# Aachen Online Statistics School

## GDC Lecture 2: Frequentist probability and confidence intervals



RWTH Aachen (online)
13-17 March 2023

https://indico.desy.de/event/37562/

Glen Cowan
Physics Department
Royal Holloway, University of London
`g.cowan@rhul.ac.uk`
`www.pp.rhul.ac.uk/~cowan`

# Outline of GDC lectures

Tue. 14.3        Probability (Bayes vs. Frequentist)
Bayesian parameter and interval estimation

→ Wed. 15.3       Frequentist confidence regions and intervals

Thu. 16.3       Python software for frequentist and Bayesian confidence regions.

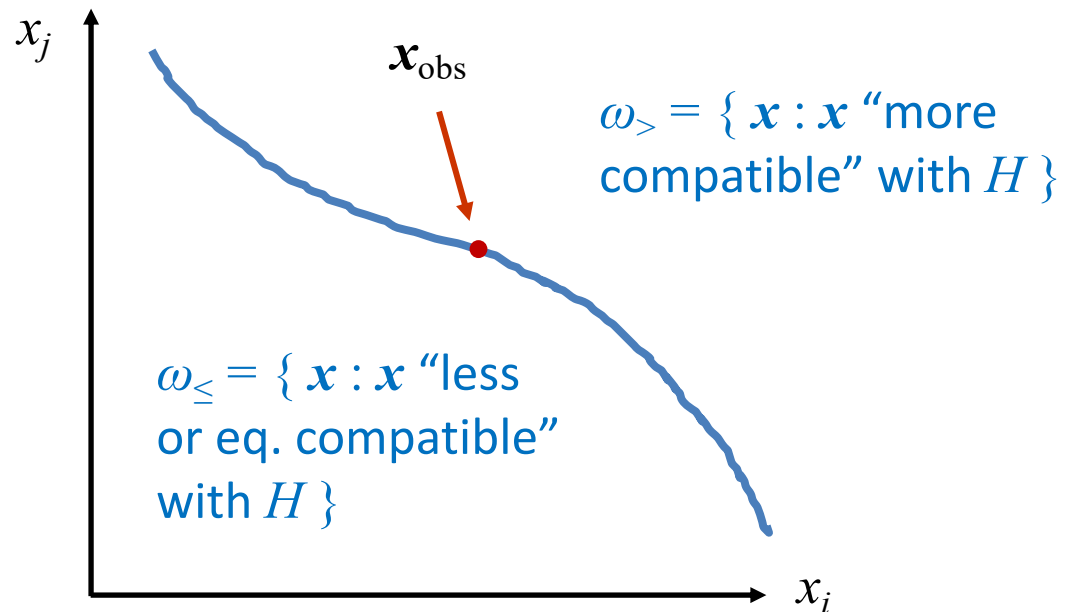Fri. 17.3        Searches and discoveries using likelihoods

# Review of *p*-values

Suppose hypothesis $H$ predicts pdf $f(\boldsymbol{x}|H)$ for a set of observations $\boldsymbol{x} = (x_1,...x_n)$.

We observe a single point in this space: $\boldsymbol{x}_{\mathrm{obs}}$.

How can we quantify the level of compatibility between the data and the predictions of $H$?

Decide what part of the data space represents equal or less compatibility with $H$ than does the point $\boldsymbol{x}_{\mathrm{obs}}$.  (Not unique!)

$x_j$

$\boldsymbol{x}_{\mathrm{obs}}$

$\omega_> = \{\, \boldsymbol{x} : \boldsymbol{x}$ "more compatible" with $H \,\}$

$\omega_\leq = \{\, \boldsymbol{x} : \boldsymbol{x}$ "less or eq. compatible" with $H \,\}$

$x_i$

# $p$-values

Express level of compatibility between data and hypothesis (sometimes 'goodness-of-fit') by giving the $p$-value for $H$:

$$p = P(\mathbf{x} \in \omega_{\leq}(\mathbf{x}_{\mathrm{obs}})|H)$$

=    probability, under assumption of $H$, to observe data with equal or lesser compatibility with $H$ relative to the data we got.

=    probability, under assumption of $H$, to observe data as discrepant with $H$ as the data we got or more so.

Basic idea: if there is only a very small probability to find data with even worse (or equal) compatibility, then $H$ is "disfavoured by the data".

If the $p$-value is below a user-defined threshold $\alpha$ (e.g. 0.05) then $H$ is rejected – equivalent to hypothesis test of size $\alpha$.

# $p$-value of $H$ is not $P(H)$

The $p$-value of H is not the probability that $H$ is true!

In frequentist statistics we don't talk about $P(H)$ (unless $H$ represents a repeatable observation).

If we do define $P(H)$, e.g., in Bayesian statistics as a degree of belief, then we need to use Bayes' theorem to obtain

$$P(H|\vec{x}) = \frac{P(\vec{x}|H)\pi(H)}{\int P(\vec{x}|H)\pi(H)\,dH}$$

where $\pi(H)$ is the prior probability for $H$.

For now stick with the frequentist approach;
result is $p$-value, regrettably easy to misinterpret as $P(H)$.

# Confidence intervals by inverting a test

In addition to a 'point estimate' of a parameter we should report an interval reflecting its statistical uncertainty.

