# Statistics 101 – a very fast introduction Lecture 2

# 

## PHYSTAT-Gamma (online) 26-30 September 2022

https://indico.cern.ch/event/1122011/



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

# Outline

Lecture 1: Probability and Bayes theorem, Frequentist and Bayesian statistics, likelihood function, parameter estimation, maximum likelihood, information inequality, properties of MLE

→ Lecture 2: Frequentist hypothesis tests, Neyman-Pearson lemma/likelihood ratio, goodness of fit, p values and significances, confidence interval from a test, Wilk's theorem and confidence regions

Almost everything is a subset of the University of London course: http://www.pp.rhul.ac.uk/~cowan/stat\_course.html

#### Frequentist hypothesis tests

Suppose a measurement produces data x; consider a hypothesis  $H_0$  we want to test and alternative  $H_1$ 

 $H_0$ ,  $H_1$  specify probability for  $\mathbf{x}$ :  $P(\mathbf{x}|H_0)$ ,  $P(\mathbf{x}|H_1)$ 

A test of  $H_0$  is defined by specifying a critical region w of the data space such that there is no more than some (small) probability  $\alpha$ , assuming  $H_0$  is correct, to observe the data there, i.e.,

 $P(\mathbf{x} \in w \mid H_0) \le \alpha$ 

Need inequality if data are discrete.

 $\alpha$  is called the size or significance level of the test.

If x is observed in the critical region, reject  $H_0$ .



### Definition of a test (2)

But in general there are an infinite number of possible critical regions that give the same size  $\alpha$ .

Use the alternative hypothesis  $H_1$  to motivate where to place the critical region.

Roughly speaking, place the critical region where there is a low probability ( $\alpha$ ) to be found if  $H_0$  is true, but high if  $H_1$  is true:



#### Classification viewed as a statistical test

Suppose events come in two possible types:

s (signal) and b (background)

For each event, test hypothesis that it is background, i.e.,  $H_0 = b$ .

Carry out test on many events, each is either of type s or b, i.e., here the hypothesis is the "true class label", which varies randomly from event to event, so we can assign to it a frequentist probability.

Select events for which where  $H_0$  is rejected as "candidate events of type s". Equivalent Physics terminology:

background efficiency 
$$arepsilon_{
m b} = \int_W f({f x}|H_0)\,d{f x} = lpha$$

$$\varepsilon_{\mathbf{s}} = \int_{W} f(\mathbf{x}|H_1) \, d\mathbf{x} = 1 - \beta = \text{power}$$

signal efficiency

G. Cowan / RHUL Physics

#### Example of a test for classification



For each event in a mixture of signal (s) and background (b) test

 $H_0$ : event is of type b

using a critical region W of the form:  $W = \{x : x \le x_c\}$ , where  $x_c$  is a constant that we choose to give a test with the desired size  $\alpha$ .

G. Cowan / RHUL Physics

#### Classification example (2)

Suppose we want  $\alpha = 10^{-4}$ . Require:

$$\alpha = P(x \le x_{c}|b) = \int_{0}^{x_{c}} f(x|b) \, dx = \frac{4x^{4}}{4} \Big|_{0}^{x_{c}} = x_{c}^{4}$$

and therefore  $x_{\rm c} = \alpha^{1/4} = 0.1$ 

For this test (i.e. this critical region W), the power with respect to the signal hypothesis (s) is

$$M = P(x \le x_{\rm c}|{\rm s}) = \int_0^{x_{\rm c}} f(x|{\rm s}) \, dx = 2x_{\rm c} - x_{\rm c}^2 = 0.19$$

Note: the optimal size and power is a separate question that will depend on goals of the subsequent analysis.

G. Cowan / RHUL Physics

#### Classification example (3)

Suppose that the prior probabilities for an event to be of type s or b are:

 $\pi_{\rm s} = 0.001$  $\pi_{\rm b} = 0.999$ 

The "purity" of the selected signal sample (events where b hypothesis rejected) is found using Bayes' theorem:

$$P(\mathbf{s}|x \le x_{\mathbf{c}}) = \frac{P(x \le x_{\mathbf{c}}|\mathbf{s})\pi_{\mathbf{s}}}{P(x \le x_{\mathbf{c}}|\mathbf{s})\pi_{\mathbf{s}} + P(x \le x_{\mathbf{c}}|\mathbf{b})\pi_{\mathbf{b}}}$$

= 0.655

G. Cowan / RHUL Physics

#### Test statistic based on likelihood ratio

How can we choose a test's critical region in an 'optimal way', in particular if the data space is multidimensional?

