



Analisi dati

Metodi di analisi statistica
dei dati di esperimenti di
fisica delle particelle

Argomenti

- Introduzione
- Dalle misure sperimentali ai risultati finali

- Variabili aleatorie e distribuzioni statistiche
- Medie e varianze
- Distribuzioni binomiali, di Poisson e di Gauss

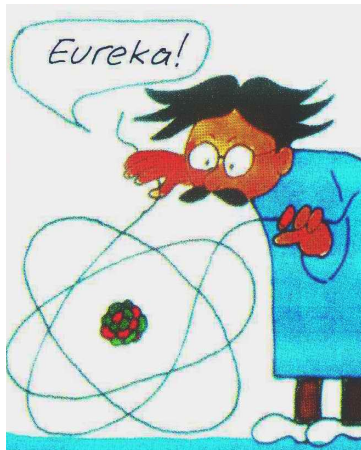
- Modelli teorici e fit sperimentali
- Misure di vita media

- Selezione dei campioni
- Livelli di confidenza e plot di esclusione
- I risultati della ricerca dell'Higgs

Introduzione

- **La fisica fondamentale e' costruita su ipotesi falsificabili**
- Le ipotesi non falsificate dagli esperimenti costituiscono la base delle teorie interpretative dei processi fisici
- Gli apparati sperimentali sono progettati per ottenere informazioni stringenti sulla natura dei processi di interesse
- La sensibilita' degli esperimenti risulta definita dalle caratteristiche dell'apparato
- L' interpretazione dei dati si basa
sullo studio e controllo della strumentazione
sulle conoscenze pregresse
su ipotesi e modelli teorici

... cosa si misura nella HEP ?

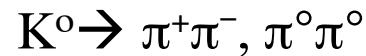


Parametri fondamentali

- Branching ratio (BR)
- vita media (τ)
- massa (m)
- costanti di accoppiamento
-

attraverso...

- quantita' di moto \vec{q} ;
- energia E rilasciata nel calorimetro ;
- angoli e direzioni delle particelle prodotte ;
- intervalli temporali ;
- efficienza del rivelatore ;
- contaminazione in un campione ;
- ecc.



New physics ! ! !

Misure e variabili aleatorie

- I processi di misura sono per natura stocastici

ERRORI STATISTICI

- Si ottengono distribuzioni di variabili relative al campione da studiare, contaminate da altri processi, “background”

VARIABILI DISCRIMINANTI, TEST DI IPOTESI

- Inoltre strumentazione e le conoscenze “pregresse” dei fenomeni da considerare possono introdurre errori sistematici

ERRORI SISTEMATICI

Distribuzioni di probabilita'

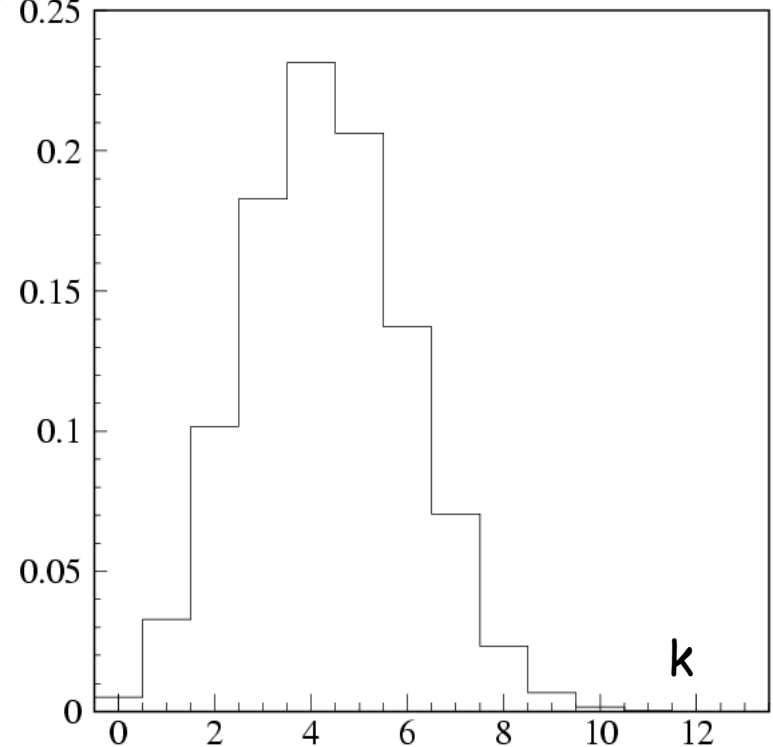
- Variabili discrete: conteggi

Se si conta classificando i conteggi come "OK" / "not-OK" questi saranno distribuiti come una binomiale

$$P(k; p, N) = \binom{N}{k} p^k (1-p)^{N-k}$$

Se i conteggi sono di un processo con frequenza μ , saranno distribuiti come una distribuzione di Poisson

$$P(k; \mu) = \frac{e^{-\mu} \mu^k}{k!}$$



$$K = \# K^0 \rightarrow \pi^0 \pi^0$$

$$BR (= p) = k/N$$

$$m = pN$$

$$\sigma = \sqrt{Np(1-p)}$$

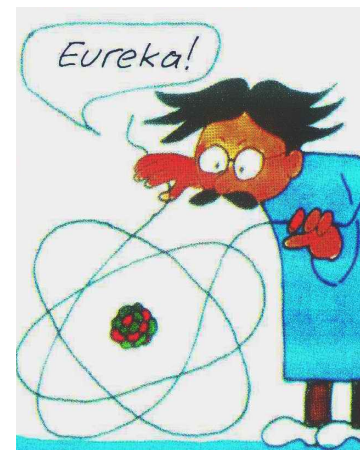
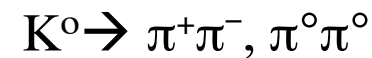
$$\frac{\sigma}{m} = \frac{\sqrt{Np(1-p)}}{pN} \approx \frac{1}{\sqrt{N}}$$