Confidence intervals for a parameter $\theta$ can be found by defining a test of the hypothesized value $\theta$ (do this for all $\theta$):

Specify values of the data that are 'disfavoured' by $\theta$ (critical region) such that $P$(data in critical region $|\theta) \leq \alpha$ for a prespecified $\alpha$, e.g., 0.05 or 0.1.

If data observed in the critical region, reject the value $\theta$.

Equivalently, define a *p*-value for $\theta$. If $p_\theta \leq \alpha$, then reject $\theta$.

# Relation between confidence interval and $p$-value

Now invert the test to define a confidence interval as:

set of $\theta$ values that are not rejected in a test of size $\alpha$ (confidence level CL is $1 - \alpha$).

or equivalently the

set of $\theta$ values for which $p_\theta > \alpha$.

E.g. an upper limit on $\theta$ is the greatest value for which $p_\theta > \alpha$.

In practice find by setting $p_\theta = \alpha$ and solve for $\theta$.

For a multidimensional parameter space $\boldsymbol{\theta} = (\theta_1, \dots \theta_M)$ *use* same idea – result is a confidence "region" with boundary determined by $p_\theta = \alpha$.

# Coverage probability of confidence interval

If the true value of $\theta$ is rejected, then it's not in the confidence interval. The probability for this is by construction (equality for continuous data):

$$P(\text{reject } \theta | \theta) \leq \alpha = \text{type-I error rate}$$

Therefore, the probability for the interval to contain or "cover" $\theta$ is

$$P(\text{conf. interval "covers" } \theta | \theta) \geq 1 - \alpha$$

This assumes that the set of $\theta$ values considered includes the true value, i.e., it assumes the composite hypothesis $P(\boldsymbol{x}|H,\theta)$.

# Frequentist upper limit on Poisson parameter

Consider again the case of observing $n \sim \text{Poisson}(s + b)$.

Suppose $b = 4.5$, $n_{\text{obs}} = 5$. Find upper limit on $s$ at 95% CL.

Relevant alternative is $s = 0$ (critical region at low $n$)

$p$-value of hypothesized $s$ is $P(n \leq n_{\text{obs}}; s, b)$

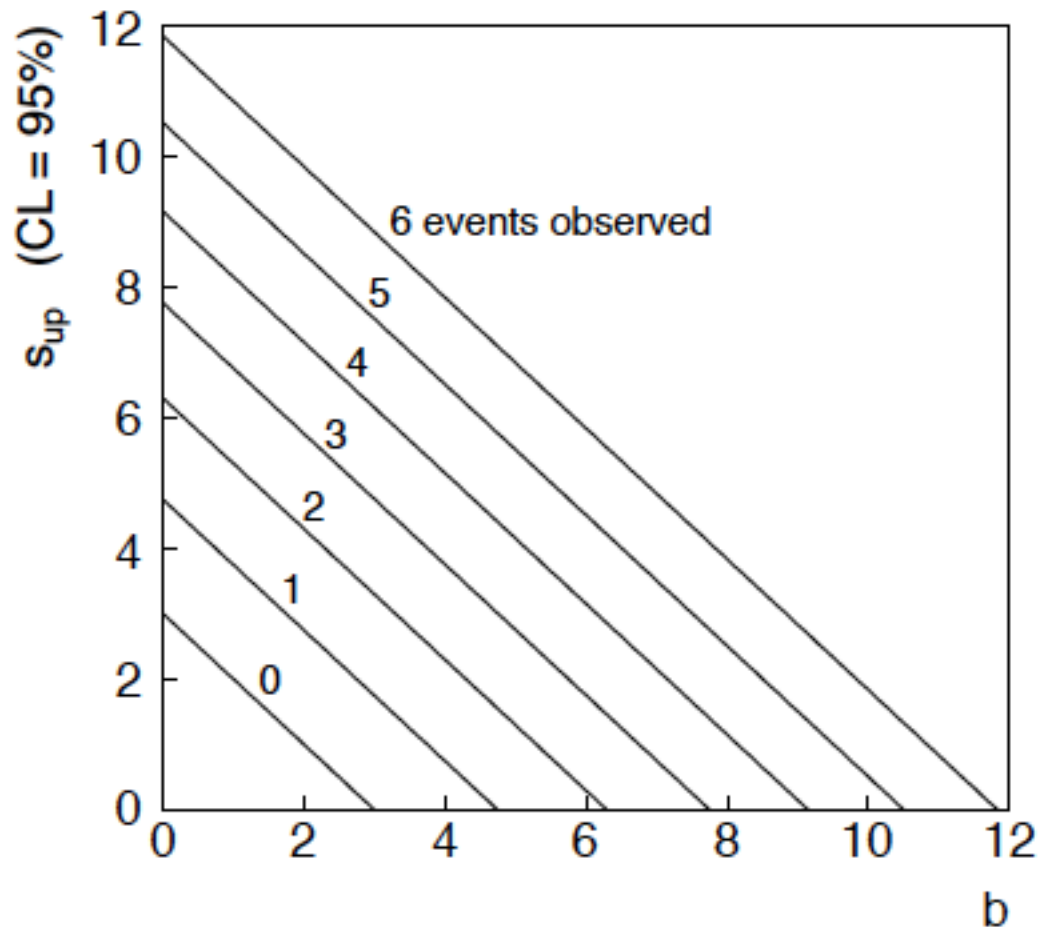Upper limit $s_{\text{up}}$ at $\text{CL} = 1 - \alpha$ found from

