



CAPITULO V

EL METODO DE LA VEROSIMILITUD MAXIMA



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 Σ 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 $(\hat{\theta} \pm \hat{\Delta}_\theta)$.

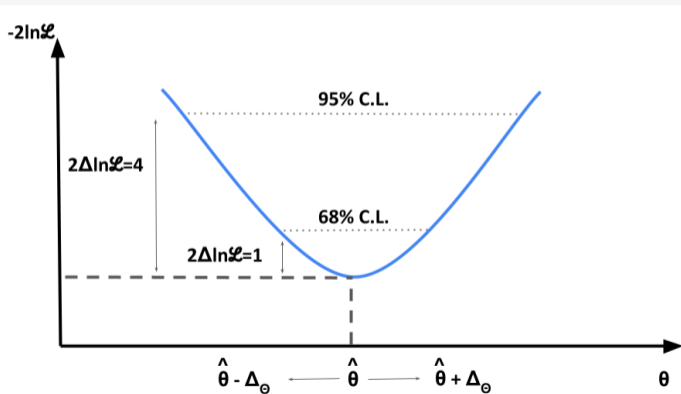


El teorema de verosimilitud máxima (III)

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 $(\hat{\theta} \pm \hat{\Delta}_{\theta})$.

(y de la misma manera, el intervalo $-2\Delta \ln \mathcal{L} < 4$ define un intervalo de confianza a 95 % corresponde a un rango $-2\Delta_{\theta} \leq \theta - \hat{\theta} \leq 2\Delta_{\theta}$, etc...)





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 μ y anchura σ . 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 μ y σ :

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

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

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

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

(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 ξ . 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 ($\hat{\xi} \pm \hat{\Delta}_{\xi}$).

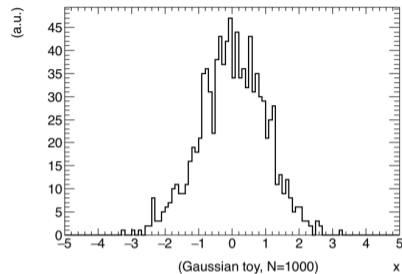
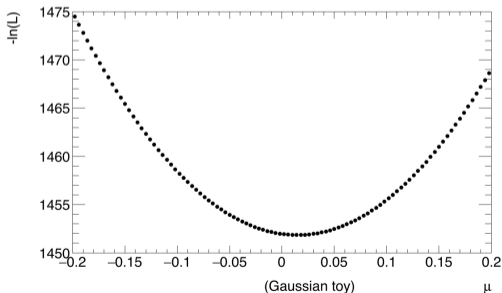
Nota: en este caso como el anterior, los errores $\hat{\Delta}$ escalan como $1/\sqrt{N}$. Esa es una propiedad muy general.



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 μ y σ (“escaneamos” los parámetros de interés)

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



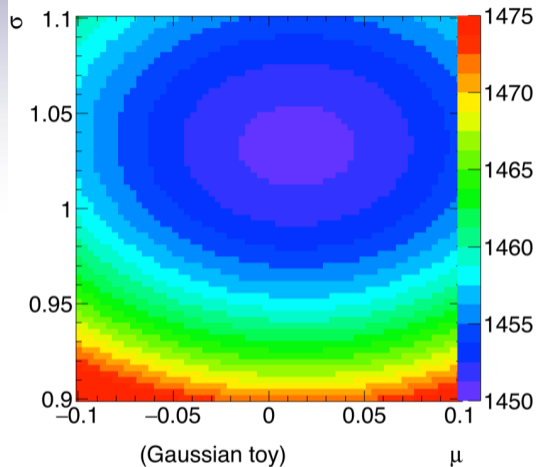
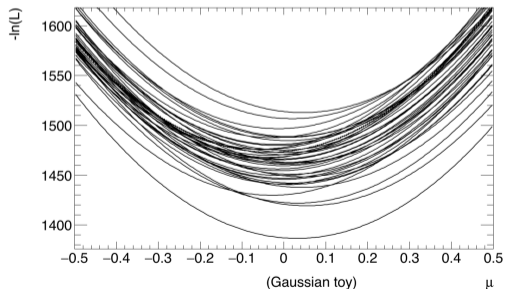
- ▶ La forma de $-\ln \mathcal{L}$ en función de μ (para σ fijo) sigue un perfil parabólico, y el mínimo coincide con el promedio empírico
- ▶ Por inspección alrededor del mínimo, se observa que el NLL aumenta en 0,5 unidades para $\Delta x \sim \pm 0,03$, que corresponde aproximadamente a σ/\sqrt{N}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