• **Precisione misura BR**

$$N=50 \quad \rightarrow \quad \sigma/m \quad \rightarrow \quad 10\%$$

$$N=10^3 \quad \rightarrow \quad \sigma/m \quad \rightarrow \quad 2.4\%$$

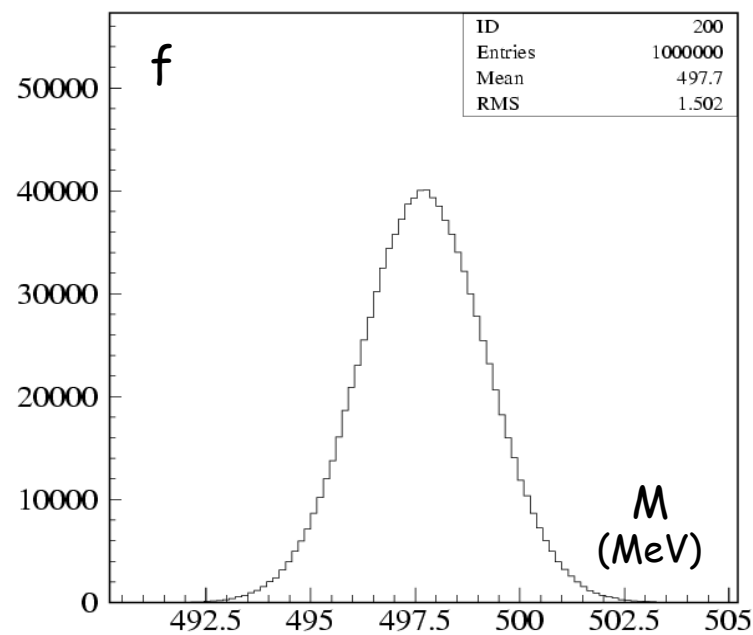
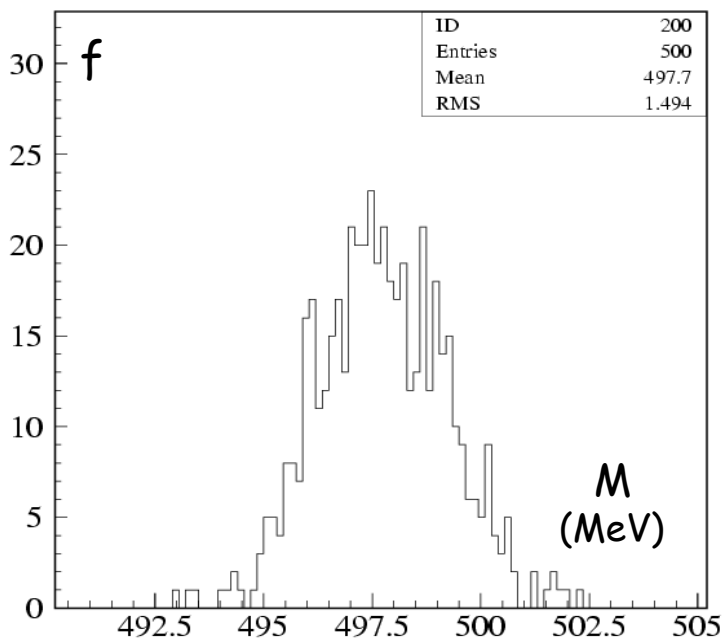
$$N=10^6 \quad \rightarrow \quad \sigma/m \quad \rightarrow \quad 0.08\%$$

