots, x_n ; \theta_1, \theta_2, \dots, \theta_k ) ,$$

para la cual queremos estimar los valores  $\theta_1, \dots, \theta_k$  a partir de una muestra de talla finita, utilizando un conjunto de estimadores  $\hat{\theta}_1, \dots, \hat{\theta}_k$ . Esos estimadores son también variables aleatorias, con sus propias medias y varianzas: sus valores diferirían al ser estimados sobre otras muestras. Esos estimadores deben satisfacer dos propiedades clave:

- ▶ *Consistencia* : asegura que, en el límite de una muestra de talla infinita, el estimador converge al verdadero valor del parámetro;
- ▶ *no sesgado* : la ausencia de sesgo asegura que el valore de expectación del estimador es el verdadero valor del parámetro, para toda talla de la muestra.

Un estimador sesgado pero consistente (también llamado asintóticamente no-sesgado) es tal que el sesgo disminuye al aumentar la talla de la muestra.



Otros criterios son útiles para caracterizar la calidad de los estimadores; por ejemplo

- ▶ eficiencia: un estimador de pequeña varianza es más eficiente que uno de mayor varianza ;
- ▶ robusteza: este criterio describe la “sensibilidad” del estimador a incertidumbres en la forma de la PDF. Por ejemplo, el promedio es robusto contra incertidumbres sobre los momentos de orden par, pero es menos robusto contra incertidumbres en los momentos de orden impar.

Nota : estos criterios son en ocasiones mutuamente contradictorios; por razones prácticas, puede ser preferible tener un estimador eficiente pero sesgado, a uno no sesgado pero de pobre convergencia.



## Estimación de parámetros (II)

El promedio empírico  $\bar{x}$  es un estimador convergente, no sesgado de la media  $\mu$  de la PDF subyacente:  $\hat{\mu} = \bar{x}$ . Esto se demuestra fácilmente, evaluando el valor de expectación y la varianza de  $\bar{x}$ :

$$E[\bar{x}] = \frac{1}{N} \sum_{i=1}^N E[x] = \mu,$$
$$V[\bar{x}] = E[(\bar{x} - \mu)^2] = \frac{\sigma^2}{N}.$$

Al contrario, el RMS empírico de una muestra es un estimador sesgado (aunque asintóticamente no-sesgado) de la varianza  $\sigma^2$ . Esto se demuestra fácilmente también, reescribiendo su cuadrado en términos de la media:

$$\text{RMS}^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2 - (\bar{x} - \mu)^2,$$

de manera que su valor de expectación es

$$E[\text{RMS}^2] = \sigma^2 - V[\bar{x}] = \frac{N-1}{N} \sigma^2,$$

que si bien converge a la verdadera varianza  $\sigma^2$  en el límite  $N \rightarrow \infty$ , subestima sistemáticamente su valor para muestras de talla finita.



Si bien converge a la verdadera varianza  $\sigma^2$  en el límite  $N \rightarrow \infty$ , el RMS subestima sistemáticamente su valor para muestras de talla finita. Pero es inmediato definir un estimador modificado

$$\frac{N}{N-1} \text{RMS}^2 = \frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})^2 ,$$

que es, para muestras de talla finita, un estimador no sesgado de la varianza.

En resumen: para una PDF desconocida, tenemos estimadores consistentes y no sesgados de su media  $\mu$  y su varianza  $\sigma^2$ , que pueden ser extraídos de muestras de talla finita:

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N x_i ,$$
$$\hat{\sigma}^2 = \frac{1}{N-1} \sum_{i=1}^N (x_i - \hat{\mu})^2 .$$

Nota: los factores  $1/N$  y  $1/(N-1)$  para  $\hat{\mu}$  y  $\hat{\sigma}^2$  se entienden intuitivamente: el promedio empírico puede medirse incluso en la muestra más pequeña posible de un solo evento, mientras que al menos dos eventos son necesarios para estimar la dispersión empírica de una muestra.



El ejemplo clásico anterior trataba de una única variable aleatoria.

En presencia de múltiples variables aleatorias  $\vec{x} = \{x_1, \dots, x_n\}$ , la generalización del resultado anterior lleva a definir la covarianza empírica, cuyos elementos  $\hat{C}_{ab}$  son estimados en una muestra de  $N$  eventos de la manera siguiente:

$$\hat{C}_{ab} = \frac{1}{N-1} \sum_{i=1}^N (x_{a,i} - \hat{\mu}_a) (x_{b,i} - \hat{\mu}_b) .$$

(los índices  $a, b$  recorren la lista de variables aleatorias,  $1 \leq a, b \leq n$ ).

Suponiendo que la verdadera covarianza es conocida, la varianza de una función arbitraria  $f(\vec{x})$  de las variables aleatorias se evalúa a partir de la expansión de Taylor alrededor de las medias  $\hat{\vec{\mu}}$  según

$$f(\vec{x}) = f(\hat{\vec{\mu}}) + \sum_{a=1}^n \left. \frac{df}{dx_a} \right|_{\vec{x}=\hat{\vec{\mu}}} (x_a - \hat{\mu}_a) ,$$

en otras palabras,  $E[f(\vec{x})] \simeq f(\hat{\vec{\mu}})$ .



De manera similar,

$$E [f^2(\vec{x})] \simeq f^2(\hat{\vec{\mu}}) + \sum_{a,b=1}^n \left. \frac{df}{dx_a} \frac{df}{dx_b} \right|_{\vec{x}=\hat{\vec{\mu}}} \hat{C}_{ab} ,$$

con lo que la varianza de  $f$  se estima como

$$\hat{\sigma}_f^2 \simeq \sum_{a,b=1}^n \left. \frac{df}{dx_a} \frac{df}{dx_b} \right|_{\vec{x}=\hat{\vec{\mu}}} \hat{C}_{ab} .$$

Esta expresión, llamada *fórmula de propagación de errores*, estima la varianza de una función genérica  $f(\vec{x})$  a partir de los estimadores de sus medias y covarianzas.

Ejemplos particulares de propagación de errores:

- ▶ cuando todas las variables aleatorias  $\{x_a\}$  son no-correlacionadas, la matriz de covarianza es diagonal,  $C_{ab} = \sigma_a^2 \delta_{ab}$  y la covarianza de  $f(\vec{x})$  se reduce a

$$\hat{\sigma}_f^2 \simeq \sum_{a=1}^n \left( \left. \frac{df}{dx_a} \right|_{\vec{x}=\hat{\vec{\mu}}} \right)^2 \hat{\sigma}_a^2 .$$



- ▶ para la suma de dos variables aleatorias  $S = x_1 + x_2$ , la varianza es

$$\sigma_S^2 = \sigma_1^2 + \sigma_2^2 + 2C_{12} = \sigma_1^2 + \sigma_2^2 + 2\sigma_1\sigma_2\rho_{12},$$

- ▶ la generalización a más de dos variables es:

$$\sigma_S^2 = \sum_{a,b} \sigma_a \sigma_b \rho_{ab} .$$

En ausencia de correlaciones, se dice que los errores absolutos se suman “en cuadratura”:

$$\sigma_S = \sigma_1 \oplus \sigma_2 = \sqrt{\sigma_1^2 + \sigma_2^2}$$





- ▶ para el producto de dos variables aleatorias  $P = x_1 x_2$ , la varianza es

$$\left(\frac{\sigma_P}{P}\right)^2 = \left(\frac{\sigma_1}{x_1}\right)^2 + \left(\frac{\sigma_2}{x_2}\right)^2 + 2\frac{\sigma_1}{x_1}\frac{\sigma_2}{x_2}\rho_{12} ,$$

- ▶ la generalización a más de dos variables:

$$\left(\frac{\sigma_P}{P}\right)^2 = \sum_{a,b} \frac{\sigma_a}{x_a} \frac{\sigma_b}{x_b} \rho_{ab} .$$

En ausencia de correlaciones, se dice que los errores relativos se suman en cuadratura:

$$\sigma_P/P = \sigma_1/x_1 \oplus \sigma_2/x_2 .$$



- ▶ para una función genérica en ley de potencia,  $Z = x_1^{n_1} x_2^{n_2} \dots$ , si todas las variables son no-correlacionadas, la varianza es

$$\frac{\sigma_Z}{Z} = n_1 \frac{\sigma_1}{x_1} \oplus n_2 \frac{\sigma_2}{x_2} \oplus \dots$$



## Lista no-exhaustiva de distribuciones de uso común

Distribución	Uso(s) en altas energías	Otro(s) nombre(s)
Binomial	Tasa de decaimiento, eficiencias	Bernouilli
Poisson	Conteo de eventos	"ley de eventos raros"
Uniforme	Integración Monte-Carlo	
Exponencial	Vida media, tiempos de relajación	
Gaussiana	Resolución	Normal
Breit-Wigner	Resonancia	Cauchy, Lorentz
$\chi^2$	"Goodness-of-fit"	"bondad de ajuste"

Lista no exhaustiva; otras de uso frecuente incluyen la distribución de Student, Galton o lognormal...



## Ejemplo de una distribución discreta: Binomial (I)

Escenario con dos únicas realizaciones posibles: “éxito” y “fracaso”, y con una probabilidad fija  $p$  de “éxito”. Nos interesamos solamente en el número  $k$  de “éxitos” después de  $n$  intentos  $0 \leq k \leq n$ ; (suponemos que la secuencia de intentos es irrelevante)

El número entero  $k$  sigue la distribución binomial  $P(k; n, p)$ :

$$P_{\text{binomial}}(k; n, p) = \frac{n!}{k!(n-k)!} p^k (1-p)^{n-k},$$

donde  $k$  es la variable aleatoria, mientras que  $n$  y  $p$  son parámetros.

Ejemplo típico: el número de eventos en una sub-categoría específica de eventos (p.e. una tasa de decaimiento)

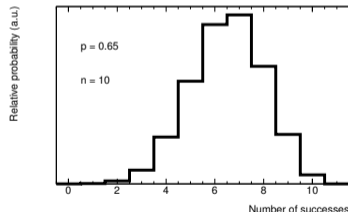
### Ejercicios :

- ▶ mostrar que la media y la varianza de una distr. binomial son

$$E[k] = \sum_{k=1}^n k P(k; n, p) = np,$$

$$V[k] = np(1-p).$$

- ▶ Sean dos binomiales con números de intentos  $n_1$  y  $n_2$  y de misma probabilidad  $p$ . Mostrar que la suma es también una binomial.





## Ejemplo de una distribución discreta: Binomial (II)

Teorema de la varianza de suma de binomiales:

La varianza de la suma de variables aleatorias que siguen distribuciones binomiales con probabilidades  $p_i$  diferentes,

$$k = \sum_i k_i ,$$

viene dada por

$$V[k] = n\bar{p}(1 - \bar{p}) - ns^2 , \text{ con } s^2 = \frac{1}{n} \sum_i (p_i - \bar{p})^2 ,$$

y es por tanto inferior o igual a la varianza de una variable binomial de probabilidad  $\bar{p}$ .

Ejercicio:

Verificar numéricamente la validez de este teorema, estimando con una simulación sencilla la varianza de la suma de dos variables aleatorias binomiales con probabilidades  $p_1$  y  $p_2$ , y comparándola con la varianza de una variable aleatoria binomial con  $p = (p_1 + p_2)/2$ . Elegir varios valores de  $p_1$  y  $p_2$  que ilustren casos extremos del teorema.



## Ejemplo de una distribución discreta: Poisson

Para la distribución binomial en el límite  $n \rightarrow \infty$ ,  $p \rightarrow 0$  (con  $\lambda = np$  finito y no-nulo) la variable aleatoria  $k$  sigue la llamada distribución de Poisson  $P(k; \lambda)$ ,

$$P_{\text{Poisson}}(k; \lambda) = \frac{\lambda^k e^{-\lambda}}{k!},$$

que tiene  $\lambda$  por único parámetro. For Poisson, the mean value and variance are the same:

$$E[k] = V[k] = \lambda.$$

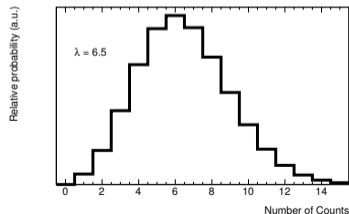
Esta distribución es también llamada “ley de los eventos raros” (debido al límite  $p \rightarrow 0$ ), describe el número de observaciones en condiciones fijas si la tasa de ocurrencia es constante. Ejemplo: el número de decaimientos en un intervalo de tiempo fijo, para una fuente de actividad constante.

### Ejercicios :

- ▶ mostrar que la media y la varianza de una Poisson son

$$E[k] = V[k] = \lambda.$$

- ▶ Sean dos distribuciones de Poisson con parámetros  $\lambda_1$  y  $\lambda_2$ . Mostrar que la distribución de la suma es también una Poisson, con parámetro  $\lambda = \lambda_1 + \lambda_2$ .





## Ejemplo de una distribución continua: Uniforme

Una variable aleatoria continua  $x$ , con densidad de probabilidad  $P(x; a, b)$  constante y no-nula únicamente dentro de un intervalo finito  $[a, b]$ :

$$P_{\text{uniform}}(x; a, b) = \begin{cases} \frac{1}{b-a} & , \quad a \leq x \leq b , \\ 0 & , \quad x < a \text{ o } x > b . \end{cases}$$

Ejercicio : mostrar que para esta distribución uniforme, la media y la varianza son

$$E[x] = \frac{a+b}{2} , \quad V[x] = \frac{(b-a)^2}{12} .$$

Uso común de la distribución uniforme : generación de distribuciones aleatorias

- ▶ “acepto-rechazo” para simular un proceso de probabilidad  $p(x)$ : utilizar dos variables aleatorias:  $x$  uniforme en el intervalo de interés, y  $x_0$ , uniforme en el intervalo  $[0, 1]$ . Aceptar el evento si  $x_0 < p(x)$ . Intuitivo y sencillo, pero usualmente poco eficaz.
- ▶ Transformada inversa: si se conoce la CDF de la distribución deseada, a partir de  $x$  uniforme en el intervalo  $[0, 1]$ , la variable aleatoria  $z = CDF(x)$  sigue la distribución correspondiente.
- ▶ Transformación de Box-Muller: a partir de dos variables aleatorias  $x_1, x_2$  uniformes en  $[0, 1]$  e independientes, las variables  $z_1, z_2$  dadas por

$$z_1 = \sqrt{-2 \ln x_1} \cos(2\pi x_2) , \quad z_2 = \sqrt{-2 \ln x_1} \sin(2\pi x_2) ,$$

son independientes y siguen cada distribuciones Gaussianas, ambas de media nula y varianza unitaria.



## Ejemplo de una distribución continua: Exponencial

Una variable aleatoria que sigue una densidad de probabilidad  $P(x; \xi)$  dada por

$$P_{\text{exponential}}(x; \xi) = \begin{cases} \frac{1}{\xi} e^{-x/\xi} & , \quad x \geq 0 , \\ 0 & , \quad x < 0 , \end{cases}$$

Y tiene por media y varianza

$$E[x] = \xi , \quad V[x] = \xi^2 .$$

Ejercicio: evaluar la media y la varianza de una distribución exponencial truncada, que es no-nula solamente en un intervalo finito  $a \leq x \leq b$ . (Nota: verificar que la PDF que usen está correctamente normalizada)

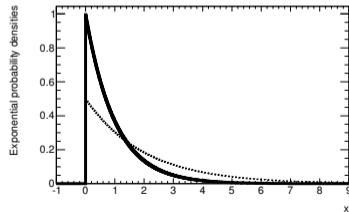
La aplicación más común de esta distribución exponencial es la descripción de fenómenos independientes que se realizan a una tasa constante, como el tiempo de vida de una partícula inestable.