Un “ejemplo de juguete” (II)

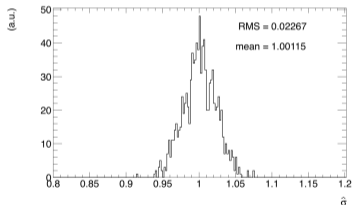
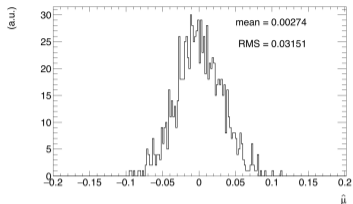
- ▶ El perfil bidimensional de $-\ln \mathcal{L}$ es el de un paraboloide, con semi-ejes diferentes y ortogonales
- ▶ El mínimo coincide con la posición del promedio y el RMS empíricos
- ▶ El semi-eje a lo largo de σ es más estrecho que el de μ , como se espera de la relación $1/\sqrt{2N}$ vs. $1/\sqrt{N}$

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



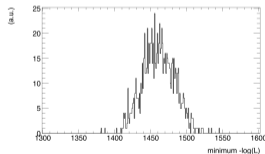


Un “ejemplo de juguete” (III)



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 σ/\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



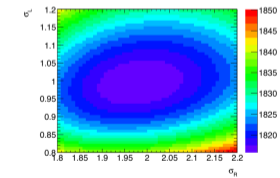
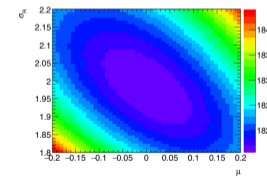
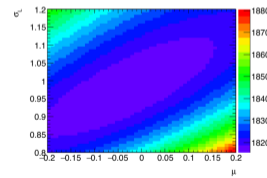
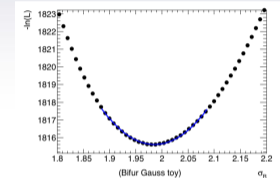
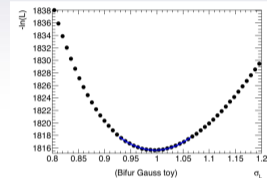
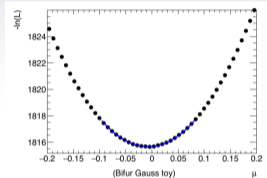
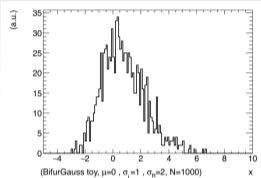
La distribución del $-\ln \mathcal{L}$ no siempre sigue una forma precisa (aquí parece bastante Gaussiana) pero permite hacer una estimación del “goddness-of-fit” :

- ▶ 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 : μ , σ_L y σ_R . Los mínimos del $-\log \mathcal{L}$ no son del todo parabólicos, y hay correlaciones importantes entre los parámetros.

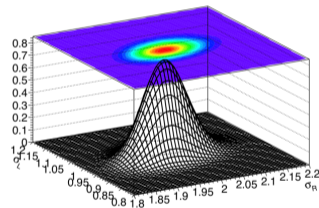
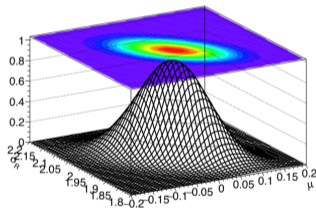
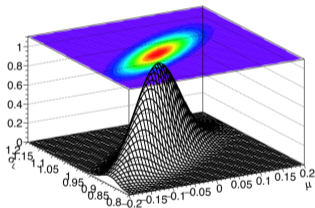




Otro "ejemplo de juguete" (II)

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 μ y σ_L : +79 %

Correlación entre μ y σ_R : -59 %

Correlación entre σ_L y σ_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](#)^[1] 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.^[2] 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.^[3] There is also a [Java](#) port^[4] as well as a [Python](#) frontend to the C++ code.^[5]

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.



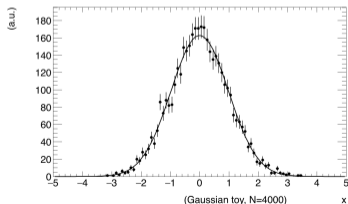
Diferencia entre un fit NLL y un fit de χ^2

Muchos paquetes y programas contienen algoritmos de minimización más sencillos, del tipo “mínimos cuadrados” o χ^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 σ_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 χ^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 $\sigma_i = \sqrt{y_i}$.



