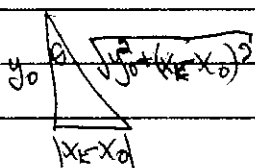


9/4/19

Physics 880S: Learning from Data: Lecture 4

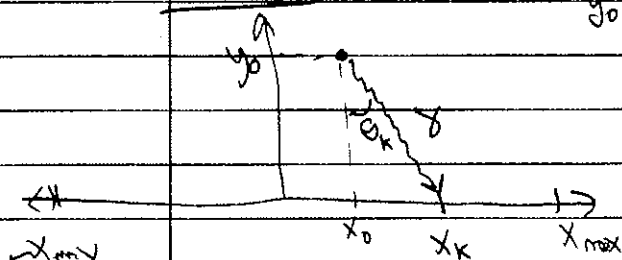
Before class? Set up notebooks exploring pdfs, ipynb, radioactive lighthouse exercise key, ipynb, parameter estimation Gaussian noise, ipynb, parameter estimation fitting straight line I, ipynb.

On board:



$$\cos^2 \theta = \frac{y_0^2}{y_0^2 + (x_k - x_0)^2}$$

$$x_k = x_0 + y_0 \tan \theta_k$$



Recap of radioactive lighthouse problem. \Rightarrow step through key.

0. review set up: geometry, trigonometry in figure.
- 1, 2, 3. \Rightarrow as in key. Think about autoscaled plots & posteriors.
4. Discussion of prior - much to follow.
5. Go over different types of independence.

Extra example: Flip two ^{fair} coins. Propositions: A = 1st is heads

Independent A, B: $p(A, B) = p(A|B)p(B) = p(A)p(B)$ B = 2nd is heads

$p(A) = p(B) = \frac{1}{2}$ and these are independent C = the two results are the same

pairwise independent $\Rightarrow p(A|C) = p(A)$; $p(B|C) = p(B)$; $p(C|A) = p(C)$; $p(C|B) = p(C)$
 \Rightarrow we don't know anything more given the proposition to the right of 1.

BUT A and B are not conditionally independent given C!

If you know A and C are true, then B is determined!

In our case, knowing x_0, y_0 tells us about x_k , but also knowing x_k tells us nothing in addition.

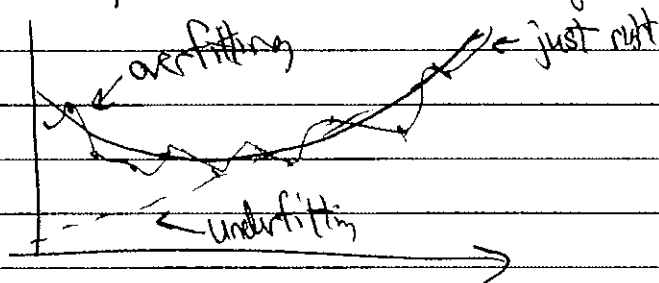
6. emphasize $p(\theta_k) d\theta_k = p(x_k) dx_k$ (suppressing known x_0, y_0)

7, 8. We can see the problem: heavy tails! Note fluctuations - always remember!
 9. Bottom line: Bayes works where naive point estimate fails. Central limit failure! (More to come)

• Go back to (20) and get a moving start.

9/14/19

As a teaser, let's ask: what can go wrong in a fit?



Bayesian methods can prevent/identify both underfitting (model is not complex enough to describe the fit data) or overfitting (model tunes to data fluctuations or terms are underdetermined, leading to them playing off each other.)

⇒ We'll see how this plays out

Let's step through parameter estimation - Gaussian noise in python
 • run using RISE. Include some "footnotes" on Python, Jupyter, etc.

• Import of modules

- Note "cell magic" %matplotlib inline (cf. %matplotlib notebook)
- Using seaborn just to make nice graphs. with interactive figures
- We'll use emcee (cf. MC → Monte Carlo) to do "sampling".
corner is used to make a particular type of plot.

