Statistics for Particle Physicists Lecture 2



https://pages.lip.pt/data-science/school/



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

LIP Data Science School 2022, Coimbra / Lecture 2

Outline

Lecture 1: Introduction

Probability

Hypothesis tests

From hypothesis tests to machine learning

→ Lecture 2: Parameter estimation

Confidence limits

Systematic uncertainties

Hypothesis, likelihood

Suppose the entire result of an experiment (set of measurements) is a collection of numbers *x*.

A (simple) hypothesis is a rule that assigns a probability to each possible data outcome:

 $P(\mathbf{x}|H)$ = the likelihood of H

Often we deal with a family of hypotheses labeled by one or more undetermined parameters (a composite hypothesis):

$$P(\mathbf{x}|oldsymbol{ heta}) = L(oldsymbol{ heta})$$
 = the "likelihood function"

Note:

- 1) For the likelihood we treat the data x as fixed.
- 2) The likelihood function $L(\theta)$ is not a pdf for θ .

The likelihood function for i.i.d.* data

* i.i.d. = independent and identically distributed

Consider *n* independent observations of *x*: $x_1, ..., x_n$, where *x* follows $f(x; \theta)$. The joint pdf for the whole data sample is:

$$f(x_1,\ldots,x_n;\theta) = \prod_{i=1}^n f(x_i;\theta)$$

In this case the likelihood function is

$$L(\vec{\theta}) = \prod_{i=1}^{n} f(x_i; \vec{\theta}) \qquad (x_i \text{ constant})$$

Parameter estimation

The parameters of a pdf are any constants that characterize it,



i.e., θ indexes a set of hypotheses.

Suppose we have a sample of observed values: $x = (x_1, ..., x_n)$

We want to find some function of the data to estimate the parameter(s):

 $\hat{\theta}(\vec{x}) \leftarrow \text{estimator written with a hat}$

Sometimes we say 'estimator' for the function of $x_1, ..., x_n$; 'estimate' for the value of the estimator with a particular data set.

Properties of estimators

If we were to repeat the entire measurement, the estimates from each would follow a pdf:



We want small (or zero) bias (systematic error): $b = E[\hat{\theta}] - \theta$

 \rightarrow average of repeated measurements should tend to true value.

And we want a small variance (statistical error): $V[\hat{\theta}]$

 \rightarrow small bias & variance are in general conflicting criteria

Maximum Likelihood Estimators (MLEs)

We *define* the maximum likelihood estimators or MLEs to be the parameter values for which the likelihood is maximum.



Could have multiple maxima (take highest).

MLEs not guaranteed to have any 'optimal' properties, (but in practice they're very good).

MLE example: parameter of exponential pdf

Consider exponential pdf,
$$f(t; \tau) = \frac{1}{\tau}e^{-t/\tau}$$

and suppose we have i.i.d. data, t_1, \ldots, t_n

The likelihood function is
$$L(\tau) = \prod_{i=1}^{n} \frac{1}{\tau} e^{-t_i/\tau}$$

The value of τ for which $L(\tau)$ is maximum also gives the maximum value of its logarithm (the log-likelihood function):

$$\ln L(\tau) = \sum_{i=1}^{n} \ln f(t_i; \tau) = \sum_{i=1}^{n} \left(\ln \frac{1}{\tau} - \frac{t_i}{\tau} \right)$$

MLE example: parameter of exponential pdf (2)

Find its maximum by setting

 $\rightarrow \hat{\tau} = \frac{1}{n} \sum_{i=1}^{n} t_i$

$$\frac{\partial \ln L(\tau)}{\partial \tau} = 0 ,$$

Monte Carlo test: generate 50 values using $\tau = 1$:

We find the ML estimate:

$$\hat{\tau} = 1.062$$



MLE example: parameter of exponential pdf (3)

For the exponential distribution one has for mean, variance:

$$E[t] = \int_0^\infty t \, \frac{1}{\tau} e^{-t/\tau} \, dt = \tau$$

$$V[t] = \int_0^\infty (t - \tau)^2 \, \frac{1}{\tau} e^{-t/\tau} \, dt = \tau^2$$
For the MLE $\hat{\tau} = \frac{1}{n} \sum_{i=1}^n t_i$ we therefore find
$$E[\hat{\tau}] = E\left[\frac{1}{n} \sum_{i=1}^n t_i\right] = \frac{1}{n} \sum_{i=1}^n E[t_i] = \tau \implies b = E[\hat{\tau}] - \tau = 0$$

$$V[\hat{\tau}] = V\left[\frac{1}{n} \sum_{i=1}^n t_i\right] = \frac{1}{n^2} \sum_{i=1}^n V[t_i] = \frac{\tau^2}{n} \implies \sigma_{\hat{\tau}} = \frac{\tau}{\sqrt{n}}$$

LIP Data Science School 2022, Coimbra / Lecture 2

Variance of estimators: Monte Carlo method

Having estimated our parameter we now need to report its 'statistical error', i.e., how widely distributed would estimates be if we were to repeat the entire measurement many times.