En vista del carácter auto-similar de la función exponencial:

$$P(t - t_0 | t > t_0) = P(t) ,$$

se dice que esta distribución es “sin memoria”.

La figura muestra dos PDF exponenciales, con parámetros  $\xi = 1$  (línea sólida) y  $\xi = 2$  (punteada).







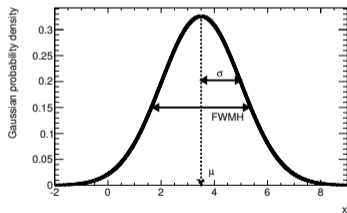
# La distribución continua de mayor ubicuidad: la Gaussiana (I)

$$P_{\text{Gauss}}(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}} .$$

Para la distribución Normal (o Gaussiana), su media y varianza vienen dadas por

$$E[x] = \mu, \quad V[x] = \sigma^2 .$$

Usamos deliberadamente los símbolos  $\mu$  y  $\sigma$  tanto para los parámetros de la PDF Gaussiana como para su media y su varianza: la Gaussiana está caracterizada unívocamente por sus primer y segundo momentos;  $\mu_i = 0 \forall i > 2$ .



La dispersión de una distribución a un solo máximo es en ocasiones caracterizada en términos de su *anchura a media altura* o FWHM (full width at half-maximum); para Gaussianas, la relación con la varianza es as  $\text{FWHM} = 2\sqrt{2 \ln 2} \simeq 2,35\sigma$ .

La Gaussiana es también la distribución límite para las binomiales y Poisson, en los límites a gran  $n$  y gran  $\lambda$ , respectivamente:

$$P_{\text{binomial}}(k; n \rightarrow \infty, p) \rightarrow P_{\text{Gauss}}(k; np, np(1-p)) ,$$

$$P_{\text{Poisson}}(k; \lambda \rightarrow \infty) \rightarrow P_{\text{Gauss}}(k; \lambda, \sqrt{\lambda}) .$$

Nota :, una corrección de continuidad es requerida: el rango de la Gaussiana se extiende a valores negativos, mientras que Binomial y Poisson están solamente definidas en el rango positivo.



## La distribución continua de mayor ubicuidad: la Gaussiana (II)

La primacía de la Gaussiana en términos de su relevancia conceptual y sus aplicaciones prácticas proviene en gran parte del *teorema del límite central* : si tenemos  $n$  variables aleatorias independendientes  $\vec{x} = \{x_1, x_2, \dots, x_n\}$ , cada una con medias y varianzas  $\mu_i$  y  $\sigma_i^2$ , la resultante  $S(\vec{x})$  de sumarlas todas

$$S = \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{x_i - \mu_i}{\sigma_i},$$

sigue una distribución que, en el límite de gran  $n$ , tiende a una distribución normal reducida.

(El caso particular  $\mu = 0$ ,  $\sigma = 1$  es llamado en ocasiones “normal reducida”).

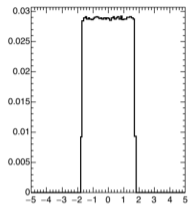
Por ello, una gran variedad de procesos, sujetos a múltiples fuentes independientes de incertidumbre, pueden describirse en buena aproximación con distribuciones Gaussianas, sin necesidad de conocer los detalles específicos de cada fuente de incertidumbre.

(ejemplo gráfico en la lámina siguiente, ilustrando la convergencia rápida del teorema del valor central para suma de distribuciones uniformes)

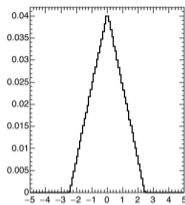
Ejercicio : verificar con una aplicación numérica la validez del teorema para otras distribuciones, por ejemplo sumas de exponenciales.



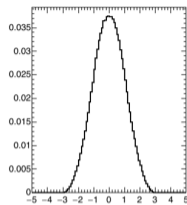
# La distribución continua de mayor ubicuidad: la Gaussiana (III)



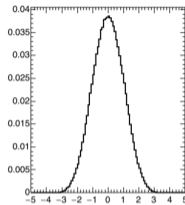
$$\sqrt{3} x_1$$



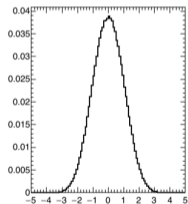
$$\sqrt{3}(x_1+x_2)/\sqrt{2}$$



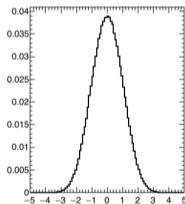
$$\sqrt{3}(x_1+x_2+x_3)/\sqrt{3}$$



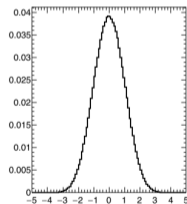
$$\sqrt{3}(x_1+x_2+x_3+x_4)/\sqrt{4}$$



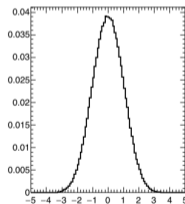
$$\sqrt{3}(x_1+x_2+x_3+x_4+x_5)/\sqrt{5}$$



$$\sqrt{3}(x_1+x_2+x_3+x_4+x_5+x_6)/\sqrt{6}$$



$$\sqrt{3}(x_1+x_2+x_3+x_4+x_5+x_6+x_7)/\sqrt{7}$$



$$\sqrt{3}(x_1+x_2+x_3+x_4+x_5+x_6+x_7+x_8)/\sqrt{8}$$



## Ejemplo de una distribución continua: $\chi^2$

$$P_{\chi^2}(x; n) = \begin{cases} \frac{x^{n/2-1} e^{-x/2}}{2^{n/2-1} \Gamma(\frac{n}{2})} & , \quad x \geq 0, n \text{ entero} \\ 0 & , \quad x < 0 \text{ o } n \text{ no - entero,} \end{cases}$$

con único parámetro  $n$ , y donde  $\Gamma(n/2)$  es la función Gamma. Su media y su varianza vienen dadas por  
 $E[x] = n$ ,  $V[x] = 2n$ .

La forma de la distribución de  $\chi^2$  depende de  $n$ .

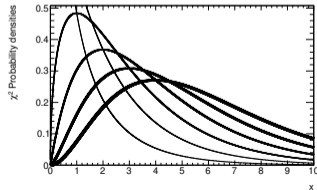
Importante: La distribución de  $\chi^2$  puede escribirse como la suma de cuadrados de  $n$  variables normales reducidas  $x_i$ , cada una de media  $\mu_i$  y varianza  $\sigma_i^2$ :

$$P_{\chi^2}(x; n) = \sum_{i=1}^n \left( \frac{x_i - \mu_i}{\sigma_i} \right)^2 .$$

El parámetro  $n$  es también llamado “número de grados de libertad”, y se refiere al comportamiento de los ajustes por el método de los cuadrados mínimos: cuando  $n_d$  puntos son utilizados para estimar  $n_p$  parámetros, el número correspondiente de grados de libertad es  $n_d - n_p$ .

Si el modelo utilizado para ajustar es adecuado, los valores de los  $\chi^2$  obtenidos deben seguir  $\chi^2$  distribution.

Este criterio es llamado *goodness-of-fit* o “bondad del ajuste”.





# Ejemplo de una distribución continua: Breit-Wigner (I)

También llamada distribución de Cauchy, o Lorentziana, la distribución

$$P_{\text{BW}}(x; \Gamma, x_0) = \frac{1}{\pi} \frac{\Gamma/2}{(x - x_0)^2 + \Gamma^2/4},$$

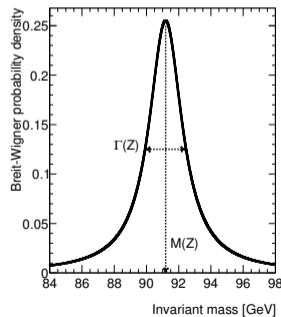
es a menudo utilizada para describir un proceso resonante (p.e. la masa invariante de productos del decaimiento de un estado intermedio resonante), con lo que  $x_0$  y  $\Gamma$  son la masa y la anchura natural, que es su FWHM.

La BW es un ejemplo de una distribución “fat-tailed”:

- ▶ promedio empírico y RMS mal definidos (idem momentos superiores)
- ▶ sus varianzas aumentan con el tamaño de la muestra
- ▶ ejercicio : verificar con una simulación numérica que el RMS empírico de la Breit-Wigner aumenta (siguiendo una ley de potencia) con el rango en el que es estimado (ver gráfica en la lámina siguiente)
- ▶ un teorema curioso (que no demostraremos) el estimador de  $x_0$  menos ineficaz es una media truncada sobre el 24 % central de la muestra : otras truncaciones son menos eficaces
- ▶ ejercicio opcional : verificar este teorema con una simulación numérica

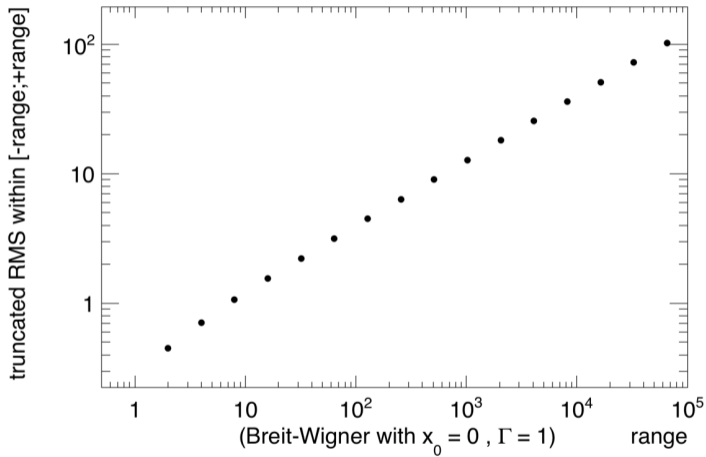
Gráfico : una Breit-Wigner con la masa y la anchura natural del bosón Z :

$$M_Z = 91,1876 \pm 0,0021 \text{ GeV}/c^2, \Gamma_Z = 2,4952 \pm 0,0023 \text{ GeV}/c^2.$$





## Ejemplo de una distribución continua: Breit-Wigner (II)





## Ejemplo de una distribución continua: Voigtiana (I)

La Voigtiana es la convolución de una Breit-Wigner con una Gaussiana,

$$P_{\text{Voigt}}(x; x_0, \Gamma, \sigma) = \int_{-\infty}^{+\infty} dx' P_{\text{Gauss}}(x'; 0, \sigma) P_{\text{BW}}(x - x'; x_0, \Gamma) ,$$

- ▶ La Voigtiana es una distribución a tres parámetros : masa  $x_0$ , anchura natural  $\Gamma$  y resolución  $\sigma$ .
- ▶ Es un modelo de medidas de procesos resonantes, suponiendo que la resolución instrumental sea Gaussiana.
- ▶ No hay una forma analítica cerrada sencilla para la Voigtiana, pero hay implementaciones numéricas precisas y eficientes, p.e. la función miembro `TMath::Voigt` en ROOT, o la clase `RooVoigtian` en RooFit.
- ▶ Para valores de  $\Gamma$  y  $\sigma$  razonablemente similares, la FWHM de la Voigtiana se aproxima a una combinación en suma y en cuadratura :

$$\text{FWHM}_{\text{Voigt}} \simeq \left[ (\Gamma/2) \oplus 2\sqrt{2 \ln 2} \sigma \right] + \Gamma/2 .$$

### Ejercicio :

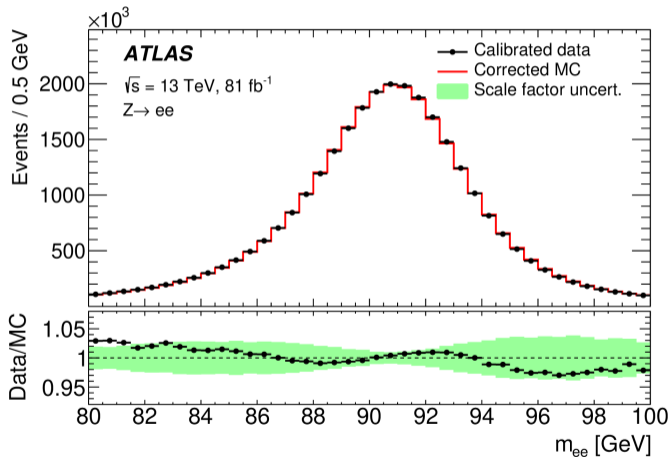
- ▶ La figura en la lámina siguiente muestra la distribución de masa invariante dielectrón alrededor de la masa del bosón Z, obtenida con datos ATLAS a 13 TeV.
- ▶ Suponiendo que la distribución puede ser descrita por una Voigtiana, determinar aproximadamente cuál es la resolución en masa dielectrón del detector ATLAS.
- ▶ La parte inferior de la figura muestra una comparación entre los datos y la simulación. Comentar.



# Ejemplo de una distribución continua: Voigtiana (II)

$$\text{FWHM}_{\text{Voigt}} \simeq \Gamma/2 + \left[ (\Gamma/2) \oplus 2\sqrt{2 \ln 2} \sigma \right]$$

$$M_Z = 91,1876 \pm 0,0021 \text{ GeV}/c^2$$
$$\Gamma_Z = 2,4952 \pm 0,0023 \text{ GeV}/c^2$$
$$\sigma(m_{ee}) = ?$$







# Otros ejemplos de distribuciones continuas

La “Gaussiana Bifurcada” : útil para describir distribuciones asimétricas

$$P_{BG}(x; \mu, \sigma_L, \sigma_R) = \sqrt{\frac{2}{\pi(\sigma_L + \sigma_R)}} \times \begin{cases} e^{-\frac{(x-\mu)^2}{2\sigma_L^2}} & , x \leq \mu \\ e^{-\frac{(x-\mu)^2}{2\sigma_R^2}} & , x \geq \mu \end{cases}$$

La “Crystal Ball” : en ocasiones utilizada para describir resoluciones con una componente de “fuga” (leakage)

$$P_{CB}(x; \mu, \sigma, \alpha, n) = N \times \begin{cases} e^{-\frac{(x-\mu)^2}{2\sigma^2}} & , \frac{x-\mu}{\sigma} \leq \alpha \\ \left( \frac{n}{n-\alpha^2 + \alpha\sigma^{-1}(x-\mu)} \right)^n e^{-\frac{\alpha^2}{2}} & , \frac{x-\mu}{\sigma} \geq \alpha \end{cases}$$

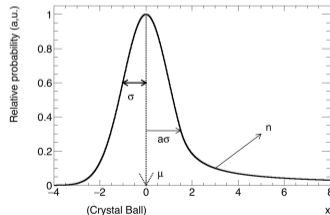
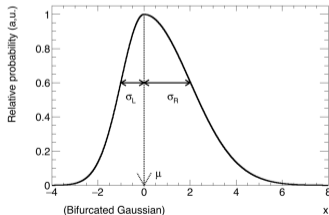
Funciones implementadas p.e. en RooFit

Interpretación “intuitiva” de los parámetros, pero

- ▶ los estimadores simples están sesgados
- ▶ correlaciones muy grandes entre ellos
- ▶ utilizar, pero con prudencia y ojo crítico

Ejercicio numérico : evaluar las matrices de correlación

Otras funciones similares : DSCB, ACB ...





