Statistics for Particle Physics Lecture 3



Helmholtz Alliance



Terascale Statistics School https://indico.desy.de/event/46667/

DESY, Hamburg 24-28 Feb 2025



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

Outline

Monday 15:30 Hypothesis testing
 Monday 16:30 Hypothesis testing tutorial
 Tuesday 10:30 Parameter estimation tutorial
 Tuesday 14:00 Setting limits
 → Wednesday 9:00 Bayesian parameter estimation

Wednesday 11:00 Errors on errors

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 $p(x|\theta)$.

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

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

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

ML and Bayesian estimators

Purist Bayesian: $p(\theta|x)$ contains all knowledge about θ .

Pragmatist Bayesian: $p(\theta|x)$ could be a complicated function,

 \rightarrow summarize using an estimator $\hat{\theta}_{\text{Bayes}}$

Take mode of $p(\theta|x)$, (could also use e.g. expectation value)

What do we use for $\pi(\theta)$? No golden rule (subjective!), often represent 'prior ignorance' by $\pi(\theta) = \text{constant}$, in which case

$$p(\theta|\mathbf{x}) \propto P(\mathbf{x}|\theta) = L(\theta)$$

$$\longrightarrow \hat{\theta}_{\text{Bayes}} = \hat{\theta}_{\text{ML}}$$



ML and Bayesian estimators (2)

Note $\Box_{(}()$ = const. cannot be normalized – "improper prior".

Can be allowed for some problems; prior always appears multiplied by likelihood, so product $L(\theta)\pi_{\theta}(\theta)$ can result in normalizable posterior probability.

But... we could have used a different parameter, e.g., $\lambda = 1/\langle$, and if prior $\Box_{(}()$ is constant, then $\pi_{\lambda}(\lambda)$ is not:

$$\pi_{\lambda}(\lambda) = \pi_{\theta}(\theta) \left| \frac{d\theta}{d\lambda} \right| \propto \frac{1}{\lambda^2}$$

Maybe we know say we nothing about λ , so take $\pi_{\lambda}(\lambda) = \text{const.}$

Then
$$\hat{\lambda}_{\text{Bayes}} = \hat{\lambda}_{ML} \neq \frac{1}{\hat{\theta}_{\text{Bayes}}}$$

'Complete prior ignorance' is not well defined.

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 $\sigma_{\hat{\theta}_{0}}.$

$$x^{2}_{0} = \frac{1}{2\theta_{0}} \frac{1}{12\theta_{0}} \frac{1}{12\theta_$$

1.32

θ

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 $\hat{\theta}_0$.



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}

G. Cowan / RHUL Physics

Terascale Statistics 2025 / Lecture 3

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

Putting the ingredients into Bayes' theorem gives:

$$p(\theta_{0},\theta_{1}|\vec{y}) \propto \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi}\sigma_{i}} e^{-(y_{i}-\mu(x_{i};\theta_{0},\theta_{1}))^{2}/2\sigma_{i}^{2}} \pi_{0} \frac{1}{\sqrt{2\pi}\sigma_{t_{1}}} e^{-(\theta_{1}-t_{1})^{2}/2\sigma_{t_{1}}^{2}}$$

$$posterior \propto likelihood \times prior$$

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 $P(\mathbf{y},t|\theta_0,\theta_1)$ and the constant "Ur-prior" for θ_1 .

Here posterior only found as a proportionality.

G. Cowan / RHUL Physics

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

Terascale Statistics 2025 / Lecture 3

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 $p(\theta)$ up to a proportionality constant, generate a sequence of points θ_1 , θ_2 , θ_3 ,...

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

Proposal density
$$q(\theta; \theta_0)$$

e.g. Gaussian centred
about θ_0

$$\alpha = \min\left[1, \frac{p(\vec{\theta})q(\vec{\theta}_0; \vec{\theta})}{p(\vec{\theta}_0)q(\vec{\theta}; \vec{\theta}_0)}\right]$$

- 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: $p(\theta|x) \propto p(x|\theta)\pi(\theta)$.

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): $\alpha = \min \left[1, \frac{p(\vec{\theta})}{p(\vec{\theta}_0)}\right]$

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.

G. Cowan / RHUL Physics

Terascale Statistics 2025 / Lecture 3

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 :



Tutorial: Bayesian parameter estimation

The exercise is described

https://www.pp.rhul.ac.uk/~cowan/stat/exercises/bayesFit/

in the file bayes_fit_exercise.pdf.

