

8/28/19

Physics 8805: Learning From Data: Lecture 3Before class: Set up projector with laptop and notebooks:

- medical example key
- Flipping coins
- exploring pdfs
- Parameter estimation: Gaussian noise
- Parameter estimation: fitting straight line I

On board:

I will be out of town on Friday but we will have a
 Jupiter notebooks lab session facilitated by Xilin and Jordan.
 You'll work on:

- * • radioactive lighthouse (quick overview today)
- exploring pdfs
- parameter estimation (two notebooks available)

Step through medical example key.
 Comments:

2. Why is it $p(H|D)$ and not $p(H, D)$?
 Recall that $p(H, D) = p(H|D) \cdot p(D) \leftarrow$ ^{if you know} $p(D)=1$, then they are the same.
 \leftarrow generally what you are interested in.

3, 4 straight forward

5. Emphasis on sum rule. Why didn't a column of the sum-product rule notebook add to 1? Because we were looking at $p(\text{tall}, \text{blue}) + p(\text{short}, \text{blue}) \neq 1$ whereas $p(\text{tall} | \text{blue}) + p(\text{short} | \text{blue}) = 1$ ✓

6. Emphasize usefulness of Bayes' rule to express $p(H|D)$ in terms of $p(D|H)$.
 7. straight forward

8. Standard, but not so obvious at first. After familiar, jump right to end.

9. Stress the analysis of the result.

8/28/19

Recap of coin flipping notebook

- Recall the names of the pdfs in Bayes' rule: posterior, likelihood, prior, evidence; and Bayesian updating: prior + data \rightarrow posterior \rightarrow new prior

Take-aways from coin flipping. \Rightarrow Note: added "New data" button (show code)

① Different priors eventually give the same posterior with enough data \Rightarrow This is called Bayesian convergence.

- How many tosses are enough?

- Hit "New data" multiple times to see fluctuations.
- Clearly depends on p_n and how close you want the posteriors to be
- $p_n = 0.4 \Rightarrow \approx 200$ tosses get you most of the way
- $p_n = 0.9 \Rightarrow$ much longer for informative prior.

② Ask class: Why does the "anti-prior" work well even though its dominant assumptions are proven wrong early on.

\Rightarrow "heavy tails" mean it is like uniform (renormalized!) after the ends are eliminated. A lesson for formulating priors: allow for deviations from your expectations.

③ Return to pages (14) and (15) and go through them, code: `y[i] = stats.beta.pdf(X, alpha_i + heads, beta_i + N - heads)`

④ Is there a difference between updating sequentially or all at once?

Do simplest problem first: two tosses

- Let results be $D = \{D_k\}$ (in practice D 's all 1's $\Rightarrow R = \sum_k D_k$)

$$\text{General } p(p_n | \{D_k\}, I) \propto p(\{D_k\} | p_n, I) p(p_n | I)$$

$$K=1 \Rightarrow p(p_n | D_1, I) \propto p(D_1 | p_n, I) p(p_n | I)$$

class: why?

Bayes' rule

Bayes' rule ($D_1 \in I$)

tosses are independent

Last line is sequential! (Prior for 2nd flip is posterior from first flip)

$K=2$ only \Rightarrow so all at once is same as sequential as function of p_n when renormalized

8/28/19

(19)

To go to 3:

$$p(p_h | D_1, D_2, D_3, I) \propto p(D_3 | p_h, I) \times p(p_h | D_1, D_2, I) \\ \propto p(D_3 | p_h, I) \times p(D_2 | p_h, I) \times p(D_1 | p_h, I) p(p_h)$$

and so on.

(5) What about "bootstrapping"?

Why can't we use the data to improve the prior and apply it (repeatedly) for the same data. Consider an extreme case!

$$p_2(p_h | D_1, I) \propto p(D_1 | p_h, I) p(p_h | I)$$

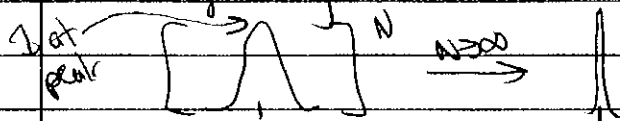
$$\Rightarrow p_2(p_h | D_1, I) \propto p(D_1 | p_h, I) p_1(p_h | D_1, I) \\ \propto [p(D_1 | p_h, I)]^2 p(p_h | I)$$

keep going?

$$p_N(p_h | D_1, I) \propto p(D_1 | p_h, I) p_{N-1}(p_h | D_1, I) \\ \propto [p(D_1 | p_h, I)]^N p(p_h | I)$$

Suppose D_1 was 0, then $[p(\text{tails} | p_h, I)]^N \propto (1-p_h)^N p(p_h | I) \rightarrow \text{only } p_h = 0$
 If D_1 was 1, then $[p(\text{heads} | p_h, I)]^N \propto p_h^N p(p_h | I) \rightarrow \text{only } p_h = 1$

More generally it would cause the posterior to get narrower and narrower.



Moral: Don't do that!

Something to come back to: Frequentist point estimates.

Maximum-likelihood means: What value of p_h maximizes $L = N p_h^R (1-p_h)^{N-R}$

$$\frac{d}{dp_h} () = N(R p_h^{R-1} (1-p_h)^{N-R} - (N-R) p_h^R (1-p_h)^{N-R-1}) = 0 \Rightarrow \boxed{p_h = \frac{R}{N}}$$

$$\text{Similarly, } \sigma = \sqrt{\frac{p_h(1-p_h)}{N}}$$

Revisit next week in context of Gaussian approximation.

8/28/19

(20)

- Set up radioactive lighthouse problem as another Bayesian updating (and parameter estimation) problem. See page (16).
 - Go through notebook on Friday.
- Interlude: bring up exploring pdfs, ipynb
 - Point out projected posterior plots and preview what we learn.
 - Look at multimodal plot and note that projections don't always show multiple modes.
 - Sampling from 1d pdfs \Rightarrow fluctuating.

Parameter estimation

Overview comments

- In general terms, "parameter estimation" in physics means obtaining values for parameters (constants) that appear in a theoretical model that describes data. (Exceptions exist, of course)
 - examples:
 - couplings in a Hamiltonian
 - coefficients of a polynomial or exponential model of data
- Conventionally this process is known as "fitting the parameters" and the goal is to find the "best fit" and maybe error bars.
- We will make particular interpretations of these phrases from our Bayesian point of view.
- Plan: set up the problem and look at how familiar ideas like "least-squares fitting" show up from a Bayesian perspective.
- As we proceed, we'll make the case that for physics a Bayesian approach works well.