

9/6/19

(29)

Physics 8805: Learning from Data: Lecture 5

Before class: Set up laptop with Jupyter notebooks:

list
on
board

- (A) radioactive_lighthouse_exercise_key.ipynb
- (B) parameter_estimation_Gaussian_noise.ipynb
- (C) parameter_estimation_fitting_straight_line_T.ipynb
- (D) assignment_01.ipynb

On board: • Office hours next week: ?? (What works?)

• Assignment 1, due next Friday, is a follow-up to the radioactive lighthouse exercise using MCMC.

Due: Complete notebook (D) by combining and extending results and code from (A) and (B) [or any other of our notebooks]

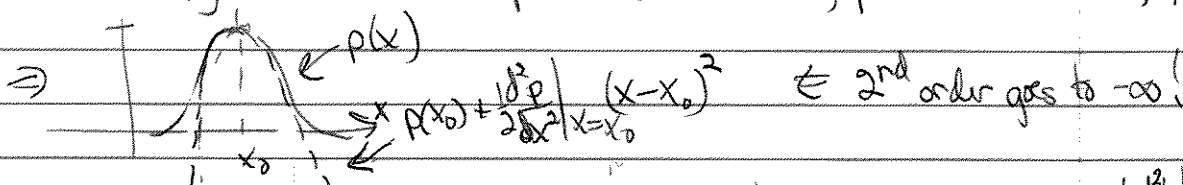
Learning goals: (remember: toss in water to learn how to swim)

note:
shouldn't need
to recalculate
anything
⇒ focus on
notebook
ingredients

- Be able to re-use md and Python from existing notebooks to perform similar tasks (even if not understanding all details)
 - ⇒ generating data, sampling via MCMC, making plots.
- Apply basic ideas of Bayesian statistics ⇒ Bayes theorem, priors, sampling posterior
- Analyze results (with hints)
- Try out markdown.

• Follow-up to discussion of Gaussian approximation to a pdf (see (27))

• Critical problem (not yet mentioned) with doing a straight Taylor series approximation to a pdf: doesn't satisfy conditions for pdf (normalizable, positive definite).



• In contrast $L(x) = \log p(x) \doteq L(x_0) + \frac{1}{2} \frac{d^2 L}{dx^2} \bigg|_{x=x_0} (x-x_0)^2 \Rightarrow p(x) \propto \mathcal{A} e^{-\frac{1}{2} \frac{d^2 L}{dx^2} \bigg|_{x=x_0} (x-x_0)^2}$
 is positive definite, renormalizable, and a higher order approximation to $p(x)$ because it includes all orders in $(x-x_0)^2$.

9/6/19

Let's recap notebook (B) in light of Assignment 1.

- Jump to Bayesian approach and come back to contrast with frequentist approach. Goal: sample a posterior

two random variables

data (x points)

data (where flashes hit (x_k))

$$p(\underbrace{\mu, \sigma}_\theta | D, I) \rightarrow p(\underbrace{x_0, y_0}_\theta | D, I) \quad \text{same general problem!}$$

$$\Rightarrow p(\theta | D, I)$$

- What do we need? From Bayes Theorem, need
 - likelihood: $p(D | \mu, \sigma, I) \rightarrow p(D | x_0, y_0, I)$
 - priors $p(\mu, \sigma | I) \rightarrow p(x_0, y_0 | I)$
 } you figure out based on lighthouse exercise

- Note that you are generalizing the functions for log pdfs and the plotting of posteriors in notebook (A)

- Note in (B) the functions for log-prior and log-likelihood,
 - Here $\theta = [\mu, \sigma]$ is a vector of parameters,
 - $\Rightarrow \theta = [x_0, y_0]$

- Step through set up for emcee \leftarrow need to install emcee and corner. Best to use environment, but email if difficulties.

• hint: nothing here will change!!

- More next week on what is happening, but basically we are doing 50 random walks in parallel to explore the posterior. Where the walkers end up will define our samples of μ, σ
 - \Rightarrow the histogram is an approximation to the joint posterior.

- Plotting is also the same, once you change labels and $\mu_{\text{true}}, \sigma_{\text{true}}$ to $x_0_{\text{true}}, y_0_{\text{true}}$. (And skip maxlike part!)

- Analysis:

- maximum likelihood here is frequentist estimate \rightarrow optimization problem
- Are μ and σ correlated or uncorrelated?
- Read off marginalized estimates for μ and σ .

- Return as time permits to: Bayesian vs. frequentist probability on (24), (28) and confidence intervals denominator $p(D | I)$ on (25) and why it is a normalization

9/6/19

General form of the central limit theorem (CLT):

The sum of n random variables that are drawn from any pdf(s) & finite variance σ^2 tends as $n \rightarrow \infty$ to be Gaussian distributed about the expectation value of the sum, with variance $n\sigma^2$.
 (So we scale the sum by $\frac{1}{\sqrt{n}}$ - see next page.)

Consequences:

1. The mean of a large number of values becomes normally distributed regardless of the probability distribution from which the values are drawn. (This fails for lighthouse!)