$$\alpha = P(n \leq n_{\text{obs}}; s_{\text{up}}, b) = \sum_{n=0}^{n_{\text{obs}}} \frac{(s_{\text{up}} + b)^n}{n!} e^{-(s_{\text{up}}+b)}$$

$$s_{\text{up}} = \frac{1}{2} F_{\chi^2}^{-1}(1 - \alpha; 2(n_{\text{obs}} + 1)) - b$$

$$= \frac{1}{2} F_{\chi^2}^{-1}(0.95; 2(5 + 1)) - 4.5 = 6.0$$

# $n \sim$ Poisson($s+b$): frequentist upper limit on $s$

For low fluctuation of $n$, formula can give negative result for $s_{up}$; i.e. confidence interval is empty; all values of $s \geq 0$ have $p_s \leq \alpha$.

# Limits near a boundary of the parameter space

Suppose e.g. $b = 2.5$ and we observe $n = 0$.

If we choose CL $= 0.9$, we find from the formula for $s_{\text{up}}$

$$s_{\text{up}} = -0.197 \quad (\text{CL} = 0.90)$$

Physicist:

We already knew $s \geq 0$ before we started; can't use negative upper limit to report result of expensive experiment!

Statistician:

The interval is designed to cover the true value only 90% of the time — this was clearly not one of those times.

Not uncommon dilemma when testing parameter values for which one has very little experimental sensitivity, e.g., very small $s$.
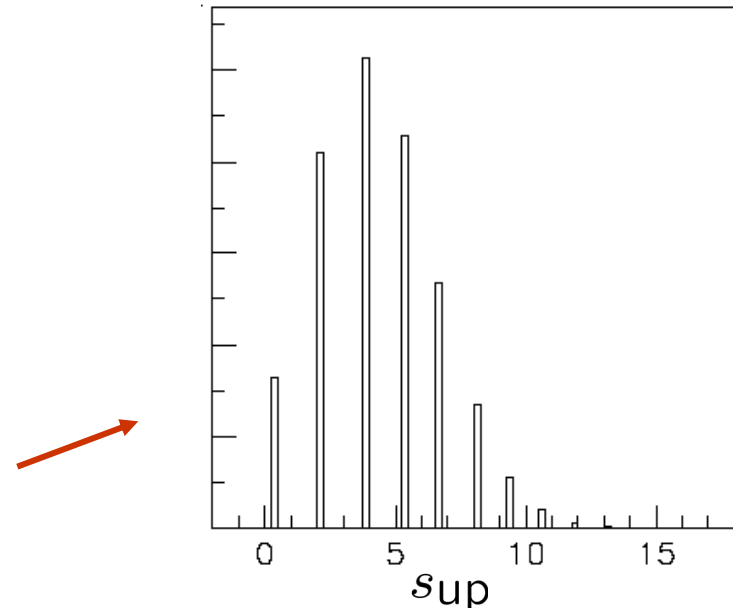
# Expected limit for $s = 0$

Physicist: I should have used CL = 0.95 — then $s_{up} = 0.496$

Even better: for CL = 0.917923 we get $s_{up} = 10^{-4}$ !

Reality check: with $b = 2.5$, typical Poisson fluctuation in $n$ is at least $\sqrt{2.5} = 1.6$. How can the limit be so low?

Look at the mean limit for the no-signal hypothesis ($s = 0$) (sensitivity).

Distribution of 95% CL limits with $b = 2.5$, $s = 0$.
Mean upper limit = 4.44

# Approximate confidence intervals/regions from the likelihood function

Suppose we test parameter value(s) $\boldsymbol{\theta} = (\theta_1, ..., \theta_N)$ using the ratio

$$\lambda(\boldsymbol{\theta}) = \frac{L(\boldsymbol{\theta})}{L(\hat{\boldsymbol{\theta}})} \qquad\qquad 0 \le \lambda(\boldsymbol{\theta}) \le 1$$

Lower $\lambda(\boldsymbol{\theta})$ means worse agreement between data and hypothesized $\boldsymbol{\theta}$. Equivalently, usually define

$$t_{\boldsymbol{\theta}} = -2 \ln \lambda(\boldsymbol{\theta})$$

so higher $t_{\boldsymbol{\theta}}$ means worse agreement between $\boldsymbol{\theta}$ and the data.

$p$-value of $\boldsymbol{\theta}$ therefore

$$p_{\boldsymbol{\theta}} = \int_{t_{\boldsymbol{\theta},\mathrm{obs}}}^{\infty} f(t_{\boldsymbol{\theta}}|\boldsymbol{\theta})\, dt_{\boldsymbol{\theta}}$$

need pdf

# Confidence region from Wilks' theorem

Wilks' theorem says (in large-sample limit and provided certain conditions hold...)

$$f(t_{\boldsymbol{\theta}}|\boldsymbol{\theta}) \sim \chi^2_N$$

chi-square dist. with # d.o.f. = # of components in $\boldsymbol{\theta} = (\theta_1, ..., \theta_N)$.