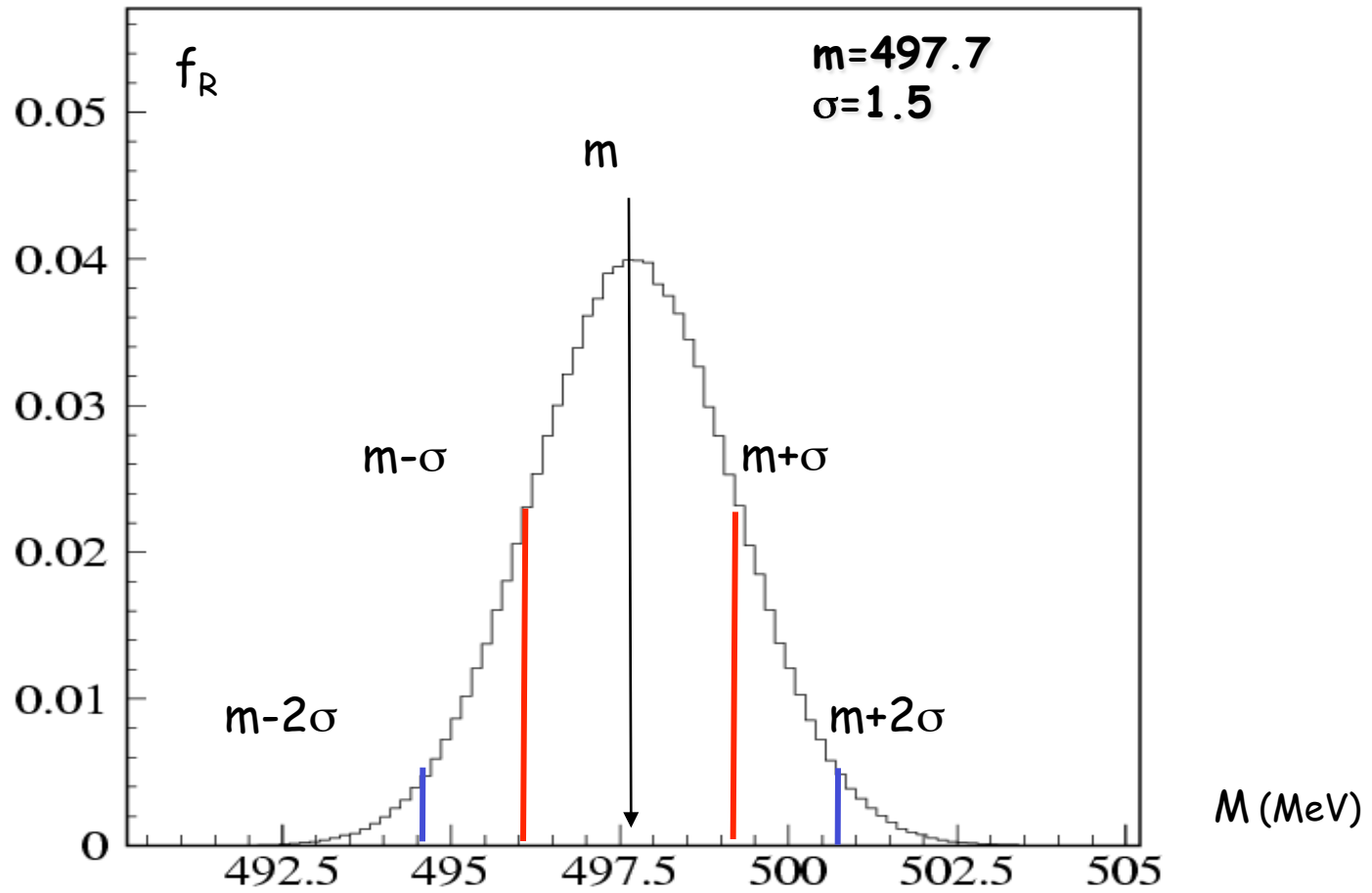


Distribuzione di Gauss

- Teorema del limite centrale

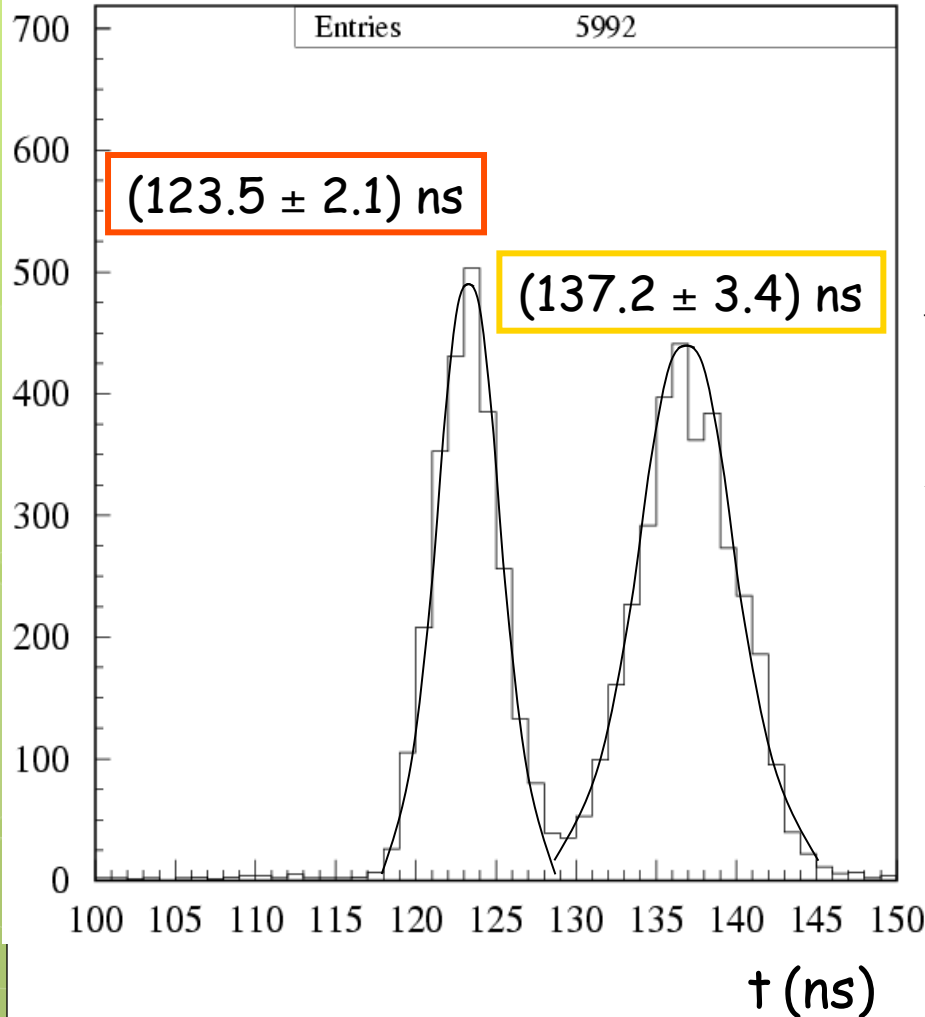
$$f(x; m, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-m)^2}{2\sigma^2}}$$





Tempi di volo

$v = ?$



$$t_{12} = \bar{t}_2 - \bar{t}_1 = 13.7 \text{ ns}$$

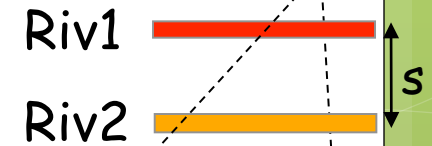
$$\Delta t_{12} = \Delta t_1 \oplus \Delta t_2 = 4.0 \text{ ns}$$

$$v = \frac{s}{t_{12}} = \frac{3.8 \text{ m}}{13.7 \times 10^{-9} \text{ s}} = 2.77 \times 10^8 \text{ m/s}$$

$$\Delta v = v \left(\frac{\Delta s}{s} \oplus \frac{\Delta t_{12}}{t_{12}} \right) = 0.83 \times 10^8 \text{ m/s}$$

$\sim 2\%$

$\sim 30\%$



Binomiale

Il numero di **successi** in N prove e' una variabile aleatoria: $0 \leq k \leq N$.



risultato lancio di un dado/moneta
conteggi / efficienza
conteggi/ selezione

Distribuzione di Poisson

campionamenti idealmente infiniti



nascite al mese
decadimenti di una
sostanza radioattiva
in un intervallo di
tempo fissato

Distribuzione di Gauss

innumerevoli fonti di errori casuali



spessori, ...
masse, impulsi, tempi, ...

