# INTERVAL ESTIMATION

• Goal. Confidence interval.

In interval estimation, we want to find the range $\theta_a \leq \theta \leq \theta_b$ which contains the true value $\theta_0$ with probability $\beta$.

Such interval is called "confidence interval" with probability content $\beta$.

Typically, we choose $\beta = 68.3\%$ and call it 1 standard deviation error. However, the 68% interval corresponds to $\pm 1$ STD only for a Gaussian distro.

Given an observation $x$ from a PDF $f(x|\theta)$, the probability content $\beta$ of the region $[a, b]$ in $x$-space is:

$$\beta = P(a \leq x \leq b) = \int_a^b f(x|\theta) \, dx$$

If $f(x|\theta)$ and the parameter $\theta$ are known, one can always compute $\beta$ given $a$ and $b$.

If the parameter $\theta$ is unknown, we need to find another variable

$$z = z(x, \theta)$$

such that the PDF of $z$ is independent of $\theta$: ~~$f(z|\theta) = f()$~~

$$f(z|x,\theta) = f(z|x)$$

If this can be found, we can find the optimal range $[\theta_a, \theta_b]$ in $\theta$ space such that:
$$P(\theta_a < \theta < \theta_b) = \beta$$

This interval $[\theta_a, \theta_b]$ is called "confidence interval".

A method which yields ~~an~~ such an interval $[\theta_a, \theta_b)$ is said to possess the property of <u>coverage</u>.

Notice that: → $\theta_0$ is an unknown constant

→ $\theta_a$ and $\theta_b$ are functions of $x$, not of $\theta$.

• Normally distributed data

Let $f(x|\theta)$ be a normal distribution: $\quad f(x|\theta) = \frac{1}{\sqrt{2\pi}\,\sigma^2} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$

→ If $\mu$ and $\sigma$ are known:

$$\beta = P(a \leq x \leq b) = \int_a^b f(x|\theta) = \int_a^b \frac{1}{\sqrt{2\pi}\,\sigma} \exp\left(-\frac{(x-\mu)^2}{\sigma^2}\right) dx$$

$$= \Phi\left(\frac{b-\mu}{\sigma}\right) - \Phi\left(\frac{a-\mu}{\sigma}\right) \qquad \text{with } \Phi = \text{cumulative of Gaussian}$$

→ If $\mu$ is unknown, we can't compute the probability content of the interval $[a,b]$. But we can compute the probability $\beta$ that $x$ lies in some interval relative to the unknown mean, say $[\mu+c, \mu+d]$.

Let's define: $\quad y = \frac{x-\mu}{\sigma}$

We have: $\quad \beta = P(\mu+c \leq x \leq \mu+d) = \int_{\mu+c}^{\mu+d} N(\mu, \sigma^2)\, dx'$

$$= \int_{c/\sigma}^{d/\sigma} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{y^2}{2}\right) dy = \Phi\left(\frac{d}{\sigma}\right) - \Phi\left(\frac{c}{\sigma}\right)$$

We can invert it to obtain: $\quad \beta = P(x-d \leq \mu \leq x-c)$

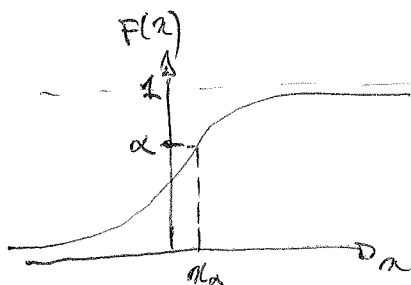$\hookrightarrow$ This is still a probability statement about $x$. $\mu$ is an unknown constant.

• $\alpha$ - point

Given a random variable $x$ with PDF $f(x)$ and CDF $F(x)$,
The $\alpha$ point is defined as:

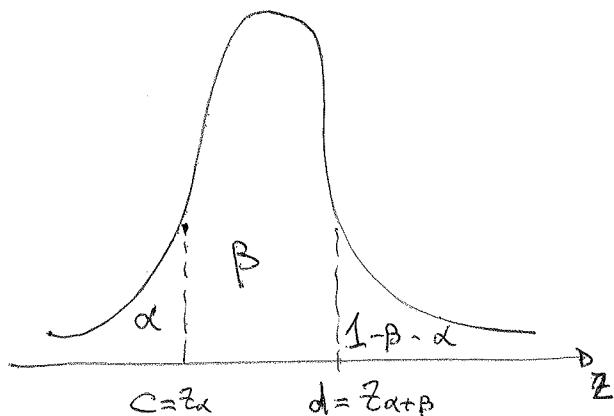$$\int_{-\infty}^{x_\alpha} f(x)dx = F(x_\alpha) = \alpha$$

• Confidence intervals for the mean of a Gaussian

For any Gaussian, we can re-define: $z = \frac{x-\mu}{\sigma}$

~~and~~ which is a standard-normal variable: $f(z) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{z^2}{2}\right)$

Estimating the interval $[c,d]$ so that $P(c \leq z \leq d) = \beta$ is equivalent to ~~the~~ finding $[z_\alpha, z_{\alpha+\beta}]$:



→ There's infinite choices of the interval!