En las previas láminas describimos una lógica “intuitiva” o “caso a caso” de la estimación de parámetros. Eso, por supuesto, no es generalizable (ni robusto) : de manera general, y más allá de los ejemplos sencillos

- ▶ los estimadores sencillos de los momentos de la función característica están sometidos a sesgos
- ▶ una descripción completa de la PDF requiere a priori un número infinito de momentos

Una descripción más general del problema es la siguiente :

- ▶ tenemos una muestra compuesta por  $N$  realizaciones independientes de variables aleatorias  $\vec{x}$
- ▶ suponemos que esas realizaciones resultan de muestrear una PDF  $n$ -paramétrica

$$P(\vec{x}; \theta_1, \dots, \theta_n) ,$$

- ▶ suponemos también que la dependencia funcional de la PDF es conocida, y que solamente ignoramos los valores numéricos de los parámetros

Fisher, 1921: el *teorema de la verosimilitud máxima* (maximum likelihood) es una herramienta poderosa para la estimación de parámetros  $\theta_1, \dots, \theta_n$  de nuestra PDF.



# El teorema de verosimilitud máxima (I)

Definimos la función de verosimilitud  $\mathcal{L}$ , evaluada sobre una muestra compuesta por  $N$  eventos :

$$\mathcal{L}(\theta_1, \dots, \theta_n) = \prod_{i=1}^N P(\vec{x}_i; \theta_1, \dots, \theta_n) .$$

Teorema : los valores  $\hat{\theta}_1, \dots, \hat{\theta}_n$  que maximizan la función  $\mathcal{L}$  son estimadores de los parámetros  $\theta_1, \dots, \theta_n$  de nuestra PDF,

$$\mathcal{L}(\hat{\theta}_1, \dots, \hat{\theta}_n) = \max_{\theta} \{ \mathcal{L}(\theta_1, \dots, \theta_n) \} ,$$

con varianzas  $\hat{\sigma}_{\theta}$  que se extraen a partir de la matriz de covarianza de  $\mathcal{L}$  alrededor de su máximo.

En palabras intuitivas: para una muestra dada, el MLE corresponde a los valores que maximizan la probabilidad de realizar esa muestra !

No es un pleonasma : la función  $\mathcal{L}$  debe satisfacer ciertas condiciones:

- ▶ ser derivable al menos dos veces con respecto a los parámetros  $\theta_1, \dots, \theta_n$ ,
- ▶ ser (asintóticamente) no sesgada y eficiente (condición llamada "Cramer-Rao bound"),
- ▶ seguir una distribución (asintóticamente) multi-normal,

$$f(\vec{\hat{\theta}}, \vec{\theta}, \Sigma) = \frac{1}{\sqrt{2\pi} |\Sigma|} \exp \left\{ -\frac{1}{2} (\vec{\hat{\theta}} - \vec{\theta}) \Sigma^{-1} (\vec{\hat{\theta}} - \vec{\theta}) \right\} .$$



## El teorema de verosimilitud máxima (II)

- ▶ seguir una distribución (asintóticamente) multi-normal,

$$f(\hat{\vec{\theta}}, \vec{\theta}, \Sigma) = \frac{1}{\sqrt{2\pi} |\Sigma|} \exp \left\{ -\frac{1}{2} (\hat{\theta}_i - \bar{\theta}_i) \Sigma_{ij}^{-1} (\hat{\theta}_j - \bar{\theta}_j) \right\}.$$

donde la matriz de covarianza  $\Sigma$  es

$$\Sigma_{ij}^{-1} = -E \left[ \frac{\partial^2 \mathcal{L}}{\partial \theta_i \partial \theta_j} \right].$$

Al alejarnos del máximo, el valor de la función  $\mathcal{L}$  disminuye, a una tasa que depende de los elementos de la matriz de covarianza :

$$-2\Delta \ln \mathcal{L} = -2 \left[ \ln \mathcal{L}(\vec{\theta}) - \ln \mathcal{L}(\hat{\vec{\theta}}) \right] = \sum_{i,j} (\theta_i - \hat{\theta}_i) \Sigma_{ij}^{-1} (\theta_i - \hat{\theta}_j).$$

En otras palabras: la matriz de covarianza define mapas de contorno alrededor de su máximo, que corresponden a *intervalos de confianza*.

En el caso de una  $\mathcal{L}$  con un parámetro único  $\mathcal{L}(\theta)$ , el intervalo contenido dentro de  $-2\Delta \ln \mathcal{L} < 1$  alrededor de  $\hat{\theta}$  define un intervalo de confianza a 68 % que corresponde a un rango  $-\Delta_\theta \leq \theta - \hat{\theta} \leq \Delta_\theta$  alrededor del punto máximo.

Por ello el resultado de MLE se escribe a menudo como  $(\hat{\theta} \pm \hat{\Delta}_\theta)$ .



## Ejemplo de estimación por verosimilitud máxima (I) : la Gaussiana

Una muestra compuesta por  $N$  realizaciones de una única variable aleatoria  $x$ , que suponemos sigue una distribución Gaussiana de media  $\mu$  y anchura  $\sigma$ . El teorema MLE nos permite estimar los parámetros así:

$$\mathcal{L}(\mu, \sigma) = \prod_{i=1}^N P_{\text{Gaus}}(x_i; \mu, \sigma) ; -\ln \mathcal{L} = N \ln \sigma + \sum_{i=1}^N \frac{(x_i - \mu)^2}{2\sigma^2} \quad (+\text{constante}).$$

Los estimadores  $\hat{\mu}$  y  $\hat{\sigma}$  son los ceros de las primeras derivadas de  $-\ln \mathcal{L}$  con respecto a  $\mu$  y  $\sigma$  :

$$\left. \frac{\partial}{\partial \mu} (-\ln \mathcal{L}) \right|_{\hat{\mu}, \hat{\sigma}} = 0 \quad \rightarrow \quad \hat{\mu} = \frac{1}{N} \sum_{i=1}^N x_i .$$

$$\left. \frac{\partial}{\partial \sigma} (-\ln \mathcal{L}) \right|_{\hat{\sigma}, \hat{\mu}} = 0 \quad \rightarrow \quad \hat{\sigma}^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2 .$$

Las segundas derivadas nos dan los elementos de la matriz de covarianza:

$$\left. \frac{\partial^2}{\partial \mu^2} (-\ln \mathcal{L}) \right|_{\hat{\mu}, \hat{\sigma}} = \frac{N}{\hat{\sigma}^2} \quad \rightarrow \quad \Sigma_{\mu\mu} = \frac{\hat{\sigma}^2}{N} \quad \rightarrow \quad \hat{\Delta}_{\mu} = \frac{\hat{\sigma}}{\sqrt{N}} ,$$
$$\left. \frac{\partial^2}{\partial \sigma^2} (-\ln \mathcal{L}) \right|_{\hat{\mu}, \hat{\sigma}} = \frac{2N}{\hat{\sigma}^2} \quad \rightarrow \quad \Sigma_{\sigma\sigma} = \frac{\hat{\sigma}^2}{2N} \quad \rightarrow \quad \hat{\Delta}_{\sigma} = \frac{\hat{\sigma}}{\sqrt{2N}} ,$$

(y los términos no diagonales de la covarianza son ambos cero).

El MLE extrae de manera formal y robusta los parámetros y errores ( $\hat{\mu} \pm \hat{\Delta}_{\mu}$ ) y ( $\hat{\sigma} \pm \hat{\Delta}_{\sigma}$ ) de una Gaussiana.



## Ejemplo de estimación por verosimilitud máxima (II) : la exponencial

Una muestra compuesta por  $N$  realizaciones de una única variable aleatoria  $x$ , que suponemos sigue una distribución exponencial con parámetro de forma  $\xi$ . El teorema MLE nos permite estimar ese parámetro  $x_i$  analíticamente, al menos para el caso en que la variable  $x$  cubre el rango  $[0; +\infty]$ .

$$\mathcal{L}(\xi) = \prod_{i=1}^N P_{\text{exponencial}}(x_i; \xi) = \prod_{i=1}^N \frac{1}{\xi} e^{-x_i/\xi} .$$

De manera análoga al ejemplo anterior, determinamos el NLL :

$$-\ln \mathcal{L} = N \ln \xi - \frac{1}{\xi} \sum_{i=1}^N x_i ,$$

y evaluamos sus primera y segunda derivadas para el valor  $\hat{\xi}$  que anula la primera derivada :

$$\begin{aligned} \left. \frac{\partial}{\partial \xi} (-\ln \mathcal{L}) \right|_{\hat{\xi}} = 0 &\longrightarrow \hat{\xi} = \frac{1}{N} \sum_{i=1}^N x_i , \\ \left. \frac{\partial^2}{\partial \xi^2} \right|_{\hat{\xi}} = \frac{N}{\hat{\xi}^2} &\longrightarrow \Sigma_{\xi\xi} = \frac{\hat{\xi}^2}{N} \longrightarrow \hat{\Delta}_{\xi} = \frac{\hat{\xi}}{\sqrt{N}} . \end{aligned}$$

Con lo que el MLE nos da en este caso también una solución analítica para la estimación del parámetro de forma exponencial y su correspondiente error ( $\hat{\xi} \pm \hat{\Delta}_{\xi}$ ).

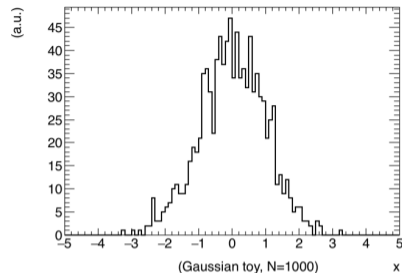
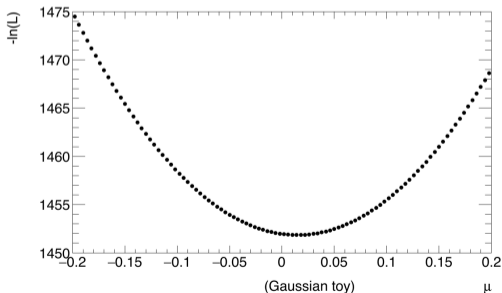
Nota: en este caso como el anterior, los errores  $\hat{\Delta}$  escalan como  $1/\sqrt{N}$ . Esa es una propiedad muy general.



# Un “ejemplo de juguete” (I)

- ▶ Generamos una realización aleatoria con  $N=1000$  eventos (“toy sample”) a partir de una PDF Gaussiana reducida, con  $\mu = 0, \sigma = 1$ .
- ▶ Evaluamos (el negativo del logaritmo de) la función de verosimilitud para esa muestra, para diferentes valores de  $\mu$  y  $\sigma$  (“escaneamos” los parámetros de interés)

$$-\ln \mathcal{L}(\mu, \sigma) = N \ln \sigma + \sum_{i=1}^N \frac{(x_i - \mu)^2}{2\sigma^2} .$$



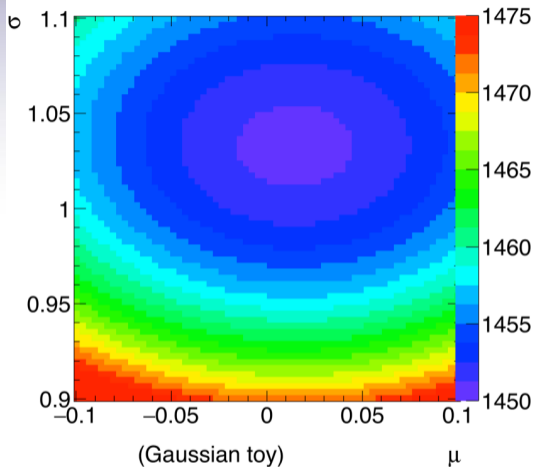
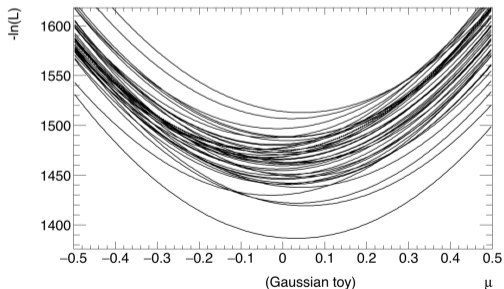
- ▶ La forma de  $-\ln \mathcal{L}$  en función de  $\mu$  sigue un perfil parabólico, y el mínimo coincide con el promedio empírico
- ▶ Por inspección alrededor del mínimo, se observa que el NLL aumenta en 0,5 unidades para  $\Delta x \sim \pm 0,03$ , que corresponde aproximadamente a  $\sigma/\sqrt{N}$  para  $N = 1000$
- ▶ La lámina siguiente muestra el gráfico del escaneo en 2D de  $\mu$  y  $\sigma$



## Un “ejemplo de juguete” (II)

- ▶ El perfil bidimensional de  $-\ln \mathcal{L}$  es el de un paraboloide, con semi-ejes diferentes y ortogonales
- ▶ El mínimo coincide con la posición del promedio y el RMS empíricos
- ▶ El semi-eje a lo largo de  $\sigma$  es más estrecho que el de  $\mu$ , como se espera de la relación  $1/\sqrt{2N}$  vs.  $1/\sqrt{N}$

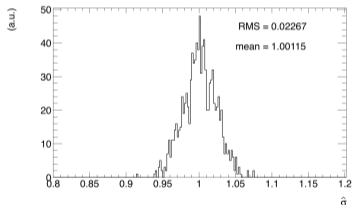
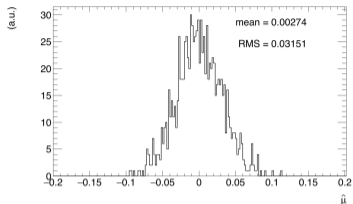
Para otras realizaciones aleatorias independientes, la posición y la anchura de los mínimos cambia :





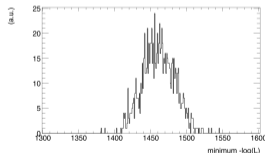


## Un “ejemplo de juguete” (III)



Si realizamos el mismo estudio sobre un ensamble de muestras generadas con la misma PDF (aquí  $N=1000$ ,  $\mu = 0$ ,  $\sigma = 1$ ) obtenemos los resultados siguientes :

- ▶ los mínimos de  $\hat{\mu}$  fluctúan alrededor de su valor verdadero  $\mu = 0$  con una dispersión de  $\pm 3,1\%$ , que corresponde a  $\sigma/\sqrt{N}$
- ▶ los mínimos de  $\hat{\sigma}$  fluctúan alrededor de su valor verdadero  $\sigma = 1$  con una dispersión de  $\pm 2,3\%$ , que corresponde a  $\sigma/\sqrt{2N}$
- ▶ los intervalos con  $-\Delta \ln \mathcal{L} < 0,5$  alrededor del mínimo corresponden bien a las regiones que cubren 68,3% de la dispersión
- ▶ el teorema MLE nos da una definición rigurosa y precisa de la incertidumbre estadística



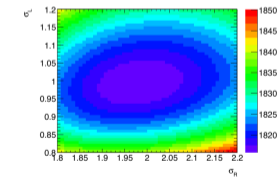
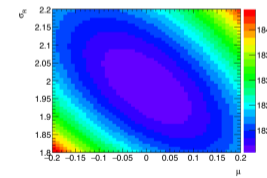
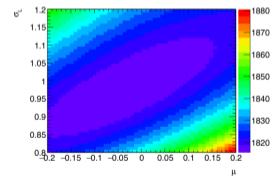
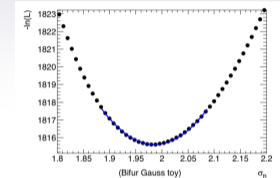
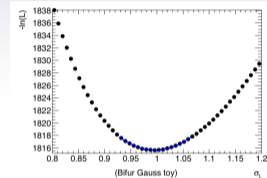
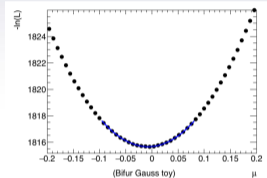
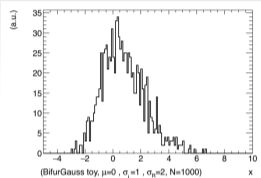
La distribución del  $-\ln \mathcal{L}$  no siempre sigue una forma precisa (aquí parece bastante Gaussiana) pero permite hacer una estimación del “goddness-of-fit” :

- ▶ si el  $-\ln \mathcal{L}$  observado en una muestra se aleja significativamente del intervalo cubierto en un “estudio de juguete”, la calidad del modelo es sospechosa...