One way to do this would be to simulate the entire experiment many times with a Monte Carlo program (use ML estimate for MC).

For exponential example, from sample variance of estimates we find:

 $\hat{\sigma}_{\hat{\tau}} = 0.151$

Note distribution of estimates is roughly Gaussian – (almost) always true for ML in large sample limit.



Variance of estimators from information inequality

The information inequality (RCF) sets a lower bound on the variance of any estimator (not only ML):

$$V[\hat{\theta}] \ge \left(1 + \frac{\partial b}{\partial \theta}\right)^2 / E\left[-\frac{\partial^2 \ln L}{\partial \theta^2}\right] \qquad \text{Bound (MVB)}$$
$$(b = E[\hat{\theta}] - \theta)$$

Often the bias *b* is small, and equality either holds exactly or is a good approximation (e.g. large data sample limit). Then,

$$V[\hat{\theta}] \approx -1 \left/ E\left[\frac{\partial^2 \ln L}{\partial \theta^2}\right] \right.$$

Estimate this using the 2nd derivative of $\ln L$ at its maximum:

$$\widehat{V}[\widehat{\theta}] = -\left(\frac{\partial^2 \ln L}{\partial \theta^2}\right)^{-1} \bigg|_{\theta = \widehat{\theta}}$$

LIP Data Science School 2022, Coimbra / Lecture 2

MVB for MLE of exponential parameter

Find MVB =
$$-\left(1 + \frac{\partial b}{\partial \tau}\right)^2 / E\left[\frac{\partial^2 \ln L}{\partial \tau^2}\right]$$

We found for the exponential parameter the MLE

$$\hat{\tau} = \frac{1}{n} \sum_{i=1}^{n} t_i$$

and we showed b = 0, hence $\partial b / \partial \tau = 0$.

We find
$$\frac{\partial^2 \ln L}{\partial \tau^2} = \sum_{i=1}^n \left(\frac{1}{\tau^2} - \frac{2t_i}{\tau^3} \right)$$

and since $E[t_i] = \tau$ for all i , $E\left[\frac{\partial^2 \ln L}{\partial \tau^2} \right] = -\frac{n}{\tau^2}$,
and therefore MVB $= \frac{\tau^2}{n} = V[\hat{\tau}]$. (Here MLE is "efficient").

LIP Data Science School 2022, Coimbra / Lecture 2

Variance of estimators: graphical method

Expand $lnL(\theta)$ about its maximum:

$$\ln L(\theta) = \ln L(\hat{\theta}) + \left[\frac{\partial \ln L}{\partial \theta}\right]_{\theta = \hat{\theta}} (\theta - \hat{\theta}) + \frac{1}{2!} \left[\frac{\partial^2 \ln L}{\partial \theta^2}\right]_{\theta = \hat{\theta}} (\theta - \hat{\theta})^2 + \dots$$

First term is $\ln L_{max}$, second term is zero, for third term use information inequality (assume equality):

$$\ln L(\theta) \approx \ln L_{\max} - \frac{(\theta - \hat{\theta})^2}{2\hat{\sigma}^2_{\hat{\theta}}}$$

i.e.,
$$\ln L(\hat{\theta} \pm \hat{\sigma}_{\hat{\theta}}) \approx \ln L_{\max} - \frac{1}{2}$$

 \rightarrow to get $\hat{\sigma}_{\hat{\theta}}$, change θ away from $\hat{\theta}$ until $\ln L$ decreases by 1/2.

Example of variance by graphical method



Not quite parabolic $\ln L$ since finite sample size (n = 50).