→ A standard choice is the central interval symmetric around zero, so that:
$$\alpha = \frac{1-\beta}{2}$$

| $\beta = \frac{1-\alpha}{2}$ | $z_\alpha$ | $z_{\alpha+\beta}$ | |
|---|---|---|---|
| 0.6827 | −1 | +1 | → ±1$\sigma$ |
| 0.9 | −1.65 | +1.65 | |
| 0.95 | −1.96 | +1.96 | |
| 0.9545 | −2 | +2 | → ±2$\sigma$ |
| 0.9973 | −3 | +3 | → ±3$\sigma$ |

• Confidence intervals for several parameters

Suppose we have an n-dimensional Gaussian:
$$f(\vec{x}\,|\vec{\theta}) = \frac{1}{(2\pi)^{\frac{n}{2}} \sqrt{|C|}} \exp\left(-\frac{1}{2}(\vec{x}-\vec{\theta})^T C^{-1} (\vec{x}-\vec{\theta})\right)$$
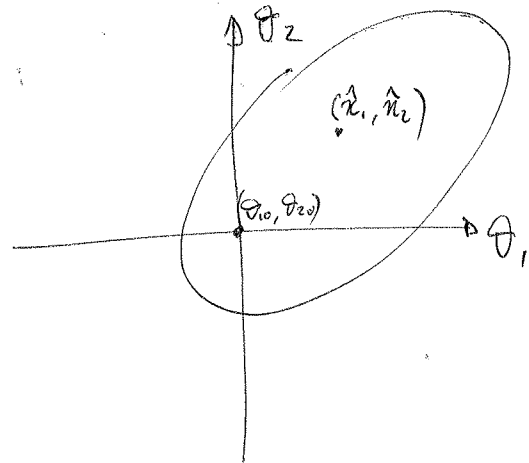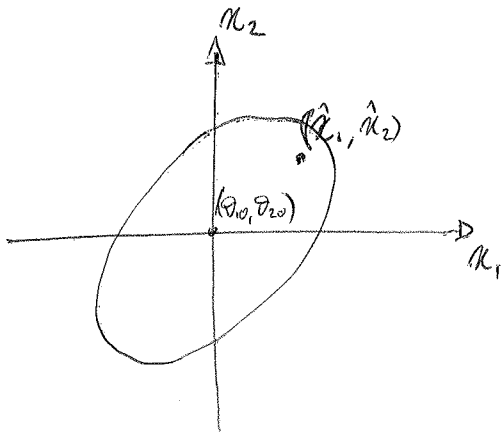
Each $x_i$ is normal, therefore $Q(\vec{x},\vec{\theta}) = (\vec{x}-\vec{\theta})^T C^{-1} (\vec{x}-\vec{\theta})$ is a $\chi^2(n)$ distribution, and does not depend on $\vec{\theta}$:
$$Q(\vec{x},\vec{\theta}) = Q(\vec{x})$$

We can write: $P\left[Q(\vec{x}, \vec{\theta}) \leq K_\beta^2\right] = \beta$

Where $K_\beta^2$ is the $\beta$-point of the $\chi^2(n)$ distribution

This region is defined in $\vec{x}$-space as: $Q(\vec{x}, \vec{\theta}) \leq K_\beta^2$
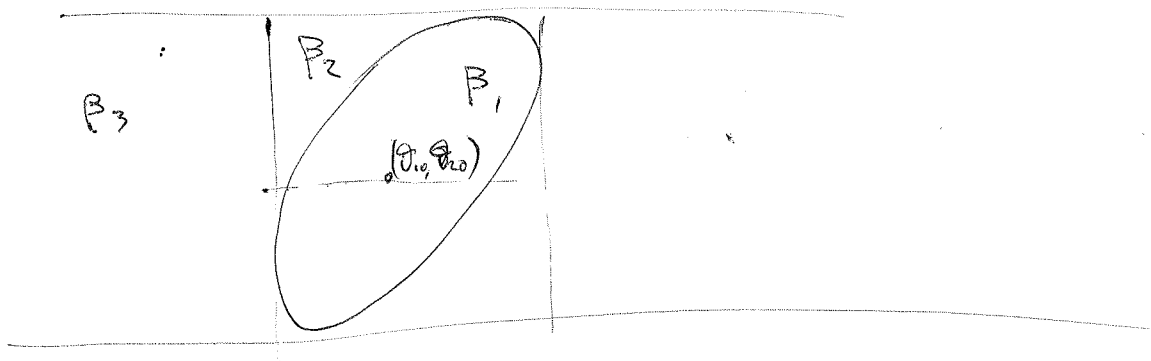and is a hyperellipsoid of constant probability:



This holds for any dimension.

Notice that we have 3 ways to quote an interval for e.g. $x_2$:

~~a) $P\left[(\theta_1 - \sigma_1 \leq x_1 \leq \theta_1 + \sigma_1) \wedge (\theta_2 - \sigma_2 \leq x_2 \leq \theta_2 + \sigma_2)\right] = \beta_2$~~

a) $P\left[Q(\vec{x}, \vec{\theta}) \leq \vec{K}_\beta\right] = \beta_a \Rightarrow$ ellipse

b) $P\left[(\theta_1 - \sigma_1 \leq x_1 \leq \theta_1 + \sigma_1) \wedge (\theta_2 - \sigma_2 \leq x_2 \leq \theta_2 + \sigma_2)\right] = \beta_b \Rightarrow$ square

c) $x_2 = \theta_{20} \pm \sigma_2 \Rightarrow P\left(\theta_2 - \sigma_2 \leq x_2 \leq \theta_2 + \sigma_2\right) = \beta_c \Rightarrow$ band



The probability content is different for the 3 cases, and corresponds to the 1-dim case just for case Ⓒ !

$\boxed{\text{INT } 4}$

- **Second derivative matrix**

Assume $\vec{x}$ has an n-dim Gaussian PDF.

One can prove that the n-dim covariance matrix $C$ can be obtained from the inverse of the 2$^{nd}$ order partial derivative matrix of $-\ln \mathcal{L}$:

$$C_{ij}^{-1} = - \frac{\partial \ln \mathcal{L}(\vec{x}|\vec{\vartheta})}{\partial \vartheta_i \, \partial \vartheta_j}$$

This covariance matrix gives an n-dim elliptic contour with the correct coverage only if the PDF is exactly Gaussian!