• Example from Sivik's book: Gaussian noise and averages.
 • Excerpts in modules.

class:
dimensions
of
pdf?
(cons. $\frac{1}{\sigma^2}$)

$$p(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \leftarrow \mu, \sigma \text{ are given. Normalized } \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx = 1$$

• Justification as theoretical model by maximum entropy, "central limit theorem" (how many know what CLT is?), or general considerations we'll get to next time.

9/4/19

- M measurements: $D \equiv \{X_k\} = (X_1, \dots, X_M)$ eg. $M=100$
distributed according to $p(x|\mu, \sigma)$.

How do we get μ and σ here? As in Exploring-pdfs.ipynb.

- "Sample" from $N(\mu, \sigma^2) \Rightarrow$ we'll see this.

Goal: Find approximate μ, σ given D

Frequentist: maximum likelihood method

Bayesian: compute posterior pdf $p(\mu, \sigma | D, I)$ ← other information
every time you repeat

- Random seed of 1 means same series of random numbers. If you put 2 or 42, then different from 1, but still the same with every run.
"normal" pdf ← random variates

- stats.norm, rvs as in Exploring-pdfs.ipynb
 - size=M is a "keyword argument" (often kw \equiv keyword)
 \Rightarrow optional and there is a default value (here 1).

- shift-tab-tab after evaluating cell.
eg. place on "norm" or "rvs"

- Output D is a numpy array. ← everything in Python is an object. So more than just a datatype \Rightarrow extra methods.

- Put cursor after D and shift-tab-tab
- $[\dots]$ when printed.

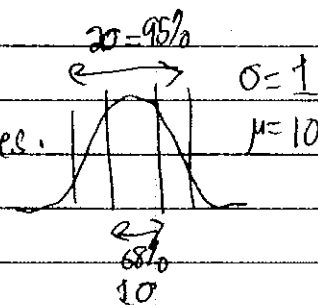
- Discuss about number of entries in tails amongst selves.

- Hint: "tail" of Gaussian, say beyond 20
 $\Rightarrow x > 12$ or $x < 8$.

- How many do you expect on average?

20 \Rightarrow 95% so about 5/100.

- Here 4 in that range. If there were zero is there a bug?
No, there is a chance that will happen.



- Note the pattern (or lack) and repeat to get different numbers. How?
Change the random seed from 1. (You are invited to try.) Always play!

9/4/19

• Questions about plotting?

- We'll repeatedly use constructions like this, so get used to it!
- ; means we put on same line. Not necessary.
- alpha=0.5 just makes the (default) color lighter.
- try color='red' on your own in scatter plot (as in vlines)
- might 'prefer side-by-side' \Rightarrow alternative code.
- An "axis" in Matplotlib means an entire subfigure, not just the x-axis or y-axis.
- If you want to know about a plotting command already there, shift-tab-tab (usually, sometimes not).
- To find vlines (vertical lines), google "matplotlib vertical line" (try it).
- fig.tight_layout() for good spacing with subplots.

• Ask questions in class and via email, etc., if you are confused by code!

• Observations on graphs?

- scatter plot shows tail \Rightarrow in this case there are 5, but rerun and it will be more or less \Rightarrow everything is a pdf
- histogram is imperfect. Problem: of Exploring pdfs at end (sampling 1D pdfs)
- tails fluctuate

• Frequentist approach

- true value for parameters μ, σ , not a pdf
- Use of \mathcal{L} is notation commonly used
- Why the product? Assumed independent. Reasonable?
- log \mathcal{L} for several reasons (note: "log" always means ln. If we want base 10, then log₁₀)

*

eg. to avoid extreme values

$\mathcal{L} \propto (\text{const}) e^{-\chi^2}$ so maximizing \mathcal{L} = maximizing log \mathcal{L} = minimizing χ^2

• You can all carry out the maximization

eg. $\frac{\partial \log \mathcal{L}}{\partial \mu} = -\frac{1}{\sigma^2} \sum_{i=1}^M \frac{x_i - \mu}{\sigma^2} \cdot -1 = \frac{1}{\sigma^2} \sum_{i=1}^M (x_i - \mu) = \frac{1}{\sigma^2} \left(\sum_{i=1}^M x_i - M\mu \right)$

\Rightarrow set to zero $\Rightarrow M\mu_0 = \sum_{i=1}^M x_i$ or $\mu_0 = \frac{1}{M} \sum_{i=1}^M x_i$ You do σ_0^2 (easier than σ_0 is σ_0^2)

9/4/19

(24)

* Do these make sense?

