Aachen Online Statistics School

GDC Lecture 2: Frequentist probability and confidence intervals





RNTHAACHEN UNIVERSITY

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. 14.3Probability (Bayes vs. Frequentist)Bayesian parameter and interval estimation
- → Wed. 15.3 Frequentist confidence regions and intervals
 - Thu. 16.3Python 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(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 α (e.g. 0.05) then *H* is rejected – equivalent to hypothesis test of size α .

Confidence interval from *p*-values

We can define a *p*-value for all hypothesized values of θ . Then the confidence region at confidence level $CL = 1 - \alpha$ is the set of θ values for which $p_{\theta} > \alpha$.

or equivalently

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

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

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

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 θ 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" θ is

P(conf. interval "covers" $\theta | \theta \ge 1 - \alpha$

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

Example: upper limit on mean of Gaussian

When we test the parameter, we should take the critical region to maximize the power with respect to the relevant alternative(s).

Example: $x \sim \text{Gauss}(\mu, \sigma)$ (take σ known)

Test $H_0: \mu = \mu_0$ versus the alternative $H_1: \mu < \mu_0$

 \rightarrow Put w_{μ} at region of x-space characteristic of low μ (i.e. at low x)



Equivalently, take the *p*-value to be

$$p_{\mu_0} = P(x \le x_{\text{obs}} | \mu_0) = \int_{-\infty}^{x_{\text{obs}}} \frac{1}{\sqrt{2\pi\sigma}} e^{-(x-\mu_0)^2/2\sigma^2} \, dx = \Phi\left(\frac{x_{\text{obs}} - \mu_0}{\sigma}\right)$$

Upper limit on Gaussian mean (2)

To find confidence interval, repeat for all μ_0 , i.e., set $p_{\mu 0} = \alpha$ and solve for μ_0 to find the interval's boundary



$$\mu_0 \to \mu_{\rm up} = x_{\rm obs} - \sigma \Phi^{-1}(\alpha) = x_{\rm obs} + \sigma \Phi^{-1}(1 - \alpha)$$

This is an upper limit on μ , i.e., higher μ have even lower p-value and are in even worse agreement with the data.

Usually use $\Phi^{-1}(\alpha) = -\Phi^{-1}(1-\alpha)$ so as to express the upper limit as x_{obs} plus a positive quantity. E.g. for $\alpha = 0.05$, $\Phi^{-1}(1-0.05) = 1.64$.

Upper limit on Gaussian mean (3)

 μ_{up} = the hypothetical value of μ such that there is only a probability α to find $x < x_{obs}$.



1-vs. 2-sided intervals

Now test: $H_0: \mu = \mu_0$ versus the alternative $H_1: \mu \neq \mu_0$

I.e. we consider the alternative to μ_0 to include higher and lower values, so take critical region on both sides:



Result is a "central" confidence interval [μ_{lo} , μ_{up}]:

$$\mu_{\rm lo} = x_{\rm obs} - \sigma \Phi^{-1} \left(1 - \frac{\alpha}{2} \right) \qquad \text{E.g. for } \alpha = 0.05$$
$$\mu_{\rm up} = x_{\rm obs} + \sigma \Phi^{-1} \left(1 - \frac{\alpha}{2} \right) \qquad \Phi^{-1} \left(1 - \frac{\alpha}{2} \right) = 1.96 \approx 2$$

Note upper edge of two-sided interval is higher (i.e. not as tight of a limit) than obtained from the one-sided test.

G. Cowan / RHUL Physics

Statistical Data Analysis / lecture week 9

Approximate confidence intervals/regions from the likelihood function

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

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

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

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

so higher t_{θ} means worse agreement between θ and the data.

p-value of θ therefore

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

RWTH Aachen online course / GDC lecture 2

Confidence region from Wilks' theorem

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

 $f(t_{\theta}|\theta) \sim \chi_N^2$ chi-square dist. with # d.o.f. = # of components in $\theta = (\theta_1, ..., \theta_N)$.

Assuming this holds, the *p*-value is

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

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

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

Recall also

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

G. Cowan / RHUL Physics

RWTH Aachen online course / GDC lecture 2

Confidence region from Wilks' theorem (cont.) i.e., boundary of confidence region in θ space is where

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

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, $Q_{\alpha} = 1$.



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}		-				
	n = 1	n = 2	n = 3	n = 4	n = 5	\leftarrow # of par.
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_{α} increases with *n* for a given $CL = 1 - \alpha$.

$1 - \alpha$						
	n = 1	n = 2	n = 3	n = 4	n = 5	\leftarrow # of par.
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	_

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 θ 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.

Example: fitting a straight line

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

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

$$\mu(x;\theta_0,\theta_1)=\theta_0+\theta_1x,$$

assume x_i and σ_i known.

Goal: estimate θ_0

Here suppose we don't care about θ_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)

θ_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 $\chi^2 \rightarrow \text{estimator } \widehat{\theta}_0$. Come up one unit from χ^2_{min} to find $\sigma_{\hat{\theta}_0}$.



ML (or LS) fit of θ_0 and θ_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

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

Correlation between $\hat{\theta}_0, \ \hat{\theta}_1$ causes errors to increase.



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

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

$$\chi^2(\theta_0, \theta_1) = \sum_{i=1}^n \frac{(y_i - \mu(x_i; \theta_0, \theta_1))^2}{\sigma_i^2} + \frac{(t_1 - \theta_1)^2}{\sigma_{t_1}^2}$$

The information on θ_1 improves accuracy of $\hat{\theta}_0$.

Here the contour corresponds to $lnL = lnL_{max} - \frac{1}{2}$, so:

The 2-parameter region corresponds to CL = 0.393.

The interval for θ_0 is a conf. interval with CL = 0.683.



Profiling

The $\ln L = \ln L_{max} - \frac{1}{2}$ contour in the (θ_0 , θ_1) plane is a confidence region at CL = 39.3%.

Furthermore if one wants to know only about, say, θ_0 , then the interval in θ_0 corresponding to $\ln L = \ln L_{\max} - \frac{1}{2}$ is a confidence interval at CL = 68.3% (i.e., ±1 std. dev.).

I.e., form the interval for θ_0 using

$$\ln L(\theta_0, \hat{\hat{\theta}}_1(\theta_0)) = \ln L_{\max} - 1/2$$

where θ_1 is replaced by its "profiled" value

$$\hat{\hat{\theta}}_1(\theta_0) = \operatorname*{argmax}_{\theta_1} L(\theta_0, \theta_1)$$



G. Cowan / RHUL Physics

Profile Likelihood Ratio – Wilks theorem

Goal is to test/reject regions of μ space (param. of interest).

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

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 μ). So *p*-value for μ 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

Statistical Data Analysis / lecture week 10

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 μ 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_{μ} .

G. Cowan / RHUL Physics

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 < α .

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



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).

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 α is

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



So we can define the critical region of a test of H_0 with size α 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 .

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_{up} at $\text{CL} = 1 - \alpha$ found from

$$\alpha = P(n \le n_{\text{obs}}; s_{\text{up}}, b) = \sum_{n=0}^{n_{\text{obs}}} \frac{(s_{\text{up}} + b)^n}{n!} e^{-(s_{\text{up}} + b)}$$
$$s_{\text{up}} = \frac{1}{2} F_{\chi^2}^{-1} (1 - \alpha; 2(n_{\text{obs}} + 1)) - b$$

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

G. Cowan / RHUL Physics

RWTH Aachen online course / GDC lecture 2

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



RWTH Aachen online course / GDC lecture 2

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?