This is the "standard" classical method used to compute uncertainties in common fitting algorithms, e.g. Migrad/Hesse of Minuit/ROOT.


- **Log-Likelihood scan**

Another common method consists in taking a scan of $-2\ln \mathcal{L}$ around its minimum value, $-2\ln \mathcal{L}_{max}$.
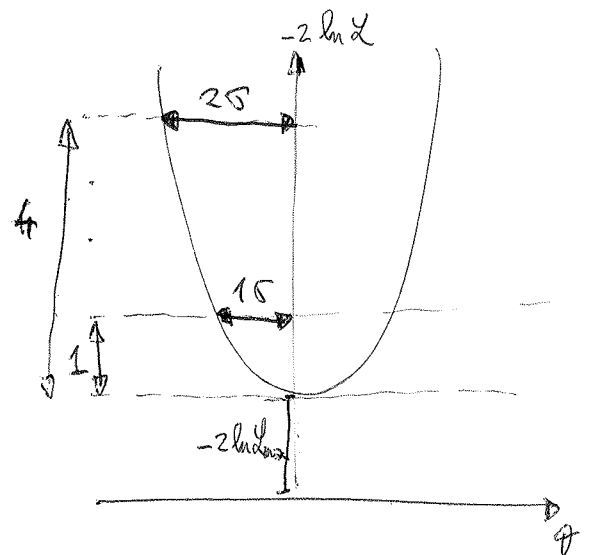
→ For a Gaussian 1-dim distribution:

$$\ln \mathcal{L}(x|\mu) = \ln C - \frac{(\mu-x)^2}{2\sigma^2} \qquad \Rightarrow \text{parabola in } \mu$$

$$-2\ln \mathcal{L} = -2\ln \mathcal{L}_{max} + \frac{(\mu-x)^2}{\sigma^2}$$

The intercept at $-2\ln \mathcal{L} = -2\ln \mathcal{L}_{max} + 1$ provides the $\pm 1\sigma$ interval.
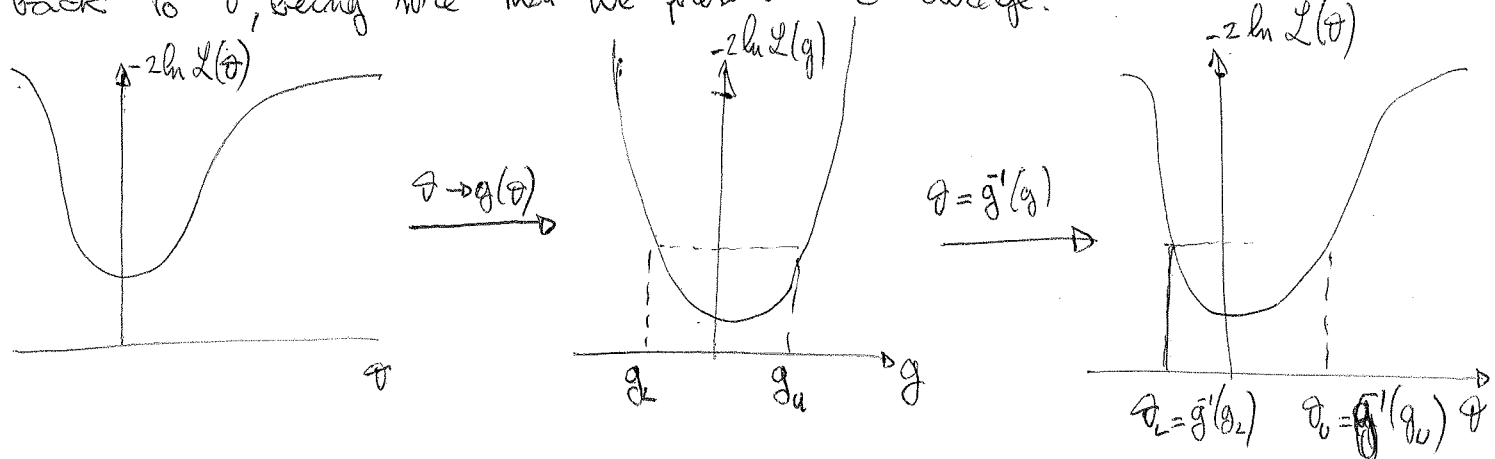
The intercept at $+4$ provides the $\pm 2\sigma$, and so on

→ For a non-Gaussian PDF, the $\ln \mathcal{L}$ is not parabolic.
However, the maximum $\mathcal{L}$ estimate is invariant, so instead of estimating $\theta$, one can estimate any (monotonic) function of $\theta$.

So one can find (in principle) a transformation $\theta \to g(\theta)$ that makes the $\ln \mathcal{L}$ parabolic, compute the confidence interval on $g$, and transform back to $\theta$, being sure that we preserve the coverage.



In practice, we can ~~use the same metho~~. find the interval $[\theta_L, \theta_u]$ without actually doing the transformations $g$ and $g^{-1}$, but simply finding the intercepts at $-2 \ln \mathcal{L}_{max} + 1$, as we did for the Gaussian case.

Note however that the uncertainties are now asymmetric.