Assuming this holds, the $p$-value is

$$p_{\boldsymbol{\theta}} = 1 - F_{\chi^2_N}(t_{\boldsymbol{\theta}}|\boldsymbol{\theta}) \quad \leftarrow \text{set equal to } \alpha$$

To find boundary of confidence region set $p_{\boldsymbol{\theta}} = \alpha$ and solve for $t_{\boldsymbol{\theta}}$:

$$t_{\boldsymbol{\theta}} = F^{-1}_{\chi^2_N}(1 - \alpha)$$

Recall also

$$t_{\boldsymbol{\theta}} = -2\ln\frac{L(\theta)}{L(\hat{\theta})}$$

# Confidence region from Wilks' theorem (cont.)

i.e., boundary of confidence region in $\boldsymbol{\theta}$ space is where

$$\ln L(\boldsymbol{\theta}) = \ln L(\hat{\boldsymbol{\theta}}) - \tfrac{1}{2} F^{-1}_{\chi^2_N}(1 - \alpha)$$

For example, for $1 - \alpha = 68.3\%$ and $n = 1$ parameter,

$$F^{-1}_{\chi^2_1}(0.683) = 1$$

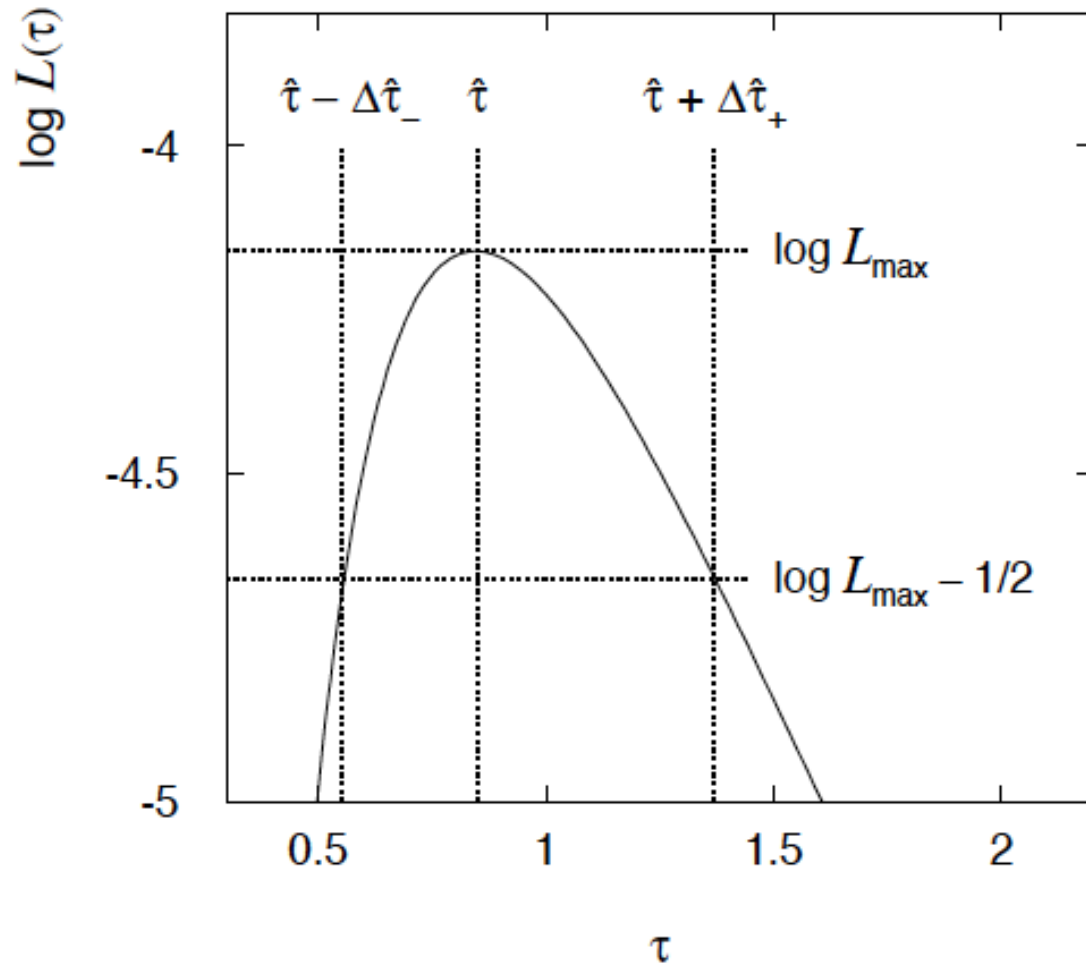and so the 68.3% confidence level interval is determined by

$$\ln L(\theta) = \ln L(\hat{\theta}) - \frac{1}{2}$$

Same as recipe for finding the estimator's standard deviation, i.e.,

$[\hat{\theta} - \sigma_{\hat{\theta}}, \hat{\theta} + \sigma_{\hat{\theta}}]$   is a 68.3% CL confidence interval.

# Example of interval from $\ln L(\theta)$

For $N=1$ parameter, CL = 0.683, $Q_\alpha = 1$.



