



CAPITULO VII

CONTRASTE DE HIPOTESIS



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 mi hipótesis, ¿cuál es la probabilidad de observar una muestra de datos menos representativa de mi hipótesis que la que realmente observé?”
- ▶ equivale a cuantificar la probabilidad matemática

$$\mathcal{P}(\text{datos} \mid \vec{\mu}) ,$$

donde $\vec{\mu}$ son nuestros parámetros de interés POI.



Digestión rápida: la interpretación *bayesiana* de la estadística

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?”
- ▶ equivale a cuantificar una probabilidad matemática diferente a la frecuentista:

$$\mathcal{P}(\vec{\mu} | \text{datos}) = \mathcal{P}(\vec{\mu}) \times \frac{\mathcal{P}(\text{datos} | \vec{\mu})}{\mathcal{P}(\text{datos})},$$

donde

- ▶ $\mathcal{P}(\text{datos} | \vec{\mu})$ es la probabilidad frecuentista (p.e. extraída de un análisis de verosimilitud máxima),
- ▶ $\mathcal{P}(\vec{\mu})$ es la *probabilidad previa*, o *prior bayesiano*,
- ▶ $\mathcal{P}(\text{datos})$ es un coeficiente de normalización sin importancia.
- ▶ $\mathcal{P}(\vec{\mu} | \text{datos})$ es la *probabilidad subjetiva* o “*grado de creencia*” (!)

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...

“Bayesians address the questions everyone is interested in by using assumptions that no one believes. Frequentists use impeccable logic to deal with an issue that is of no interest to anyone.”
(L. Lyons)



El test del χ^2

Dado un conjunto de n medidas independientes x_i con varianzas σ_i^2 , y un conjunto de predicciones μ_i , se define un test estadístico llamado χ^2 de la manera siguiente:

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

El test de χ^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 χ^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 aleatorio, y del cual extraemos una muestra de datos. El 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 .

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.

Caso particular : física de altas energías

Allí, 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?” .



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 (“falso negativo”);
- ▶ se incurre en un “Error de Tipo-II” si \mathcal{H}_0 es aceptada aún siendo falsa (“falso positivo”).

Las tasa de errores de Tipo-I y Tipo-II se llaman usualmente α y β 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 α es llamada *tamaño del contraste* (*size of the test* en inglés), puesto que fijar α 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_λ :

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

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_λ , y la p -value observada se obtiene al integrarla con respecto al valor observado de q_λ .



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 hábitos 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”.

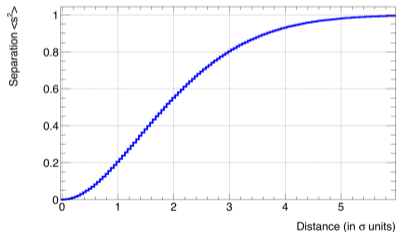


Separación

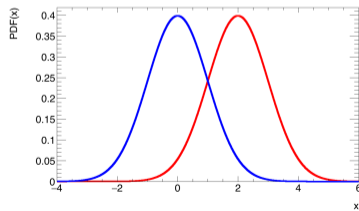
Las definiciones generales de *tamaño* y *potencia* pueden ser complementadas con otras definiciones más sencillas. Por ejemplo la “separación”, que se determina a partir de las PDFs de las dos especies (señal S y fondo B) para las cuales se quiere cuantificar el poder de discriminación:

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

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 σ , en función de la distancia entre sus picos (en unidades de σ). Ejemplo para una distancia de 2σ :



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



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_{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 λ , derivado de un cociente de verosimilitudes, es

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

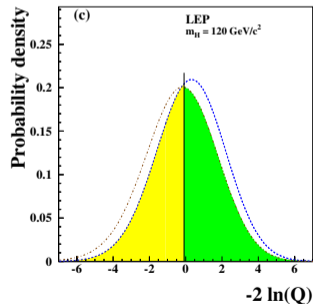
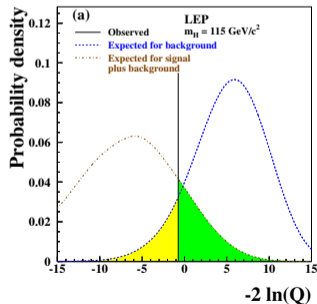
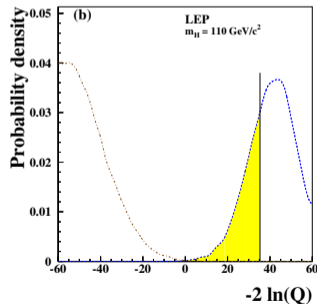
de manera que valores positivos de λ favorezcan un escenario “background-like”, y valores negativos estén más tono con un escenario “señal-más-fondo”; valores cercanos a cero indicando una sensibilidad pobre para distinguir entre ambos.