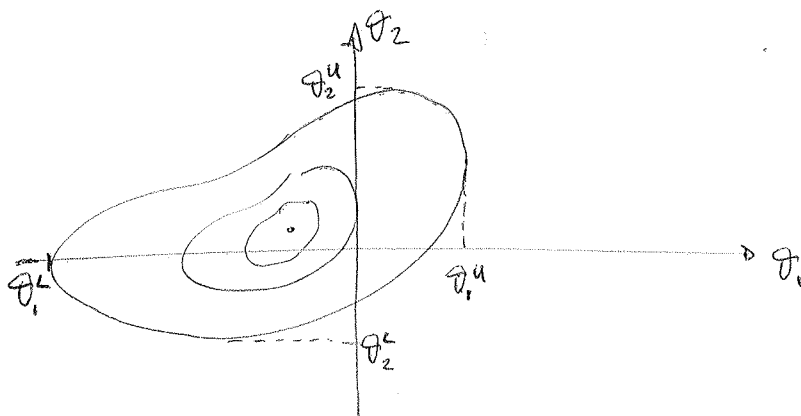This is the method used by MINOS in ROOT/Minuit.

Notice that we have made the experimental $\mathcal{L}$ distribution Gaussian in the parameter instead of finding a transformation that would make the Theoretical distribution Gaussian (i.e. using the true value $\theta_o$). To order $1/N$ this is the same, but for small samples it is not exact!

- $\mathcal{L}$ scan in dim > 1 : Profile likelihood

For dim > 1 we have :
$$\ln \mathcal{L}_{\vec{\theta}} (\vec{n} | \vec{\theta}) = \ln \mathcal{L}_{max} - \frac{1}{2} \chi^2_P (n)$$
$$\hookleftarrow -2 \ln \mathcal{L} (\vec{n} | \vec{\theta}) = -2 \ln \mathcal{L}_{max} + \chi^2_P (n)$$

In principle, we can compute the contours :



Notice that :
→ The inner contour is more nearly elliptical than the outer ones
→ The coverage is improved with respect to the Gaussian approximation
→ This is still an approximation valid for large n.

- Variance of Transformed variables, aka Error propagation

Suppose we have a variable $\vec{n}$ with a given PDF$_n$ and mean $\vec{\mu}$ an variance $V(\vec{n})$
Suppose we want to compute the variance of the transformed variable $\vec{y} = \vec{y}(\vec{n})$,
and that the function can be expanded in Taylor series around $\vec{\mu}$ :
$$\vec{y}(\vec{n}) = \vec{y}(\vec{\mu}) + \sum_i (x_i - \mu_i) \frac{\partial y}{\partial \mu_i} + \dots$$

The expectation value of $y$ is : $\bar{y} = y(\vec{\mu})$

The variance is :
$$V(y) = E[y - E(y)]^2$$
$$\simeq E\left[ \sum_i (x_i - \mu_i) \frac{\partial y}{\partial x_i} \right]^2$$
$$\simeq \sum_i \sum_j \frac{\partial y}{\partial x_i} \frac{\partial y}{\partial x_j} E[(x_i - \mu_i)(x_j - \mu_j)]$$
$$\simeq \sum_{i,j} \frac{\partial y}{\partial x_i} \frac{\partial y}{\partial x_j} \text{Cov}(x_i, x_j)$$

Int. 7

If The variables $x_i$ are independent: $\quad \sigma_y^2 \simeq \sum_i \left(\frac{\partial y}{\partial x_i}\right)^2 \sigma_i^2 \quad$ with $\sigma_i^2 = V(x_i)$

Example: linear combination of variables:

$$z = ax + by$$

$$\sigma_z^2 = \left(\frac{\partial z}{\partial x}\right)^2 \sigma_x^2 + \left(\frac{\partial z}{\partial y}\right)^2 \sigma_y^2 + 2\sigma_x\sigma_y \; 2 \frac{\partial z}{\partial x}\frac{\partial z}{\partial y} cov(x,y)$$

$$= a^2 \sigma_x^2 + b^2 \sigma_y^2 + 2ab \, \rho \sigma_x \sigma_y$$

Example: Product of variables: $\quad \left(\frac{\sigma_{xy}}{xy}\right)^2 = \frac{\sigma_x^2}{x^2} + \frac{\sigma_y^2}{y^2} + \frac{2\rho\sigma_x\sigma_y}{xy}$

→ Product and ratio of uncorrelated variables:

$$\left(\frac{\sigma_{xy}}{xy}\right)^2 = \left(\frac{\sigma_{x/y}}{x/y}\right)^2 = \frac{\sigma_x^2}{x^2} + \frac{\sigma_y^2}{y^2}$$

→ Logarithm of variable: $\quad \sigma_{\ln x} = \frac{\sigma_x}{x}$

Caveats: → Valid only if The Transformation of variable can be expanded
in Taylor series

Alternatively, one can use numerical methods, to extract The PDF of The
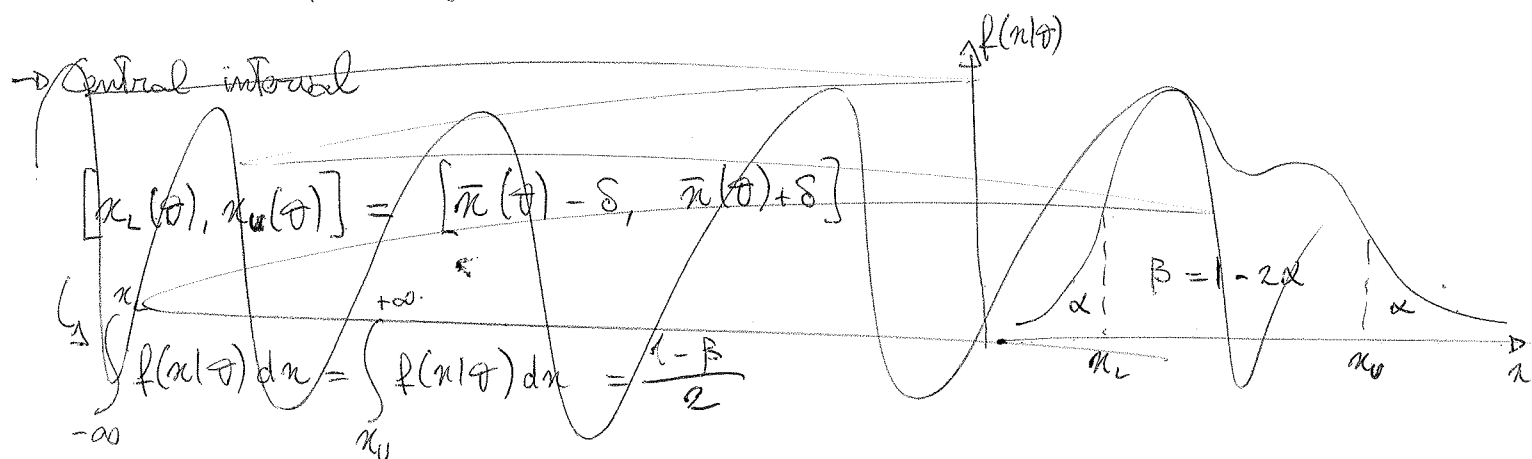Transformed variable, and compute the confidence interval on it.

