

9/11/19

Physics 8805: Learning from Data: Lecture 6

On board:

Notebooks for today:

- parameter estimation fitting straight line I.ipynb
- amplitude in presence of background.ipynb
- Metropolis Poisson example.ipynb
- MCMC - random-walk-and-sampling.ipynb

Websites for MCMC visualization (MCMC visualization page on Coursera):

<http://elevaroth.org/blog/2017/11/28/build-a-better-markov-chain/>
<https://chi-feng.github.io/mcmc-demo/>

** Questions or problems for Assignment 1?

For today: go back first to (26), (33), (34) [fitting straight line 1 notebook]

• Follow-up to priors: where does $p(m) \propto 1/(1+m^2)^{3/2}$ come from?• We can consider the line as $y = mx + b$ or $x = m'y + b'$
where $m' = 1/m$ and $b' = -b/m$ gives the same result.• Let us require that the priors on m, b have the same functional form as m', b' , because the labeling is arbitrary.• Then for the probabilities we must have

$$p(m, b) dm db = p(m', b') dm' db' = p\left(\frac{1}{m}, -\frac{b}{m}\right) \left| \frac{\partial m'}{\partial m} \frac{\partial b'}{\partial b} \right| dm db$$
• Evaluating the Jacobian gives $1/m^2$, so we need to solve the functional equation:

$$p\left(\frac{1}{m}, -\frac{b}{m}\right) = m^2 p(m, b)$$

• A solution is $p(m, b) = \frac{C}{(1+m^2)^{3/2}}$, as claimed,
 $C \leftarrow \text{arbitrary}$

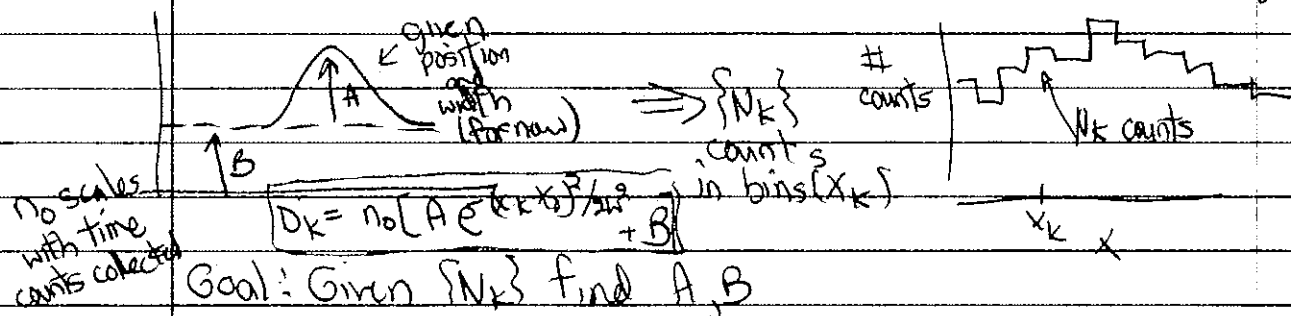
• How would you solve this without guessing the answer?

• Continue with the Sinia example on signal on top of background

9/11/19

(36)

Sina example: Amplitude of a signal in the presence of background



So what is the posterior we want?

$p(A, B | \{N_k\}, I)$ with $I = x_0, \sigma, \text{Gaussian, Flat background}$

Actual counts we get will be integers, and we can expect a Poisson distribution

$$\Rightarrow p(n|\mu) = \frac{\mu^n e^{-\mu}}{n!} \text{ for } n \geq 0 \text{ integer}$$

or

$$n \rightarrow N_k, \mu \rightarrow D_k \quad p(N_k | D_k) = \frac{D_k^{N_k} e^{-D_k}}{N_k!} \text{ for } k^{\text{th}} \text{ bin at } x_k$$

* What do we learn from the plots of the Poisson distribution?

$$p(A, B | \{N_k\}, I) \propto p(\{N_k\} | A, B, I) \times p(A, B | I)$$

posterior \propto likelihood \times prior

$$\Rightarrow L = \log[p(A, B | \{N_k\}, I)] = \text{constant} + \sum_{k=1}^M [N_k \log(D_k) - D_k]$$

• Choose constant for convenience: independent of A, B .