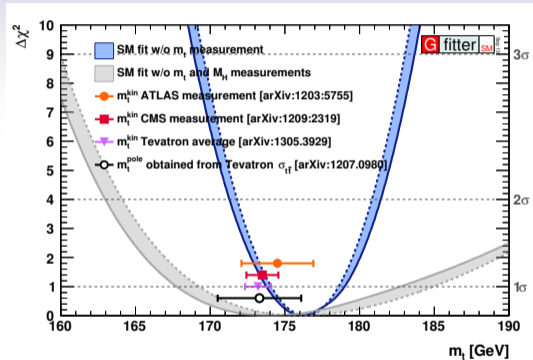
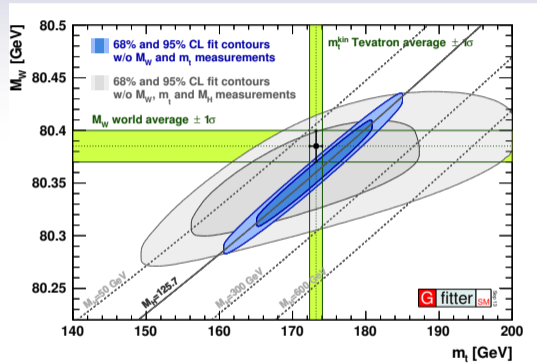
El ajuste de este histograma a una función Gaussiana nos da $\mu = (0,008 \pm 0,016)$ y $\sigma = (1,007 \pm 0,011)$.

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 RMS/\sqrt{n} y $\text{RMS}/\sqrt{2n}$.

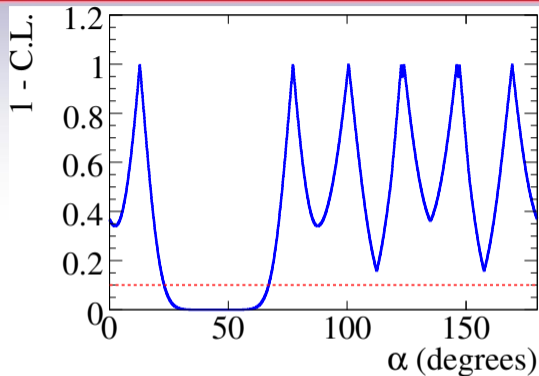


Los intervalos de confianza de Gfitter

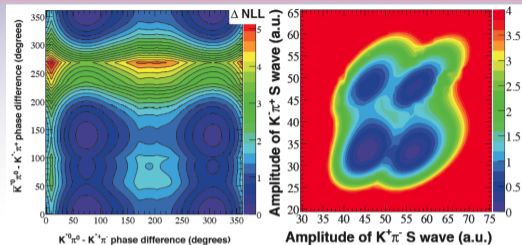




Ejemplos de funciones de verosimilitud complicadas



Un elemento importante del programa científico del experimento *BABAR* (y en general, de la *física de sabores*): la medida del ángulo α 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