- Confidence intervals for any PDF: Ordering rules

The approximation of $-2\ln L$ with a Gaussian or its excursion around its minimum guarantee an exact coverage only for a small set of cases, and in particular for large $n$.

Here we'll see a general approach.
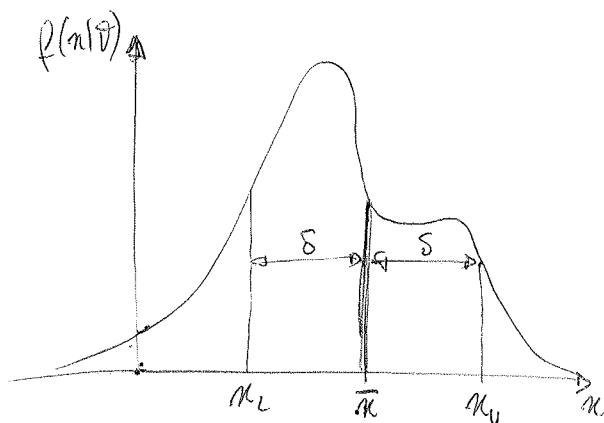
Suppose we have a variable $x$ with PDF $f(x|\theta)$.

In general, $f$ is not symmetric, so we need to decide how to compute an interval corresponding to some predefined probability content $\beta$.
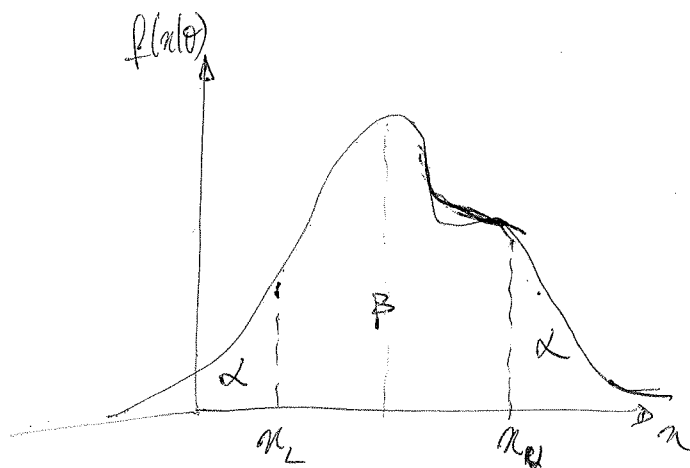
→ Central interval

$$[x_L(\theta), x_U(\theta)] = [\bar{x}(\theta) - \delta, \ \bar{x}(\theta) + \delta]$$

$$\int_{-\infty}^{x_L} f(x|\theta)\,dx = \int_{x_U}^{+\infty} f(x|\theta)\,dx = \frac{1-\beta}{2}$$

$\beta = 1 - 2\alpha$

→ Central interval

$$[x_L(\theta), x_U(\theta)] = [\bar{x}(\theta) - \delta, \ \bar{x}(\theta + \delta)]$$
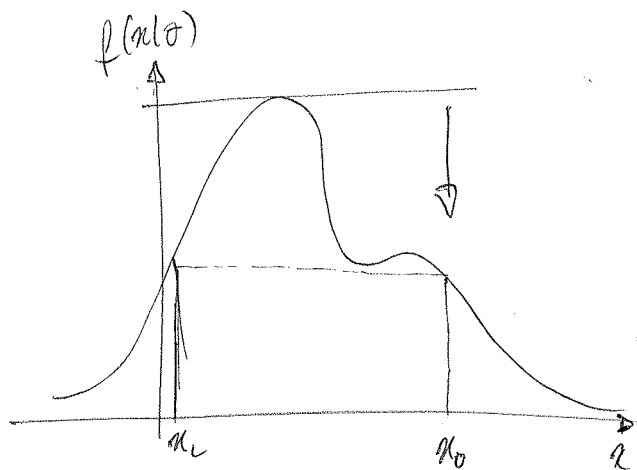
$\bar{x}$ could be the mean or the mode

→ Equal areas

$$\int_{-\infty}^{x_L} f(x|\theta)\,dx = \int_{x_U}^{+\infty} f(x|\theta)\,dx = \frac{1-\beta}{2}$$

→ Shortest interval

$$f(x_L | \vartheta) = f(x_U | \vartheta) \wedge \int_{x_L}^{x_R} f(x | \vartheta) \, dx = \beta$$

↳ Start with a horizontal line $y = f_{max}(x | \vartheta)$, and lower until you satisfy such conditions.
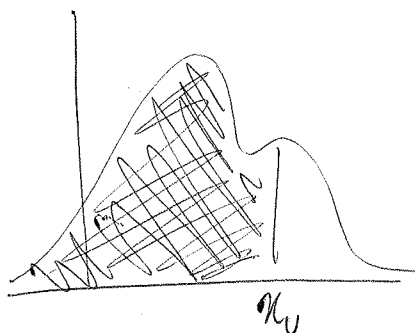