Information inequality for N parameters Suppose we have estimated N parameters $\theta = (\theta_1, ..., \theta_N)$ The Fisher information matrix is

$$I_{ij} = -E\left[\frac{\partial^2 \ln L}{\partial \theta_i \partial \theta_j}\right] = -\int \frac{\partial^2 \ln P(\mathbf{x}|\boldsymbol{\theta})}{\partial \theta_i \partial \theta_j} P(\mathbf{x}|\boldsymbol{\theta}) \, d\mathbf{x}$$

and the covariance matrix of estimators $\hat{\theta}$ is $V_{ij} = ext{cov}[\hat{ heta}_i, \hat{ heta}_j]$

The information inequality states that the matrix

$$M_{ij} = V_{ij} - \sum_{k,l} \left(\delta_{ik} + \frac{\partial b_i}{\partial \theta_k} \right) I_{kl}^{-1} \left(\delta_{lj} + \frac{\partial b_l}{\partial \theta_j} \right)$$

is positive semi-definite:

 $z^{\mathrm{T}}Mz \ge 0$ for all $z \ne 0$, diagonal elements ≥ 0

Information inequality for N parameters (2)

In practice the inequality is ~always used in the large-sample limit: bias $\rightarrow 0$ inequality \rightarrow equality, i.e, M = 0, and therefore $V^{-1} = I$

That is,
$$V_{ij}^{-1} = -E\left[\frac{\partial^2 \ln L}{\partial \theta_i \partial \theta_j}\right]$$

This can be estimated from data using $\hat{V}_{ij}^{-1} = -\frac{\partial^2 \ln L}{\partial \theta_i \partial \theta_j}\Big|_{\hat{\theta}}$

Find the matrix V^{-1} numerically (or with automatic differentiation), then invert to get the covariance matrix of the estimators

$$\widehat{V}_{ij} = \widehat{\text{cov}}[\widehat{\theta}_i, \widehat{\theta}_j]$$

Example of ML with 2 parameters

Consider a scattering angle distribution with $x = \cos \theta$,

$$f(x;\alpha,\beta) = \frac{1+\alpha x + \beta x^2}{2+2\beta/3}$$



or if $x_{\min} < x < x_{\max}$, need to normalize so that

$$\int_{x_{\mathsf{min}}}^{x_{\mathsf{max}}} f(x; lpha, eta) \, dx = \mathbf{1} \; .$$

Example: $\alpha = 0.5$, $\beta = 0.5$, $x_{\min} = -0.95$, $x_{\max} = 0.95$, generate n = 2000 events with Monte Carlo.

$$\ln L(\alpha,\beta) = \sum_{i=1}^{n} \ln f(x_i;\alpha,\beta) \quad \longleftarrow \quad \text{need to find maximum} \\ \text{numerically}$$

G. Cowan / RHUL Physics

LIP Data Science School 2022, Coimbra / Lecture 2

$$\hat{\alpha} = 0.508$$

$$\hat{\beta} = 0.47$$

N.B. No binning of data for fit, but can compare to histogram for goodness-of-fit (e.g. 'visual' or χ^2).



х

(Co)variances from
$$(\widehat{V^{-1}})_{ij} = -\frac{\partial^2 \ln L}{\partial \theta_i \partial \theta_j}\Big|_{\vec{\theta} = \hat{\vec{\theta}}}$$

 $\hat{\sigma}_{\hat{\alpha}} = 0.052 \quad \operatorname{cov}[\hat{\alpha}, \hat{\beta}] = 0.0026$

0.11 r = 0.46 = correlation coefficient

 \widehat{x}

LIP Data Science School 2022, Coimbra / Lecture 2

Two-parameter fit: MC study

Repeat ML fit with 500 experiments, all with n = 2000 events:



G. Cowan / RHUL Physics

Multiparameter graphical method for variances

Expand $\ln L(\theta)$ to 2nd order about MLE:

$$\ln L(\boldsymbol{\theta}) \approx \ln L(\hat{\boldsymbol{\theta}}) + \sum_{i} \frac{\partial \ln L}{\partial \theta_{i}} \Big|_{\hat{\boldsymbol{\theta}}} (\theta_{i} - \hat{\theta}_{i}) + \frac{1}{2!} \sum_{i,j} \frac{\partial^{2} \ln L}{\partial \theta_{i} \partial \theta_{j}} \Big|_{\hat{\boldsymbol{\theta}}} (\theta_{i} - \hat{\theta}_{i})(\theta_{j} - \hat{\theta}_{j})$$

$$\int_{\ln L_{\max}} zero relate to covariance matrix of$$

relate to covariance matrix of MLEs using information (in)equality.

Result:
$$\ln L(\boldsymbol{\theta}) = \ln L_{\max} - \frac{1}{2} \sum_{i,j} (\theta_i - \hat{\theta}_i) V_{ij}^{-1} (\theta_j - \hat{\theta}_j)$$

So the surface $\ln L(\theta) = \ln L_{\max} - \frac{1}{2}$ corresponds to

 $(\theta - \hat{\theta})^T V^{-1} (\theta - \hat{\theta}) = 1$, which is the equation of a (hyper-) ellipse.