# Otro "ejemplo de juguete" (I)

Consideremos ahora una Gaussiana bifurcada, PDF con tres parámetros :  $\mu$ ,  $\sigma_L$  y  $\sigma_R$ . Los mínimos del  $-\log \mathcal{L}$  no son del todo parabólicos, y hay correlaciones importantes entre los parámetros.

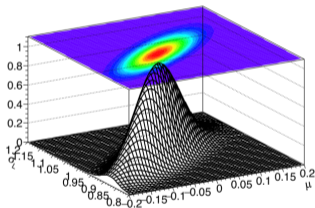




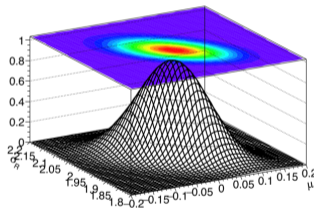
## Otro "ejemplo de juguete" (II)

Los parámetros de la PDF son ellos mismos variables aleatorias:

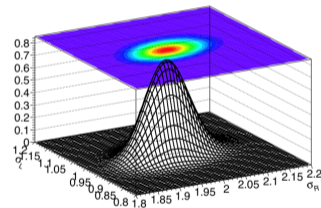
- ▶ si efectuamos otras realizaciones aleatorias a partir de la misma PDF, obtendremos valores diferentes de los parámetros
- ▶ y éstos se distribuirán siguiendo la matriz de covariancia de los parámetros (y por tanto tomando en cuenta las correlaciones)



Correlación entre  $\mu$  y  $\sigma_L$  : +79 %



Correlación entre  $\mu$  y  $\sigma_R$  : -59 %



Correlación entre  $\sigma_L$  y  $\sigma_R$  : +19 %

Nota: se trata por supuesto de una PDF en 3 dimensiones: por eso se muestra 3 proyecciones bidimensionales.



Los ejemplos de juguete discutidos previamente son casos sencillos, con 2 o 3 parámetros, y pudimos explorar el espacio bi- o tri-dimensional con bucles sencillas.

En situaciones más generales, el número de parámetros puede ser significativamente superior, así que la aproximación de “escanear” los parámetros para identificar el mínimo del  $-\log \mathcal{L}$  no es eficaz (y se vuelve rápidamente imposible).

Por ello se utilizan algoritmos de minimización numérica, que optimizan la búsqueda del mínimo de una función. El proceso se llama *ajuste numérico* o “fit”. El algoritmo más usado en altas energías es MINUIT diseñado en el CERN en los años 1970. MINUIT está implementado en ROOT, en la clase TMinuit.

## MINUIT

From Wikipedia, the free encyclopedia

**MINUIT**, now **MINUIT2**, is a [numerical minimization computer program](#) originally written in the [FORTRAN programming language](#)<sup>[1]</sup> by CERN staff physicist Fred James in the 1970s. The program searches for a minimum in a user-defined [function](#) with respect to one or more [parameters](#) using several different methods as specified by the user. In addition to that it can compute confidence intervals for the parameters by scanning the function around the minimum.

The original FORTRAN code was later ported to [C++](#) by the [ROOT](#) project; both the FORTRAN and C++ versions are in use today. The program is very widely used in [particle physics](#), and thousands of published papers cite use of MINUIT.<sup>[2]</sup> In the early 2000s, Fred James started a project to implement MINUIT in C++ using [object-oriented programming](#). The new MINUIT is an optional package (minuit2) in the ROOT release. As of October 2014 the latest version is 5.34.14, released on 24 January 2014.<sup>[3]</sup> There is also a [Java](#) port<sup>[4]</sup> as well as a [Python](#) frontend to the C++ code.<sup>[5]</sup>

MINUIT is not a program that can be distributed as an [executable](#) binary to be run by a relatively unskilled user: the user must write and [compile](#) a subroutine defining the function to be optimized, and oversee the optimization process.



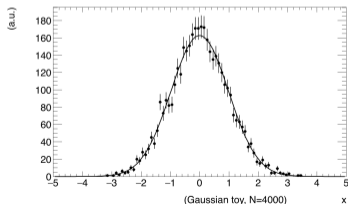
## Diferencia entre un fit NLL y un fit de $\chi^2$

Muchos paquetes y programas contienen algoritmos de minimización más sencillos, del tipo “mínimos cuadrados” o  $\chi^2$ . Estos difieren de un ajuste de verosimilitud máxima: aquí la muestra se compone de un conjunto de  $n$  puntos  $y_i$  con sus incertidumbres  $\sigma_i$ , y la función a minimizar es:

$$\chi^2(\theta) = \sum_{i=1}^n \left( \frac{y_i - f(x_i; \theta)}{\sigma_i} \right)^2,$$

donde  $f(x; \theta)$  es la función que pretende describir una dependencia funcional  $y = f(x)$ . La minimización del  $\chi^2$  provee los valores  $\hat{\theta}$  estimados por el ajuste.

El conjunto de puntos puede provenir de medidas individuales, o ser una reducción por histogramado de una muestra completa: en el ejemplo aquí se tiene una muestra compuesta de 4000 valores de una variable aleatoria  $x$ , histogramada con 100 *bines* de anchura  $\delta x = 0,1$  en el intervalo  $-5 \leq x \leq +5$ . El equivalente a  $y_i$  es el número de eventos con valores contenidos en el bin  $i$ , y la incertidumbre asociada es  $\sigma_i = \sqrt{y_i}$ .



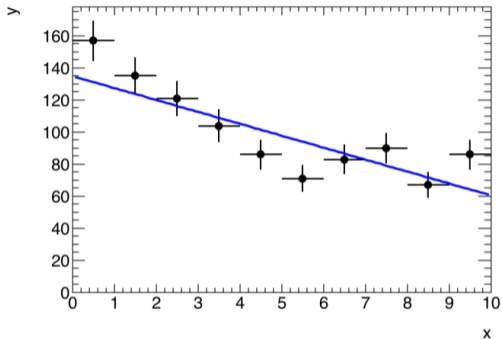
El ajuste de este histograma a una función Gaussiana nos da  $\mu = (0,008 \pm 0,016)$  y  $\sigma = (1,007 \pm 0,011)$ .

Nota: si se hubiera usado otro “binning” el resultado puede ser un poco diferente!

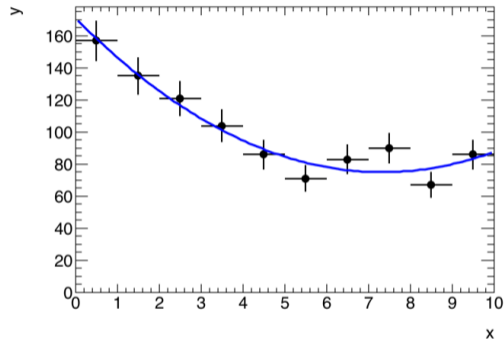
Nota: en cambio, el ajuste por verosimilitud máxima a esta misma muestra nos dará exactamente los resultados teóricos: la media empírica y el RMS empírico, con sus incertidumbres  $\text{RMS}/\sqrt{n}$  y  $\text{RMS}/\sqrt{2n}$ .



## Otro ejemplo de un fit de $\chi^2$



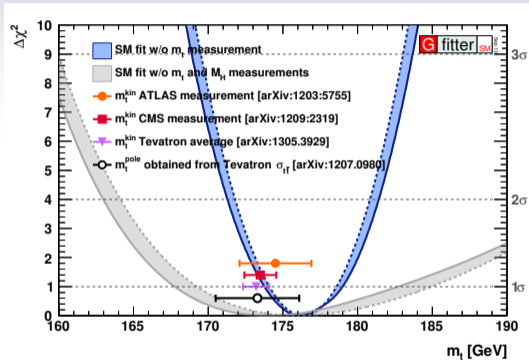
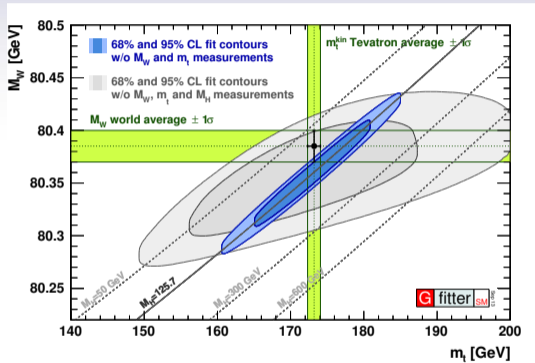
Un fit a un polinomio de orden 2 da  $\chi^2 = 6,5$  para 7 grados de libertad.



Mientras que un fit a un polinomio de orden 1 da  $\chi^2 = 21,2$  para 8 grados de libertad.

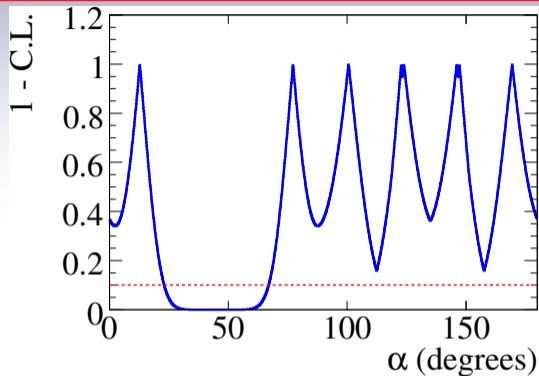


# Los intervalos de confianza de Gfitter

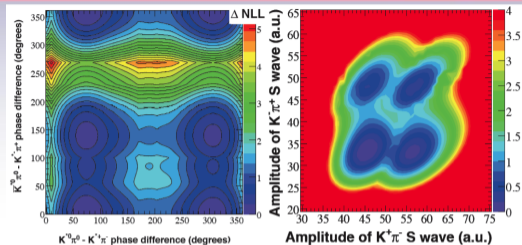




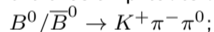
# Ejemplos de funciones de verosimilitud complicadas



Un elemento importante del programa científico del experimento *BABAR* (y en general, de la *física de sabores*): la medida del ángulo  $\alpha$  de la matriz CKM:  $\alpha \neq 0 \rightarrow$  violación de la simetría CP.  
 Problema : el observable físico es la asimetría dependiente del tiempo  $B^0/\bar{B}^0 \rightarrow \pi^+\pi^-$ , que es una función de  $\sin 2(\alpha - \delta)$ : ambigüedad octuple.



Situación más complicada: el perfil de interferencias entre las amplitudes de decaimiento



- ▶ ambigüedades múltiples, sin forma analítica precisa...
- ▶ numerosas amplitudes resonantes intermedias:  $K^{*\pm}\pi^\mp$ ,  $K^{*0}\pi^0$ ,  $\rho^\pm\pi^\mp$ , varios estados excitados  $K^*$ ,  $\rho$ , otras resonancias... todo para  $B^0$  y para  $\bar{B}^0$
- ▶ la función de verosimilitud contenía unos 60 parámetros de interés: amplitudes y fases...





# Ejemplos de funciones de verosimilitud complicadas

B. AUBERT *et al.*

PHYSICAL REVIEW D **86**, 112001 (2009)

TABLE VIII. Full correlation matrix for the isobar parameters of solution I. The entries are given in percent. Since the matrix is symmetric, all elements above the diagonal are omitted.

[i]	[j]																	
	$\rho^0$	$K^*$	$S$	$f_1$	$f_2$	$f_x$	$NR$	$x$	$f_0$	$\rho^0$	$K^*$	$S$	$f_1$	$f_2$	$f_x$	$NR$	$x$	
[i]	$\rho^0$	100.0																
	$K^*$	51.9	100.0															
	$S$	54.0	65.0	100.0														
	$f_1$	8.4	2.8	21.0	100.0													
	$f_2$	14.9	23.2	32.2	22.7	100.0												
	$NR$	5.2	35.0	24.4	12.6	39.3	100.0											
	$x$	6.4	9.9	7.8	2.0	7.4	6.1	100.0										
[i]	$f_0$	31.3	30.3	39.9	25.2	36.7	31.3	8.0	100.0									
	$\rho^0$	20.6	48.6	51.2	8.0	27.7	27.5	5.6	17.3	100.0								
	$K^*$	44.7	73.5	56.3	-4.8	26.9	22.0	9.5	22.6	43.4	100.0							
	$S$	59.6	79.9	79.7	21.8	39.3	26.9	11.3	35.2	49.4	57.7	100.0						
	$f_1$	2.4	-30.1	6.3	-56.1	-1.5	3.9	-0.3	10.7	-6.2	-21.5	5.0	100.0					
	$f_2$	14.5	34.1	12.5	16.1	-23.0	12.4	2.5	34.5	7.3	8.3	12.9	-6.2	100.0				
	$NR$	17.8	57.6	41.7	13.7	30.1	49.7	2.4	40.0	32.1	25.0	31.7	7.5	46.2	100.0			
	$x$	18.9	27.0	20.6	5.8	11.8	9.5	-84.2	21.5	17.8	28.1	27.8	0.8	8.1	20.2	100.0		
ang(c)	$\rho^0$	-11.2	13.3	4.0	-16.1	-2.9	-2.1	-0.5	-0.2	24.1	16.3	3.2	-3.3	8.9	2.1	4.2		
	$K^*$	25.0	8.6	-3.2	-0.2	-15.7	-9.7	6.3	-10.4	-3.9	5.5	16.0	3.8	6.3	-6.5	-3.2		
	$S$	33.0	10.6	3.4	-4.7	-17.3	-16.5	6.2	-9.6	1.0	18.7	21.3	-4.2	9.6	-4.2	1.3		
	$f_1$	12.1	-0.6	-9.8	-2.6	-23.1	-27.4	0.9	-16.7	-7.2	2.2	1.1	-10.6	7.2	-14.1	-2.6		
	$f_2$	25.0	10.2	5.4	-0.5	-11.4	-11.8	1.0	-0.8	2.6	8.5	11.8	-3.8	15.6	2.4	0.4		
	$NR$	31.6	17.0	39.3	1.0	-27.1	-31.7	-6.7	11.3	12.8	14.5	19.0	3.3	21.5	19.6	14.2		
	$x$	8.6	1.8	9.8	0.6	-9.9	-8.9	-7.9	2.8	3.8	1.3	4.2	3.5	7.3	8.9	12.4		
ang(f)	$\rho^0$	22.2	11.7	18.9	3.5	-20.3	-26.2	-1.6	-3.6	-6.9	7.3	18.2	1.8	20.9	-7.1	4.3		
	$K^*$	14.5	18.0	14.6	-17.3	-13.4	-21.0	-0.7	-8.7	14.3	19.8	13.4	1.7	7.2	-4.4	5.4		
	$S$	17.1	7.1	22.0	5.2	-13.5	-17.3	-2.1	5.0	7.2	6.5	13.8	8.1	12.8	29.5	9.6		
	$f_1$	22.5	15.9	25.2	-3.2	-16.9	-21.6	-0.5	4.2	10.6	17.7	16.1	1.7	15.7	28.8	30.8		
	$f_2$	15.1	4.9	15.5	-5.0	-15.5	-17.9	-2.1	10.0	-2.5	3.9	3.9	11.1	15.7	18.6	7.5		
	$NR$	8.1	2.7	12.3	-0.6	16.5	-20.4	-0.9	12.2	6.1	3.4	4.8	1.4	-14.6	4.7	6.5		
	$x$	15.3	4.1	14.5	-3.0	-22.6	-20.8	0.8	1.7	8.2	1.8	5.2	2.6	20.0	15.1	3.2		
	$y$	10.9	1.1	12.8	0.7	-13.9	-18.0	-4.7	2.1	3.3	0.6	3.9	5.9	9.8	13.4	8.2		

112001-26

TIME-DEPENDENT AMPLITUDE ANALYSIS OF ...