→ Lower limit

$$\int_{-\infty}^{x_L} f(x | \vartheta) \, dx = 1 - \beta$$



→ Upper limit

$$\int_{x_U}^{+\infty} f(x | \vartheta) \, dx = 1 - \beta$$



Example : Ordering rules

- Neyman confidence belt
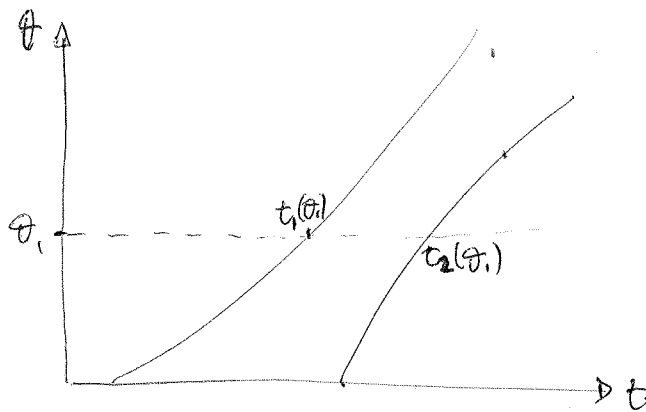
Take a variable $x$ with PDF $f(x|\theta)$ and $\theta$ unknown.

~~Assume $x$ could be an estimator of the parameter $\theta$.~~

Suppose $t(x)$ is some function of the data.

We can write: $\beta = P(t_1 \leq t \leq t_2) = P(t_1(\theta) \leq t \leq t_2(\theta))$

$$\overset{!}{=} \int_{t_1}^{t_2} f(t|\theta) \, dt$$

Assume that we have a way to determine $t_1$ and $t_2$ for each value of $\theta$.
Such values form two curves in the $(t, \theta)$ space: ~~Then~~
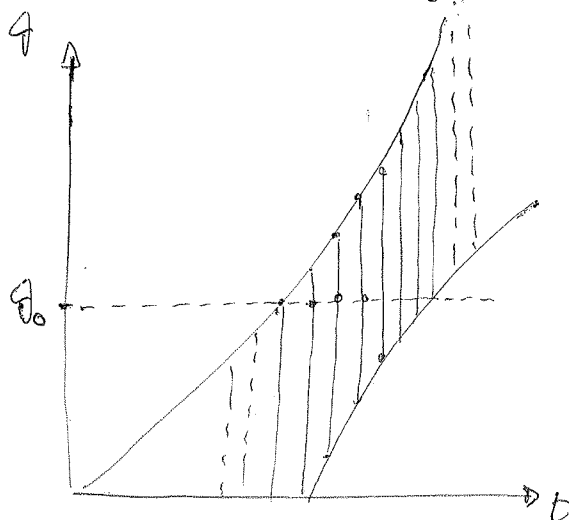
From the PDF of $t$:

$f(t(x)|\theta)$ we get $t_1$ and $t_2$



The ~~space~~ space between these curves is called the "Neyman confidence belt".

From this belt, ~~we want to~~ given a specific measured value $t_0$, we want to ~~extract~~ extract an interval on the parameter $\theta$.

Suppose $\theta_0$ is the true, unknown value of $\theta$.

If we repeat the measurement many times, a fraction $\beta$ of the measurements will fall ~~also~~ in $[t_1(\theta_0), t_2(\theta_0)]$ by ~~our~~ definition:

By construction, if I can "invert" the belt and compute $\theta_1(t)$ and $\theta_2(t)$ for any of the measured values of $t$, the true value $\theta_0$ will fall in $[\theta_1, \theta_2]$ exactly a fraction $\beta$ of the times:

$$P(\theta_1 \leq \theta \leq \theta_2) = \beta$$

Notice that: → $\theta_1$ and $\theta_2$ depend only on the data.

→ If the observable $t$ is discrete, we need to construct the belt so that:

$$P(t_1 \leq t \leq t_2) \geq \beta$$

Therefore $P(\theta_1 \leq \theta \leq \theta_2) \geq \beta$

In such a case, we say that the method "overcovers".

Examples: → Neyman belt for Gaussian case
→ Neyman belt for binomial interval

• Flip-flopping problem

Suppose a variable $x$ has a PDF $f(x|\mu) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$ with known $\sigma$.

Suppose we have a physical theory that limits $\mu$ to be non-negative: $\mu \geq 0$.
This is for example the case of $\mu$ representing a mass.

