

10/15/19

Physics 8805: Learning From Data: Lecture 15

On board:

- Mini-project 2a is due Friday (really by end of Monday).

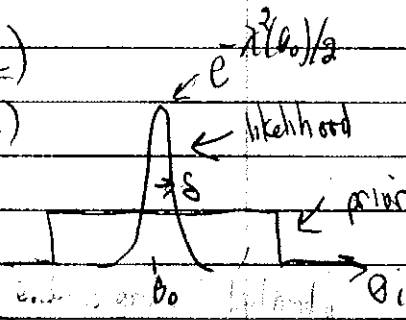
Notebooks for today:

model-selection-mini-project-IIa.ipynb

Evidence for model-EFT-coefficients.ipynb

- Review evidence ratio calculation \Rightarrow top of (78):

$$\frac{p(A|D,I)}{p(B|D,I)} = \frac{p(D|A,I) p(A|I)}{p(D|B,I) p(B|I)}$$



If model A depends on parameters $\vec{\theta}$, then

$$p(D|A,I) = \int p(D|\vec{\theta}, A, I) p(\vec{\theta}|A, I) d\vec{\theta}$$

usual likelihood
prior for these particular parameters
integrate over all parameters

model A will have greater evidence when peak of likelihood increases more than the ratio of the prior data to pre-data volume of the parameter space

- We did an approximate evaluation of $p(A|D,I)/p(B|D,I)$ in a special case. We'll see a similar argument today for the toy model in Mini-project-I.

- Do comparison to parameter estimation on (81)

- Do evaluation with Laplace's method on (81)

- Another look at model-selection-mini-project-IIa.ipynb. (82)

prior $p(\theta_i|I) = \frac{1}{\theta_{\max} - \theta_{\min}}$ if $\theta_{\min} \leq \theta_i \leq \theta_{\max}$, otherwise 0.

$$\int d\vec{\theta} p(D|\vec{\theta}, A, I) p(\vec{\theta}|I) \approx \frac{1}{(\theta_{\max} - \theta_{\min})} e^{-\chi^2(\theta_0)/2} \sqrt{\frac{(2\pi)^N}{\text{det}(Z)}}$$

shrinkage \rightarrow

Does improvement in likelihood out-weigh the penalty from shrinkage of the parameter phase space?

10/15/19

(84)

Problem: In the toy model example, the evidence is highly sensitive to the prior \Rightarrow need criterion to determine prior.

Plan for us: Apply an informative prior from EFT naturalness.

\Rightarrow look at Evidence for model EFT coefficients: i.e. μ

- Illustrate with toy model as in mini-project I.

- von Neumann quote \Rightarrow we will constrain use of higher-order parameters to prevent elephant fitting.

- Look at some EFT slides for additional motivation of expansion (request for more physics).

\Rightarrow EFT

- Look at classical analog in extra slides. $9/1/19 \rightarrow 2/1/19$

- Note on "models" EFT is said to be model independent because it uses the most general form consistent with symmetries of underlying physics. \Rightarrow no extra assumptions.

- "model" is any theoretical construct for computing an observable.

- Model selection in EFT context could be with models having different dots (nucleons only, nucleons plus pions, nucleons + Δ s + pions) or for different orders in the same EFT (cf model problem).
 \Rightarrow first is a frontier, 2nd demonstrated in paper with model
 \Rightarrow do this here.

- Look at punch line in Evidence for model EFT coefficients and then try to explain. Come back to details.

- Class: play with the effects of changing the range of data, the relative error, and the # of points.

10/15/91

Revisit two model discussion

- Models M_1, M_2 with same data set D .

- Evidence: $p(M_1|D, I)$ vs. $p(M_2|D, I)$ no reference to a particular parameter set \Rightarrow comparison between two models, not two fits.

- Note: In Bayesian model selection, only a comparison makes sense. One does not deal with the hypothesis like "Model M_2 is correct."