Neyman-Pearson lemma states:

For a test of  $H_0$  of size  $\alpha$ , to get the highest power with respect to the alternative  $H_1$  we need for all x in the critical region W

"likelihood ratio (LR)" 
$$\frac{P(\mathbf{x}|H_1)}{P(\mathbf{x}|H_0)} \ge c_{\alpha}$$

inside W and  $\leq c_{\alpha}$  outside, where  $c_{\alpha}$  is a constant chosen to give a test of the desired size.

Equivalently, optimal scalar test statistic is

$$t(\mathbf{x}) = \frac{P(\mathbf{x}|H_1)}{P(\mathbf{x}|H_0)}$$

N.B. any monotonic function of this is leads to the same test.

G. Cowan / RHUL Physics

#### Neyman-Pearson doesn't usually help

We usually don't have explicit formulae for the pdfs f(x|s), f(x|b), so for a given x we can't evaluate the likelihood ratio

$$t(\mathbf{x}) = \frac{f(\mathbf{x}|s)}{f(\mathbf{x}|b)}$$

Instead we may have Monte Carlo models for signal and background processes, so we can produce simulated data:

generate 
$$\boldsymbol{x} \sim f(\boldsymbol{x}|\mathbf{s}) \rightarrow \boldsymbol{x}_1, \dots, \boldsymbol{x}_N$$

generate 
$$\mathbf{x} \sim f(\mathbf{x}|\mathbf{b}) \longrightarrow \mathbf{x}_1, ..., \mathbf{x}_N$$

This gives samples of "training data" with events of known type.

Use these to construct a statistic that is as close as possible to the optimal likelihood ratio (→ Machine Learning).

G. Cowan / RHUL Physics

#### Testing significance / goodness-of-fit

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

We observe a single point in this space:  $x_{obs}$ .

 $X_i$ 

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  $x_{obs}$ . (Not unique!)



#### *p*-values

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

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

- 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 as discussed previously).



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).

#### *p*-value from test statistic



If e.g. we define the region of less or eq. compatibility to be  $t(x) \ge t_{obs}$  then the *p*-value of *H* is

$$p_H = \int_{t_{\text{obs}}}^{\infty} f(t|H) \, dt = \int_{\{\mathbf{x}: t(\mathbf{x}) \ge t_{\text{obs}}\}} f(\mathbf{x}|H) \, d\mathbf{x}$$

G. Cowan / RHUL Physics

The Poisson counting experiment Suppose we do a counting experiment and observe *n* events.

Events could be from *signal* process or from *background* – we only count the total number.

Poisson model:

$$P(n|s,b) = \frac{(s+b)^n}{n!}e^{-(s+b)}$$

*s* = mean (i.e., expected) # of signal events

*b* = mean # of background events

Goal is to make inference about *s*, e.g.,

test s = 0 (rejecting  $H_0 \approx$  "discovery of signal process")

test all non-zero *s* (values not rejected = confidence interval)

In both cases need to ask what is relevant alternative hypothesis.

Poisson counting experiment: discovery *p*-value Suppose b = 0.5 (known), and we observe  $n_{obs} = 5$ . Should we claim evidence for a new discovery?

Give *p*-value for hypothesis s = 0:

$$p$$
-value =  $P(n \ge 5; b = 0.5, s = 0)$   
=  $1.7 \times 10^{-4} \ne P(s = 0)!$ 



G. Cowan / RHUL Physics

PHYSTAT-Gamma 2022 / Lecture 2

#### Significance from *p*-value

Often define significance Z as the number of standard deviations that a Gaussian variable would fluctuate in one direction to give the same p-value.



$$p = \int_{Z}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx = 1 - \Phi(Z)$$
  
 $Z = \Phi^{-1}(1-p)$ 

in ROOT: in pyt
p = 1 - TMath::Freq(Z) p = 1
Z = TMath::NormQuantile(1-p) Z = no

in python (scipy.stats): p = 1 - norm.cdf(Z) = norm.sf(Z) Z = norm.ppf(1-p)

Result Z is a "number of sigmas". Note this does not mean that the original data was Gaussian distributed.