The program is in bayesFit.py or bayesFit.ipynb.

This exercise treats the same fitting problem as seen with maximum likelihood, here using the Bayesian approach.

Bayes' theorem is used to find the posterior pdf for the parameters, and these are summarized using the posterior mode (MAP estimators).

The posterior pdf is marginalized over the nuisance parameters using Markov Chain Monte Carlo.

Gaussian signal on exponential background

Same pdf as from mlFit.py (see tutorial 1) with n = 400 independent values of x from

$$f(x|\boldsymbol{\lambda}) = \theta \frac{1}{\sqrt{2\pi\sigma}} e^{-(x-\mu)^2/2\sigma^2} + (1-\theta) \frac{1}{\xi} e^{-x/\xi}$$

Posterior pdf for parameters $\lambda = (\theta, \mu, \sigma, \xi)$ from Bayes theorem,

$$p(oldsymbol{\lambda}|\mathbf{x}) \propto p(\mathbf{x}|oldsymbol{\lambda})\pi(oldsymbol{\lambda}) \;, \;\;\; ext{where} \;\;\; p(\mathbf{x}|oldsymbol{\lambda}) = \prod_{i=1}^n f(x_i|oldsymbol{\lambda})$$

At first take prior pdf constant for all parameters subject to $0 \le \theta \le 1, \sigma > 0, \xi > 0$ (later try different priors).

G. Cowan / RHUL Physics

Data and MAP estimates

Maximize posterior with minuit (minimize $-\ln p(\lambda | \mathbf{x})$).



Standard deviations from minuit correspond to approximating posterior as Gaussian near its peak.

Here priors constant so MAP estimates same as MLE, covariance matrix $V_{ij} = \operatorname{cov}[\theta_i, \theta_j]$ also same.

A look at bayesFit.py

Find maximum of posterior with iminuit (minimize $-\ln p(\lambda | \mathbf{x})$), similar to maximum likelihood:

```
# Negative log-likelihood
def negLogL(par):
  fx = f(xData, par)
  return -np.sum(np.log(fx))
```

```
# Prior pdf
def prior(par):
    theta = par[0]
    mu = par[1]
    sigma = par[2]
    xi = par[3]
    pi_theta = 1. if theta >= 0. and theta <= 1. else 0.
    pi_mu = 1. if mu >= 0. else 0.
    pi_sigma = 1. if sigma > 0. else 0.
    pi_xi = 1. if sigma > 0. else 0.
    pi_xi = 1. if xi > 0. else 0.
    piArr = np.array([pi_theta, pi_mu, pi_sigma, pi_xi])
    pi = np.product(piArr[np.array(parfix) == False]) # exclude fixed par
    return pi
```

```
# Negative log of posterior pdf
def negLogPost(par):
    return negLogL(par) - np.log(prior(par))
```

minimize with iminuit

Metropolis-Hastings algorithm in bayesFit.py

```
# Iterate with Metropolis-Hastings algorithm
chain = [np.array(MAP)]
                            # start point is MAP estimate
numlterate = 10000
numBurn = 100
numAccept = 0
print("Start MCMC iterations: ", end="")
while len(chain) < numlterate:
  par = chain[-1]
  log_post = -negLogL(par) + np.log(prior(par))
  par prop = np.random.multivariate normal(par, cov prop)
  if prior(par prop) <= 0:
    chain.append(chain[-1]) # never accept if prob<=0.
  else:
    log_post_prop = -negLogL(par_prop) + np.log(prior(par_prop))
    alpha = np.exp(log post prop - log post)
    u = np.random.uniform(0, 1)
    if u \leq alpha:
      chain.append(par prop)
      numAccept += 1
    else:
      chain.append(chain[-1])
    if len(chain)%(numlterate/100) == 0:
      print(".", end="", flush=True)
chain = np.array(chain)
```

Try increasing number of iterations (10k runs in about 20 s).