Media e varianza

Valor medio

$$m = \frac{\sum x_i}{N}$$

$$\bar{x} = \int x \cdot f(x) dx$$

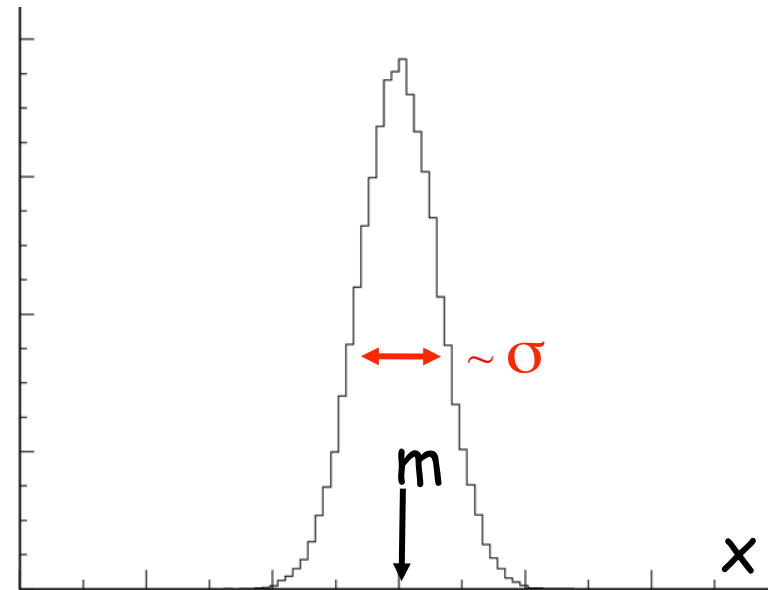
Deviazione standard

$$\sigma = \sqrt{\frac{\sum (m - x_i)^2}{N - 1}}$$

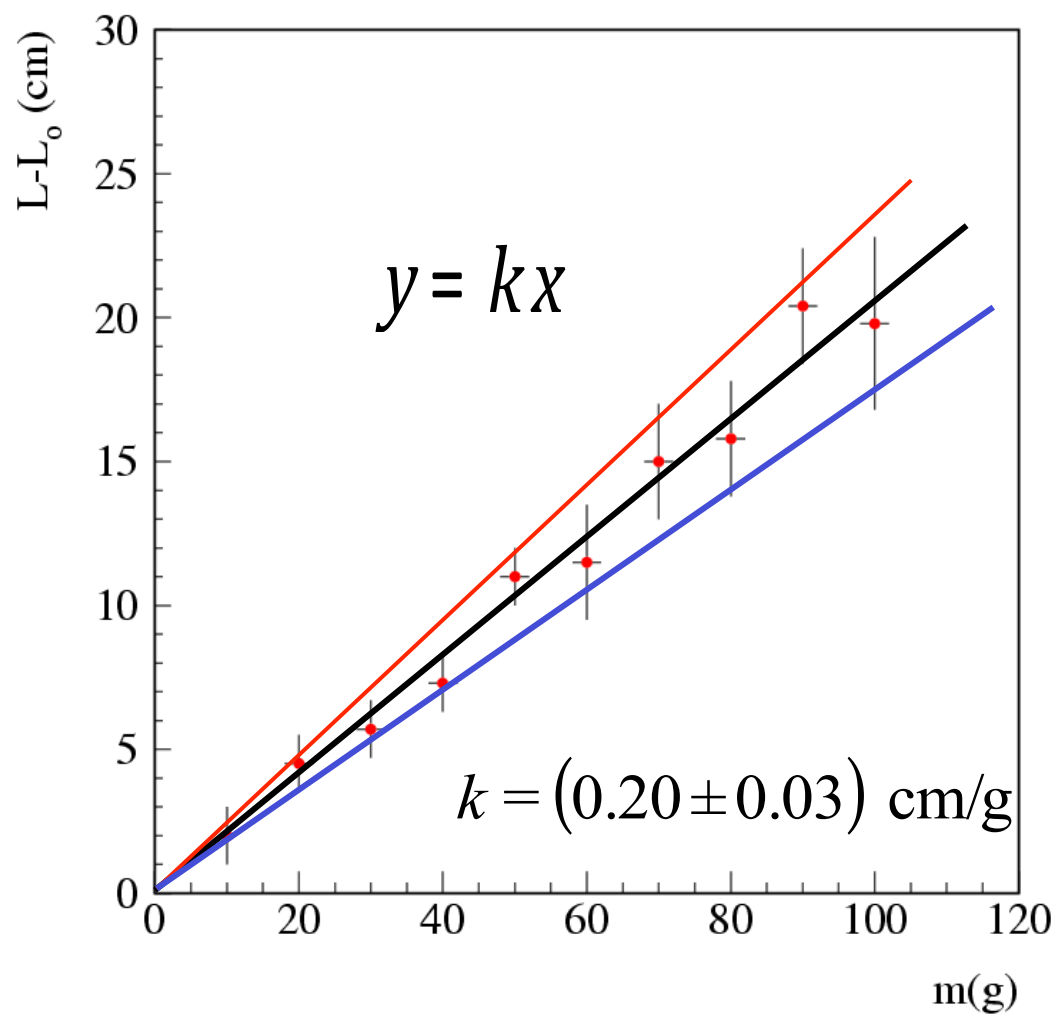
$$\sigma^2 = \int (x - \bar{x})^2 \cdot f(x) dx$$