- μ_0 is mean of data \rightarrow estimator for "true mean"
- σ_0 gives spread about μ_0 .

• Note use of `.sum` to add up D array elements

• Printing with f strings

f'...' or f'''...''' \leftarrow multiline string

- `.2f` means float with 2 decimal points

• Note comment on "unbiased estimator"

- an accurate statistic

• Here compare μ_0 from $\frac{1}{M}$ and $\frac{1}{M-1}$

• If you do this many times, you'll find $\frac{1}{M}$ doesn't quite give μ_{true} correctly (take mean of μ_0 's from many trials), but $\frac{1}{M-1}$ does. (Try it!)

- The difference is $O(\frac{1}{M})$, so small for large M.

• Compare estimates to true. Are they good estimates?

How can you tell? E.g. should they be within .1, .01, or what?
 \Rightarrow more as we go!

• Bayesian approach

$p(\mu, \sigma | D, I)$ is posterior: probability (density) of finding some μ, σ given data D and what else we know (I).

"I could be that $\sigma > 0$ or μ should be near zero.

Frequentist probability: long-run frequency of (real or imagined) trials.

\Rightarrow data is probabilistic (repeat experiment and get different result)

but model parameters are not (universe stays the same with more observations)

Bayesian probability: quantification of information (what you know, often said "what you believe"). Data are fixed (it's what you found) but knowledge of true model parameters is fuzzy (and gets update with more trials - coin flipping).

8/4/19

(25)

Bayes' Theorem

class: you label each term

likelihood

$$\text{posterior} \rightarrow P(\mu, \sigma | D, I) = \frac{p(D | \mu, \sigma, I) p(\mu, \sigma | I)}{p(D | I)}$$

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

It will become intuitive!

• tells you how to flip $p(\mu, \sigma | D, I) \leftrightarrow p(D | \mu, \sigma, I)$
hard easy

Aside on denominator

general vector
of parameters

$$p(D | I) = \int p(D | \theta, I) p(\theta) d\theta$$

first step
 $\int p(D | \theta, I) d\theta$
then

so integrate ("marginalize") over all values of θ . Numerically costly \Rightarrow more later on how to do it.
(parameter)

• For model fitting, we don't need $p(D | I)$ calculated. Find the posterior and just normalize that function (or we might only need relative probabilities).

If $p(\mu, \sigma | I) \propto 1 \Rightarrow$ "flat prior" (more later)
then

$$p(\mu, \sigma | D, I) \propto \mathcal{L}(D | \mu, \sigma)$$

then F and B get same answer for most likely values μ_0, σ_0 (called "point estimates" as opposed to a full pdf)

• Back to the prior, \Rightarrow include additional information. What you know before a measurement.

• We will talk much more.

• F says it is nonsense; subjective, individual

\Rightarrow discuss this amongst yourselves.

• How to compute $p(\mu, \sigma | D, I)$ in practice? Often with MCMC. Just look now and we'll discuss later.

9/4/19

(50)

Now turn to parameter estimation fitting straight line - I. ipynb
Fitting a straight line.