Our exponential example, now with only $n = 5$ events.

Can report ML estimate with approx. confidence interval from $\ln L_{max} - 1/2$ as "asymmetric error bar":

$$\hat{\tau} = 0.85^{+0.52}_{-0.30}$$

# Multiparameter case

For increasing number of parameters, CL $= 1 - \alpha$ decreases for confidence region determined by a given

$$Q_\alpha = F^{-1}_{\chi^2_n}(1 - \alpha)$$

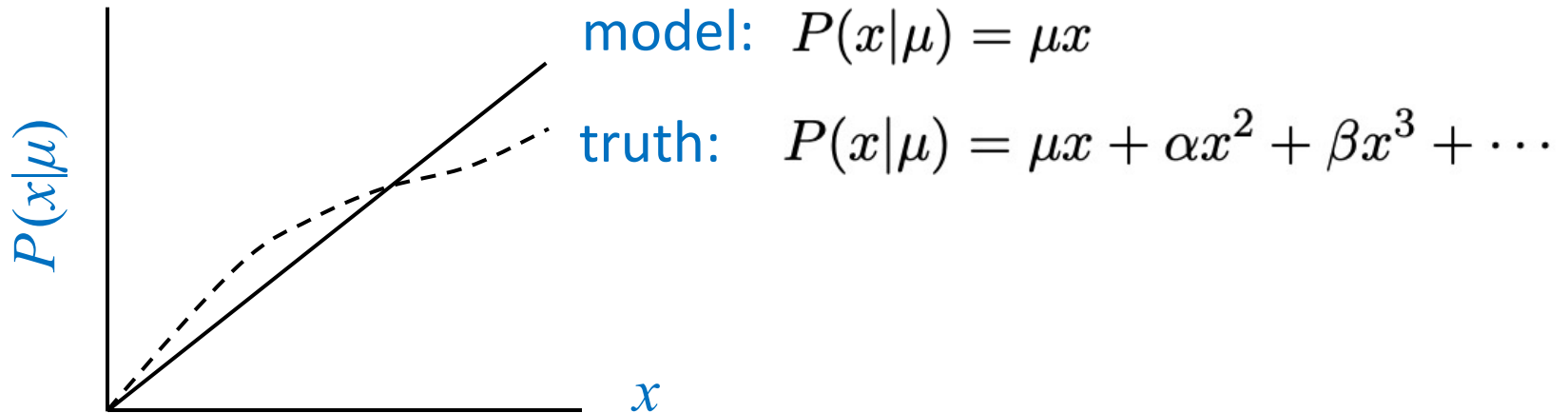| $Q_\alpha$ | $1 - \alpha$ | | | | |
|------------|-------------|-------------|-------------|-------------|-------------|
|            | $n = 1$ | $n = 2$ | $n = 3$ | $n = 4$ | $n = 5$ |
| 1.0 | 0.683 | 0.393 | 0.199 | 0.090 | 0.037 |
| 2.0 | 0.843 | 0.632 | 0.428 | 0.264 | 0.151 |
| 4.0 | 0.954 | 0.865 | 0.739 | 0.594 | 0.451 |
| 9.0 | 0.997 | 0.989 | 0.971 | 0.939 | 0.891 |

← # of par.

# Multiparameter case (cont.)

Equivalently, $Q_\alpha$ increases with $n$ for a given CL $= 1 - \alpha$.

| $1 - \alpha$ | $\hat{Q}_\alpha$ | | | | |
|---|---|---|---|---|---|
| | $n = 1$ | $n = 2$ | $n = 3$ | $n = 4$ | $n = 5$ |
| 0.683 | 1.00 | 2.30 | 3.53 | 4.72 | 5.89 |
| 0.90 | 2.71 | 4.61 | 6.25 | 7.78 | 9.24 |
| 0.95 | 3.84 | 5.99 | 7.82 | 9.49 | 11.1 |
| 0.99 | 6.63 | 9.21 | 11.3 | 13.3 | 15.1 |

← # of par.

# Systematic uncertainties and nuisance parameters

In general, our model of the data is not perfect:

model: $P(x|\mu) = \mu x$

truth: $P(x|\mu) = \mu x + \alpha x^2 + \beta x^3 + \cdots$

Can improve model by including additional adjustable parameters.

$$P(x|\mu) \to P(x|\mu, \boldsymbol{\theta})$$

Nuisance parameter $\leftrightarrow$ systematic uncertainty. Some point in the parameter space of the enlarged model should be "true".

Presence of nuisance parameter decreases sensitivity of analysis to the parameter of interest (e.g., increases variance of estimate).

# Profile Likelihood

Suppose we have a likelihood $L(\boldsymbol{\mu}, \boldsymbol{\theta}) = P(\boldsymbol{x}|\boldsymbol{\mu}, \boldsymbol{\theta})$ with $N$ parameters of interest $\boldsymbol{\mu} = (\mu_1, ..., \mu_N)$ and $M$ nuisance parameters $\boldsymbol{\theta} = (\theta_1, ..., \theta_M)$. The "profiled" (or "constrained") values of $\boldsymbol{\theta}$ are:

$$\hat{\hat{\boldsymbol{\theta}}}(\boldsymbol{\mu}) = \underset{\boldsymbol{\theta}}{\mathrm{argmax}}\, L(\boldsymbol{\mu}, \boldsymbol{\theta})$$

and the profile likelihood is: $\quad L_{\mathrm{p}}(\boldsymbol{\mu}) = L(\boldsymbol{\mu}, \hat{\hat{\boldsymbol{\theta}}})$

The profile likelihood depends only on the parameters of interest; the nuisance parameters are replaced by their profiled values.