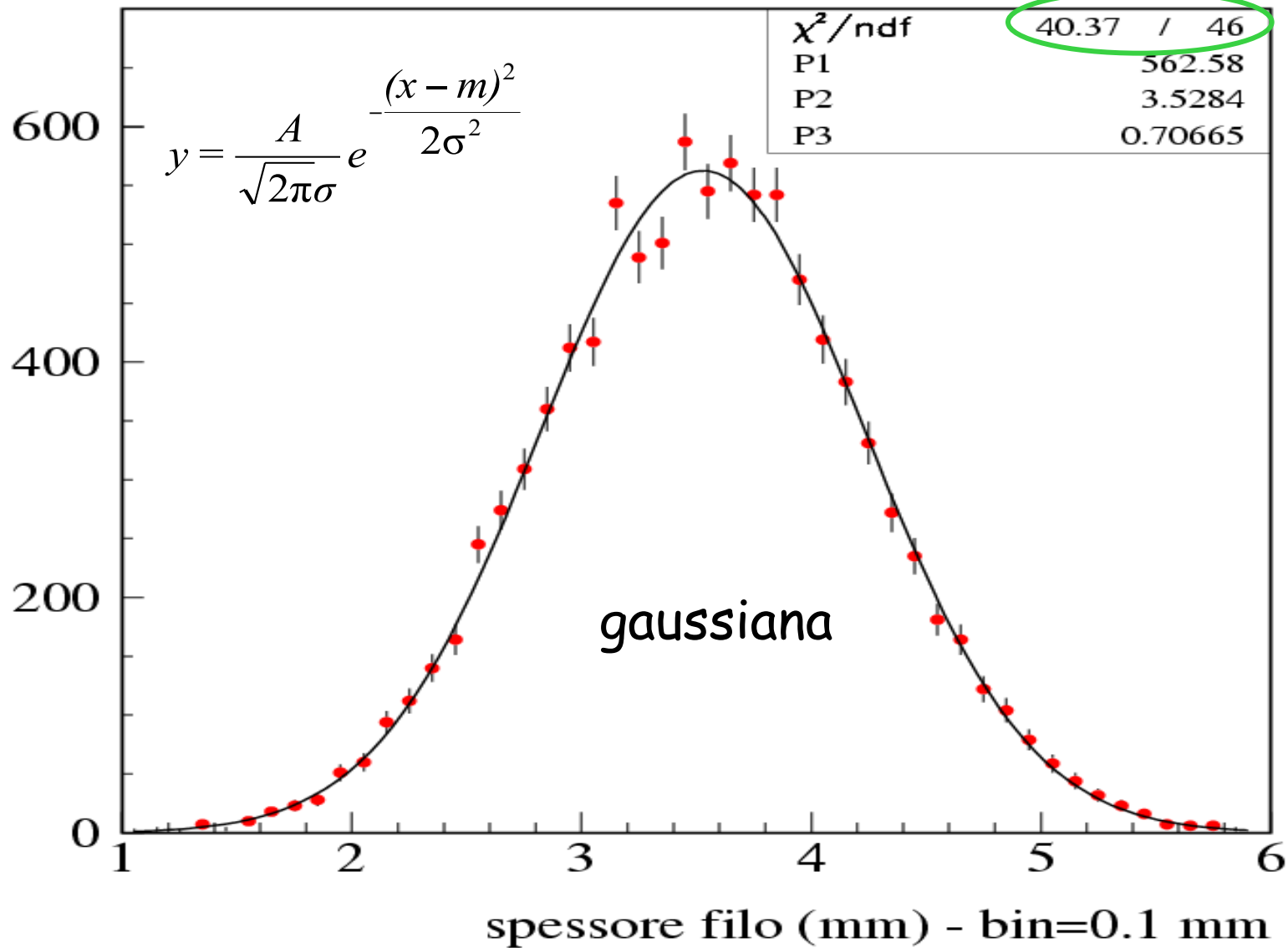
$$x = m \pm \sigma$$

$$e_x = \sigma/m$$

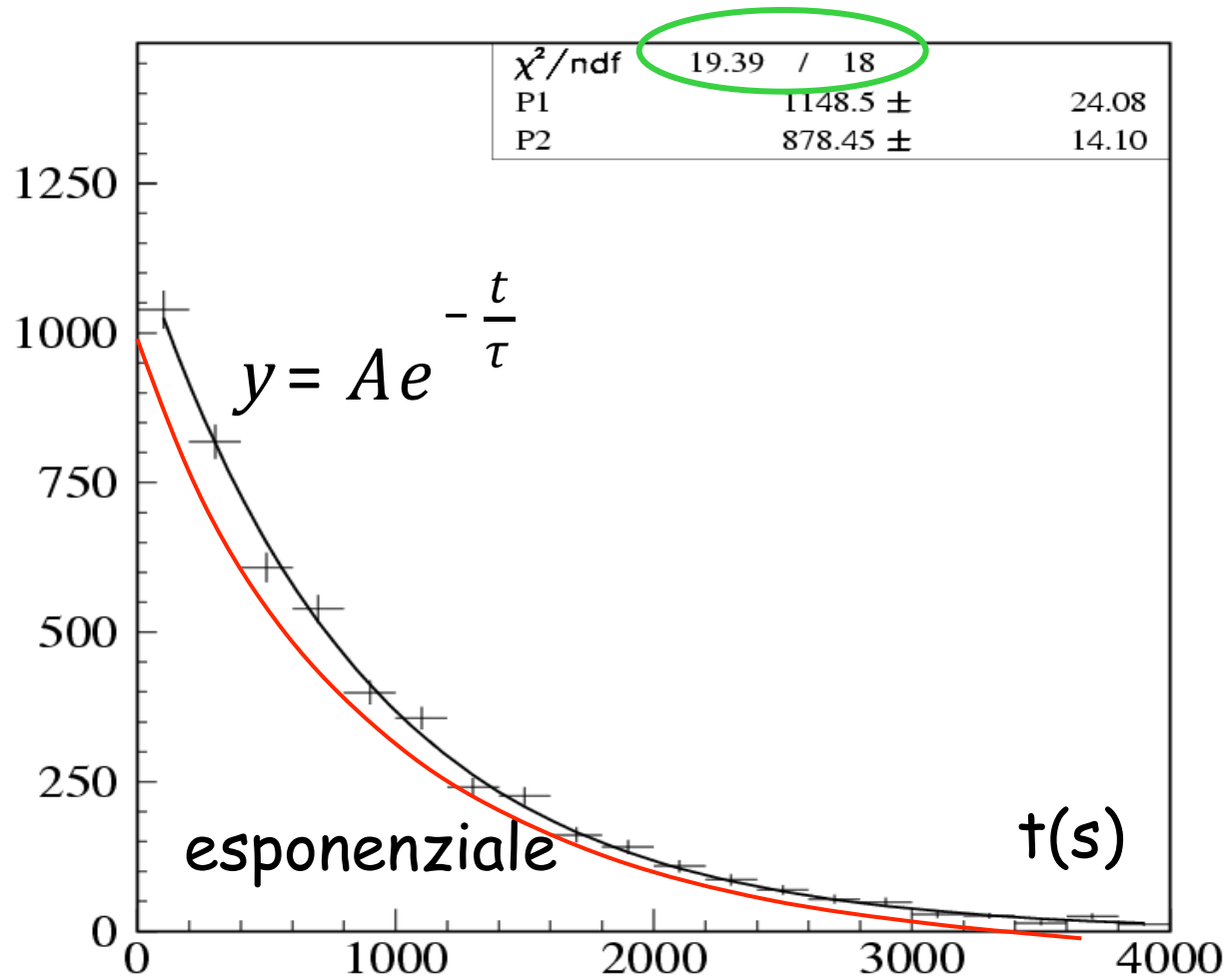


Test di ipotesi, Fit





Misura vita media



Funzione del χ^2

• Se il fenomeno in studio e' ben noto, la funzione $y = f(x; \alpha, \beta)$ e' gia' decisa. Lo scopo del *fit* e' allora di determinare i valori dei parametri (α e β in questo