

8/23/19

Physics 8805: Learning from Data Lecture 2

Handouts: Extra copies of simple sum-product rule.ipynb
Copies of medical example by Bayes.ipynb

Before class: Set up projector with laptop and notebooks
for modules 0 and 1. Put up Carmen 8805 Modules.

On board: • Give feedback!!

- Questions on Anaconda installation?
- Questions on Python or Jupyter notebooks?
- Questions on course logistics?
- Have you used git? Google "Furnstahl github 8805" for repository

Overview:

• We took a quick look at continuous pdfs, but before looking further (finishing ⑧), we'll put sum and product rules on the table.

• Do ⑨ and ⑩

• Class: fill out simple sum-product rule.ipynb in note book or hardcopy.

• Brief guidance on how you might fill in the table

• `fstring print(f'ratio = {ratio:.3f}')` or `np.around(number, digits)`

• How do you use numpy arrays?

• experts: write function to work with either ints or floats or numpy arrays.

• Recap of simple sum-product rule.

• medical example by Bayes for homework.

• Return to exploring pdfs.ipynb and finish ⑧.

Take-aways: Bayesian confidence intervals, various "point estimates", characteristics of different pdfs (eg. symmetry, heavy tails, ...), what "sampling" means, projected posteriors, ...

identify these
↓

8/23/19

Bayesian updating via Bayes' Theorem

$$\overbrace{p(\theta | \text{data}, I)}^{\text{posterior}} = \frac{\overbrace{p(\text{data} | \theta, I)}^{\text{likelihood}} \times \overbrace{p(\theta | I)}^{\text{prior}}}{\underbrace{p(\text{data} | I)}_{\text{data probability (or "fully marginalized likelihood" or evidence or ...)}}}$$

↑
general
vector of
parameters

data probability
(or "fully marginalized likelihood"
or evidence or ...)

⇒ more on denominator later, but often just a normalization

• The prior pdf is what information we have (or believe) about θ before we observe the data.

• The posterior pdf is our new pdf, given that we have observed the data.

• So Bayes' Theorem tells us how to update our expectations.

* Coin tossing example to illustrate updating:
Bayesian updating coinflip interactive.ipynb

Storyline: We are observing successive flips of a coin (or any binary process). There is a definite true probability of getting heads, (p_h) true but we don't know what it is, other than it is between 0 and 1.

- We characterize our information about p_h as a pdf.
- Before any flips, we start with a preconceived notion of the probability, this is the prior pdf $p(p_h | I)$ [I is only info we have]
- With each flip of the coin, we gain additional information, so we update our expectation of p_h , so we find the posterior:

$$p(p_h | \overset{N}{\# \text{ losses}}, \overset{R}{\# \text{ heads}}, I)$$

• note that outcome is discrete but p_h is continuous as $p_h \in [0, 1]$

(13)

8/23/19

Let's first play a bit with the simulation and then come back and think of the details.

- Note a few of the Python Features *class structure*
 - class for data called Data. Easy compared to C++!
 - function to make a plot that is made repeatedly
 - elaborate widget \Rightarrow use as guide for making your own!
(read from bottom up to understand the structure)

• User interface Features

- tabs to control parameters or look at documentation
- set the true p_0 by the slider
- press "Next" to flip "jump" # of times, Reset to go back to beginning
- plot shows updating from three different initial prior pdfs.

- Class: tell your neighbor how to interpret each of the priors.
 - uniform prior: any probability is equally likely. Is this uninformative?
 - informative prior: we have reason to believe the coin is fair.
 - anti-prior: Could be anything but most likely a two-headed or two-tailed coin.

What is minimal common info?
 $\Rightarrow 0 \leq p_0 \leq 1$
 and 'Separable'?

Things to try:

- First one flip at a time. How do you understand the changes intuitively?
- What happens with more and more tosses?
- Try different values of the true p_0 .
- What happens with enough data? All converge to narrow pdf including p_0 true (but not really)
- Which prior(s) get to the correct conclusion fastest for $p_0 = .4, .9, .5$?
- Does it matter if you update after every toss or all at once?
 - We'll need to look at details!

What about the difference between a particular ordering of data for R heads versus any order?

(14)

8/23/19

Suppose we had a fair coin $\Rightarrow p_h = 0.5$

$$\Rightarrow p(R \text{ heads out of } N \text{ tosses} | \text{fair coin}) = \binom{N}{R} (0.5)^R (0.5)^{N-R}$$
$$\Rightarrow p(R, N | p_h = 0.5)$$

combinations = $\frac{N!}{R!(N-R)!}$

Is the sum rule obeyed?

$$\sum_{R=0}^N p(R, N | p_h = \frac{1}{2}) = \sum_{R=0}^N \binom{N}{R} (\frac{1}{2})^N = (\frac{1}{2})^N \cdot 2^N = 1 \checkmark$$

$$[\text{proof: } (x+y)^N = \sum_{R=0}^N \binom{N}{R} x^R y^{N-R} \xrightarrow{x=y=1} 2^N = \sum_{R=0}^N \binom{N}{R}]$$

More general p_h !