G. Cowan / RHUL Physics

# Poisson counting experiment: discovery significance Equivalent significance for $p = 1.7 \times 10^{-4}$ : $Z = \Phi^{-1}(1-p) = 3.6$ Often claim discovery if Z > 5 ( $p < 2.9 \times 10^{-7}$ , i.e., a "5-sigma effect")



In fact this tradition should be revisited: *p*-value intended to quantify probability of a signallike fluctuation assuming background only; not intended to cover, e.g., hidden systematics, plausibility signal model, compatibility of data with signal, "look-elsewhere effect" (~multiple testing), etc.

#### 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  $P(\text{data in critical region} | \theta) \le \alpha$ for a prespecified  $\alpha$ , e.g., 0.05 or 0.1.

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

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  $1 - \alpha$ ).

Relation between confidence interval and *p*-value

Equivalently we can consider a significance test for each hypothesized value of  $\theta$ , resulting in a *p*-value,  $p_{\theta}$ .

If  $p_{\theta} \leq \alpha$ , then we reject  $\theta$ .

The confidence interval at  $CL = 1 - \alpha$  consists of those values of  $\theta$  that are not rejected.

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  $\theta = (\theta_1, \dots, \theta_M)$  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*(conf. interval "covers"  $\theta | \theta \ge 1 - \alpha$ 

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

#### 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  $P(n \le n_{\text{obs}}; s, b)$ Upper limit  $s_{\text{up}}$  at  $\text{CL} = 1 - \alpha$  found from

$$\begin{aligned} \alpha &= P(n \le n_{\text{obs}}; s_{\text{up}}, b) = \sum_{n=0}^{n_{\text{obs}}} \frac{(s_{\text{up}} + b)^n}{n!} e^{-(s_{\text{up}} + b)} \\ s_{\text{up}} &= \frac{1}{2} F_{\chi^2}^{-1} (1 - \alpha; 2(n_{\text{obs}} + 1)) - b \\ &= \frac{1}{2} F_{\chi^2}^{-1} (0.95; 2(5 + 1)) - 4.5 = 6.0 \end{aligned}$$

G. Cowan / RHUL Physics

### Approximate confidence intervals/regions from the likelihood function

Suppose we test parameter value(s)  $\theta = (\theta_1, ..., \theta_n)$  using the ratio

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

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

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

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

*p*-value of  $\theta$  therefore

$$p_{\theta} = \int_{t_{\theta,\text{obs}}}^{\infty} f(t_{\theta}|\theta) \, dt_{\theta}$$
need pdf

G. Cowan / RHUL Physics

#### Confidence region from Wilks' theorem

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

 $f(t_{\theta}|\theta) \sim \chi_n^2 \qquad \begin{array}{l} \text{chi-square dist. with $\#$ d.o.f. =} \\ \# \text{ of components in $\theta = (\theta_1, ..., \theta_n)$.} \end{array}$ 

Assuming this holds, the *p*-value is

$$p_{m{ heta}} = 1 - F_{\chi^2_n}(t_{m{ heta}}) \quad \leftarrow \text{set equal to } lpha$$

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

$$t_{\theta} = F_{\chi_n^2}^{-1}(1-\alpha)$$

Recall also

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

G. Cowan / RHUL Physics

Confidence region from Wilks' theorem (cont.)

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

$$\ln L(\theta) = \ln L(\hat{\theta}) - \frac{1}{2}F_{\chi_n^2}^{-1}(1-\alpha)$$

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

$$F_{\chi_1^2}^{-1}(0.683) = 1$$

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,  $Q_{\alpha} = 1$ .



moved to tutorial

#### Multiparameter case

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

$$Q_{\alpha} = F_{\chi_n^2}^{-1}(1-\alpha)$$

$Q_{lpha}$	$1-\alpha$					
	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	

#### Multiparameter case (cont.)

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

$1 - \alpha$	$\widehat{Q}_{lpha}$						
	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		

### Finally

#### Two lectures only enough for a brief introduction to:

Parameter estimation

Hypothesis tests ( $\rightarrow$  path to Machine Learning)

Limits (confidence intervals/regions)

No time today for many important things, e.g.,

Systematics (nuisance parameters)

**Experimental sensitivity** 

Final thought: once the basic formalism is fixed, most of the work focuses on writing down the likelihood, e.g.,  $P(x|\theta)$ , and including in it enough parameters to adequately describe the data (true for both Bayesian and frequentist approaches) so often best to invest most of your time with it.

