Confidence Interval Estimation Part I

<u>Roman Kogler</u> (roman.kogler@desy.de)

Terascale Statistics School 2023 DESY Hamburg, July 2023





What do we learn?



What have we learned?

Now?



What have we learned?

Now? Does the data favour H1?



What have we learned?

Now?

Does the data favour H1?

We can't say unless we know the uncertainty (confidence intervals)!

Confidence Intervals

Start with a simple question

We know a true value x and the corresponding pdf, and we want to know in which interval a certain amount of measurements \hat{x}_i will fall.

The value of the measurement \hat{x}_i lies in the interval [X-, X+] in "CL" % of the time.

Corresponds to the statement "the interval [X-, X+] has CL% confidence level"



$$p(X_{-} \ge x_i \ge X_{+}) = \int_{X_{-}}^{X_{+}} f(x)dx = CL$$

More Realistic Problem - Gaussian Case

In a real experiment, x is unknown - need to estimate it

• Gaussian distributed estimator, cumulative probability (CLT):

$$G(\hat{x}, x, \sigma_{\hat{x}}) = \int_{-\infty}^{\hat{x}} \frac{1}{\sqrt{2\pi\sigma_{\hat{x}}}} \exp\left(-\frac{(x'-x)^2}{2\sigma_{\hat{x}}}\right) dx'$$

with \hat{x} being the measured value and σ_x the standard deviation (resolution) of the measurement

• Confidence interval [a,b] can be obtained through solving

$$\begin{split} \alpha &= 1 - G(\hat{x}, a, \sigma_{\hat{x}}) = 1 - \phi \left(\frac{\hat{x} - a}{\sigma_{\hat{x}}}\right) & \text{lower CL is } I - \alpha \\ \beta &= G(\hat{x}, b, \sigma_{\hat{x}}) = \phi \left(\frac{\hat{x} - a}{\sigma_{\hat{x}}}\right) & \text{upper CL is } I - \beta \end{split}$$

Gaussian Case

We obtain

$$a = \hat{x} - \sigma_{\hat{x}} \phi^{-1} (1 - \alpha) \qquad \Phi^{-1} \text{ are quantiles:}$$
$$b = \hat{x} + \sigma_{\hat{x}} \phi^{-1} (1 - \beta)$$

Results in the typical "I σ error bar":

$$[a,b] = [\hat{x} - \sigma_{\hat{x}}, \hat{x} + \sigma_{\hat{x}}]$$



Note:

- More complicated if $\sigma_{\hat{\chi}}$ is unknown, have to rely on $\hat{\sigma}_{\hat{\chi}}$
- If \hat{x} does not follow a Gaussian pdf

Coverage of Uncertainties

- Parameter μ_m estimated from data sample **x**, errors σ_1 and σ_2 (example: fitting procedure)
- The data follow a probability density $p(x|\mu)$ with given value of parameter μ (fixed, true)
- The result $\mu_{m-\sigma_1}^{+\sigma_2}$ determines an interval $[\mu_m \sigma_1, \mu_m + \sigma_2]$ (the region between error bars).
- What ist the coverage probability of the interval? (short: coverage)

Coverage of Uncertainties

- Parameter μ_m estimated from data sample **x**, errors σ_1 and σ_2 (example: fitting procedure)
- The data follow a probability density $p(x|\mu)$ with given value of parameter μ (fixed, true)
- The result $\mu_{m-\sigma_1}^{+\sigma_2}$ determines an interval $[\mu_m \sigma_1, \mu_m + \sigma_2]$ (the region between error bars).
- What ist the coverage probability of the interval? (short: coverage)
- Coverage $C(\mu)$ is a function of the (unknown) parameter μ , defined as the probability that, with $\mu_m \equiv \mu(n)$, $\sigma_1 \equiv \sigma_1(n)$ and $\sigma_2 \equiv \sigma_2(n)$

$$\mu_m - \sigma_1 \le \mu \le \mu_m + \sigma_2$$

 $C(\mu)$ is the probability that an experiment will obtain an interval that includes, or "covers", the true value μ .

Frequentist require $C(\mu) = C_0$ or at least $C(\mu) \ge C_0$ with $C_0 = 0.6827$, corresponding to 1σ of a Gaussian distribution

Roman Kogler

Poisson Distributed Data

- The origin of a statistical error in particle physics is most often the observation of Poisson distributed data.
- The probability of observing n events, if the mean value is $\mu,$ is given by

$$p(n,\mu) = \frac{\mu^n e^{-\mu}}{n!}$$

- n is a random variable and μ is fixed.
- There are several different methods for the error-bar scheme.

Joel G. Heinrich, Coverage of error bars for Poisson data, <u>CDF/MEMO/STATISTICS/PUBLIC/6438</u> (2003)

Roman Kogler

Choice of $\sigma = \sqrt{n}$

- Best estimate of µ: Observation n
- Calculate the uncertainty as \sqrt{n} :

$$n - \sqrt{n} < \mu < n + \sqrt{n}$$

• Common choice, motivated by Poisson distribution,

$$\sqrt{V[x]} = \sqrt{\mu} = \sqrt{n}$$

- Study the results and implications in the exercise
- Serious undercoverage

• Pearson's χ^2 is given by

$$\chi^2(\mu, n) = \frac{(n-\mu)^2}{\mu}$$

• Pearson's χ^2 is given by

$$\chi^2(\mu, n) = \frac{(n-\mu)^2}{\mu}$$

• After observing n events, the estimate is obtained by minimising χ^2 ,

i.e. $\mu_m = n$ with $\chi^2 = 0$

 Uncertainties are obtained from the interval such that

 $\chi^2(\mu_m, n) < \Delta$

and we get

$$\sigma_1 = \sqrt{n\Delta + \Delta^2/4} - \Delta/2$$

$$\sigma_2 = \sqrt{n\Delta + \Delta^2/4} + \Delta/2$$

• Pearson's χ^2 is given by

$$\chi^2(\mu, n) = \frac{(n-\mu)^2}{\mu}$$

• After observing n events, the estimate is obtained by minimising χ^2 ,

i.e. $\mu_m = n$ with $\chi^2 = 0$

 Uncertainties are obtained from the interval such that

 $\chi^2(\mu_m, n) < \Delta$

and we get

$$\sigma_1 = \sqrt{n\Delta + \Delta^2/4 - \Delta/2}$$

$$\sigma_2 = \sqrt{n\Delta + \Delta^2/4} + \Delta/2$$



many discontinuities, values between 0 (no coverage) and 1 (overcoverage)