• Best point estimate: maximize $L(A, B)$ to find A_0, B_0 .

$$p(A, B | I) = \begin{cases} \text{constant} & 0 \leq A \leq A_{\max} \\ & 0 \leq B \leq B_{\max} \\ 0 & \text{otherwise} \end{cases}$$

• Look at code for likelihood and prior

• Uniform flat prior for $0 \leq A \leq A_{\max}, 0 \leq B \leq B_{\max}$

• Not sensitive to A_{\max}, B_{\max} if larger than support of likelihood

| Fig # | data bins | Δx | $(x_k)_{\max}$ | D_{\max} |
|-------|-----------|------------|----------------|------------|
| 1 | 15 | 1 | 7 | 100 |
| 2 | 15 | 1 | 7 | 10 |
| 3 | 31 | 1 | 15 | 100 |
| 4 | 7 | 1 | 3 | 100 |

9/11/19

(137)

Comments on Figures.

Fig 1: 15 bins and $D_{\max} = 100$

- Contours are at 20% intervals showing height.
- Read off best estimates and compare to Γ_{true} .
 - does find signal is about half background
- Marginalization of B
 - what if we don't care about B? "nuisance parameter"

$$p(A | \{N_k\}, I) = \int_0^{\infty} p(A, B | \{N_k\}, I) dB$$

compare to $p(A | \{N_k\}, B_{\text{true}}, I) \Rightarrow$ plotted on graph

- Also can marginalize over A

$$p(B | \{N_k\}, I) = \int_0^{\infty} p(A, B | \{N_k\}, I) dA$$

- See how these are done in code: B-marginalized
B true fixed
- note the normalization at the end.

- Set extra plots to true

• different representations of same info and contours in
first 3. Last one is attempt at 68%, 95%, 99.7% but looks wrong.

- * note difference between contours showing pdf height and showing integrated volume.

- Look at the other Figures and draw conclusions!

- How should you design your experiments?

Eg, how should you bin data, how many counts are needed, what $(x_k)_{\max}$, and so on.

9/11/19

(38)

Why MCMC? (Based on Gregory, Chap 12)

We have been emphasizing that in the Bayesian approach, everything is a pdf. One type of pdf is for the parameters of theory, which we'll denote by the vector $\vec{\theta}$, given data D :

$$p(\vec{\theta} | D, I)$$

and information I . Suppose we have a theoretical model for this.

- Maybe these are the LECs for an effective field theory Hamiltonian. And now we want to calculate the expectation value of a function of $\vec{\theta}$: $\langle f(\vec{\theta}) \rangle$. Or $\vec{\theta}$ characterizes a signal and background.
- As we discussed in doing the central limit theorem:

$$\langle f(\vec{\theta}) \rangle = \int f(\vec{\theta}) p(\vec{\theta} | D, I) d\vec{\theta} = \int g(\vec{\theta}) d\vec{\theta}$$

Note that this is more than traditional calculations, in which we would have single values of $\vec{\theta}$, e.g. denoted $\hat{\vec{\theta}}$, that we might have found by minimizing a χ^2 . E.g. we identified the particular values of $\theta_1, \theta_2, \dots, \theta_n$ that best reproduced scattering data. Then we would calculate $f(\hat{\vec{\theta}})$, which might be the binding energy of a nucleus.

- But $\langle f(\vec{\theta}) \rangle$ means we do a multidimensional integral over the full range of possible $\vec{\theta}$ values, weighted by the probability density function $p(\vec{\theta} | D, I)$, which we have worked out.

- This is a lot more work!

- We frequently also have a situation where we want to integrate (marginalize) over a subset of parameters $\vec{\theta}_1$ to find a probability for the rest $\vec{\theta}_2$. E.g. over parameters for the width of a signal and other parameters characterizing our model for the Higgs mass.

- These multidimensional integrals then become a necessity to do, but conventional methods for low dimension (Gaussian quadrature or Simpson's rule) become inadequate rapidly with the increase of dimension.