### Extra slides

#### Proof of Neyman-Pearson Lemma

Consider a critical region W and suppose the LR satisfies the criterion of the Neyman-Pearson lemma:

 $P(\mathbf{x}|H_1)/P(\mathbf{x}|H_0) \ge c_{\alpha} \text{ for all } \mathbf{x} \text{ in } W,$  $P(\mathbf{x}|H_1)/P(\mathbf{x}|H_0) \le c_{\alpha} \text{ for all } \mathbf{x} \text{ not in } W.$ 

Try to change this into a different critical region W' retaining the same size  $\alpha$ , i.e.,

$$P(\mathbf{x} \in W'|H_0) = P(\mathbf{x} \in W|H_0) = \alpha$$

To do so add a part  $\delta W_+$ , but to keep the size  $\alpha$ , we need to remove a part  $\delta W_-$ , i.e.,

$$W \to W' = W + \delta W_+ - \delta W_-$$

$$P(\mathbf{x} \in \delta W_+ | H_0) = P(\mathbf{x} \in \delta W_- | H_0)$$





G. Cowan / RHUL Physics

#### Proof of Neyman-Pearson Lemma (2)

But we are supposing the LR is higher for all x in  $\delta W_{-}$  removed than for the x in  $\delta W_{+}$  added, and therefore

$$P(\mathbf{x} \in \delta W_+ | H_1) \le P(\mathbf{x} \in \delta W_+ | H_0) c_\alpha$$

$$\frac{\delta W_{+}}{\delta W_{-}}$$

 $P(\mathbf{x} \in \delta W_{-}|H_{1}) \ge P(\mathbf{x} \in \delta W_{-}|H_{0})c_{\alpha}$ 

The right-hand sides are equal and therefore

 $P(\mathbf{x} \in \delta W_+ | H_1) \le P(\mathbf{x} \in \delta W_- | H_1)$ 

#### Proof of Neyman-Pearson Lemma (3)

#### We have

$$W \cup W' = W \cup \delta W_+ = W' \cup \delta W_-$$

Note W and  $\delta W_+$  are disjoint, and W' and  $\delta W_-$  are disjoint, so by Kolmogorov's 3<sup>rd</sup> axiom,



$$P(\mathbf{x} \in W') + P(\mathbf{x} \in \delta W_{-}) = P(\mathbf{x} \in W) + P(\mathbf{x} \in \delta W_{+})$$

#### Therefore

$$P(\mathbf{x} \in W'|H_1) = P(\mathbf{x} \in W|H_1) + P(\mathbf{x} \in \delta W_+|H_1) - P(\mathbf{x} \in \delta W_-|H_1)$$

G. Cowan / RHUL Physics

#### Proof of Neyman-Pearson Lemma (4)

And therefore

$$P(\mathbf{x} \in W'|H_1) \le P(\mathbf{x} \in W|H_1)$$

i.e. the deformed critical region W' cannot have higher power than the original one that satisfied the LR criterion of the Neyman-Pearson lemma.

#### Example of *p*-value: exponential decay time

A nuclear sample contains two radioactive isotopes with mean lifetimes  $\tau = 0.2$  s and  $\tau = 1.0$  s.

For either isotope we expect the decay time to follow  $f(t|\tau) = \frac{1}{\tau}e^{-t/\tau}$ 

A nucleus is observed to decay after a time  $t_{obs}$  = 0.6 s.

The *p*-value of the hypothesis *H* that the nucleus is of the type with  $\tau = 0.2$  s is

$$p_H = P(t \ge t_{\rm obs} | \tau = 0.2 \,\mathrm{s}) = 0.0498$$

Here we take  $t \ge t_{obs}$  as being less compatible with  $\tau = 0.2$  s , because greater t is more characteristic of  $\tau = 1.0$  s.