- Here M_2 is M_1 with one extra order (one more parameter) eventually

- Apply Bayes' Theorem

$$\frac{p(M_2|D, I)}{p(M_1|D, I)} = \frac{p(D|M_2, I) p(M_2|I)}{p(D|M_1, I) p(M_1|I)} \quad \text{cancel}$$

Bayes Factor

we'll take $=1$ for our example \Rightarrow no a priori preference for best order

$$(I): \frac{p(D|M_2, I)}{p(D|M_1, I)} = \frac{\int d\vec{a}_2 p(D|\vec{a}_2, M_2, I) p(\vec{a}_2|M_2, I)}{\int d\vec{a}_1 p(D|\vec{a}_1, M_1, I) p(\vec{a}_1|M_1, I)} \quad \leftarrow \begin{array}{l} \text{recall steps} \\ p(D|M_2, I) \\ \rightarrow p(D, \vec{a}_2|M_2, I) \\ \rightarrow p(D|M_2, \vec{a}_2, I) p(\vec{a}_2|M_2) \end{array}$$

so, integration over the entire parameter space.

\Rightarrow Difficult numerically since likelihoods usually peaked but can have long tails that contribute to integrals (cf. Averaging over likelihood vs. finding peak)

Easiest example $\left. \begin{array}{l} M_1 \rightarrow M_k \\ M_2 \rightarrow M_{k+1} \end{array} \right\} \begin{array}{l} \text{Is going to a higher-order} \\ \text{favored by the given data?} \\ \text{order in EFT expansion} \end{array}$

Simplify: M_{k+1} has an additional parameter a' and assume powers factorize

eg. $e^{-a^2/2a^2} \rightarrow e^{-a^2/2a^2} e^{-a'^2/2a'^2} \dots e^{-a'^2/2a'^2}$ for Gaussian

10/15/19

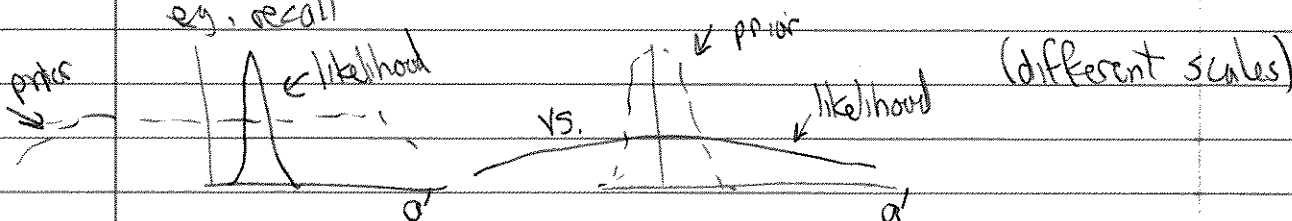
(86)

Then $p(\vec{a}_2 | M_{k+1}, I) = p(\vec{a}_2, a' | M_{k+1}, I) = p(\vec{a}_2 | M_{k+1}, I) p(a' | M_{k+1}, I)$

Consider cases...

i) values of a' that contribute to integrand in numerator of (I) are determined by the likelihood peaked region.

eg. recall



How can we approximate? $p(a' | M_{k+1}, I)$ $\left[\text{rectangle of width } \Delta a' \text{ and height } \frac{1}{\Delta a'} \right]$

Call the value of the likelihood peak \hat{a} and the width $\Delta a'$.

So two different widths: before data $\Delta a'$ (prior) and after $\Delta a'$ (likelihood)

$$\Rightarrow \frac{p(D | M_{k+1}, I)}{p(D | M_k)} = \frac{\Delta a'}{\Delta a'} \underbrace{\int_{\Delta a'} d\vec{a}_2}_{\leftarrow \text{from integral over } a'} p(D | \vec{a}_2, \hat{a}', M_{k+1}, I) \times \underbrace{p(\vec{a}_2 | M_{k+1}, I)}_{\leftarrow \text{peak value}} \underbrace{p(a' | M_{k+1}, I)}_{\rightarrow \text{posterior}}$$

* \Rightarrow The ratio of the integrals is the gain in the likelihood from an extra parameter with value \hat{a} (cf. $M_{k+1}(\hat{a} \neq 0) = M_k$)

• But also "Occam factor" or "Occam penalty" $\frac{\Delta a'}{\Delta a'}$

\Rightarrow How much parameter space collapses in face of data, we thought initially that a' could be anywhere in $\Delta a'$, but find after data it is only in $\Delta a'$. What a waste (less predictive) if $\Delta a' \ll \Delta a'$.