2. Functions such as the Binomial and Poisson distribution all tend to look like Gaussian distributions in the limit of a large number of drawings

$$\text{P.e.g., } P_n = \frac{x^n e^{-x}}{n!} \text{ (n integer)} \xrightarrow[n \rightarrow x \rightarrow \text{large}]{x \rightarrow \text{large}} p(x) = \frac{e^{-(x-1)^2/2}}{\sqrt{2\pi x}} \quad \leftarrow \begin{matrix} \mu=1 \\ \sigma^2=1 \end{matrix}$$

(Class: How would you verify this in a Jupyter notebook?
 How would you prove it analytically?)

Start with independent random variables x_1, \dots, x_n drawn from a distribution with mean $\langle x \rangle = \int x p(x) dx = 0 \leftarrow \text{generalize later}$ and $\langle x^2 \rangle = \sigma^2$ [$\langle x^n \rangle = \int x^n p(x) dx$]

$$\text{Let } X = \frac{1}{\sqrt{n}}(x_1 + x_2 + \dots + x_n) = \sum_{j=1}^n \frac{x_j}{\sqrt{n}} \quad \left[\text{need to scale by } \frac{1}{\sqrt{n}} \text{ for finite } X \text{ in } n \rightarrow \infty \text{ limit} \right]$$

What is the distribution of X ?

\Rightarrow call it $p(X|I)$ where I is the information about how X is drawn.

Plan: Use the sum and product rule and their consequences to relate $p(X)$ to what we know of $p(x_j)$. [Be careful of large X vs. small x_j !]

9/6/19

Class: fill in the rule used to justify the following steps:
 (Note: we'll suppress I to keep from getting too cluttered.)

introduce x_i
 into problem

$$p(X) = \int_{-\infty}^{\infty} dx_1 \dots dx_n p(X, x_1, \dots, x_n)$$

why?

marginalization

$$= \int_{-\infty}^{\infty} dx_1 \dots dx_n p(X | x_1, \dots, x_n) p(x_1, \dots, x_n)$$

product rule

$$= \int_{-\infty}^{\infty} dx_1 \dots dx_n p(X | x_1, \dots, x_n) p(x_1) \dots p(x_n)$$

independence

What is $p(X | x_1, \dots, x_n)$? $\Rightarrow \delta(X - \frac{1}{\sqrt{n}}(x_1 + \dots + x_n))$

Rather than use the δ -function to evaluate one integral, use a Fourier representation:

$$\delta(X - \frac{1}{\sqrt{n}}(x_1 + \dots + x_n)) = \frac{1}{2\pi} \int_{-\infty}^{\infty} dw e^{i w (X - \frac{1}{\sqrt{n}} \sum_{j=1}^n x_j)} = \frac{1}{2\pi} \int_{-\infty}^{\infty} dw e^{-i w X} \prod_{j=1}^n \left[\int_{-\infty}^{\infty} dx_j e^{\frac{i w x_j}{\sqrt{n}}} p(x_j) \right]$$

Observe that the terms in $[\]$ are all the same and they have factored!

Suppose we Taylor expand $e^{\frac{i w x_j}{\sqrt{n}}}$, assuming that the Fourier integral is dominated by small x as $n \rightarrow \infty$: (when does this fail?)

$$e^{\frac{i w x}{\sqrt{n}}} = 1 + \frac{i w x}{\sqrt{n}} + \frac{(i w)^2 x^2}{2n} + O\left(\frac{w^3 x^3}{n^{3/2}}\right) \Rightarrow \int_{-\infty}^{\infty} dx p(x) \left[1 + \frac{i w x}{\sqrt{n}} + \frac{(i w)^2 x^2}{2n} + \dots \right]$$

$$\Rightarrow = 1 + \frac{i w}{\sqrt{n}} \langle x \rangle - \frac{w^2}{2n} \langle x^2 \rangle + \langle x^3 \rangle O\left(\frac{w^3}{n^{3/2}}\right) \leftarrow \sigma^2 \text{ assumed to be finite}$$

$$\Rightarrow p(X) = \frac{1}{2\pi} \int_{-\infty}^{\infty} dw e^{-i w X} \left[1 - \frac{w^2}{2n} \sigma^2 + O\left(\frac{w^3}{n^{3/2}}\right) \right]^n \text{ but } \lim_{n \rightarrow \infty} \left(1 + \frac{a}{n}\right)^n = e^a$$

$$\Rightarrow \frac{1}{2\pi} \int_{-\infty}^{\infty} dw e^{-i w X} e^{-\frac{w^2 \sigma^2}{2}} = \frac{1}{\sqrt{2\pi} \sigma} e^{-\frac{X^2}{2\sigma^2}} \quad (Q.E.D.)$$

To generalize to $\langle x \rangle \neq 0$, consider $X = [x_1 + \dots + x_n - n\mu] / \sqrt{n}$ and change to $y_i = x_i - \mu \Rightarrow$ sum of y_i 's applies.

Why does this work?