caso) che corrispondono al miglior adattamento della curva ai dati sperimentali.

• I valori che corrispondono al miglior adattamento sono quelli che rendono minima la quantita' seguente :

valore sperimentale

valore teorico

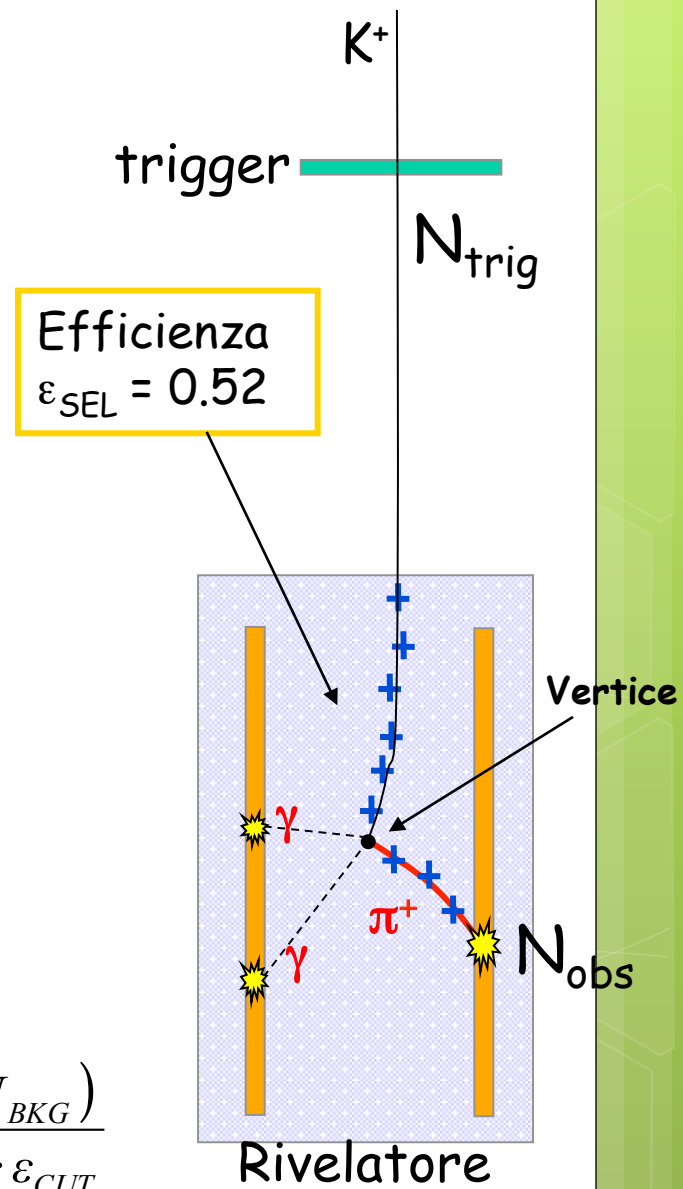
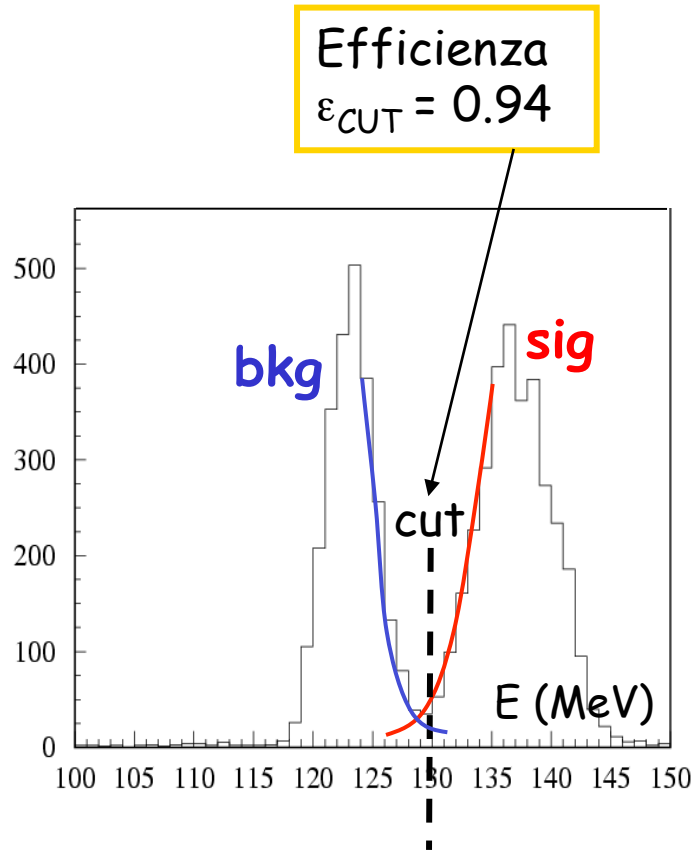
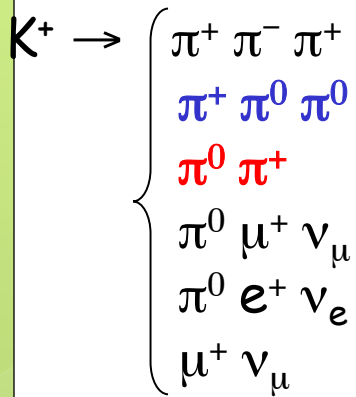
Van der Waals