Multiparameter graphical method (2)



Distance from MLE to tangent planes gives standard deviations.

The $\ln L_{\rm max} - 1/2$ contour for two parameters

For large n, $\ln L$ takes on quadratic form near maximum:

$$\ln L(\alpha,\beta) \approx \ln L_{\max}$$
$$-\frac{1}{2(1-\rho^2)} \left[\left(\frac{\alpha - \hat{\alpha}}{\sigma_{\hat{\alpha}}} \right)^2 + \left(\frac{\beta - \hat{\beta}}{\sigma_{\hat{\beta}}} \right)^2 - 2\rho \left(\frac{\alpha - \hat{\alpha}}{\sigma_{\hat{\alpha}}} \right) \left(\frac{\beta - \hat{\beta}}{\sigma_{\hat{\beta}}} \right) \right]$$

The contour $\ln L(\alpha,\beta) = \ln L_{\max} - 1/2$ is an ellipse:

$$\frac{1}{(1-\rho^2)}\left[\left(\frac{\alpha-\widehat{\alpha}}{\sigma_{\widehat{\alpha}}}\right)^2 + \left(\frac{\beta-\widehat{\beta}}{\sigma_{\widehat{\beta}}}\right)^2 - 2\rho\left(\frac{\alpha-\widehat{\alpha}}{\sigma_{\widehat{\alpha}}}\right)\left(\frac{\beta-\widehat{\beta}}{\sigma_{\widehat{\beta}}}\right)\right] = 1$$

(Co)variances from In L contour



 \rightarrow Tangent lines to contours give standard deviations.

 \rightarrow Angle of ellipse φ related to correlation: $\tan 2\phi = \frac{2\rho\sigma_{\hat{\alpha}}\sigma_{\hat{\beta}}}{\sigma_{\hat{\alpha}}^2 - \sigma_{\hat{\beta}}^2}$

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 θ can be found by defining a test of the hypothesized value θ (do this for all θ):

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

If data observed in the critical region, reject the value θ .

Now invert the test to define a confidence interval as:

set of θ values that are not rejected in a test of size α (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 θ , resulting in a *p*-value, p_{θ} .

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

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

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 $P(\mathbf{x}|H,\theta)$.

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

LIP Data Science School 2022, Coimbra / 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 $p_s \le \alpha$.



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.

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

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).

Example: fitting a straight line

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

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



If we have a 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 $\hat{\theta}_0$.



Reminder of Bayesian approach

In Bayesian statistics we can associate a probability with a hypothesis, e.g., a parameter value θ .

Interpret probability of θ as 'degree of belief' (subjective). Need to start with 'prior pdf' $\pi(\theta)$, this reflects degree of belief about θ before doing the experiment.

Our experiment has data x, \rightarrow likelihood $L(x|\theta)$.

Bayes' theorem tells how our beliefs should be updated in light of the data *x*:

$$p(\theta|x) = \frac{L(x|\theta)\pi(\theta)}{\int L(x|\theta')\pi(\theta') d\theta'} \propto L(x|\theta)\pi(\theta)$$

Posterior pdf $p(\theta|x)$ contains all our knowledge about θ .

Bayesian approach: $y_i \sim \text{Gauss}(\mu(x_i; \theta_0, \theta_1), \sigma_i)$ We need to associate prior probabilities with θ_0 and θ_1 , e.g.,

 $\pi(\theta_0, \theta_1) = \pi_0(\theta_0)\pi_1(\theta_1) \quad \leftarrow \text{suppose knowledge of } \theta_0 \text{ has}$ no influence on knowledge of θ_1

$$\pi_0(\theta_0) = \text{const.} \qquad \leftarrow \text{`non-informative', in any} \\ \text{case much broader than } L(\theta_0)$$

$$\pi_{1}(\theta_{1}) = p(\theta_{1}|t_{1}) \propto p(t_{1}|\theta_{1})\pi_{\mathrm{Ur}}(\theta_{1}) = \frac{1}{\sqrt{2\pi}\sigma_{t}}e^{-(t_{1}-\theta_{1})^{2}/2\sigma_{t}^{2}} \times \mathrm{const.}$$
prior after t_{1} , Ur = "primordial" Likelihood for control before y prior measurement t_{1}

Bayesian example: $y_i \sim \text{Gauss}(\mu(x_i; \theta_0, \theta_1), \sigma_i)$

Putting the ingredients into Bayes' theorem gives:



Note here the likelihood only reflects the measurements *y*.