The profile likelihood can be used to obtain confidence intervals/regions for the parameters of interest in the same way as one would for all of the parameters from the full likelihood.

# Example: fitting a straight line

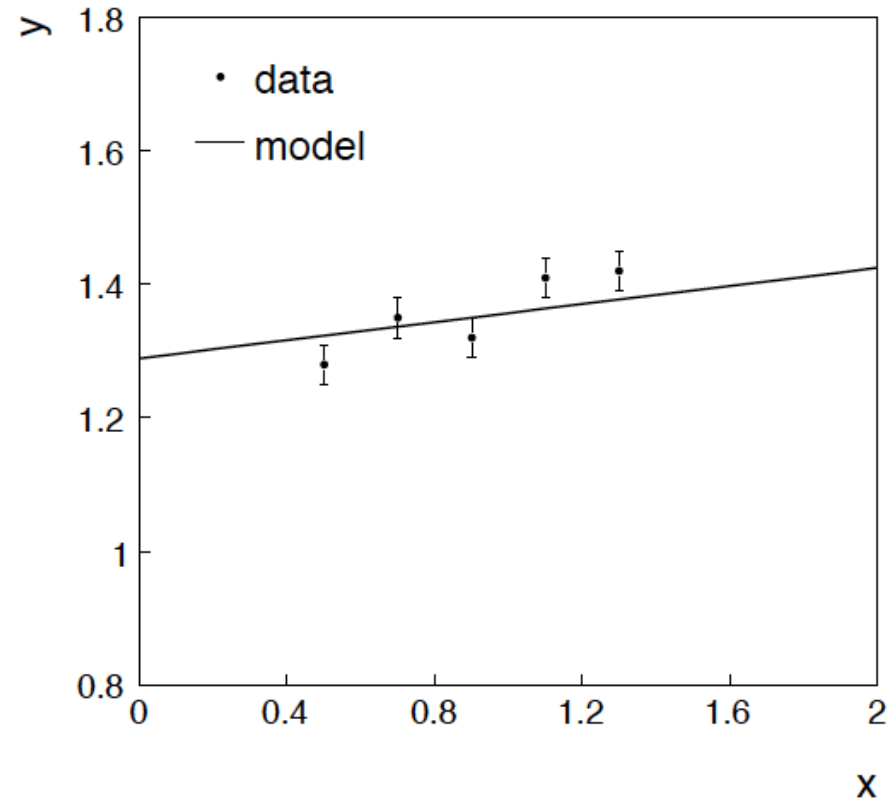Data: $(x_i, y_i, \sigma_i)$ , $i = 1, \ldots, n$ .

Model: $y_i$ independent and all follow $y_i \sim \text{Gauss}(\mu(x_i), \sigma_i)$

$$\mu(x; \theta_0, \theta_1) = \theta_0 + \theta_1 x \,,$$

assume $x_i$ and $\sigma_i$ known.

Goal: estimate $\theta_0$

Here suppose we don't care about $\theta_1$ (example of a "nuisance parameter")

# Maximum likelihood fit with Gaussian data

In this example, the $y_i$ are assumed independent, so the likelihood function is a product of Gaussians:

$$L(\theta_0, \theta_1) = \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi}\sigma_i} \exp\left[-\frac{1}{2}\frac{(y_i - \mu(x_i; \theta_0, \theta_1))^2}{\sigma_i^2}\right] ,$$

Maximizing the likelihood is here equivalent to minimizing

$$\chi^2(\theta_0, \theta_1) = -2\ln L(\theta_0, \theta_1) + \text{const} = \sum_{i=1}^{n} \frac{(y_i - \mu(x_i; \theta_0, \theta_1))^2}{\sigma_i^2} .$$

i.e., for Gaussian data, ML same as Method of Least Squares (LS)