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

$\sim 1 \rightarrow \text{O.K.}!$

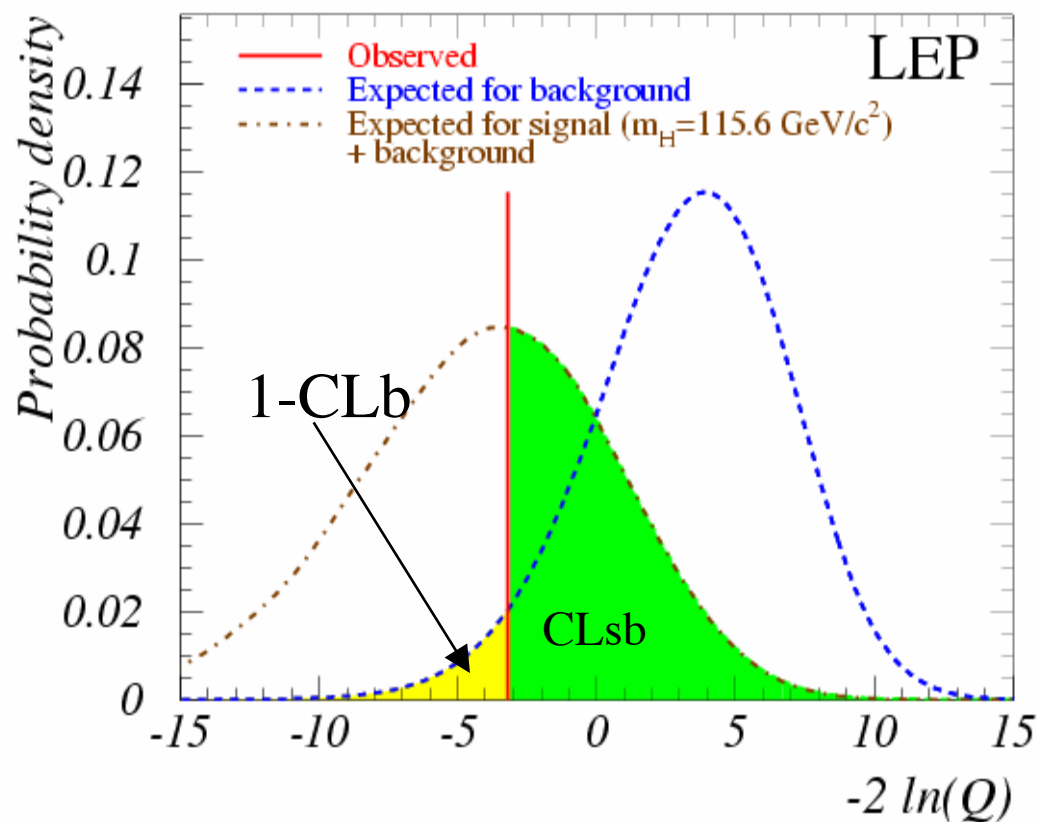
$> 1 \rightarrow \text{K.O.}!$

Branching ratio



$$BR(K^+ \rightarrow \pi^0 \pi^+) = \frac{N_{SIG}}{N_{trig}} = \frac{\epsilon_{SEL} \epsilon_{CUT}}{N_{trig}} = \frac{(N_{OBS} - N_{BKG})}{N_{trig} \cdot \epsilon_{SEL} \cdot \epsilon_{CUT}}$$

Variabili discriminanti, Likelihood



Confronto tra due ipotesi: H_0 e H_1

$$Q = L(H_1 | \text{dati}) / L(H_0 | \text{dati})$$

L'ipotesi H_0 generalmente non è associata a parametri liberi

La distribuzione nell'ipotesi H_1 cambia quando cambiano i parametri liberi della teoria