The information from the control measurement t_1 has been put into the prior for θ_1 .

We would get the same result using the likelihood $P(y,t|\theta_0,\theta_1)$ and the constant "Ur-prior" for θ_1 .

Marginalizing the posterior pdf

We then integrate (marginalize) $p(\theta_0, \theta_1 | \mathbf{y})$ to find $p(\theta_0 | \mathbf{y})$:

$$p(\theta_0|\mathbf{y}) = \int p(\theta_0, \theta_1|\mathbf{y}) \, d\theta_1$$

In this example we can do the integral (rare). We find

$$p(\theta_0|\mathbf{y}) = \frac{1}{\sqrt{2\pi}\sigma_{\theta_0}} e^{-(\theta_0 - \hat{\theta}_0)^2/2\sigma_{\theta_0^2}}$$

 $\hat{\theta}_0 = \text{same as MLE}$

 $\sigma_{\theta_0} = \sigma_{\hat{\theta}_0}$ (same as for MLE)

For this example, numbers come out same as in frequentist approach, but interpretation different.

G. Cowan / RHUL Physics

INFN 2022, Paestum / Parameter Estimation 2

Marginalization with MCMC

Bayesian computations involve integrals like

$$p(\theta_0|x) = \int p(\theta_0, \theta_1|x) d\theta_1$$
.

often high dimensionality and impossible in closed form, also impossible with 'normal' acceptance-rejection Monte Carlo.

Markov Chain Monte Carlo (MCMC) has revolutionized Bayesian computation.

MCMC (e.g., Metropolis-Hastings algorithm) generates correlated sequence of random numbers:

cannot use for many applications, e.g., detector MC; effective stat. error greater than if all values independent .

Basic idea: sample multidimensional θ but look only at distribution of parameters of interest.

MCMC basics: Metropolis-Hastings algorithm Goal: given an *n*-dimensional pdf $p(\theta)$, generate a sequence of points $\theta_1, \theta_2, \theta_3, \dots$

- 1) Start at some point $\vec{\theta}_0$
- 2) Generate $\vec{\theta} \sim q(\vec{\theta}; \vec{\theta}_0)$

Proposal density $q(\theta; \theta_0)$ e.g. Gaussian centred about θ_0

3) Form Hastings test ratio $\alpha = \min \left| 1, \frac{\pi}{n} \right|$

$$1, \frac{p(\vec{\theta})q(\vec{\theta}_{0}; \vec{\theta})}{p(\vec{\theta}_{0})q(\vec{\theta}; \vec{\theta}_{0})} \bigg]$$

- 4) Generate $u \sim \text{Uniform}[0, 1]$
- 5) If $u \le \alpha$, $\vec{\theta_1} = \vec{\theta}$, \leftarrow move to proposed point else $\vec{\theta_1} = \vec{\theta_0} \leftarrow$ old point repeated 6) Iterate

Metropolis-Hastings (continued)

This rule produces a *correlated* sequence of points (note how each new point depends on the previous one).

Still works if $p(\theta)$ is known only as a proportionality, which is usually what we have from Bayes' theorem: $p(\theta|\mathbf{x}) \propto p(\mathbf{x}|\theta)\pi(\theta)$.

The proposal density can be (almost) anything, but choose so as to minimize autocorrelation. Often take proposal density symmetric: $q(\theta; \theta_0) = q(\theta_0; \theta)$

Test ratio is (*Metropolis*-Hastings): $\alpha = \min \left[1, \frac{p(\theta)}{p(\vec{\theta}_0)} \right]$

I.e. if the proposed step is to a point of higher $p(\theta)$, take it; if not, only take the step with probability $p(\theta)/p(\theta_0)$. If proposed step rejected, repeat the current point.

Example: posterior pdf from MCMC

Sample the posterior pdf from previous example with MCMC:



Bayesian method with alternative priors

Suppose we don't have a previous measurement of θ_1 but rather, an "expert" says it should be positive and not too much greater than 0.1 or so, i.e., something like

$$\pi_1(\theta_1) = \frac{1}{\tau} e^{-\theta_1/\tau} , \quad \theta_1 \ge 0 , \quad \tau = 0.1 .$$

From this we obtain (numerically) the posterior pdf for θ_0 :



Finally

Two lectures only enough for a brief introduction to:

- Parameter estimation
- Hypothesis tests (\rightarrow path to Machine Learning)
- Limits (confidence intervals/regions)
- Systematics (nuisance parameters)
- A bit beyond... (Bayesian methods, MCMC)

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