Algunas convenciones estadísticas en física de partículas: LEP

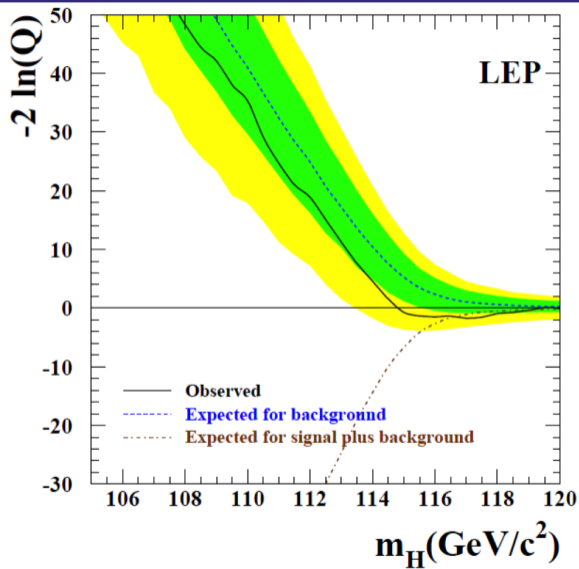
- ▶ Bajo la hipótesis “background-only”, $CL(b)$ es la probabilidad de tener $-2 \ln Q$ más pequeño que el observado (amarillo);
- ▶ bajo la hipótesis “señal-más-fondo”, $CL(s + b)$ es la probabilidad de tener $-2 \ln Q$ más grande que el observado (verde).

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 $1 - CL(b)$.





La lógica de exclusión en LEP





El estimador modificado $CL(s)$

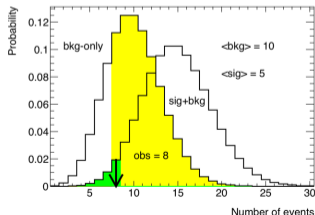
Tomemos un experimento de conteo de eventos: esperamos 10 eventos de tipo fondo, y 5 eventos de tipo señal

- ▶ pero... observamos 8 eventos en total...
- ▶ lo razonable es suponer es que tanto señal como fondo experimentaron una fluctuación negativa...
- ▶ pero en la interpretación estándar, ¡ se hubiera asignado una exclusión a 95 % C.L. !
- ▶ (ello incluso si 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 como estadístico de prueba el llamado *cociente de perfil de verosimilitud* (profiled likelihood ratio) definido así:

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

- ▶ el parámetro de interés es $\mu = \sigma/\sigma_{\text{SM}}$, la “intensidad de señal”: la tasa de conteo de candidatos de señal comparada a la predicción (p.e. la sección eficaz de producción del Higgs vs. la predicción teórica),
- ▶ $\hat{\hat{\theta}}$ son los valores de los NPs obtenidos en un ajuste a intensidad de señal μ fija,
- ▶ $\hat{\mu}$ y $\hat{\theta}$ son los valores ajustados cuando tanto μ 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).

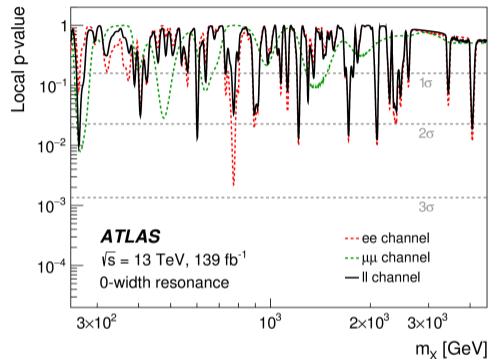
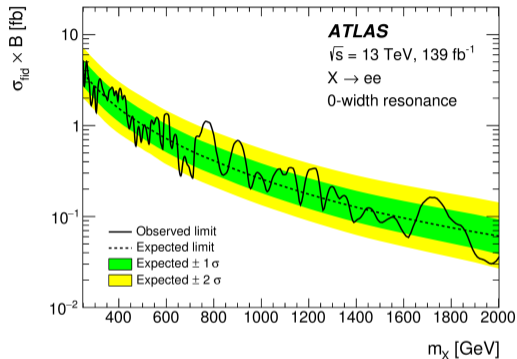
Para un valor observado del estadístico $\hat{\hat{q}}_\mu$, las p -values asignadas a las hipótesis de signal-plus-background y background-only, $p(s+b)$ y $p(b)$, son

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

- ▶ los resultados de exclusión se muestran bajo forma de un “Brazil-plot”,
- ▶ los resultados 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 $X \rightarrow \mu^+\mu^-$...)