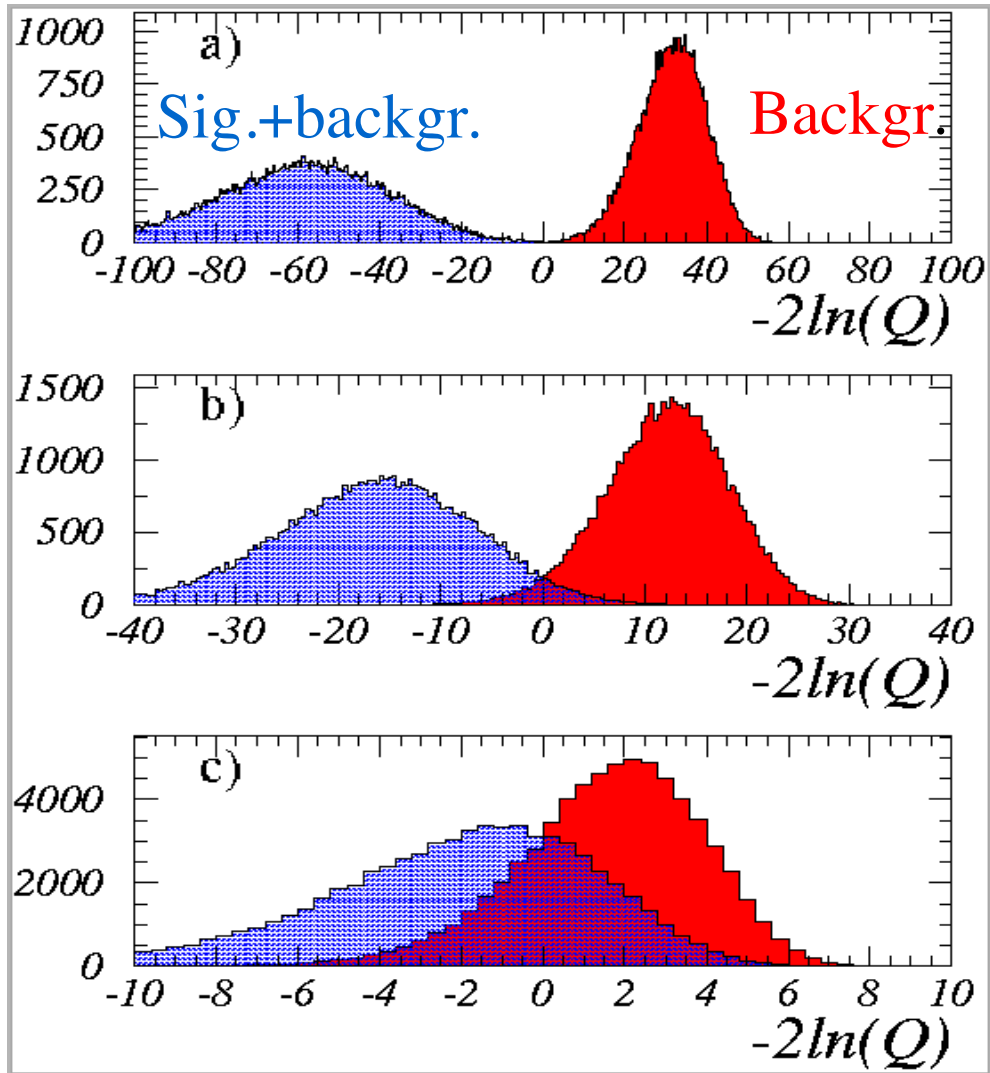
CLb: conf. nel background

1-CLb: significatività del test

CLsb: conf nel sig.+bg.

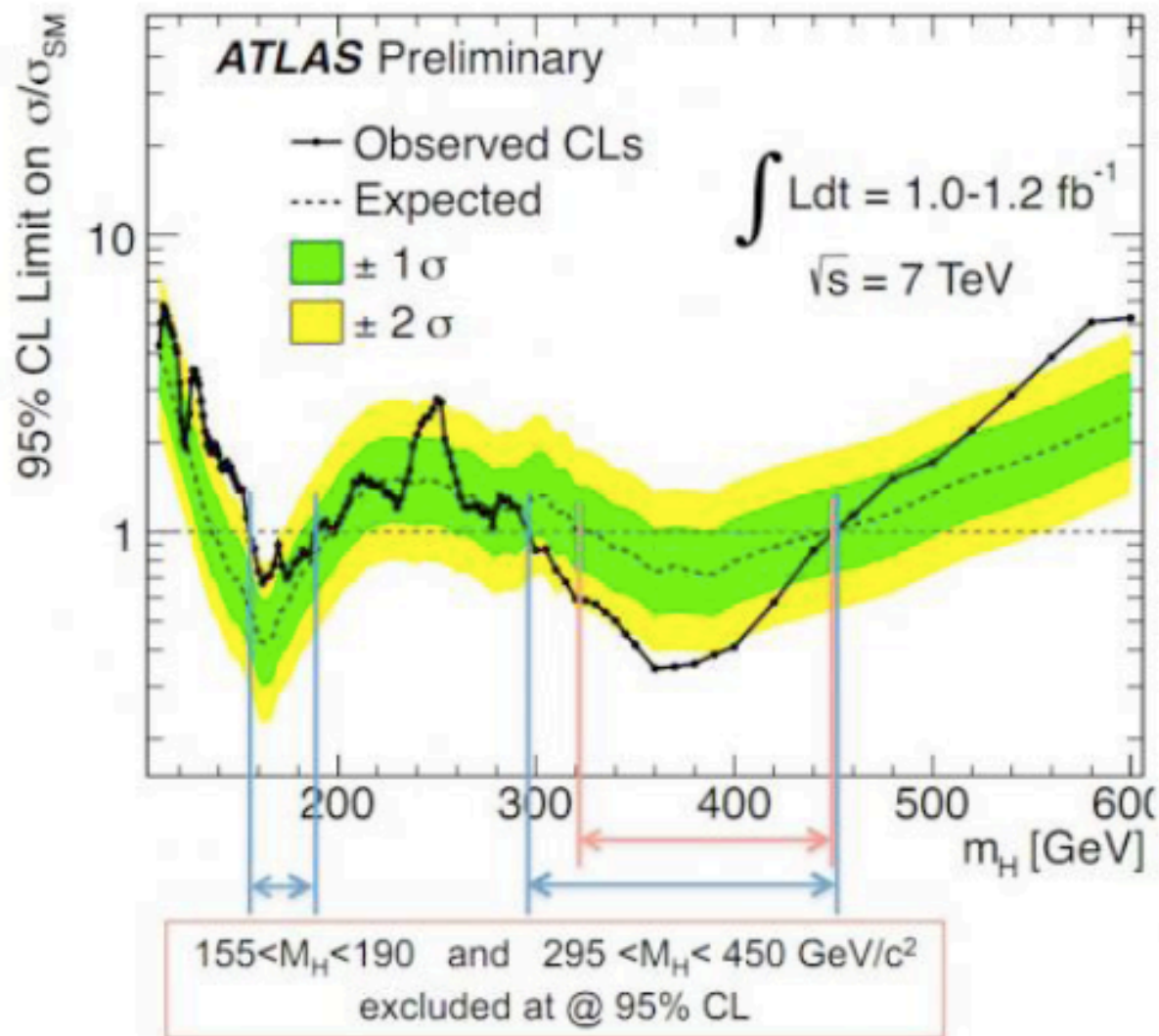
CLs=CLsb/CLb: approssimazione della conf. nel sig.

Sensibilita' sperimentale



- Ottima sensibilita'. Ottime possibilita' di scoperta o in caso di assenza del fenomeno, esclusione del segnale o limiti molto stringenti
- Buona sensibilita'. Probabile un risultato di esclusione con buoni livelli di confidenza
- **Scarsa sensibilita', probabili risultati ambigui, importante la scelta di metodi statistici raffinati**

Ricerca dell'Higgs: risultati di ATLAS



Conclusioni

- Metodi avanzati per il trattamento statistico dei dati sono necessari per l'interpretazione delle misure sperimentali
- La sensibilità ai fenomeni di interesse stabilisce quanto critica sia la modalità d'analisi sui risultati

- Restano comunque tre gli step necessari alla misura:
 - studio della componente stocastica delle distribuzioni
 - valutazione e trattamento del background
 - valutazione e trattamento degli errori sistematici