9/6/19

The pdfs we've seen for $p(\mu, \sigma | D, I)$ were characterized by elliptical contours of equal probability density whose major axes are aligned with the μ and σ axes.

We have commented that this is a signal of independent random variables.

Let's look at a case where this is not true and then look analytically at what we should expect with correlations.

So return to notebook (C) on fitting a straight line.

Step through (D), renewing the statistical model.

What are we trying to find? $p(\theta | D, I)$, just as for (B) and (D), with now $\theta = [b, m]$.

intercept \rightarrow slope

Comments on notebook:

- note that x_i is also randomly distributed uniformly
- log likelihood gives 'fluctuating' results whose size depend on # of data points N and standard deviation of noise σ_y .
- \Rightarrow if time, explore in exercise session how size varies with N .

intercept
 $b = 25.0$
slope
 $m = 0.5$
 $\sigma = 5.0$

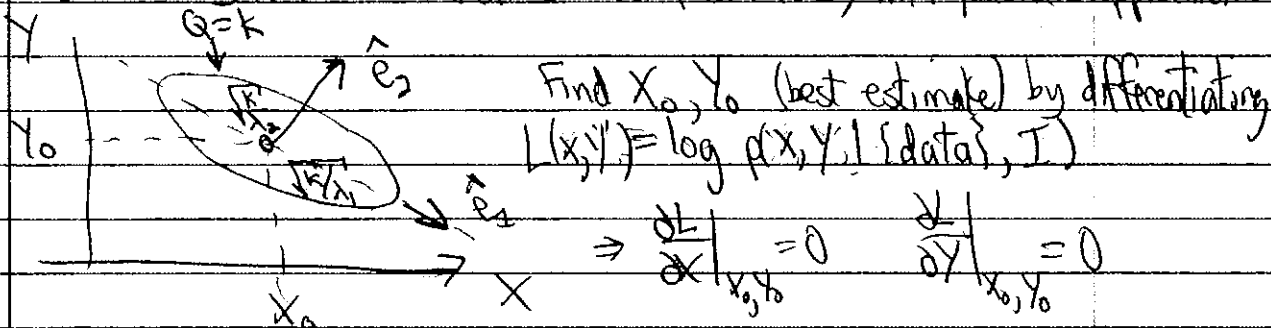
- Compare priors on slope \Rightarrow uniform in m vs. uniform in angle
- implementation of plots comparing priors - class comments
 - with first set of data with $N=20$ points, does prior matter? No!
 - with second set " " " $N=3$ points, " " " ? Yes!
- note $\log \text{posterior} = \log(\text{likelihood}) + \log(\text{prior})$
 - maximum taken to be 1 for plotting
 - exponentiate: $\text{posterior} = \exp(\log \text{posterior})$

* What does it mean that the ellipses are slanted? (coming soon!)

2nd set of data: flat gives $b = -50 \pm 75$, $m = 1.5 \pm 1$ so badly in σ
symmetric gives $b = 25 \pm 50$, $m = .5 \pm .75$ so much better!
Switch to corner!

9/6/19

Likelihoods with two variables (or posteriors) with quadratic approximation



• To check reliability, Taylor expand around $L(x_0, y_0)$:

$$L = L(x_0, y_0) + \frac{1}{2} \left[\left. \frac{\partial^2 L}{\partial x^2} \right|_{x_0, y_0} (x - x_0)^2 + \left. \frac{\partial^2 L}{\partial y^2} \right|_{x_0, y_0} (y - y_0)^2 + 2 \left. \frac{\partial^2 L}{\partial x \partial y} \right|_{x_0, y_0} (x - x_0)(y - y_0) + \dots \right] + \dots \equiv L(x_0, y_0) + \frac{1}{2} Q + \dots$$

Makes sense to do this in matrix notation

$$Q = \begin{pmatrix} x - x_0 & y - y_0 \end{pmatrix} \begin{pmatrix} A & C \\ C & B \end{pmatrix} \begin{pmatrix} x - x_0 \\ y - y_0 \end{pmatrix} \quad \text{symmetric}$$

$$A = \left. \frac{\partial^2 L}{\partial x^2} \right|_{x_0, y_0} \quad B = \left. \frac{\partial^2 L}{\partial y^2} \right|_{x_0, y_0} \quad C = \left. \frac{\partial^2 L}{\partial x \partial y} \right|_{x_0, y_0}$$

• So in quadratic approximation, the contour $Q=k$ for some k is an ellipse centered at x_0, y_0 . Orientation and eccentricity determined by A, B , and C .

• Principal axes found from eigenvectors of $\begin{pmatrix} A & C \\ C & B \end{pmatrix}$ (Hessian matrix)

$$\begin{pmatrix} A & C \\ C & B \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} = \lambda \begin{pmatrix} x \\ y \end{pmatrix} \Rightarrow \lambda_1, \lambda_2 < 0 \text{ (so } x_0, y_0 \text{ is a maximum)}$$

What if ellipse is skewed? We can marginalize (discussed in future lecture)