$\downarrow x=p_h, y=1-p_h$
 \Rightarrow sum rule works in general!

$$p(R, N | p_h) = \binom{N}{R} p_h^R (1-p_h)^{N-R}$$

But we want the other way around: $p(p_h | R, N)$

$$\text{Bayes says: } p(p_h | R, N) \propto p(R, N | p_h) \cdot p(p_h)$$

denominator
just normalizes,
doesn't depend on p_h .

\leftarrow So we don't need to distinguish particular ordering.

Claim: Can do sequentially or all at once. When is this true?
(When the tosses are independent)

What would happen if you tried to update using the same results over and over again?

8/23/19

So how are we doing the calculation of the updated posterior?
 In this case, we can do analytic calculations.

Case I: uniform (flat) prior

$$\Rightarrow p(p_n | R, N, I) = N \underset{\text{suppress}}{p(R, N | p_n)} \overset{1}{p(p_n)} = N p_n^R (1-p_n)^{N-R}$$

$$\text{But } \int_0^1 dp_n p(p_n | R, N) = 1 \Rightarrow N \frac{\Gamma(1+N-R) \Gamma(1+R)}{\Gamma(2+N)}$$

$$[\text{Recall Beta function } B(x, y) = \int_0^1 t^{x-1} (1-t)^{y-1} dt = \frac{\Gamma(x) \Gamma(y)}{\Gamma(x+y)}]$$

and $\Gamma(x) = (x-1)!$ for integer x

$$\Rightarrow N = \frac{\Gamma(2+N)}{\Gamma(1+N-R) \Gamma(1+R)} \quad \text{and we can trivially update.}$$

Case II: conjugate prior: posterior has same form as prior.
 Here if we pick a beta distribution as prior, it is conjugate.
 From scipy.stats.beta documentation:

$$\text{ia. } p(x|a, b) = \frac{\Gamma(a+b) x^{a-1} (1-x)^{b-1}}{\Gamma(a) \Gamma(b)} \Rightarrow \text{likelihood is } f(p_n, 1+R, 1+N-R)!$$

If prior $p(p_n | I) = f(p_n, \alpha, \beta)$, then normalized posterior

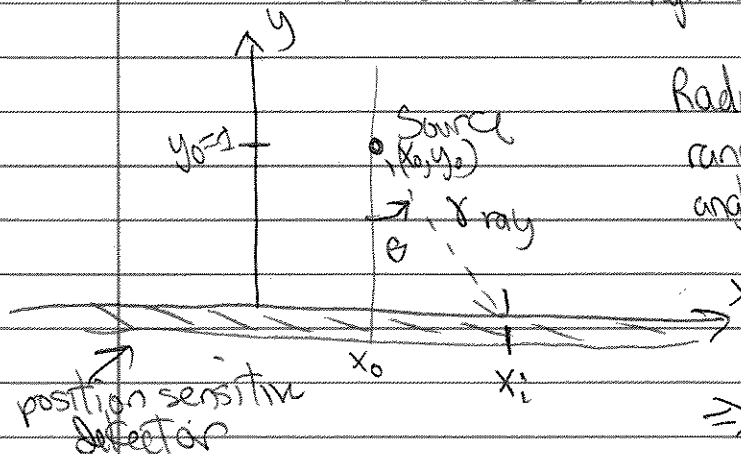
$$p(p_n | R, N) \propto p(R, N | p_n) p(p_n) = f(p_n, \alpha+R, \beta+N-R)$$

So update by $\alpha \Rightarrow \alpha+R$, $\beta \Rightarrow \beta+N-R$!

Check against code! Valuable to use conjugate prior if it is

8/23/19

First look at the radioactive lighthouse problem
 \Rightarrow radioactive_lighthouses_exercise.ipynb



Radioactive source emits gamma rays randomly in time but uniformly in angle. The source is at (x_0, y_0) .

Gamma rays are detected on the x-axis and positions recorded
 $\Rightarrow x_1, x_2, x_3, \dots, x_N \Rightarrow \{x_k\}$

Goal: Given the $\{x_k\}$, estimate (x_0, y_0) .

1st pass: take $y_0=1$ as given, so only estimate x_0 .

Naively, how would you estimate x_0 given x_1, \dots, x_N ? (Take average)
 • Let's compare the naive approach to a Bayesian approach.

Follow the notebook leading questions for the Bayesian estimate!