Suppose we do some measurement of $x$, and that our instrument is subject to fluctuations so that $x$ can be negative.
This is the case of an analogic scale.
The quoted ~~central interval must~~ central value for $\mu$ must always be $\geq 0$.
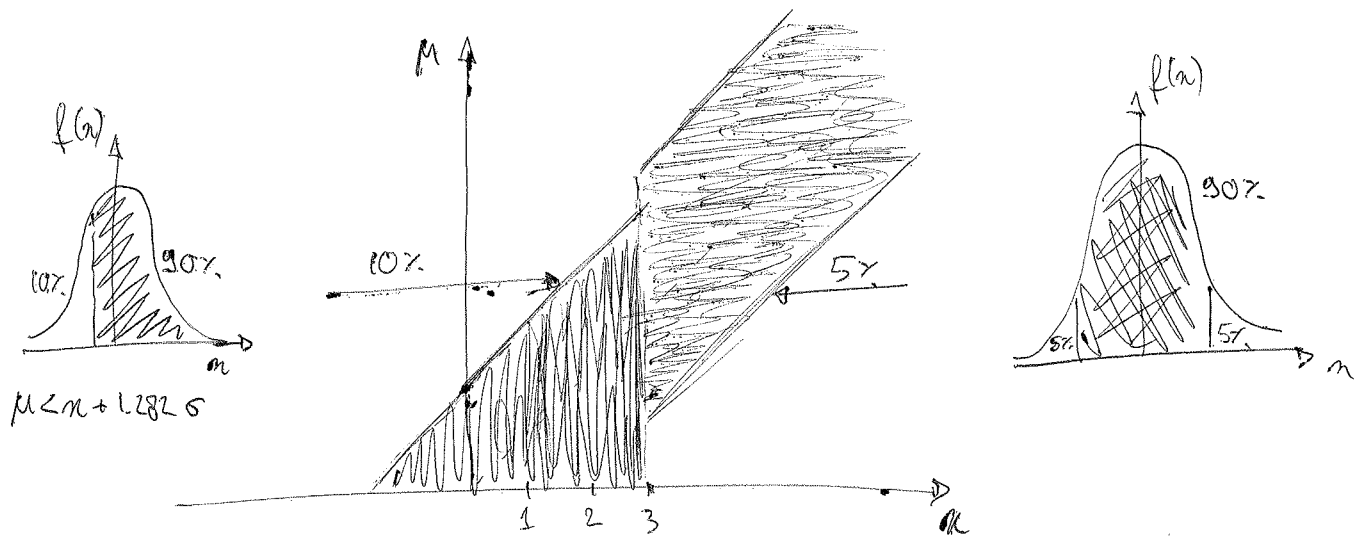Imagine we decide to quote:

$$\hat{\mu}(x) = \begin{cases} x & \text{if } \frac{x}{\sigma} \geq 3 \\ 0 & \text{if } \frac{x}{\sigma} < 3 \end{cases}$$

The corresponding 90% interval will be:

$$[\mu_1, \mu_2] = \begin{cases} [\hat{\mu} - 1.65\sigma, \hat{\mu} + 1.65\sigma] & \text{if } \frac{x}{\sigma} \geq 3 \\ [0, \hat{\mu} + 1.282\sigma] & \text{if } \frac{x}{\sigma} < 3 \end{cases}$$

The corresponding belt would be:



$$\mu < x + 1.282\,\sigma$$

For some values of $\mu$, e.g. $\mu = 2.5$, we have a coverage of 85% only!

Notice that: → The coverage is a property of the method, not of the a particular interval

→ The flip-flopping issue arises from the fact that our ordering rule depends on the outcome of the measurement.
Feldman and Cousin showed that this should not be done!

Example: Flip-flopping

• Unified Feldman - Cousins approach

Recall That we can build a method with exact coverage in an infinite number of ways, just by changing the ordering rule.

So we can look for an ordering rule with a smooth Transition from a central interval To an upper limit To avoid the flip-flopping issue.

The Feldman-Cousins approach uses The $R$-ratio as an ordering rule:

$$\lambda(x|\theta_0) = \frac{\mathcal{L}(x|\theta_0)}{\mathcal{L}(x|\hat{\theta})}$$, where $\hat{\theta}$ is the value of $\theta$ That maximizes $\mathcal{L}$.
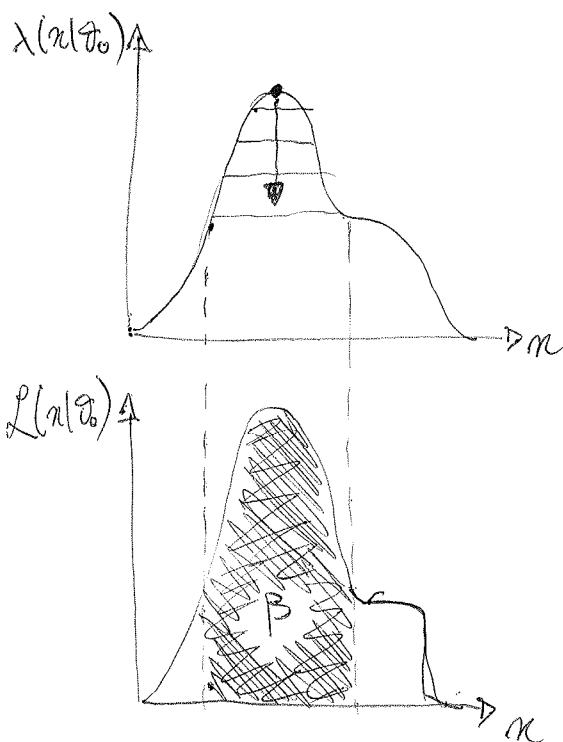
To build The FC confidence belt we proceed as follows:

1) Find $\hat{\theta}$

~~2) Start from belt For each~~

2) Fix $\theta$ To some value $\theta_0$.
   Then start from $\hat{x} = \underset{x}{\text{Max}} \ \lambda(x|\theta_0)$ and move To The left and right until we get a coverage $\beta_0$ on $\mathcal{L}(x|\theta_0)$



3) Repeat for all values of $\theta$ (in The physical range).