# $\theta_1$ known a priori

$$L(\theta_0) = \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi}\sigma_i} \exp\left[-\frac{1}{2}\frac{(y_i - \mu(x_i; \theta_0, \theta_1))^2}{\sigma_i^2}\right] \, .$$
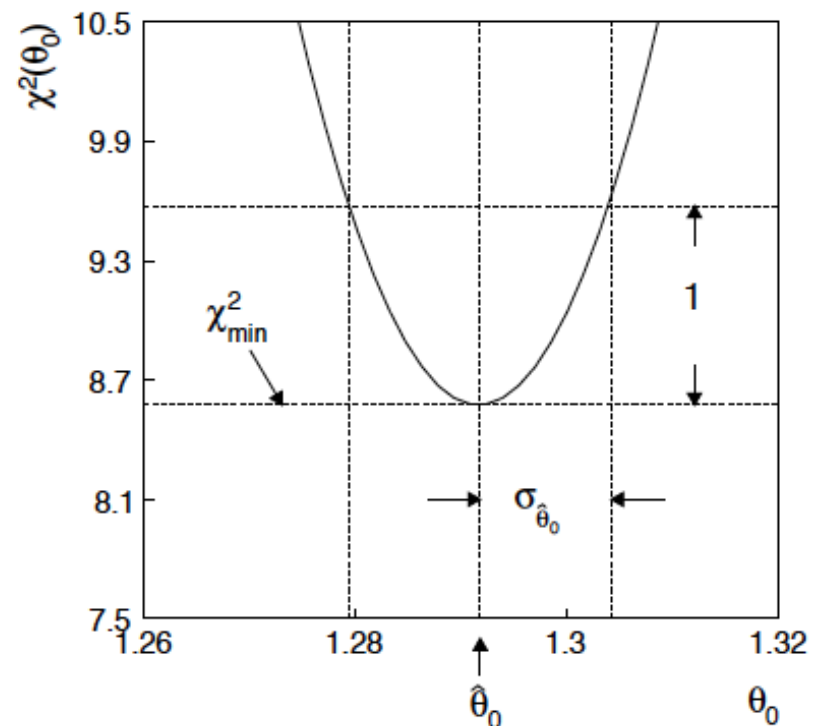
$$\chi^2(\theta_0) = -2\ln L(\theta_0) + \text{const} = \sum_{i=1}^{n} \frac{(y_i - \mu(x_i; \theta_0, \theta_1))^2}{\sigma_i^2} \, .$$

For Gaussian $y_i$, ML same as LS

Minimize $\chi^2 \rightarrow$ estimator $\hat{\theta}_0$ .

Come up one unit from $\chi^2_{\min}$

to find $\sigma_{\hat{\theta}_0}$ .

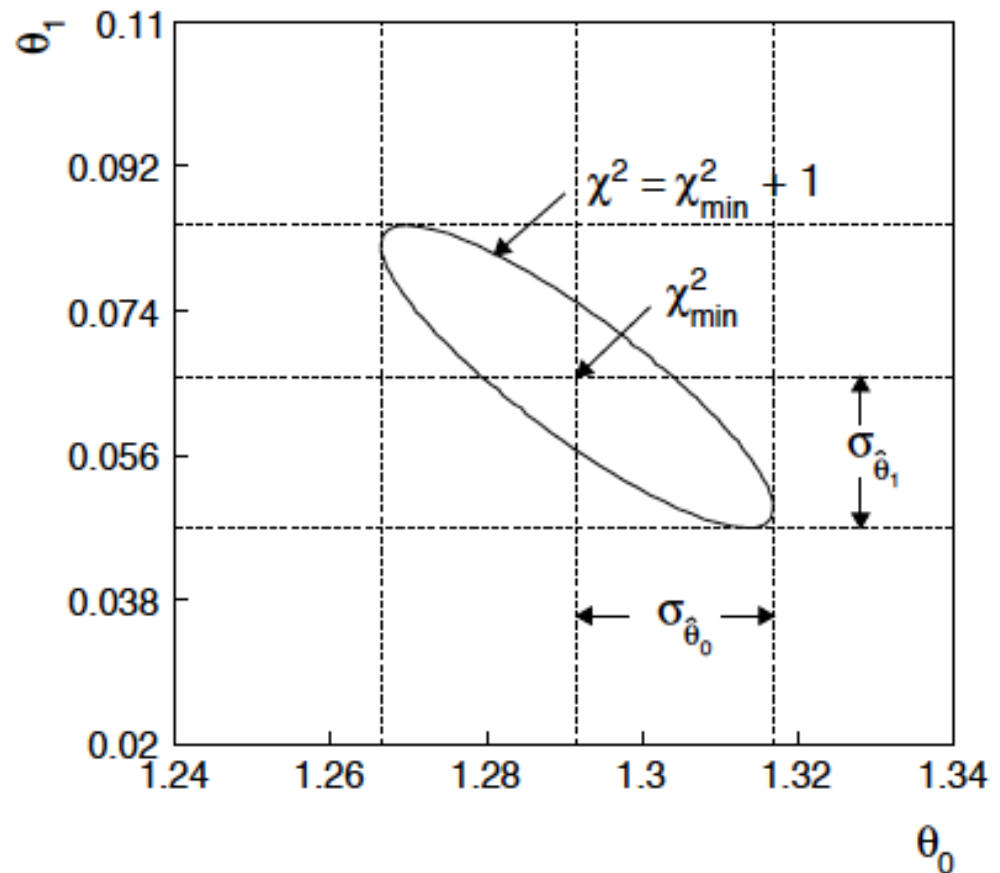# ML (or LS) fit of $\theta_0$ and $\theta_1$

$$\chi^2(\theta_0, \theta_1) = -2 \ln L(\theta_0, \theta_1) + \text{const} = \sum_{i=1}^{n} \frac{(y_i - \mu(x_i; \theta_0, \theta_1))^2}{\sigma_i^2} \, .$$

Standard deviations from

tangent lines to contour

$$\chi^2 = \chi^2_{\min} + 1 \, .$$

Correlation between

$\hat{\theta}_0, \ \hat{\theta}_1$ causes errors

to increase.

