

10/9/19

Physics 880S: Learning from Data: Lecture 14On board: Winter break this week \rightarrow no class Friday.

- Mini-project 2a is released. There will be a mini-project 2b.
- Both are Model Selection. 2a is due Friday, Oct. 18.

- Notebooks for today:

- model-selection-mini-project-IIa.ipynb
- Evidence for model EFT coefficients.ipynb

- Comments on Assignment 2 and mini-project I

- Review Mini-project I via a key notebook.
 - Using pandas to read data into a dataframe
 - Starting guesses histogram
 - What do traces tell us? Should check Gelman-Rubin ratio.
 - Explain uniform correlations.
 - Explain Gaussian correlations (not bad that they exist!)
 - Table III features
 - Figure 3/4 and bands \Rightarrow see how xfit, yfit then mean and std calculated
 - notice failure to extrapolate with good error.

emcee:
offline-instant
sampler can
handle.

10/9/19

Bayesian Model Selection

The discussion here is based heavily on Sivik, Chapter 4.

We've mostly focused on parameter estimation: given a model with parameters, what is the joint posterior for those parameters given some data - this is what we Bayesians mean by fitting the parameters: $p(\theta | D, T)$.

Now we turn to an analysis of the model itself, or, more precisely, the comparison of models.

Remember that: model selection will always be about comparisons.

We can think of many possibilities: (all given data)

• Is the signal a Gaussian or Lorentzian line shape?

• When fitting to data, what order polynomial is best?

• Given two types of Hamiltonian, which is better?

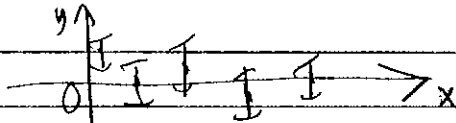
• In an EFT, which power counting (organization of Feynman diagrams) is favored?

If you decided by how well polynomials fit the data, by having the least residual, then higher order will always be better (or at least equal, as you can set coefficients to zero).

So let's think how a Bayesian would proceed. As is usually best, we'll start with the simplest possible example, dating back to Jeffreys (1939) \rightarrow Gull (1988) \rightarrow Sivik (2006) \rightarrow 8805 (2019):

The story of Dr. A and Prof. B:

"Dr. A has a theory with no adjustable parameters. Prof. B also has a theory, but with an adjustable parameter λ , whose theory should we prefer on the basis of data D ?"

E.g. could be data D : 

Dr. A thinks $y=0$, Prof. B thinks $y=\lambda$, with λ to be determined.

(70)

10/9/19

Bayesian: consider ratio of posteriors

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

again, not prob. for particular instance of B , like $x=0.2$, but prob. that the theory is correct. x doesn't appear yet.

Ok, let's proceed with Bayes' Theorem:

$$\frac{p(A|D, I)}{p(B|D, I)} = \frac{\overset{\text{likelihood}}{p(D|A, I)} \overset{\text{prior}}{p(A|I)} / \cancel{p(D|I)}}{\overset{\text{likelihood}}{p(D|B, I)} \overset{\text{prior}}{p(B|I)} / \cancel{p(D|I)}} \quad \text{so denominator cancels in the ratio.}$$

- The ratio $\frac{p(A|I)}{p(B|I)}$ might be given by our opinion of the two scientists based on their track records. But it is more typically taken to be 1 (before seeing data, no preference).

Now $p(D|A, I)$ would seem to be straightforward, but what do we do about the x in B ? Marginalize!

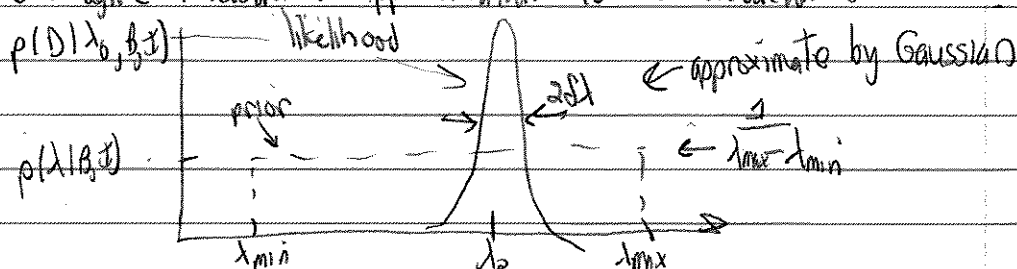
$$p(D|B, I) = \int p(D, x|B, I) dx = \int \underset{\text{ordinary likelihood}}{p(D|x, B, I)} \overset{\text{prior for } x}{p(x|B, I)} dx$$

integrates over all x

Let's suppose a uniform prior for $p(x|B, I)$:

$$p(x|B, I) = \frac{1}{x_{\max} - x_{\min}} \text{ for } x_{\min} \leq x \leq x_{\max} \text{ and zero otherwise.}$$

Suppose also we can find x_0 that maximizes the likelihood $p(D|x, B, I)$. So we imagine a reasonable approximation to the situation is



10/9/19
Plan

value at $\lambda = \lambda_0$
↓

$$p(D|\lambda, B, I) \approx p(D|\lambda_0, B, I) e^{-\frac{(\lambda - \lambda_0)^2}{2S\lambda^2}}$$

Note that $p(\lambda|B, I)$ is normalized wrt λ , but $p(D|\lambda, B, I)$ is not $\rightarrow p(D|\lambda_0, B, I)$ is not equal to $\frac{1}{\text{range}}$ (in general).
[because λ is to the right of the 1]