Standard choice of

 $\Delta = 1.0$

- minimum coverage is
 I/I.5 = 0.5518
 obtained at µ≈ I
- mean value around 0.68
- example: n = 6 then we would say

 $\mu = 6^{+3}_{-2}$

but μ is still unknown, it could be close to one we only know with certainty C \geq 0.5518



Standard choice of

 $\Delta = 1.0$

- minimum coverage is
 I/I.5 = 0.5518
 obtained at µ≈ I
- mean value around 0.68
- example: n = 6 then we would say

 $\mu = 6^{+3}_{-2}$

but μ is still unknown, it could be close to one we only know with certainty C \geq 0.5518



undercoverage can be fixed by choosing $\Delta = 1.5$, but then the mean value is around C = 0.78 (overestimation of uncertainties)

Likelihood Intervals

 Error estimate based on the value of the likelihood (case of a maximum likelihood fit)

 $-2\ln\lambda(\mu, n) =$ $2[(\mu - n) + n\ln(n/\mu)]$

• With the uncertainty defined by

 $-2\ln\lambda(\mu,n) \le \Delta$

Likelihood Intervals

 Error estimate based on the value of the likelihood (case of a maximum likelihood fit)

 $-2\ln\lambda(\mu, n) =$ $2[(\mu - n) + n\ln(n/\mu)]$

- With the uncertainty defined by $-2\ln\lambda(\mu,n) \leq \Delta$
- Minimum is 0.3033 in the vicinity of $\mu = 0.5$

С $\Delta = 1.0$ 1.0 0.9 0.8 0.7 0.6 0.5 0.4 0.3 0.2 0.1 0.0 2 3 0 9 10 13 14 15 16 17 18 19 11 12 Coverage (C) vs μ : $-2\ln\lambda < 1$ (C $\rightarrow 0.6827$ as $\mu \rightarrow \infty$)

undercoverage can be fixed by choosing $\Delta = 2.581$, but then the mean value is around C = 0.89 (even worse overestimation of uncertainties)

5 7 8 9 10 11 12 13 14 15 16 17 18 19 20 μ C) vs μ : w < 1 (C $\rightarrow 0.6827$ as $\mu \Rightarrow \infty$) ISt's Central Intervals

 Coverage achieved by the frequentist's 68.27% interval choice:

$$\sum_{k=0}^{n-1} \frac{e^{-\mu} \mu^k}{k!} \ge \frac{1-C_0}{2}$$
$$\sum_{n+1}^{\infty} \frac{e^{-\mu} \mu^k}{k!} \ge \frac{1-C_0}{2}$$

- Minimum coverage guaranteed to be larger than $C_0 = 0.6827$
- Coverage larger on average
- Special case for n = 0:

$$\mu \le \ln \frac{2}{1 - C_0}$$

C) vs μ: w < 1 (C \rightarrow 0.6827 as μ \Rightarrow ∞) ist's Central Intervals

 Coverage achieved by the frequentist's 68.27% interval choice:

$$\sum_{k=0}^{n-1} \frac{e^{-\mu} \mu^k}{k!} \ge \frac{1-C_0}{2}$$
$$\sum_{n+1}^{\infty} \frac{e^{-\mu} \mu^k}{k!} \ge \frac{1-C_0}{2}$$

- Minimum coverage guaranteed to be larger than $C_0 = 0.6827$
- Coverage larger on average
- Special case for n = 0:

 $\mu \le \ln \frac{2}{1 - C_0}$



more complicated constructions possible, but not widely used exact coverage $C(\mu) = C_0$ never possible

Back to the original problem: Quote confidence interval for true parameter μ. Get confidence belt via Neyman construction:

For a given true parameter μ the distribution of x is known

$$P(X \in [x - \sigma, x + \sigma]) = CL$$

so called "acceptance interval"



measurement

Back to the original problem: Quote confidence interval for true parameter μ. Get confidence belt via Neyman construction:

- For a given true parameter μ the distribution of x is known
- 2. For each μ we calculate the mean x



Back to the original problem: Quote confidence interval for true parameter μ. Get confidence belt via Neyman construction:

- For a given true parameter μ the distribution of x is known
- 2. For each μ we calculate the mean x
- 3. For each μ and a given CL we can calculate X_1 and X_2



measurement

Back to the original problem: Quote confidence interval for true parameter μ. Get confidence belt via Neyman construction:

- For a given true parameter μ the distribution of x is known
- 2. For each μ we calculate the mean x
- 3. For each μ and a given CL we can calculate X_{-} and X_{+}
- For a measured x we can get the confidence interval [μ₊, μ₋]



The confidence belt contains CL % of the expected measurements.

Meaning: If μ would be $\mu_+(\mu_-)$ than the probability to observe x or less (x or more) is CL/2 %. This is not a statement about μ , but about X_- and X_+ .

(wrong: construct the confidence interval by a Gaussian around x)

Neyman Construction: Poisson Distribution

We observed k events, what are the CL limits on μ ?



Roman Kogler