# Including the measurement $t_1 \sim$ Gauss $(\theta_1, \sigma_{t1})$

$$L(\theta_0, \theta_1) = \frac{1}{\sqrt{2\pi}\sigma_t} e^{-(t_1-\theta_1)^2/2\sigma_{t_1}^2} \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi}\sigma_i} \exp\left[-\frac{1}{2}\frac{(y_i - \mu(x_i; \theta_0, \theta_1))^2}{\sigma_i^2}\right]$$
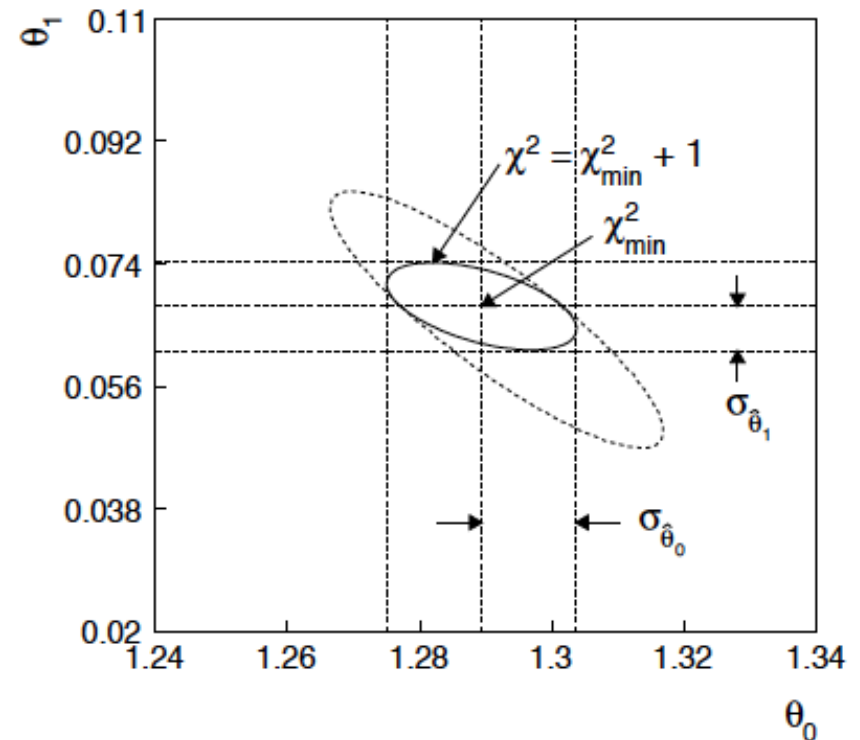
$$\chi^2(\theta_0, \theta_1) = \sum_{i=1}^{n} \frac{(y_i - \mu(x_i; \theta_0, \theta_1))^2}{\sigma_i^2} + \frac{(t_1 - \theta_1)^2}{\sigma_{t_1}^2}$$

The information on $\theta_1$ improves accuracy of $\hat{\theta}_0$ .

Here the contour corresponds to $\ln L = \ln L_{\max} - \frac{1}{2}$ , so:

The interval for $\theta_0$ is a conf. interval with CL = 0.683.

The 2-parameter region corresponds to CL = 0.393.

# Summing up...

In Frequentist statistics, probability only associated with data (and functions thereof).

Parameter estimation boils down to finding functions of the data (estimators) that are themselves random variables, having a mean, standard deviation, etc.

$p$-value of $H$ = $P$(data as "extreme" as what we saw or more so | $H$)

Confidence intervals = set of parameter values with $p$-value < $\alpha$.

Designed to "cover" a parameter with a given probability (the intervals are random, not the parameter).

Wilks' theorem allows one to find approximate confidence intervals (and multi-param. regions) directly from the likelihood.
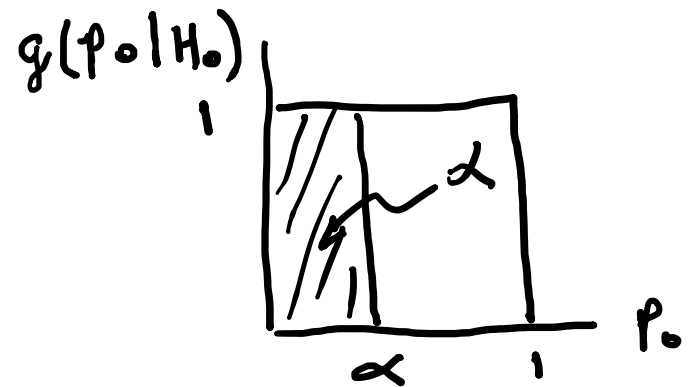
# Extra slides

# Using a $p$-value to define test of $H_0$

One can show that under assumption of a hypothesis $H_0$, its $p$-value, $p_0$, follows a uniform distribution in [0,1].

So the probability to find $p_0$ less than a given $\alpha$ is

$$P(p_0 \leq \alpha | H_0) = \alpha$$



So we can define the critical region of a test of $H_0$ with size $\alpha$ as the set of data space where $p_0 \leq \alpha$ .

Formally the $p$-value relates only to $H_0$, but the resulting test will have a given power with respect to a given alternative $H_1$.