If the relevant alternative had been  $\tau = 0.1$  s, then one would define the *p*-value as

$$p_H = P(t \le t_{\rm obs} | \tau = 0.2 \,\mathrm{s}) = 0.9502$$



#### Distribution of the *p*-value

The *p*-value is a function of the data, and is thus itself a random variable with a given distribution. Suppose the *p*-value of *H* is found from a test statistic t(x) as

$$p_H = \int_t^\infty f(t'|H)dt'$$

The pdf of  $p_H$  under assumption of H is

$$g(p_H|H) = \frac{f(t|H)}{|\partial p_H/\partial t|} = \frac{f(t|H)}{f(t|H)} = 1 \quad (0 \le p_H \le 1)$$

In general for continuous data, under assumption of H,  $p_H \sim \text{Uniform}[0,1]$  and is concentrated toward zero for some (broad) class of alternatives.



G. Cowan / RHUL Physics

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

So the probability to find the *p*-value of  $H_0$ ,  $p_0$ , less than  $\alpha$  is



We started by defining critical region in the original data space (x), then reformulated this in terms of a scalar test statistic t(x).

We can take this one step further and define the critical region of a test of  $H_0$  with size  $\alpha$  as the set of data space where  $p_0 \le \alpha$ .

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

#### $n \sim \text{Poisson}(s+b)$ : 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  $s \ge 0$  have  $p_s \le \alpha$ .



#### 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  $s_{up}$ 

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

Physicist:

We already knew  $s \ge 0$  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?



### Systematic uncertainties and nuisance parameters In general, our model of the data is not perfect:



Can improve model by including additional adjustable parameters.

$$P(x|\mu) \to P(x|\mu, \theta)$$

Nuisance parameter ↔ 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  $L(\mu, \theta) = P(x|\mu, \theta)$  with Nparameters of interest  $\mu = (\mu_1, ..., \mu_N)$  and M nuisance parameters  $\theta = (\theta_1, ..., \theta_M)$ . The "profiled" (or "constrained") values of  $\theta$  are:

$$\hat{\hat{\boldsymbol{\theta}}}(\boldsymbol{\mu}) = \operatorname*{argmax}_{\boldsymbol{\theta}} L(\boldsymbol{\mu}, \boldsymbol{\theta})$$

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

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.

#### Profile Likelihood Ratio – Wilks theorem

Goal is to test/reject regions of  $\mu$  space (param. of interest).

Rejecting a point  $\mu$  should mean  $p_{\mu} \leq \alpha$  for all possible values of the nuisance parameters  $\theta$ .

Test  $\boldsymbol{\mu}$  using the "profile likelihood ratio":  $\lambda(\boldsymbol{\mu}) = \frac{L(\boldsymbol{\mu}, \hat{\boldsymbol{\theta}})}{L(\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\theta}})}$ 

Let  $t_{\mu} = -2 \ln \lambda(\mu)$ . Wilks' theorem says in large-sample limit:  $t_{\mu} \sim \text{chi-square}(N)$ 

where the number of degrees of freedom is the number of parameters of interest (components of  $\mu$ ). So *p*-value for  $\mu$  is

$$p_{\boldsymbol{\mu}} = \int_{t_{\boldsymbol{\mu},\text{obs}}}^{\infty} f(t_{\boldsymbol{\mu}} | \boldsymbol{\mu}, \boldsymbol{\theta}) \, dt_{\boldsymbol{\mu}} = 1 - F_{\chi_N^2}(t_{\boldsymbol{\mu},\text{obs}})$$

G. Cowan / RHUL Physics

#### Profile Likelihood Ratio – Wilks theorem (2)

If we have a large enough data sample to justify use of the asymptotic chi-square pdf, then if  $\mu$  is rejected, it is rejected for any values of the nuisance parameters.

The recipe to get confidence regions/intervals for the parameters of interest at  $CL = 1 - \alpha$  is thus the same as before, simply use the profile likelihood:

$$\ln L_{\rm p}(\boldsymbol{\mu}) = \ln L_{\rm max} - \frac{1}{2} F_{\chi_N^2}^{-1} (1 - \alpha)$$

where the number of degrees of freedom N for the chi-square quantile is equal to the number of parameters of interest.

If the large-sample limit is not justified, then use e.g. Monte Carlo to get distribution of  $t_{\mu}$ .

G. Cowan / RHUL Physics