9/11/19

(40)

Bottom line: it's not feasible to draw a series of independent random samples from $p(\tilde{\theta} | D, I)$ for larger $\tilde{\theta}$.

- remember, independent means if $\tilde{\theta}_1, \tilde{\theta}_2, \dots$ is the series, knowing $\tilde{\theta}_1$ doesn't tell us anything about $\tilde{\theta}_2$.

*** But the samples don't need to be independent, they just need to generate $p(\tilde{\theta} | D, I)$ in the correct proportions (e.g. as indicated by histogramming the samples, it approximates $p(\tilde{\theta} | D, I)$).

* \Rightarrow Do a random walk in the parameter space of $\tilde{\theta}$ so that the probability for being in a region is proportional to $p(\tilde{\theta} | D, I)$ for that region.

- $\tilde{\theta}_{i+1}$ follows from $\tilde{\theta}_i$ by a transition probability (kernel)
 $\Rightarrow p(\tilde{\theta}_{i+1} | \tilde{\theta}_i)$

- assumed to be "time independent", so same $p(\tilde{\theta}_{i+1} | \tilde{\theta}_i)$ no matter when you do it.

\Rightarrow Markov chain and method is Markov chain Monte Carlo.

Basic structure of algorithm:

① Given $\tilde{\theta}_i$, propose a value for $\tilde{\theta}_{i+1}$, call it $\tilde{\phi}$, sampled from $q(\tilde{\phi} | \tilde{\theta}_i)$. This q could take many forms, so for concreteness imagine it as a multivariate normal with mean given by $\tilde{\theta}_i$ and variance $\tilde{\sigma}^2$.

- decreased probability as you get away from current sample
- $\tilde{\sigma}$ determines the step size.

② Decide whether or not to accept candidate $\tilde{\phi}$ for $\tilde{\theta}_{i+1}$. Here we'll use a Metropolis condition (later we'll see other ways but may be better).

- This dates from the 1950's in physics but didn't become widespread in statistics until almost 1980.
- Enabled Bayesian methods to take off.

9/11/19

(41) 6)

Calculate Metropolis ratio:

$$r = \frac{\overset{\text{proposed}}{p(\vec{\phi} | D, \pm)} q(\vec{\theta}_0 | \vec{\phi})}{\underset{\text{current}}{p(\vec{\theta}_i | D, \pm) q(\vec{\phi} | \vec{\theta}_i)}} \quad \left. \begin{array}{l} q \text{ may be symmetric } q(\vec{\theta}_1 | \vec{\theta}_2) \\ \text{if so } \Rightarrow \text{"Metropolis"} = q(\vec{\theta}_2 | \vec{\theta}_1) \\ \text{If not, then "Metropolis-Hastings."} \end{array} \right\}$$

Decision:

If $r \geq 1$, set $\vec{\theta}_{i+1} = \vec{\phi}$ accept

if $r < 1$, so less probable, don't always reject!
accept with probability r (remember $0 \leq r \leq 1$)
by sampling a uniform $(0, 1)$ distribution.

If $U \sim \text{Uniform}(0, 1)$ is $U \leq r$, then $\vec{\theta}_{i+1} = \vec{\phi}$
else $\vec{\theta}_{i+1} = \vec{\theta}_i$

Note that the last case means you do have a $\vec{\theta}_{i+1}$, but it is the same as $\vec{\theta}_i$ (so the chain continues to grow).

Acceptance probability is the minimum of 1, r

Algorithm pseudo code:

1. initialize $\vec{\theta}_i$, set $i \leftarrow 0$

2. Repeat {
 Obtain new candidate $\vec{\phi}$ from $q(\vec{\phi} | \vec{\theta}_i)$
 Sample $U \sim \text{uniform}(0, 1)$
 If $U \leq r$ set $\vec{\theta}_{i+1} = \vec{\phi}$, else set $\vec{\theta}_{i+1} = \vec{\theta}_i$
 $i++$
}

Plan: ① look at visualizations

② look at a basic example for Poisson distribution

③ consider MCMC-random-walk and sampling, notebook.

④ look at encee example from Assignment 1.