$$B^0/\bar{B}^0 \rightarrow K^+\pi^-\pi^0;$$

- ▶ 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^* , ρ , 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]																																			
	ρ^0	K^*	S	f_1	f_2	f_x	NR	x	f_0	ρ^0	K^*	S	f_1	f_2	f_x	NR	x																			
[i]	ρ^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	f_x 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]	ρ^0 31.3	K^* 20.6	48.6	S 44.7	73.5	56.3	f_1 2.4	-30.1	6.3	-56.1	f_2 14.5	34.1	12.5	16.1	-23.0	f_x 18.9	27.0	20.6	5.8	11.8	9.5	x 18.9	27.0	20.6	5.8	11.8	9.5	-84.2	21.5	178	28.1	27.8	0.8	8.1	20.2	100.0
ang(c)	ρ^0 -11.2	K^* 25.0	8.6	S 33.0	10.6	3.4	f_1 12.1	-0.6	-9.8	-2.6	f_2 25.0	10.2	5.4	-0.5	-11.4	f_x NR 31.6	17.0	39.3	1.0	-27.1	-31.7	x 8.6	11.7	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)	ρ^0 14.5	K^* 17.1	7.1	S 22.5	15.9	25.2	f_1 15.1	4.9	15.5	-5.0	f_2 8.1	2.7	12.3	-0.6	36.5	f_x NR 15.3	4.1	14.5	-3.0	-22.6	-20.8	x 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
ang(i)	ρ^0 18.9	K^* 17.8	8.6	S 33.0	10.6	3.4	f_1 12.1	-0.6	-9.8	-2.6	f_2 25.0	10.2	5.4	-0.5	-11.4	f_x NR 31.6	17.0	39.3	1.0	-27.1	-31.7	x 8.6	11.7	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(j)	ρ^0 14.5	K^* 17.1	7.1	S 22.5	15.9	25.2	f_1 15.1	4.9	15.5	-5.0	f_2 8.1	2.7	12.3	-0.6	36.5	f_x NR 15.3	4.1	14.5	-3.0	-22.6	-20.8	x 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
ang(k)	ρ^0 18.9	K^* 17.8	8.6	S 33.0	10.6	3.4	f_1 12.1	-0.6	-9.8	-2.6	f_2 25.0	10.2	5.4	-0.5	-11.4	f_x NR 31.6	17.0	39.3	1.0	-27.1	-31.7	x 8.6	11.7	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(l)	ρ^0 14.5	K^* 17.1	7.1	S 22.5	15.9	25.2	f_1 15.1	4.9	15.5	-5.0	f_2 8.1	2.7	12.3	-0.6	36.5	f_x NR 15.3	4.1	14.5	-3.0	-22.6	-20.8	x 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]																																			
	ρ^0	K^*	S	f_1	f_2	f_x	NR	x	f_0	ρ^0	K^*	S	f_1	f_2	f_x	NR	x																			
[i]	ρ^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	f_x 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]	ρ^0 31.3	K^* 20.6	48.6	S 44.7	73.5	56.3	f_1 2.4	-30.1	6.3	-56.1	f_2 14.5	34.1	12.5	16.1	-23.0	f_x NR 17.8	57.6	41.7	13.7	30.1	49.7	x 18.9	27.0	20.6	5.8	11.8	9.5	-84.2	21.5	178	28.1	27.8	0.8	8.1	20.2	100.0
ang(c)	ρ^0 -11.2	K^* 25.0	8.6	S 33.0	10.6	3.4	f_1 12.1	-0.6	-9.8	-2.6	f_2 25.0	10.2	5.4	-0.5	-11.4	f_x NR 31.6	17.0	39.3	1.0	-27.1	-31.7	x 8.6	11.7	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)	ρ^0 14.5	K^* 17.1	7.1	S 22.5	15.9	25.2	f_1 15.1	4.9	15.5	-5.0	f_2 8.1	2.7	12.3	-0.6	36.5	f_x NR 15.3	4.1	14.5	-3.0	-22.6	-20.8	x 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
ang(i)	ρ^0 18.9	K^* 17.8	8.6	S 33.0	10.6	3.4	f_1 12.1	-0.6	-9.8	-2.6	f_2 25.0	10.2	5.4	-0.5	-11.4	f_x NR 31.6	17.0	39.3	1.0	-27.1	-31.7	x 8.6	11.7	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(j)	ρ^0 14.5	K^* 17.1	7.1	S 22.5	15.9	25.2	f_1 15.1	4.9	15.5	-5.0	f_2 8.1	2.7	12.3	-0.6	36.5	f_x NR 15.3	4.1	14.5	-3.0	-22.6	-20.8	x 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
ang(k)	ρ^0 18.9	K^* 17.8	8.6	S 33.0	10.6	3.4	f_1 12.1	-0.6	-9.8	-2.6	f_2 25.0	10.2	5.4	-0.5	-11.4	f_x NR 31.6	17.0	39.3	1.0	-27.1	-31.7	x 8.6	11.7	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(l)	ρ^0 14.5	K^* 17.1	7.1	S 22.5	15.9	25.2	f_1 15.1	4.9	15.5	-5.0	f_2 8.1	2.7	12.3	-0.6	36.5	f_x NR 15.3	4.1	14.5	-3.0	-22.6	-20.8	x 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

112001-27



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_{sig} y P_{bkg} son las PDFs de señal y fondo, respectivamente, y la fracción de señal f_{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 :