Annotations:

- same imports as before
- assume we create data from underlying model

$$y_{\text{exp}}(x) = m_{\text{true}} x + b_{\text{true}} + \text{gaussian noise}$$

$$\theta_{\text{true}} = [b_{\text{true}}, m_{\text{true}}] = [\text{intercept}, \text{slope}]_{\text{true}} \quad \text{mean zero} \quad \text{fixed } \sigma = dy \text{ in code}$$

- x_i points are also chosen randomly according to uniform distribution $\Rightarrow \text{rand.rand}(N)$
- errors are normal $\Rightarrow y += dy * \text{rand.randn}(N)$
- These use the numpy random number generators rand while we will mostly use `scipy.stats` (see other codes)

• Theoretical model: $y_{\text{th}}(x) = mx + b$ with $\theta = [b, m]$

true in a
distributional
sense

$$y_{\text{exp}} = y_{\text{th}} + \underbrace{\delta y_{\text{exp}}}_{\text{normally distributed}} + \underbrace{\delta y_{\text{th}}}_{\text{critically important but has often (mostly?) been neglected, more later.}}$$

$$\Rightarrow y_i \sim N(y_{\text{th}}(x_i; \theta), \sigma^2) \quad \leftarrow \text{should be } \epsilon_i^2$$

mean \uparrow usually squared

• Is independent a good assumption?

• Priors - quick run through

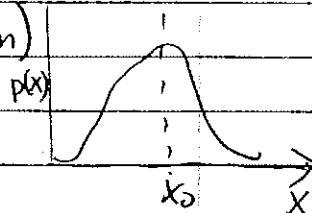
• Major point: when does prior matter?

9/4/19

General reason why Gaussians may show up:

Given $p(x|D, I)$, then our "best estimate" from

$$\frac{dp}{dx}\bigg|_{x_0} = 0 \text{ with } \frac{d^2p}{dx^2}\bigg|_{x_0} < 0 \text{ (maximum)}$$

Look nearby to characterize posterior $p(x)$. $p(x)$ varies too fast, so characterize $\log p$

$$\Rightarrow L(x) = \log p(x|D, I) = L(x_0) + \frac{dL}{dx}\bigg|_{x_0} + \frac{1}{2} \frac{d^2L}{dx^2}\bigg|_{x_0} (x-x_0)^2 + \dots$$

If we can neglect higher order terms, then

$$p(x|D, I) \hat{=} A e^{\frac{1}{2} \frac{d^2L}{dx^2}\bigg|_{x_0} (x-x_0)^2}$$

↖ normalization

 \Rightarrow very generally looks like Gaussian.

$$p(x|D, I) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \Rightarrow \mu = x_0, \sigma = \left(-\frac{d^2L}{dx^2}\bigg|_{x_0}\right)^{-1/2}$$

we usually quote $x = x_0 \pm \sigma$, because if Gaussian, this is sufficient to tell us the entire distribution.

For Bayesian: full posterior $p(x|D, I)$ for $\forall x$ is general result, and $x = x_0 \pm \sigma$ may be an approximate characterization

What if asymmetric $p(x|D, I)$? Multimodal?

9/4/19

95%

(98)

Bayesian vs. Frequentist, confidence interval

- Bayesian is easy: a credible interval or Bayesian confidence interval or degree-of-belief (DOB) interval is: given some data, 95% chance (probability) that the interval contains the true parameter.

• Frequentist 95% confidence interval

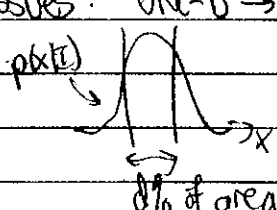
- If large # of repeat samples, 95% of these intervals include the true value of the parameter

- So the parameter is fixed (no pdf) and the confidence interval depends on data, (random sampling)

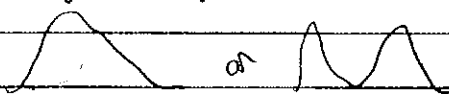
"There is a 95% probability that when I compute a confidence interval from data of this sort that the true value of θ will fall within the (hypothetical) space of observations."

• What?

• One key difference: Bayesian includes prior.

- Issues: One-D \rightarrow if symmetric pdf, then clear how to define ^{df} confidence interval,
  Algorithm: start from center, step outward adding area, stop at df
Two-D: need a way to integrate from top.

• What if asymmetric or multimodal?



Two of the possible choices:

- Equal-tailed interval (central interval): area above and below interval are equal
- Highest posterior density (HPD) region: posterior density for every point is higher than the posterior density for any point outside the interval.