PHYSICAL REVIEW D **86**, 112001 (2009)

TABLE IX. Full correlation matrix for the isobar parameters of solution II. The entries are given in percent. Since the matrix is symmetric, all elements above the diagonal are omitted.

[i]	[j]																	
	$\rho^0$	$K^*$	$S$	$f_1$	$f_2$	$f_x$	$NR$	$x$	$f_0$	$\rho^0$	$K^*$	$S$	$f_1$	$f_2$	$f_x$	$NR$	$x$	
[i]	$\rho^0$	100.0																
	$K^*$	46.9	100.0															
	$S$	49.1	66.2	100.0														
	$f_1$	8.7	7.7	25.4	100.0													
	$f_2$	16.8	40.3	38.5	26.6	100.0												
	$NR$	-8.4	30.2	21.2	9.4	49.9	100.0											
	$x$	5.5	11.7	9.3	3.4	12.1	9.1	100.0										
[i]	$f_0$	29.2	42.1	30.2	31.5	57.9	34.1	10.0	100.0									
	$\rho^0$	61.5	68.1	40.4	6.9	20.6	6.4	6.0	31.6	100.0								
	$K^*$	39.8	75.7	59.8	0.3	33.1	25.3	10.9	33.2	36.3	100.0							
	$S$	50.6	75.2	83.2	25.4	49.9	33.4	13.1	51.6	46.0	61.4	100.0						
	$f_1$	0.8	-6.1	9.6	-53.9	6.0	13.3	0.2	14.7	5.3	-18.5	10.4	100.0					
	$f_2$	10.0	-3.3	-0.9	-10.6	-68.7	-17.8	-5.2	-18.4	6.3	-4.0	-4.9	2.2	100.0				
	$NR$	23.1	68.8	44.7	13.5	39.3	34.4	5.8	45.6	58.3	32.8	45.4	14.7	-13.8	100.0			
	$x$	22.3	33.5	37.8	9.8	19.3	9.9	-79.2	31.3	20.7	30.2	36.1	3.3	-2.6	23.9	100.0		
ang(c)	$\rho^0$	-23.1	13.7	5.5	-11.4	8.0	5.2	0.0	9.0	-11.9	14.5	6.3	-0.2	0.3	3.8	6.9		
	$K^*$	30.6	2.0	-2.2	-6.3	-16.1	-28.1	-0.0	-15.2	14.3	-1.4	6.1	2.0	19.4	-10.3	0.5		
	$S$	38.1	8.9	1.8	-10.1	-17.9	-39.5	-0.1	-15.8	17.4	9.4	7.7	-8.2	19.7	-12.1	3.7		
	$f_1$	18.1	-10.0	-13.7	-7.4	-15.4	-41.3	-2.4	-18.6	1.0	-6.2	-10.2	-12.7	10.7	-21.6	-2.7		
	$f_2$	26.2	-7.8	-12.2	-5.9	-7.7	-35.9	-1.7	-14.5	7.8	-5.7	-8.8	-9.9	12.2	-15.2	-3.6		
	$NR$	32.4	-0.4	21.4	0.5	-29.5	-65.2	-10.4	-4.2	12.0	0.2	0.4	-6.4	21.2	-8.1	10.2		
	$x$	15.4	-2.2	0.2	-1.6	-9.9	-18.3	-4.9	-5.6	5.6	-3.0	-0.2	-0.8	9.2	-5.6	4.0		
ang(f)	$\rho^0$	30.1	-8.0	-2.3	-0.9	-11.2	-43.0	-2.8	-16.7	12.1	-7.2	-5.5	-4.9	10.4	-18.7	-1.6		
	$K^*$	7.6	11.4	5.8	-7.5	-1.8	-24.7	0.6	-7.5	4.1	15.1	5.5	-12.6	1.3	-7.0	4.0		
	$S$	27.0	8.0	7.6	5.6	2.8	-27.8	0.6	-2.0	9.1	1.5	7.1	3.2	-6.9	13.9	4.1		
	$f_1$	32.6	8.8	4.4	-1.1	0.6	-31.3	2.1	-4.1	12.6	12.1	7.6	-5.6	-4.4	12.2	4.7		
	$f_2$	18.7	1.7	6.6	10.1	9.8	-22.9	0.7	7.6	8.6	3.5	0.6	-5.6	-2.6	9.3	4.6		
	$NR$	21.9	1.8	4.4	9.6	-0.7	-30.2	0.1	-5.0	8.1	2.8	4.0	-17.3	1.0	-6.6	-0.2		
	$x$	27.7	-1.9	-3.0	3.9	-0.5	-30.7	2.8	-13.3	7.8	-1.5	-1.2	-13.8	-7.2	-3.7	-5.0		
	$y$	19.7	-5.0	-0.5	2.3	-4.4	-27.6	2.7	-6.1	6.2	-4.1	-2.5	-1.6	-0.2	-0.1	-2.9		

112001-27





- ▶ ¿Preguntas sobre los ejercicios?
- ▶ ¿Preguntas sobre las clases anteriores?
- ▶ Seguiremos explorando las propiedades del MLE y sus aplicaciones en física

Un pequeño acertijo :

- ▶ En el seminario del lunes 22, Xavier Bertou dijo que la colaboración Auger había liberado 10 % de sus datos
- ▶ El experimento Auger ha detectado un total de 15 eventos con energías superiores al corte GZK
- ▶ Curiosamente, 5 de esos eventos se encuentran en el lote reducido de 10 % que es de acceso público

Pregunta: ¿Cuán probable (o improbable) es esa situación?



# MLE en situaciones más elaboradas (I)

En un escenario típico, un proceso aleatorio puede tener contribuciones de origen diferente.

Para ser específicos, consideremos que los eventos que componen la muestra provienen de dos “especies”, llamadas de manera genérica “señal” y “fondo” (la generalización a más de dos especies es sencilla).

Cada especie se realiza a partir de su propia densidad de probabilidad.

Si los rangos de las variables aleatorias no son totalmente disyuntos, es imposible saber evento a evento a cuál de las especies pertenece. Pero el MLE permite efectuar una *separación estadística*: la PDF subjacente es a combinación de mas PDFs de señal y fondo,

$$\mathcal{L}(f_{\text{sig}}, \theta; \vec{x}) = \prod_{i=1}^N [f_{\text{sig}} P_{\text{sig}}(\vec{x}; \theta) + (1 - f_{\text{sig}}) P_{\text{bkg}}(\vec{x}; \theta)] ,$$

donde  $P_{\text{sig}}$  y  $P_{\text{bkg}}$  son las PDFs de señal y fondo, respectivamente, y la fracción de señal  $f_{\text{sig}}$  es el parámetro que cuantifica la pureza de la muestra :  $0 \leq f_{\text{sig}} \leq 1$ .

Ejemplo inspirado de la búsqueda del bosón de Higgs en el canal difotón :

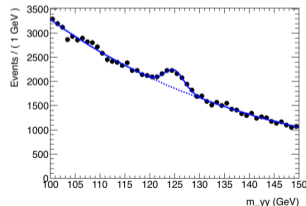
Con ROOT instalado, la macro H\_yy.cc debe correr sin problema, haciendo

```
prompt> root -l H_yy.cc
```

```
RoofitResult: minimized FCN value: 386054, estimated distance to minimum: 2.29788e-05
covariance matrix quality: Full, accurate covariance matrix
Status : MIGRAD=0 HESSE=0

Constant Parameter      Value
-----
sig_s                    2.0000e+00

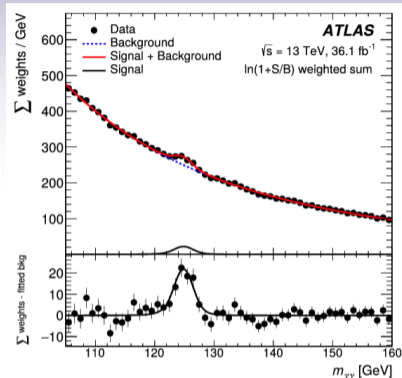
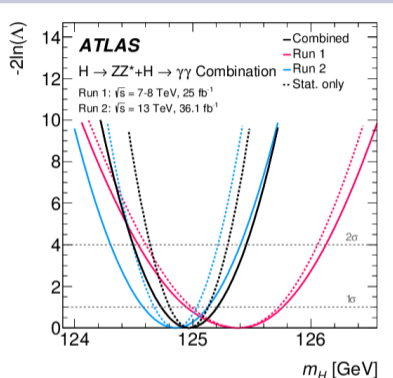
Floating Parameter      InitialValue      FinalValue +/-      Error      GblCorr.
-----
bkg_a                    -2.2500e-02      -2.2780e-02 +/-      2.31e-04      <none>
sig_f                    2.0000e-02      1.8732e-02 +/-      1.33e-03      <none>
sig_n                    1.2500e+02      1.2494e+02 +/-      1.92e-01      <none>
```





# MLE en situaciones más elaboradas (II)

Figuras tomadas de M. Aaboud *et al*, the ATLAS Collaboration, Phys.Lett.B 784 (2018) 345-366



La figura que representa  $-2\Delta(\ln \mathcal{L})$  en función de  $m_H$  ilustra claramente la interpretación del MLE en términos de *intervalos de confianza*:

- ▶ el rango de  $m_H$  que corresponde a  $-2\Delta(\ln \mathcal{L}) < 1$ , “un sigma”, cubre 68 % de los resultados que se obtendrían repitiendo el experimento ATLAS numerosas veces;
- ▶ idem para el rango de  $m_H$  correspondiendo a  $-2\Delta(\ln \mathcal{L}) < 4, 9, \dots$  (“dos sigma”, “tres sigma”, etc...)



En experimentos de conteo de eventos, el número de eventos observados puede ser un parámetro de interés. Para el caso de una especie única, esto corresponde a “extender” la verosimilitud,

$$\mathcal{L}(\lambda, \theta; \vec{x}) = \frac{\lambda^N e^{-\lambda}}{N!} \prod_{i=1}^N P(\vec{x}_i; \theta) .$$

dónde el término multiplicativo adicional, corresponde a la distribución de Poisson (el término  $N!$  en el denominador es irrelevante; es un factor global sin impacto sobre la forma de la verosimilitud)

Es fácil verificar que la verosimilitud es máxima cuando  $\hat{\lambda} = N$ , tal como se espera; ahora, si algunas de las PDFs dependen también de  $\lambda$ , el valor  $\hat{\lambda}$  que maximiza  $\mathcal{L}$  puede diferir.

La generalización a más de una especie es sencilla; para cada especie, un término multiplicativo de Poisson se incluye en la verosimilitud extendida, y las PDFs de cada especie son ponderadas por su fracción relativa de eventos.

Para el caso de dos especies, la versión extendida de la verosimilitud es

$$\mathcal{L}(N_{\text{sig}}, N_{\text{bkg}}; \theta; \vec{x}) = (N_{\text{sig}} + N_{\text{bkg}})^N e^{-(N_{\text{sig}} + N_{\text{bkg}})} \prod_{i=1}^N [N_{\text{sig}} P_{\text{sig}}(\vec{x}; \theta) + N_{\text{bkg}} P_{\text{bkg}}(\vec{x}; \theta)] .$$



# Estimar eficiencias a partir de ajustes MLE

Consideremos de nuevo el caso de un proceso aleatorio con dos resultados posibles: "yes" y "no". El estimador intuitivo de la eficiencia  $\varepsilon$  es el cociente entre el número de realizaciones de cada tipo,  $n_{\text{yes}}$  y  $n_{\text{no}}$ , y su varianza  $V[\hat{\varepsilon}]$  viene dada por :

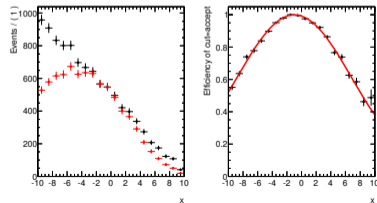
$$\hat{\varepsilon} = \frac{n_{\text{yes}}}{n_{\text{yes}} + n_{\text{no}}} , \quad V[\hat{\varepsilon}] = \frac{\hat{\varepsilon}(1 - \hat{\varepsilon})}{n} ,$$

donde  $n = n_{\text{yes}} + n_{\text{no}}$  es el número total de realizaciones. (ejercicio: reproducir este resultado)  
Este estimador ingenuo  $\hat{\varepsilon}$  claramente falla para pequeños valores de  $n$ , y en las situaciones de gran (in-)eficiencia. La técnica del MLE provee una solución robusta para la estimación de eficiencias: si la muestra contiene una variable  $x$  sensible a la eficiencia (es decir  $\varepsilon(x)$ ), que sigue la PDF  $P(x; \theta)$ , entonces la inclusión de una nueva variable aleatoria discreta y bivariada  $c = \{\text{yes}, \text{no}\}$ , nos da un modelo más elaborado:

$$P(x, c; \theta) = \delta(c - \text{yes})\varepsilon(x, \theta) + \delta(c - \text{no}) [1 - \varepsilon(x, \theta)] .$$

- ▶ la función  $\varepsilon(x)$  ha sido correctamente normalizada, para ser también una PDF
- ▶ la eficiencia ya no es un valor único, sino una función de  $x$
- ▶ (y de otros parámetros  $\theta$  que sean necesarios para caracterizar su forma)

`($ROOTSYS/tutorials/roofit/rf701_efficiencyfit.C)`





En el MLE, la matriz de covariancia es el estimador de las incertidumbres estadísticas.

Pero otras fuentes de incertidumbre contribuyen también a diluir la precisión de una medida. En lenguaje de física, estas a menudo se llaman “errores sistemáticos”. Para discutir de estas incertidumbres en el contexto MLE, modificamos ligeramente la notación, y reescribimos la función de verosimilitud como sigue:

$$\mathcal{L}(\mu_1, \dots, \mu_p, \theta_1, \dots, \theta_k; \vec{x}) ,$$

donde hemos explícitamente separado el conjunto de parámetros de  $\mathcal{L}$  en dos subconjuntos:

- ▶ los *parámetros de interés*  $\mu_1, \dots, \mu_p$ , (POI en inglés), que corresponden a las cantidades que deseamos estimar;
- ▶ los *parámetros de molestia*  $\theta_1, \dots, \theta_k$  (NP en inglés, por nuisance parameters) que representan fuentes potenciales de sesgos sistemáticos;

La lógica es la siguiente: si le asignásemos valores imprecisos o errados a algunos NPs, la formas de las PDFs resultantes pueden distorsionarse, y los estimadores de los POIs pueden estar sesgados.

A menudo, la estimación de las incertidumbres sistemáticas se lleva parte muy importante del trabajo de análisis de datos: energía, tiempo... La calidad de un resultado científico releva en gran parte de la calidad de la estimación de los sistemáticos.





## Incertidumbres sistemáticas (II)

Las incertidumbres sistemáticas debidas a los NPs se clasifican a menudo en dos categorías:

- ▶ Errores de “Tipo-I” : nuestra muestra (u otras muestras de control) pueden en principio proveer información sobre el NP considerado, y la incertidumbre proveniente de este NP debe en principio disminuir con el tamaño de las muestras utilizadas.
- ▶ Errores de “Tipo-II”, que provienen de suposiciones incorrectas del modelo (p.e. el uso de funciones inadecuadas para las PDFs), u aspectos mal controlados en los datos, como cambios en las condiciones de adquisición de las muestras, o la presencia de especies no tomadas en cuenta.