Exercises on Bayesian parameter estimation (a)

1a) Run bayesFit.py, look at the plots

1(a) Run the program and examine the plots. These include:

- 1. The data values as ticks on the x axis together with the fitted curve evaluated with MAP estimators (Fig. 1 below). The uncertainties on the parameters correspond to the covariance $V_{ij} = \operatorname{cov}[\lambda_i, \lambda_j]$ that iminuit finds by approximating the posterior as a multivariate Gaussian near its maximum (similar to finding the covariance matrix of the MLEs).
- 2. Trace plots of each of the parameters (Fig. 2). In some problems it can be useful to discard a subset of the points (called "burn-in") if the starting point λ_0 is too far from the main concentration of the target density's probability; this is indicated in the trace plots with a vertical yellow bar.
- 3. Marginal distributions of the individual parameters (Fig. 3). The histograms are normalized to unit area and the MAP estimates are indicated with the vertical bars.
- 4. The autocorrelation function for the parameters (Fig. 4).

Exercises on Bayesian parameter estimation (b,c)

1b) Investigate effect of data sample size, fixing parameters and length of MCMC chains.

1(b) Change the data sample size from n = 400 to 200 and 1000 and note the changes in the results.

Using again n = 400, fix the parameters μ and σ (by changing the corresponding elements in the array **parfix** from False to True) and note the changes in the results. When finished, go back to having all four parameters free.

Change the number of MCMC iterations from 10000 to 100000 and note the change in the results, particularly in the structures you see in the trace plots. (This probably takes some time to run; for the rest of the exercises it is probably best to change back to 10000 iterations.

1c) Investigate changing the prior

1(c) Change the prior pdfs for ξ and σ to be $\pi(\xi) \propto 1/\xi$ and $\pi(\sigma) \propto 1/\sigma$ and note the change in the results. When finished, go back to constant priors.

Exercises on Bayesian parameter estimation (d)

1d) Include auxiliary measurement to constrain ξ

1(d) Suppose that one has an independent estimate u of the parameter ξ in addition to the n = 400 values of x. Treat u as Gaussian distributed with a mean ξ and standard deviation $\sigma_u = 0.5$ and take the observed value u = 5. Find the log-likelihood function that includes both the primary measurements (x_1, \ldots, x_n) and the auxiliary measurement u and modify the fitting program accordingly. Investigate how the results are affected by including u.

1e) Investigate point and interval estimates for θ

1(e) Using the functions cc_interval and HPD_interval provided in bayesFit.py, compute the central credible interval and HPD (highest probability density) interval for the parameter of interest θ using a credibility level of 68.3%. Compare these to the intervals one obtains from a point estimate (the MAP estimate, posterior median or posterior mean) plus or minus one standard deviation. For the standard deviation, try using both the sample standard deviation from the MCMC values and the standard deviation found by iminuit, which is based on a Gaussian approximation to the peak of the posterior. Find the estimates and intervales both with and without the auxiliary measurement of ξ as in (d) above and note how this effects the results.

MCMC trace plots

Take θ as parameter of interest, rest are nuisance parameters.

Marginalize by sampling posterior pdf with Metropolis-Hastings.



Gaussian proposal pdf, covariance U = sV, $s = (2.38)^2/N_{par} = 1.41$, gives acceptance probability ~ 0.24.

Here 10000 iterations (should use more).

Marginal distributions

MAP estimates shown with vertical bars



Note long tails.

Interpretation: data distribution can be approximated by Gaussian term only, (θ large, μ small) with large width ($\sigma \sim 4-8$) and a narrow exponential ($\zeta \sim 1-3$).



Autocorrelation versus lag

MCMC samples are not independent, autocorrelation function = correlation coefficient of sample x_i with x_{i+l} as a function of the lag, l, where x = any of θ , μ , σ , ξ minus its mean:



G. Cowan / RHUL Physics

Terascale Statistics 2025 / Lecture 3

Ways to summarize the posterior

Point estimates:

Posterior mode (MAP, coincides with MLE for constant prior).

Posterior median (invariant under monotonic transformation of parameter).

Posterior mean; coincides with above in large-sample limit. Intervals:

Highest Probability Density (HPD) interval, shortest for a given probability content, not invariant under param. trans. Central credible intervals, equal upper and lower tail areas, e.g., $\alpha/2$ for CL = $1 - \alpha$.

Point estimate +/- standard deviation, std. dev. from MCMC sample or by approximating core of posterior as Gaussian (from minuit); coincides with above in large-sample limit.

Types of intervals



G. Casella and R. Berger, Statistical Inference, 2002

Equal tail -