* Examples → FC Belt for Gaussian
    → FC belt for electron neutrino mass

• ~~Analytical~~ Numerical FC belt calculation for Gaussian

Recall the flip-flopping case, where we had $x$ with a PDF:

$$f(x|\mu) = \frac{1}{\sqrt{2\pi}} \exp\left(- \frac{(x-\mu)^2}{2}\right)$$

The value ~~that~~ $\hat{\mu}$ that maximizes $f(x|\mu)$ given some measured $x$ is:

$$\hat{\mu}(x) = \max\{x, 0\}$$

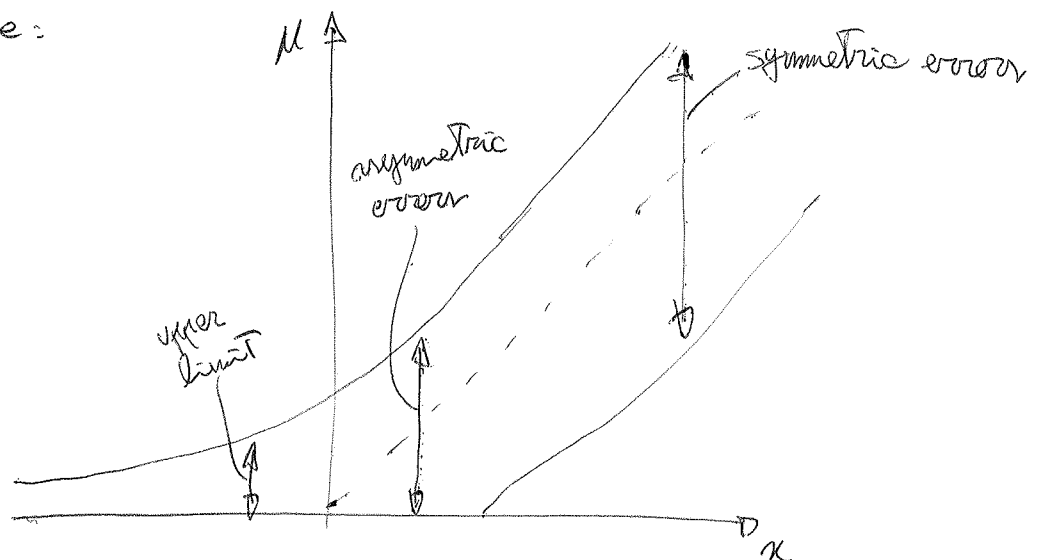The PDF for $x$, using the max-$\mathcal{L}$ estimate for $\mu$ is:

$$f(x|\hat{\mu}(x)) = \begin{cases} \dfrac{1}{\sqrt{2\pi}} & \text{if } x \geq 0 \\[3mm] \dfrac{1}{\sqrt{2\pi}} \exp\left(-\dfrac{x^2}{2}\right) & \text{if } x < 0 \end{cases}$$

The likelihood ratio becomes:

$$\lambda(x|\mu) = \frac{f(x|\mu)}{f(x|\hat{\mu}(x))} = \begin{cases} \exp\left(- \dfrac{(x-\mu)^2}{2}\right) & \text{if } x \geq 0 \\[3mm] \exp\left(x\mu - \dfrac{\mu^2}{2}\right) & \text{if } x < 0 \end{cases}$$

At this point, we can find the interval $[\mu_1, \mu_2]$ numerically for any value of $\mu$.
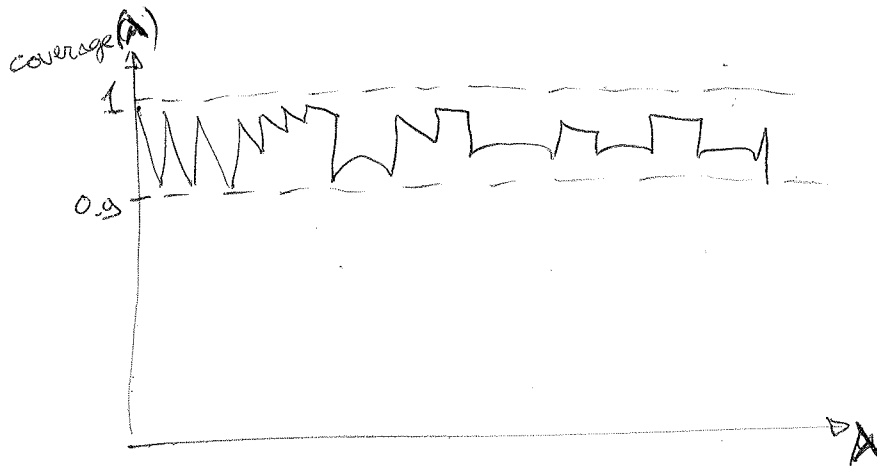
The result will be:

- FC for discrete data

  For discrete data, the requirement on the coverage must be changed:

  $$\int_{x_1}^{x_2} f(x|\theta)\,dx = \beta \quad \longrightarrow \quad \sum_{i=1}^{U} P(x_i|\theta) \geq \beta \quad .$$

  This can lead to overcoverage, but never to under coverage.

  For example, if $\lambda$ is the mean of a Poisson distribution, the coverage will look like:

  

- Summary : Which method should I use?

  → If the uncertainties in the parameter estimates are small compared to
    a) the non-linearities of the model and
    b) the distance to the nearest boundary of the physics region
    ⟹ The Normal Theory is sufficient, and all methods will provide the same numerical results for the confidence interval.

  → Empirical method to test if Normal Theory is sufficient:
    1) Compute confidence interval with Normal Theory method
    2) Compute confidence interval with profile likelihood (available in many programs)
    If (1) and (2) give the same results, the Normal Theory is sufficient.

If (1) and (2) give different results:

→ If the number of parameters is small ($\sim 2$), Feldman-Cousins is easy to implement.

→ If the number of parameters is $\gtrsim 3$, then FC might become very complicated or CPU intensive, but in this situations typically the profile-$\mathcal{L}$ method provides good coverage.

→ Coverage calculation

In any case, one should make sure the method provides the desired coverage! This can be done as follows:

1) Assume values for each parameter, generate $\gtrsim 10^4$ Toy-MC experiments, count how many times the confidence interval covers the true value of each parameter.

2) Repeat for different values of the parameters.