Algunos comentarios:

- ▶ Ciertos sistemáticos son una mezcla de ambos tipos: por ejemplo, si las muestras de control para controlar un NP no son completamente representativas de las propiedades de ese NP.
- ▶ Un NP enteramente de Tipo-I es una variable aleatoria, y por tanto se pueden estimar intervalos de confianza usando los métodos antes descritos.
- ▶ Claramente, los errores de Tipo-II son los más difíciles de manejar correctamente, y no siempre hay un procedimiento bien definido y consensual para estimar el impacto de esos errores. En los peores casos, no siempre se puede saber de qué manera un error de Tipo-II impacta los intervalos de confianza sobre los POIs.
- ▶ Frente a los casos ambiguos, hay cierto consenso (parcial) en que es mejor tener estimaciones “conservadoras” de los sistemáticos, en particular cuando contribuyen de manera subdominante al “error budget”. Pero en esas situaciones, las discusiones suelen ser controversiales, incluso amargas...



## Incertidumbres sistemáticas (III)

El método de *profile-likelihood*, (verosimilitud de perfil?) permite manejar de manera elegante los sistemáticos de Tipo-I.

Consiste en asignar una verosimilitud específica a los NPs “profileables”, de tal manera que la  $\mathcal{L}$  original se modifica en dos componentes:

$$\mathcal{L}(\mu, \theta) = \mathcal{L}_\mu(\mu, \theta) \mathcal{L}_\theta(\theta) .$$

Entonces, para un valor fijo de  $\mu$ ,  $\mathcal{L}$  es maximizada con respecto al NP  $\theta$ . Si se recorre una secuencia de valores de  $\mu$ ,  $\mathcal{L}$  es una función que solamente depende de  $\mu$ ; si dice que la molestia ha sido “profiled-out”.

Ejemplo, consideremos la medida de una sección eficaz de un proceso  $\sigma$  (initial  $\rightarrow$  final). Si solamente una fracción de los procesos de ese tipo son detectados (p.e. debido a efectos de *aceptancia geométrica* del detector, u otras fuentes de ineficiencia), se requiere por tanto conocer la eficiencia de reconstrucción  $\varepsilon$  para convertir el número observado de procesos  $\hat{N}_{\text{evt}}$  en una medida de  $\hat{\sigma}$ . La eficiencia es claramente un parámetro de molestia: un valor incorrecto de  $\varepsilon$  distorsiona directamente la medida de  $\hat{\sigma}$ , independientemente de la precisión con la cual se haya medido  $\hat{N}_{\text{evt}}$ .

Si  $\hat{\varepsilon}$  puede estimarse sobre una muestra de control de calidad (por ejemplo, una simulación detallada con copiosa estadística, o una muestra de control de alta pureza), el impacto de la molestia se atenúa. Un análisis *elegant* produciría un ajuste simultáneo a las muestra de señal y control, de tal manera que los valores y las incertidumbres de los NPs se extraen directamente de la covariancia del ajuste ML, lo cual asegura que las correlaciones con los POIs se propagan correctamente, y que los intervalos de confianza heredan de todas esas relaciones entre POIs y NPs.



Otro ejemplo de “profile-likelihood” :

Consideremos la búsqueda de una resonancia (un “bump”) sobre un fondo uniforme. Si la fracción de señal es muy pequeña, la anchura  $\Gamma$  del bump no puede estimarse directamente sobre la muestra, así que el valor de anchura a utilizar en la PDF de señal debe provenir de fuentes externas (p.e. una simulación detallada de la función de resolución del detector).

Esa anchura es claramente un parámetro de molestia: Si su valor se sobreestima, eso se traduciría en una subestimación del cociente señal/fondo, y por tanto en un aumento de la varianza de los POIs de la señal, así como tal vez posibles sesgos en sus valores centrales (p.e. si el fondo es asimétrico, el impacto será diferente a ambos lados del pico de la señal).

Si las muestras de control permiten acceder a una estimación independiente de la anchura  $\hat{\Gamma} \pm \hat{\sigma}_{\Gamma}$  de la anchura del pico de señal, esta información puede implementarse utilizando una PDF Guassiana, de media  $\hat{\Gamma}$  y de anchura  $\hat{\sigma}_{\Gamma}$ , en la componente  $\mathcal{L}_{\Gamma}$  de la verosimilitud. Este término adicional cumple el rol de una penalidad en el MLE, y por lo tanto restringe los valores de  $\Gamma$  dentro del intervalo dado por  $\pm \hat{\sigma}_{\Gamma}$ , y de ese manera su impacto en los POIs se ve controlado.



# Incertidumbres sistemáticas (V)

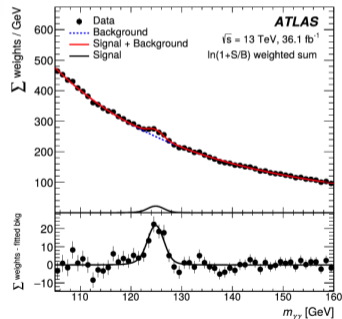
Un ejemplo práctico del “profile-likelihood” : el análisis del bosón de Higgs en el canal difotón en el LHC. Este canal se caracteriza por un cociente  $S/B$  extremadamente bajo.

Los parámetros de interés son :

- ▶ la masa del Higgs,
- ▶ el producto  $\sigma(pp \rightarrow H + X) \times \mathcal{BR}(H \rightarrow \gamma\gamma)$  de la sección eficaz de producción del Higgs en colisiones protón-protón, multiplicada por la tasa de decaimiento del Higgs en el canal difotón.

Mientras que la (larga) lista de parámetros de molestia incluye, entre otros:

- ▶ La eficiencia  $\varepsilon(H \rightarrow \gamma\gamma)$  de reconstrucción de la señal, es decir el cociente entre el número de eventos  $H(\rightarrow \gamma\gamma)_{\text{RECO}}$  registrados, reconstruidos, identificados y seleccionados, y el número de eventos  $(H \rightarrow \gamma\gamma)_{\text{TRUE}}$  realmente producidos,
- ▶ Los parámetros necesarios para caracterizar la resolución en masa invariante difotón, es decir la distribución de la diferencia entre la masa reconstruida  $m(\gamma\gamma)_{\text{RECO}}$  y la verdadera masa  $m(\gamma\gamma)_{\text{TRUE}}$ ,
- ▶ Los parámetros necesarios para caracterizar la forma del fondo. Para simplificar, supondremos que una función uni-paramétrica es suficiente (p.e. una función exponencial decreciente)





## Contraste de hipótesis

Las discusiones pasadas se centraron mayoritariamente en extraer información numérica a partir de muestras de datos: efectuar mediciones de parámetros, y reportar el resultado de esas mediciones bajo forma de

- ▶ valores centrales e incertidumbres, en particular cuando se trata de un sólo parámetro de interés;
- ▶ matrices de covarianza, en particular para medidas simultáneas de varios parámetros, y si las correlaciones no pueden ser despreciadas;
- ▶ perfiles completos de verosimilitud, cuando la aproximación “parabólica” no es suficiente;

La etapa siguiente en un análisis es producir información cualitativa a partir de los datos disponibles: se habla entonces de efectuar un *contraste estadístico de hipótesis* (hypothesis testing en inglés).

La herramienta cuantitativa para declarar el acuerdo entre una hipótesis y las observaciones (los datos) se llama un *estadístico de prueba*. El resultado de una prueba se da en términos de una “*p-value*”: la probabilidad, bajo la hipótesis en consideración, de observar un estadístico de prueba similar o “peor” que el observado en la muestra. En términos intuitivos: suponemos que los datos son una realización aleatoria de la hipótesis sometida a prueba, y comparamos cuantitativamente el estadístico observado sobre los datos con el ensamble de realizaciones aleatorias de estadísticos. Es lo que se llama la *interpretación frecuentista* de la estadística:

- ▶ “dada la hipótesis, ¿cuál es la probabilidad de observar una muestra de datos como la que observé?”

Existe una interpretación diferente de la estadística, llamada *bayesiana*, que busca resolver el *problema inverso*:

- ▶ “dada mi muestra de datos, ¿cuál es la probabilidad que mi hipótesis sea verdadera?”

Las diferencias entre ambas interpretaciones son profundas, y tocan a la esencia del proyecto de inferencia científica. Ahora, los debates entre adeptos de una u otra interpretación son en ocasiones amargos y poco estimulantes...



## El test del $\chi^2$

Dado un conjunto de  $n$  medidas independientes  $x_i$  con varianzas  $\sigma_i^2$ , y un conjunto de predicciones  $\mu_i$ , se define un test estadístico llamado  $\chi^2$  de la manera siguiente:

$$\chi^2 = \sum_{i=1}^n \frac{(x_i - \mu_i)^2}{\sigma_i^2} .$$

El test de  $\chi^2$  es una variable aleatoria que, como vimos previamente, sigue la distribución de  $P_{\chi^2}(x; n)$  para  $n$  *grados de libertad*. Su valor de expectación es  $n$  y su varianza  $2n$ .

Por ello uno espera que el valor de  $\chi^2$  observado sobre una muestra no debe alejarse mucho del número de grados de libertad, y por tanto ese valor permite sondear el acuerdo entre la observación y la predicción.

Para ser más preciso, se espera que 68% de los tests se encuentren contenidos dentro de un intervalo  $n \pm \sqrt{2n}$ .

La  $p$ -value, o probabilidad de observar un test con valores mayores viene dada por

$$p = \int_{\chi^2}^{+\infty} dx P_{\chi^2}(x; n) .$$

Intuitivamente hablando, uno debe sospechar de pruebas que arrojen pequeñas  $p$ -values, dado que esas podrían indicar un problema. Este puede provenir del lado de las predicciones, o reflejar la calidad de los datos registrados, o ser solamente resultado de la “mala suerte”.

La interpretación de las  $p$ -values (i.e. para decidir qué valores son demasiado pequeños o suficientemente grandes) es un tópico importante, que requiere un marco de análisis que apunte a separar las partes objetivas y subjetivas.



# Propiedades generales del contraste de hipótesis (I)

Consideremos dos hipótesis mutuamente excluyentes,  $\mathcal{H}_0$  y  $\mathcal{H}_1$ , que proveen ambas una descripción de un proceso, y del cual extraemos una muestra de datos. El procedimiento de contraste de hipótesis apunta a evaluar:

- ▶ cuán robusta es la *hipótesis nula*  $\mathcal{H}_0$ , en describir los datos, y
- ▶ cuán incompatible con esos mismos datos es la *hipótesis alternativa*  $\mathcal{H}_1$ .

En lenguaje de física de altas energías, un ejemplo común concierne la búsqueda de una señal (aún) desconocida, que implica dos casos:

- ▶ en la *lógica de descubrimiento*, la hipótesis nula corresponde al escenario *background-only*, mientras que la hipótesis alternativa sería *signal-plus-background*: “¿cuán probable es que una fluctuación del fondo explique el exceso que estoy observando?”;
- ▶ en la *lógica de exclusión*, las dos hipótesis se invierten: “¿cuál es la cantidad máxima de señal compatible con mi observación?”.

La dinámica del contraste de hipótesis se puede resumir así:

- ▶ construir un estadístico de prueba  $q$ , una función que reduce una muestra a un valor numérico único;
- ▶ definir un intervalo de confianza  $W \rightarrow [q_{lo} : q_{hi}]$ ;
- ▶ medir  $\hat{q}$  sobre la muestra en estudio;
- ▶ si  $\hat{q}$  está contenido en el intervalo  $W$ , se declara que la hipótesis nula es aceptada, y rechazada en caso contrario.



## Propiedades generales del contraste de hipótesis (II)

Para caracterizar el resultado de la secuencia antes descrita, se definen dos criterios:

- ▶ se incurre en un “Error de Tipo-I” si  $\mathcal{H}_0$  es rechazada aún siendo cierta;
- ▶ se incurre en un “Error de Tipo-II” si  $\mathcal{H}_0$  es aceptada aún siendo falsa.

Las tasa de errores de Tipo-I y Tipo-II se llaman usualmente  $\alpha$  y  $\beta$  respectivamente, y se determinan por integración de las densidades de probabilidad asociadas a las hipótesis  $\mathcal{H}_0$  y  $\mathcal{H}_1$  sobre el intervalo  $W$ :

$$1 - \alpha = \int_W dq \mathcal{P}(q|H_0) ,$$
$$\beta = \int_W dq \mathcal{P}(q|H_1) .$$

La tasa  $\alpha$  es llamada *tamaño del contraste* (*size of the test* en inglés), puesto que fijar  $\alpha$  determina el tamaño del intervalo  $W$ . De manera análoga,  $1 - \beta$  es llamado *potencia del contraste* (*power*).

Juntos, tamaño y potencia caracterizan el performance de un estadístico: el lema de Neyman-Pearson afirma que a *tamaño fijo*, el estadístico óptimo viene dado por el cociente de verosimilitudes  $q_\lambda$ :

$$q_\lambda(\text{data}) = \frac{\mathcal{L}(\text{data}|H_0)}{\mathcal{L}(\text{data}|H_1)} .$$

En la práctica, las distribuciones de  $\mathcal{H}_0$  y  $\mathcal{H}_1$  son obtenidas a partir de simulaciones o muestras de control. De esas distribuciones se se obtiene la distribución esperada del estadístico  $q_\lambda$ , y la  $p$ -value observada se obtiene al integrarla con respecto al valor observado de  $q_\lambda$ .





## Propiedades generales del contraste de hipótesis (III)

La significación del estadístico viene dada por su  $p$ -value,

$$p = \int_{\hat{q}}^{+\infty} dq \mathcal{P}(q|H_0) .$$

que es a menudo reportado en unidades de “sigmas”,

$$p = \int_{n\sigma}^{+\infty} dz \frac{1}{\sqrt{2\pi}} e^{-z^2/2} = 1 - \frac{1}{2} \operatorname{erf} \left( \frac{n}{\sqrt{2}} \right) ,$$

de manera que por ejemplo una  $p$ -value  $p < 0,0228$  se reporta como un “efecto a dos sigma”. Igualmente común es reportar la  $p$ -value bajo forma de un intervalo de confianza (C.L.).

Esta definición de la  $p$ -value es clara y sin ambigüedades. Pero su interpretación es parcialmente subjetiva: la conveniencia de un *umbral de tolerancia* puede depender del tipo de hipótesis bajo prueba, o de h'abitos en cada disciplina. En física de altas energías:

- ▶ en la *lógica de exclusión*, el umbral se sitúa a 95 % C.L. para declarar la exclusión de la hipótesis de *señal-más-fondo*;
- ▶ en la *lógica de descubrimiento*, el umbral se sitúa a tres sigma ( $p < 1,35 \times 10^{-3}$ ) sobre la hipótesis *solamente-fondo* para afirmar que hay “evidencia”;
- ▶ y un umbral a cinco sigma ( $p < 2,87 \times 10^{-7}$ ) es requerido para afirmar que hay “observación”.