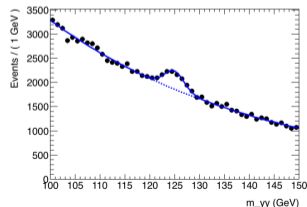
Con ROOT instalado, la macro H_yy.cc debe correr sin problema, haciendo

```
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

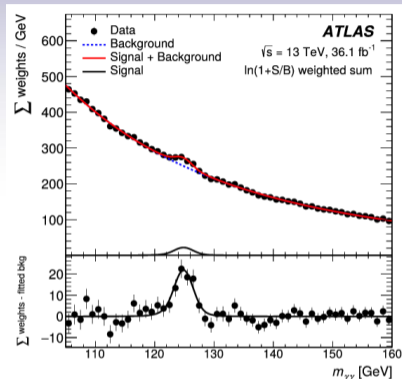
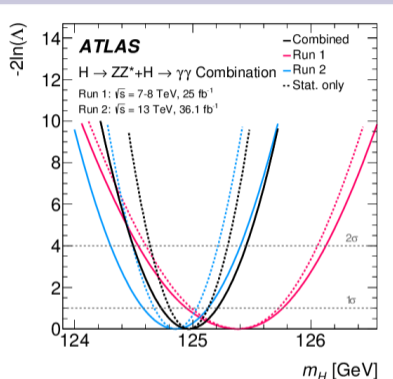
Floating Parameter      InitialValue      FinalValue +/-      Error      GblCorr.
-----
bkg_a                   -2.2500e-02      -2.2780e-02 +/-      2.31e-04      <none>
sig_f                   2.0000e-02      1.8732e-02 +/-      1.33e-03      <none>
sig_n                   1.2500e+02      1.2494e+02 +/-      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 λ , 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 ε es el cociente entre el número de realizaciones de cada tipo, n_{yes} y n_{no} , y su varianza $V[\hat{\varepsilon}]$ viene dada por :

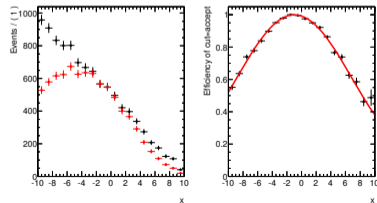
$$\hat{\varepsilon} = \frac{n_{\text{yes}}}{n_{\text{yes}} + n_{\text{no}}} , \quad V[\hat{\varepsilon}] = \frac{\hat{\varepsilon}(1 - \hat{\varepsilon})}{n} ,$$

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 biviada $c = \{\text{yes}, \text{no}\}$, 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 θ que sean necesarios para caracterizar su forma)

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



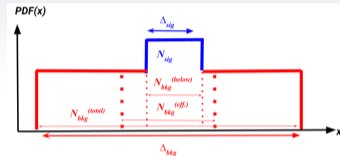


Sobre el fondo efectivo

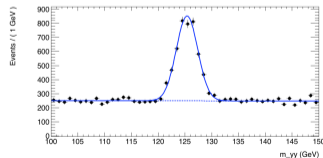
De manera general, la zona que contiene la señal también está contaminada por fondo(s). El impacto de esos fondos se traduce en una *dilución* o disminución de la precisión en la medida de los parámetros de interés.

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_{sig} y N_{bkg} .
- ▶ El parámetro de interés es el número de eventos de señal N_{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_{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.

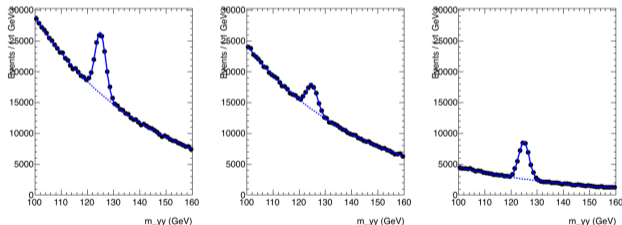




En ocasiones, la muestra de análisis puede descomponerse en dos o más submuestras (“categorías”), cada una de ellas con sus propias PDFs y purezas.

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: $\sigma(m_H) = \pm 2,25 \%$

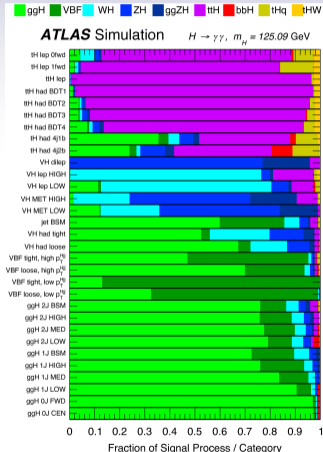
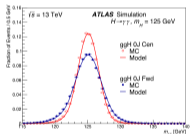
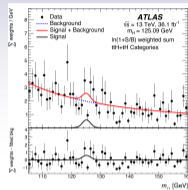
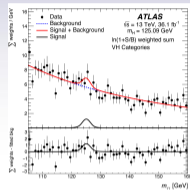
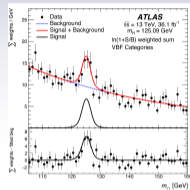
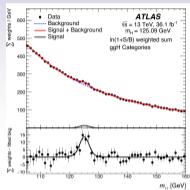
Errores en las categorías sucia y limpia : $\pm 3,70 \%$ y $\pm 1,85 \%$, respectivamente.

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 5σ de significancia en la observación del Higgs...