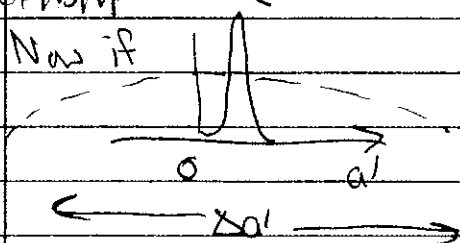
• These factors play off each other: if we add a parameter to a nested model, we expect to gain because \hat{a}' is more information (it could be $a' = 0$ instead)

10/15/19

$\Delta a'$

(87) 6

Now if



Then $a'=0$ likelihood is $\ll a'=\hat{a}'$ likelihood
 \Rightarrow evidence ratio $\gg 1$ and inclusion of this parameter is highly favored.

Unless you put flat prior from near $-\infty$ to near $+\infty$.
 But we have a naturalness prior, so $\Delta a'$ restricted.

Now suppose

$$\frac{p(D|M_{k+1}, t)}{p(D|M_k, t)} \approx \frac{\int da' \int d\tilde{a}' p(D|\tilde{a}', a', M_{k+1}, t) p(\tilde{a}'|M_{k+1}, t)}{\int da' \int d\tilde{a}' p(D|\tilde{a}', a', M_k, t) p(\tilde{a}'|M_k, t)}$$

Turn analysis on its head
 replace dependence on a' because weak
 normalization

normalization integral, dominated by prior so $\hat{a}' \approx 0$ can be used

- But M_{k+1} with $\hat{a}'=0$ is $M_k \Rightarrow$ Bayes ratio $\rightarrow 1$ (not decrease)
- Same argument for $k+1 \rightarrow k+2 \rightarrow \dots$
- \Rightarrow we have saturation of a_k 's

Summary: naturalness prior cuts down on wasted space in the parameter phase space that might be ruled out by data.

Thus EFT is a simpler model (in the model selection sense) than the same functional form with unconstrained or only weakly constrained LECs.

Seen in Fig. 8 \Rightarrow see notebook.

Ratio is about 5, so quadratic is moderately more favorable (compare logs)
 \Rightarrow returning to prior

Predict: Based on your experience, how does this behavior change if we have more data (higher energy) or more certain data? So depends on data.

(88)

10/15/19

Return to notebook to look at calculation of evidence with linear algebra.

• Integrals to calculate are Gaussians in multiple variables!
 $\vec{a} = (a_0, \dots, a_k)$ plus \hat{a} .

• We can write them with matrices.

E.g. see Lecture 10 and 11 notes,

$$\chi^2 = (Y - A\theta)^T \Sigma^{-1} (Y - A\theta)$$

← k →

where $Y = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}$ Data

$A = \begin{bmatrix} 1 & x_1 & x_1^2 & x_1^k \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_N & x_N^2 & x_N^k \end{bmatrix}$ N

$\theta = \begin{bmatrix} a_0 \\ a_1 \\ \vdots \\ a_k \end{bmatrix}$

$\Sigma = \begin{pmatrix} \sigma_1^2 & \rho_{12} \sigma_1 \sigma_2 & \dots \\ \vdots & \ddots & \vdots \\ \vdots & \vdots & \sigma_N^2 \end{pmatrix}$ $N \times N$

N data points
 we've taken these ≥ 0 mostly

$$\chi^2_{\text{MLE}} \text{ when } \hat{\theta} = (A^T \Sigma^{-1} A)^{-1} (A^T \Sigma^{-1} Y) \quad (\text{from lecture 11 notes})$$

Here we have a couple of options:

i) Use $\int e^{-\frac{1}{2} x^T A x + b^T x} dx = \sqrt{\det(A^{-1})} e^{\frac{1}{2} b^T A^{-1} b}$ (different A here, sorry!)

complete the square
 generic square matrix A and vector b

ii) Use conjugacy. See "conjugate prior" entry in Wikipedia.

$$p(\theta|D) = \frac{p(D|\theta) p(\theta)}{p(D)} \quad \text{if } p(\theta) \text{ Gaussian and } p(D|\theta) \text{ Gaussian, SD is } p(\theta|D)$$

(μ_0, σ_0) (μ, σ)

$$\hat{\mu} = \frac{1}{\frac{1}{\sigma_0^2} + \frac{N}{\sigma^2}} \left(\frac{\mu_0}{\sigma_0^2} + \frac{\sum_{i=1}^N y_i}{\sigma^2} \right) \quad \hat{\sigma}^2 = \frac{1}{\frac{1}{\sigma_0^2} + \frac{N}{\sigma^2} + 1} \quad (\text{check } N \rightarrow \infty, \mu_0 \text{ and } \sigma_0 \text{ are irrelevant})$$

→ Generalized formula for multivariate Gaussians