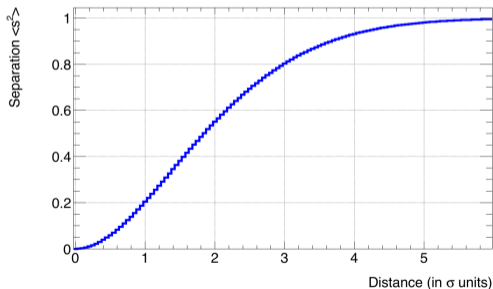


# Separación

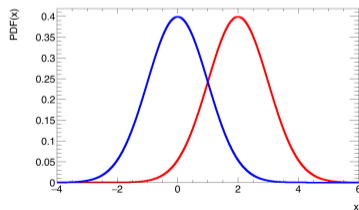
Las definiciones generales de *tamaño* y *potencia* pueden ser complementadas con otras definiciones más sencillas. Por ejemplo la "separación", que se determina a partir de las PDFs de las dos especies (señal  $S$  y fondo  $B$ ) para las cuales se quiere cuantificar el poder de discriminación:

$$\langle s^2 \rangle = \frac{1}{2} \int d\vec{x} \frac{[S(\vec{x}) - B(\vec{x})]^2}{S(\vec{x}) + B(\vec{x})} .$$

Por construcción,  $0 \leq \langle s^2 \rangle \leq 1$ , valores límites que corresponden a los casos extremos en que el poder de discriminación es nulo (señal y fondo son imposibles de distinguir) o total (no hay ningún solapamiento entre las dos especies).



La gráfica a la izquierda muestra la separación entre dos Gaussianas de misma anchura  $\sigma$ , en función de la distancia entre sus picos (en unidades de  $\sigma$ ). Ejemplo para una distancia de  $2\sigma$ :



(a menudo se habla de "separación a  $n\sigma$ ")

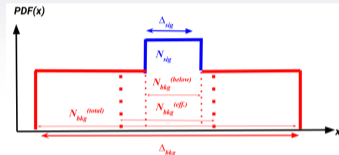


# Sobre el fondo efectivo

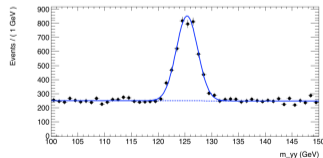
De manera general, la zona que contiene la señal también está contaminada por fondo(s). El impacto de esos fondos se traduce en una *dilución* o disminución de la precisión en la medida de los parámetros de interés.

Ejercicio : Consideremos un escenario compuesto por dos especies: una señal, distribuida de manera uniforme en un intervalo  $\Delta(\text{sig})$ , y un fondo distribuido de manera uniforme en un intervalo  $\Delta(\text{bkg})$  más amplio que cubre ambos lados del intervalo de señal ("sidebands").

- ▶ Generar una realización aleatoria, eligiendo valores particulares de los intervalos  $\Delta(\text{sig})$  y  $\Delta(\text{bkg})$ , y del número de eventos de señal y fondo  $N_{\text{sig}}$  y  $N_{\text{bkg}}$ .
- ▶ El parámetro de interés es el número de eventos de señal  $N_{\text{sig}}$ . Estimar su valor y error  $\hat{N}_s \pm \hat{\sigma}_s$  por verosimilitud máxima.
- ▶ Si el fondo fuera nulo, se tendría  $\hat{\sigma}_s = \sqrt{N_s}$ .
- ▶ Definir el "fondo efectivo" como el causante del aumento en el error,  $\hat{\sigma}_s = \sqrt{N_s + N_{\text{bkg}}^{\text{eff}}}$ .
- ▶ Comparar  $N_{\text{bkg}}^{\text{eff}}$  al fondo "debajo" de la señal,  $N_{\text{bkg}}^{\text{below}}$ .
- ▶ Repitiendo el ejercicio para diferentes valores de  $N_{\text{sig}}$ , verificar que el fondo efectivo es siempre el mismo.
- ▶ Repitiendo el ejercicio para varios valores crecientes de  $\Delta(\text{bkg})$ , verificar que el fondo efectivo tiende a  $N_{\text{bkg}}^{\text{below}}$ .
- ▶ Interpretar.



Ejercicio: Para una señal Gaussiana, realizar un ejercicio similar, con dos parámetros de interés adicionales: la posición del pico y su anchura.

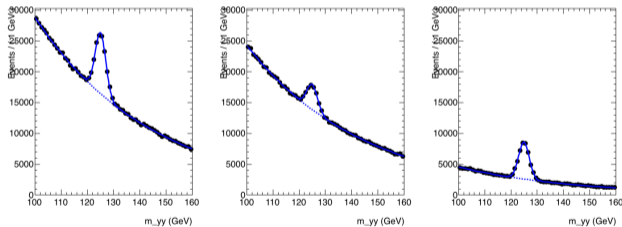




En ocasiones, la muestra de análisis puede descomponerse en dos o más submuestras (“categorías”), cada una de ellas con sus propias PDFs y purezas.

Cuando las características de cada especie son razonablemente diferentes, puede ser de interés descomponer la función de verosimilitud de tal manera que cada categoría utilice sus propias PDFs y purezas. El resultado combinado sobre un parámetro común de interés tendrá una significación superior a la que se obtendría de un análisis “inclusivo”, es decir usando PDFs y purezas promedio sobre la muestra completa.

Ejemplo sencillo (y ejercicio): supongamos que la muestra  $H \rightarrow \gamma\gamma$  se descompone en dos categorías: una “limpia” con excelente cociente señal/fondo, y una “sucias” en la que el fondo es ampliamente dominante. Si el parámetro de interés es la masa del Higgs, la ventaja de realizar un análisis en categorías puede ser significativo.



Aquí las dos categorías difieren en el cociente señal/fondo, grande para la “limpia”, y pequeño para la “sucias”. Otra posibilidad es aprovechar diferencias en resolución.

Error en el análisis inclusivo:  $\sigma(m_H) = \pm 2,25\%$

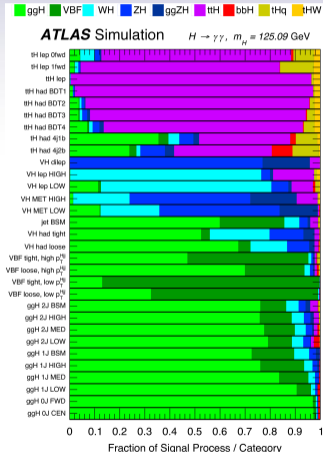
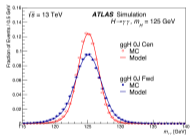
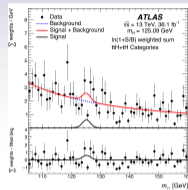
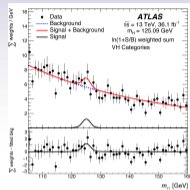
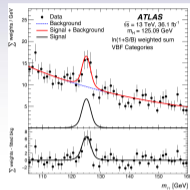
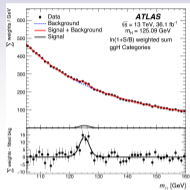
Errores en las categorías sucias y limpia :  $\pm 3,70\%$  y  $\pm 1,85\%$ , respectivamente.

Error en la combinación de ambas categorías:  $\pm 1,72\%$  ! Equivale a aumentar en 70 % la estadística inclusiva !



# El uso de categorías en el análisis $H \rightarrow \gamma\gamma$

ATLAS separa su muestra de candidatos difotón en 31 categorías, definidas en función de varios criterios: el modo de producción del Higgs, diferencias en la resolución experimental en masa, diferencias en la pureza.



En el 2012, la categorización fue crucial para alcanzar los  $5 \sigma$  de significancia en la observación del Higgs...



## Algunas convenciones estadísticas en física de partículas: LEP

En la HEP experimental, hay tradición de definir por consenso la elección de los estadísticos de prueba, para simplificar las combinaciones de resultados de diferentes experimentos, de manera que las componentes relacionadas con los detectores (específicas a cada experimento) se factoricen con respecto a los observables físicos (que son en principio universales).

Ejemplo: en el contexto de la búsqueda del Bosón de Higgs del Modelo Estándar, los cuatro experimentos en el LEP (ALEPH, DELPHI, OPAL, L3) decidieron describir sus datos utilizando las siguientes verosimilitudes:

$$\mathcal{L}(H_1) = \prod_{a=1}^{N_{\text{ch}}} \mathcal{P}_{\text{Poisson}}(n_a, s_a + b_a) \prod_{j=1}^{n_a} \frac{s_a \mathcal{S}_a(\vec{x}_j) + b_a \mathcal{B}_a(\vec{x}_j)}{s_a + b_a},$$
$$\mathcal{L}(H_0) = \prod_{a=1}^{N_{\text{ch}}} \mathcal{P}_{\text{Poisson}}(n_a, b_a) \prod_{j=1}^{n_a} \mathcal{B}_a(\vec{x}_j).$$

donde  $N_{\text{ch}}$  es el número de “canales de decaimiento” del Higgs estudiados,  $n_a$  es el número observado de candidatos en cada canal  $a$ ,  $\mathcal{S}_a$  y  $s_a$  ( $\mathcal{B}_a$  y  $b_a$ ) son las PDFs y los números de eventos para las especies de señal (fondo) de cada canal. El estadístico de prueba  $\lambda$ , derivado de un cociente de verosimilitudes, es

$$\lambda = -2 \ln Q, \text{ con } Q = \frac{\mathcal{L}(H_1)}{\mathcal{L}(H_0)};$$

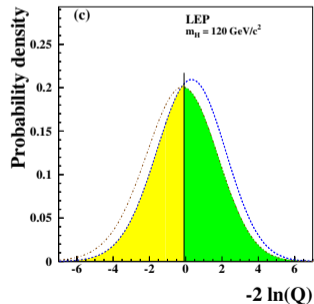
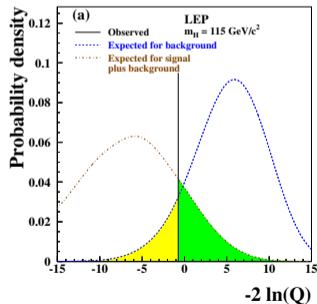
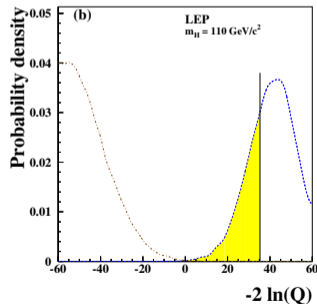
de manera que valores positivos de  $\lambda$  favorezcan un escenario “background-like”, y valores negativos estén más tono con un escenario “señal-más-fondo”; valores cercanos a cero indicando una sensibilidad pobre para distinguir entre ambos.

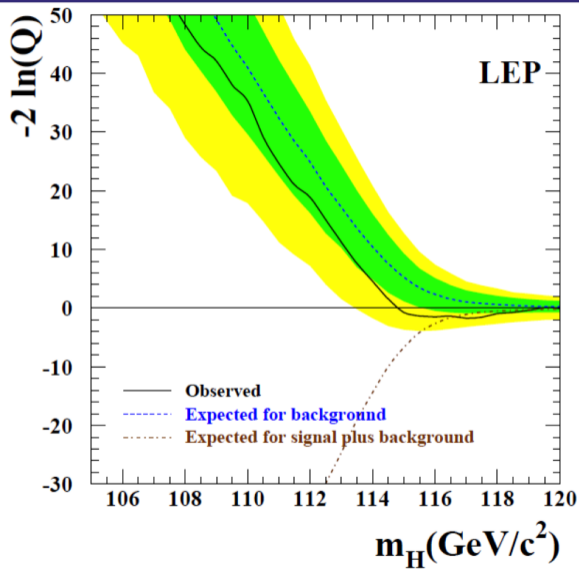


# Algunas convenciones estadísticas en física de partículas: LEP

- ▶ Bajo la hipótesis “background-only”,  $CL(b)$  es la probabilidad de tener  $-2 \ln Q$  más pequeño que el observado (amarillo);
- ▶ bajo la hipótesis “señal-más-fondo”,  $CL(s + b)$  es la probabilidad de tener  $-2 \ln Q$  más grande que el observado (verde).

Las figuras bajo muestra, para tres hipótesis diferentes de la masa del Higgs, los valores de  $-2 \ln Q$  obtenidos al combinar los resultados de los cuatro experimentos LEP. También se muestran las distribuciones de  $CL(s + b)$  y  $1 - CL(b)$ .









## El estimador modificado $CL(s)$

Tomemos un experimento de conteo de eventos. Esperamos 10 eventos de tipo fondo, y 5 eventos de tipo señal. Pero observamos 8 eventos en total... lo más probable es que tanto la señal como el fondo hayan experimentado una fluctuación negativa... pero en la interpretación estándar, se hubiera asignado una exclusión a 95 % C.L. ! (ello incluso cuando el experimento no tiene sensibilidad alguna a la señal)  
Para evitar esta situación, se define un nivel de confianza modificado,  $CL(s)$ , definido como

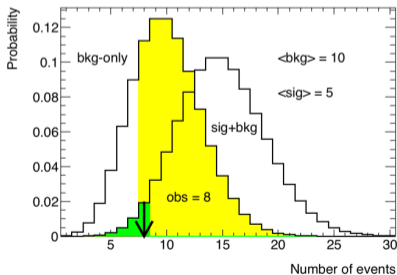
$$CL(s) = \frac{CL(s + b)}{1 - CL(b)}$$

que si bien *stricto sensu* no es una  $p$ -value (un cociente de probabilidades no es una probabilidad) por lo menos tiene la propiedad de proteger contra fluctuaciones negativas del fondo.

(cuidado con las convenciones de color amarillo/verde...)

- ▶  $CL(s + b) = 3,7 \%$
- ▶  $1 - CL(b) = 33 \%$
- ▶  $CL(s) = 11 \%$

No sin controversia, el estimador  $CL(s)$  ha sido sin embargo adoptado por varias colaboraciones internacionales, incluyendo los experimentos del Tevatron (Fermilab) y ATLAS y CMS en el LHC.





## Los cocientes de perfil de verosimilitud

Los experimentos ATLAS y CMS usan un estadístico de prueba, llamado cociente de perfil de verosimilitud (profiled likelihood ratio) definido así:

$$\tilde{q}_\mu(\mu) = -2 \ln \frac{\mathcal{L}(\mu, \hat{\hat{\theta}})}{\mathcal{L}(\hat{\mu}, \hat{\theta})}, \text{ con } 0 \leq \hat{\mu} \leq \mu,$$

Aquí el parámetro de interés es  $\mu = \sigma/\sigma_{\text{SM}}$ , la “intensidad de señal”, que es la tasa de conteo de candidatos de señal comparada a la predicción (p.e. la sección eficaz de producción del Higgs vs. la predicción del Modelo Estándar),  $\hat{\theta}$  son los valores de los NPs obtenidos en un ajuste a intensidad de señal  $\mu$  fija,  $\hat{\mu}$  y  $\hat{\theta}$  son los valores ajustados cuando tanto  $\mu$  como los NPs son libres en el ajuste. (el límite inferior en  $0 \leq \hat{\mu} \leq \mu$  es para asegurar tener una intensidad de señal positiva, y el límite superior es para evitar que una fluctuación hacia arriba no desfavorezca la hipótesis de señal).

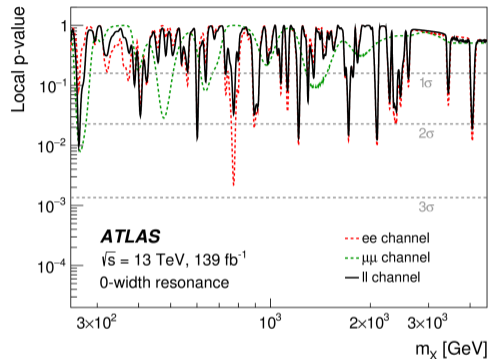
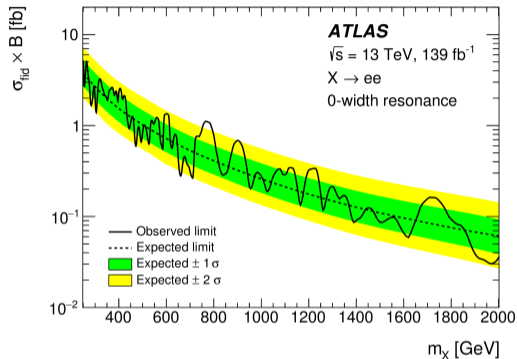
Para un valor observado del estadístico  $\hat{q}_\mu$ , las  $p$ -values asignadas a las hipótesis de signal-plus-background y background-only,  $p(s+b)$  y  $p(b)$ , son

$$p(s+b) = \int_{\hat{q}_\mu}^{\infty} dq P(q; \mu = \hat{\mu}, \hat{\theta}), \quad 1 - p(b) = \int_{\hat{q}_\mu}^{\infty} dq P(q; \mu = 0, \hat{\theta}).$$

Los resultados de exclusión se muestran bajo forma de un “Brazil-plot”, y los de observación bajo forma de un “local- $p_0$ -plot”.

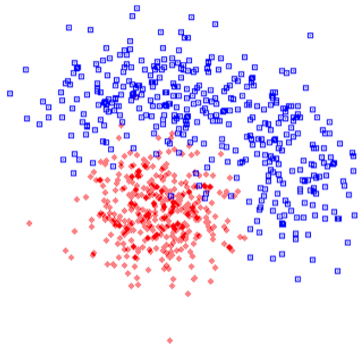


# El "look-elsewhere-effect"



Cuando varias "regiones" estadísticamente independientes son probadas en un mismo análisis (aquí la búsqueda de otras posibles resonancias en dileptón), hay que tomar en cuenta el Look-Elsewhere-Effect (o trials factors en la literatura) para evaluar la  $p$ -value...

(Nota: a alta energía la resolución en impulso es superior para los electrones que para los muones, por eso hay más "regiones" en el canal  $X \rightarrow e^+e^-$  que para el canal  $X \rightarrow \mu^+\mu^-$ ...)

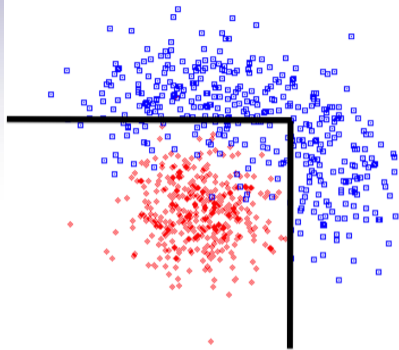


Puntos azules : muestra de control finita,  
distribuída como el fondo  
Puntos rojos : muestra de control finita,  
distribuída como la señal

A menudo, hay grandes regiones del espacio de la muestra en las cuales los fondos son ampliamente dominantes, o donde la densidad de la señal es nula o despreciable.

Si se reduce la muestra a subconjuntos “enriquecidos en señal” del espacio completo, la pérdida de información puede ser mínima, y otras ventajas pueden compensar esas posibles pérdidas:

- ▶ para muestras multidimensionales, puede ser difícil caracterizar las formas de las PDFs en las regiones de baja densidad de eventos
- ▶ el reducir el tamaño de la muestra puede aliciar el consumo de memoria y CPU en las partes numéricas del análisis (p.e. la minimización)



Las líneas negras indican una selección “cut-based”, definida de manera de conservar cerca de 100 % de la señal.

(para ser más precisos, la selección es 100 % eficaz sobre la muestra de control de señal, pero la caracterización precisa de la eficiencia de selección requiere de estimar la densidad de probabilidad de la señal fuera de la región seleccionada (y para caracterizar con precisión el nivel de fondo se requiere estimar la densidad de probabilidad del fondo dentro de ella)

El método más sencillo de reducción de una muestra es restringiendo las variables, una a una, a intervalos finitos. En la práctica, esas selecciones “cut-based” aparecen a varios niveles de la definición del espacio de muestreo: umbrales en las decisiones en línea (triggers), filtros a varios niveles posteriores del proceso de adquisición, eliminación de datos a partir de criterios de calidad...

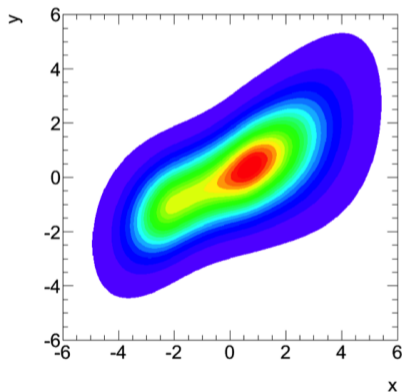
Pero ya en etapas más avanzadas del análisis de datos, esas selecciones “accept-reject” conviene ser reemplazadas por procedimientos más sofisticados. Estos son llamados de manera genérica *técnicas multivariadas*.



## Análisis multidimensional (III)

Consideremos un conjunto de  $n$  variables aleatorias  $\vec{x} = \{x_1, x_2 \dots, x_n\}$ . Si todas las variables son no-correlacionadas, Las PDFs  $n$ -dimensionales están completamente determinadas por el producto directo de sus  $n$  PDFs uni-dimensionales.

si algunas de las variables están correlacionadas, y si sus patrones de correlación son completamente lineales, es posible definir un nuevo conjunto de variables  $\vec{y}$ , que son combinaciones lineales de  $\vec{x}$ , obtenidas diagonalizando el inverso de la matriz de covarianza.



En ocasiones, cuando los patrones de correlación son no-lineales, es tal vez posible en algunos casos definir una descripción analítica: por ejemplo, el perfil de correlación que incluye una (ligera) componente no lineal representado aquí, fue producido con el paquete RooFit aplicando la opción `Conditional` en RooProdPdf para producir un producto de PDFs: la anchura de  $y$  varía de manera no-lineal con  $x$ .

En la práctica, ésta solución elegante no puede extenderse fácilmente a más de dos dimensiones, y no hay garantía que se puedan reproducir patrones no lineales complicados.

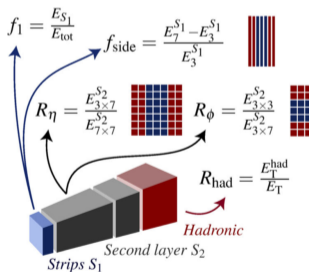
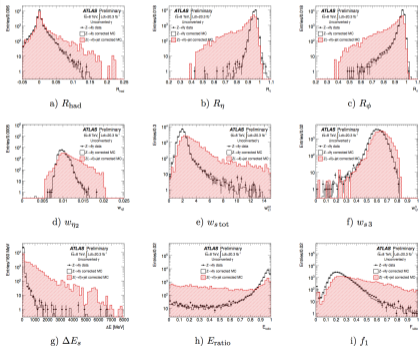
Frente a tales escenarios, un protocolo de *reducción dimensional* puede a menudo producir resultados más efectivos que un intento de descripción analítica de correlaciones.



# Reducción dimensional

Un escenario típico para considerar la reducción dimensional es cuando varias variables arrastran en gran parte información común (y por tanto exhiben correlaciones fuertes), pero contienen también algunos elementos (diluídos pero importantes) de información independiente.

Un ejemplo: la caracterización de cascadas en calorímetros segmentados. Las señales depositadas en celdas vecinas están fuertemente correlacionadas (tienen un origen común), pero permiten reconstruir detalles precisos del desarrollo de la cascada. Esas correlaciones son utilizadas para caracterizar las “formas de cascadas”, y permiten discriminar (por ejemplo) entre cascadas electromagnéticas y hadrónicas.



$$w_{\eta 2} = \sqrt{\frac{\sum E_i \eta_i^2}{\sum E_i} - \left(\frac{\sum E_i \eta_i}{\sum E_i}\right)^2}$$

width in a  $3 \times 5$  ( $\Delta\eta \times \Delta\phi$ ) region of cells in  $S_2$

$$w_s = \sqrt{\frac{\sum E_i (i - i_{max})^2}{\sum E_i}}$$

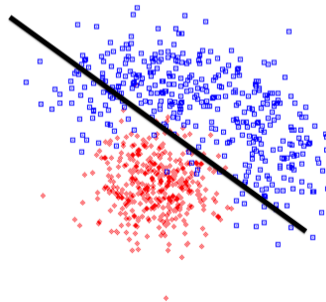
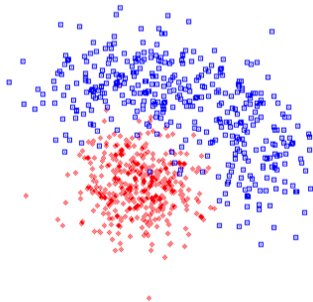
$w_{s3}$  uses  $3 \times 2$  strips ( $\eta \times \phi$ )

$w_{s\ tot}$  is defined similarly but uses  $20 \times 2$  strips



## Discriminantes lineales

El algoritmo más sencillo de reducción dimensional es el discriminante de Fisher: es una función lineal de las variables, con coeficientes definidos a partir de un criterio de optimización de la separación entre especies, que viene dado por el cociente de las varianzas *inter-especies* y las varianzas *intra-especies*. Fisher es un caso particular de los llamados PCA (principal component analysis). Un análisis MLE es también un PCA: reduce un problema multidimensional a un problema en 1 dimensión: el comportamiento del estadístico de test  $\lambda = \frac{\mathcal{L}(H_1)}{\mathcal{L}(H_0)}$ . Por construcción, los discriminantes lineales son óptimos para variables multinormales, por tanto con correlaciones perfectamente lineales.







## El discriminante de Fisher

Tenemos dos especies, “señal” y “fondo”, que distinguimos con un índice  $c = 1, 2$ . Tenemos  $n$  variables discriminantes  $\vec{x} = x_1, x_2, \dots, x_n$ , y podemos caracterizar las PDFs de cada especie usando muestras de control compuestas por  $N_c$  eventos cada una, formando las llamadas “n-tuplas”  $\mathbf{x}^c$ . Algunas definiciones:

- ▶ la “tupla-media”  $\bar{\mathbf{x}}^c$  de cada especie, y la tupla especie-promediada  $\bar{\mathbf{x}}$ :

$$\bar{\mathbf{x}}^c = \frac{1}{N_c} \sum_{k=1}^{N_c} \mathbf{x}_k^c, \quad \bar{\mathbf{x}} = \frac{1}{N_1 + N_2} \sum_{c=1}^2 \sum_{k=1}^{N_c} \mathbf{x}_k^c,$$

- ▶ las matrices de covarianza *intra-especie*  $W_{ij}$  (“within”) e *inter-especie*  $B_{ij}$  (“between”):

$$W_{ij} = \frac{1}{N_1 + N_2} \sum_{c=1}^2 \sum_{k=1}^{N_c} (x_i^{ck} - \bar{x}_i^c) (x_j^{ck} - \bar{x}_j^c), \quad B_{ij} = \frac{1}{N_1 + N_2} \sum_{c=1}^2 N_c (\bar{x}_i^c - \bar{x}_i) (\bar{x}_j^c - \bar{x}_j),$$

- ▶ y la matriz de covariancia especie-promediada  $T_{ij} = B_{ij} + W_{ij}$ . Los índices  $i, j = 1, 2, \dots, n$ .
- ▶ los coeficientes de Fisher-Mahalanobis  $f_i$ , asignados a cada variable  $i = 1, 2, \dots, n$ , y el offset  $f_0$ :

$$f_i = \frac{\sqrt{N_1 N_2}}{N_1 + N_2} \sum_{i=1}^n C_{ij}^{-1} (\bar{x}_i^1 - \bar{x}_i^2), \quad f_0 = \sum_{i=1}^n f_i (\bar{x}_i^1 + \bar{x}_i^2),$$

- (aquí la matriz  $C_{ij}^{-1}$  corresponde a  $W_{ij}^{-1}$  para Fisher, y a  $T_{ij}^{-1}$  para Mahalanobis)
- ▶ finalmente, el *discriminante de Fisher*  $\mathcal{F}$ , combinación lineal de las variables aleatorias  $\vec{x} = x_1, x_2, \dots, x_n$ :

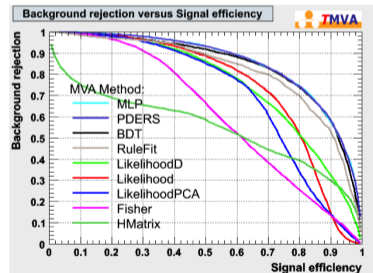
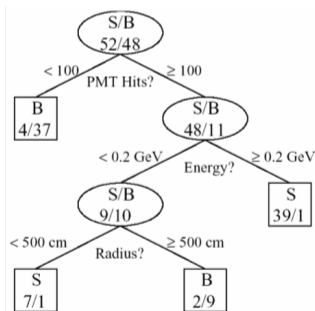
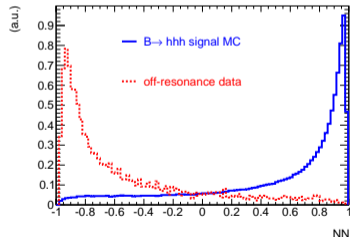
$$\mathcal{F} = f_0 + \sum_{i=1}^n f_i x_i.$$



# Discriminantes no lineales

Para tratar correlaciones no lineales y con perfiles complejos, existe una variedad de técnicas y herramientas. El paquete TMVA es una implementación popular en HEP de varios algoritmos de reducción dimensional: además de una biblioteca de discriminantes lineales y basados en likelihood, incluye métodos de entrenamiento y prueba con redes de neuronas artificiales y árboles de decisión (boosted decision trees), que forman parte de los más utilizados en HEP.

Idea general: un análisis multivariado utiliza una colección de variables de entrada, que se combinand de manera optimizada, a partir de un algoritmo “entrenado” sobre dos muestras independientes (correspondientes a señal y fondo), con dos protocolos: entrenamiento (training) y prueba (test). El performance del algoritmo entrenado se evalúa sobre las muestra de prueba, para evitar el efecto de “over-training”.



“ROC-curve”:  
Receiver Operating Characteristic



## The famous last words

En resumen, un *análisis multivariado* produce una reducción dimensional, proyectando un espacio de  $n$  variables aleatorias  $\vec{x} = \{x_1, x_2 \dots, x_n\}$ , sobre una variable final  $\mathcal{Z}$ . Esta variable puede ser

- ▶ una combinación lineal de las  $\vec{x}$  (discriminante de Fisher) ;
- ▶ un estadístico de prueba más elaborado (por ejemplo el cociente de verosimilitudes  $\lambda = \frac{\mathcal{L}(H_1)}{\mathcal{L}(H_0)}$ ) ;
- ▶ la salida de un algoritmo de *Machine Learning*. Ejemplos: *Multilayer Perceptron* (red de neuronas artificiales en sus múltiples apelaciones, tipo NN, CNN, ANN, GNN...), *Boosted Decision Tree* (árbol de decisiones), *algoritmos genéticos*, y otros más.

De manera general, los algoritmos de ML obtienen *performances* superiores que los algoritmos más sencillos, pero es crucial verificar la robustez de esos *performances* con respecto a la calidad de las muestras de control, o con respecto al algoritmo de entrenamiento. Dependiendo del objetivo del análisis se podrá preferir la robustez al *performance*, o se privilegiará una combinación entre ambas...

Dos citas sacadas de Wikipedia:

- ▶ Trees can be very non-robust. A small change in the training data can result in a large change in the tree and consequently the final predictions.
- ▶ Decision-tree learners can create over-complex trees that do not generalize well from the training data. (This is known as overfitting.)

Estos tópicos serán tratados desde varias perspectivas complementarias en el curso siguiente: “Tópicos avanzados en ciencia de datos”.

¡ Feliz continuación !