Now the prior doesn't depend on λ , so pull it out of the integral

$$p(D|B, I) = \frac{1}{\lambda_{\max} - \lambda_{\min}} \int_{\lambda_{\min}}^{\lambda_{\max}} d\lambda p(D|\lambda, B, I)$$

take $\lambda_{\max} \rightarrow \infty$
and $\lambda_{\min} \rightarrow -\infty$
with negligible error

$$= \frac{1}{\lambda_{\max} - \lambda_{\min}} p(D|\lambda_0, B, I) \cdot S\lambda\sqrt{2\pi}$$

← integral over the Gaussian

Now put this together

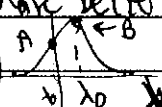
$$\frac{p(A|D, I)}{p(B|D, I)} = \frac{p(A|I)}{p(B|I)} \times \overset{i)}{\left(\frac{p(D|A, I)}{p(D|\lambda_0, B, I)} \right)} \times \overset{ii)}{\left(\frac{\lambda_{\max} - \lambda_{\min}}{S\lambda\sqrt{2\pi}} \right)}$$

a priori ratio (usually unity) ratio of likelihoods "Occam factor" (Occam's razor)

\Rightarrow competing factors

i) B with λ_0 will usually have better likelihood than A.

(always if nested theory)



ii) Another term! Penalize B for additional parameter since $\lambda_{\max} - \lambda_{\min} > S\lambda\sqrt{2\pi}$ (usually)

\Rightarrow This is a formalization of Occam's razor: add parameters as long as the gain in likelihood beats the cost in complexity.

• Occam factor is the ratio of volumes before and after data is known.

A greater collapse means a greater penalty - cf. many parameters: full prior volume vs. effective likelihood volume.

• Doesn't always get interpreted so easily: think Gaussian vs. Lorentzian.

(80)

10/9/19

Jeffreys first worried about the prior: if the new parameter has infinite range, isn't the Occam factor an infinite penalty? Indeed an issue if there are no constraints.

A more reasonable case is the naturalness prior in mini-project-I.

In most physical cases, there are limits prescribed by some a priori knowledge.

What happens if likelihood is wider than prior?

caution: recalculate the integral rather than using our formula

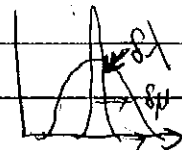
⇒ in model-selection mini-project-IIa you will take an exploration of evidence for a simple case, and test intuition.

Suppose Dr. A now also has an adjustable parameter μ .

$$\Rightarrow \frac{p(A|D, I)}{p(B|D, I)} = \frac{p(A|I)}{p(B|I)} \frac{p(D|\mu_0, A, I)}{p(D|\lambda_0, B, I)} \times \frac{\int_{\mu_{\min}}^{\mu_{\max}} d\mu}{\int_{\lambda_{\min}}^{\lambda_{\max}} d\lambda}$$

Suppose same prior range and equal a priori theories

$$\Rightarrow \frac{p(A|D, I)}{p(B|D, I)} \approx \frac{p(D|\mu_0, A, I)}{p(D|\lambda_0, B, I)} \times \frac{\int_{\mu_{\min}}^{\mu_{\max}} d\mu}{\int_{\lambda_{\min}}^{\lambda_{\max}} d\lambda}$$



If data is good, likelihood ratio is most likely to dominate.
If comparable, then shape with larger error bar for its parameter will be favored.
Why? Because in model selection, this means more parameter values are consistent with a good fit.

Finally: same theory but different prior range. Then $\frac{p(A|D, I)}{p(B|D, I)} \approx \frac{\mu_{\max} - \mu_{\min}}{\lambda_{\max} - \lambda_{\min}}$

⇒ prefer narrower prior range ⇒ must have more insight to give a narrower range for the parameter.

(81)

10/9/19

Compare to parameter estimation...

$$p(\lambda | D, B, \pm) = \frac{p(D | \lambda, B, \pm) p(\lambda | B, \pm)}{p(D | B, \pm)}$$

↑
This is what we've been calculating
for model selection.

Called evidence for B or marginal likelihood ← (marginal \Rightarrow integrate over all parameters in likelihood)
or global likelihood or prior predictive.

Parameter estimation focuses on the maximum of the likelihood
while model selection calculates an average of it.

Evidence calculations: Laplace's method

Suppose an unnormalized probability density $P^*(\theta)$ has a peak at θ_0 .
Let θ be K dimensional. Then the evidence is

$$Z_p = \int P^*(\theta) d^K \theta$$

If it is ok to expand $\log P^*(\theta)$ around its peak:

$$\log P^*(\theta) \approx \log P^*(\theta_0) - \frac{1}{2}(\theta - \theta_0)^T \Sigma^{-1}(\theta - \theta_0) + \dots$$

Saddle point approx.

where $\Sigma^{-1} = H$ is the Hessian matrix $H_{ij} = -\frac{\partial^2}{\partial \theta_i \partial \theta_j} \log P^*(\theta) \big|_{\theta=\theta_0}$

$$\Rightarrow \text{approximate } P^*(\theta) \approx P^*(\theta_0) e^{-\frac{1}{2}(\theta - \theta_0)^T \Sigma^{-1}(\theta - \theta_0)}$$

$$\text{and approximate } Z_p \hat{=} Z_0 = P^*(\theta_0) \sqrt{\frac{(2\pi)^K}{\det(\Sigma^{-1})}}$$

$$\text{if } P^*(\theta) = e^{-\frac{1}{2}\chi^2(\theta)} \Rightarrow Z_p \hat{=} e^{-\frac{1}{2}\chi^2(\theta_0)} \sqrt{\frac{(2\pi)^K}{\det(\Sigma^{-1})}}$$

10/9/19

(8)

Quick look at model-selection-mini-project-IIa.ipynb

- We have a non-polynomial function
 - Given data with errors, what polynomial is "best" to use.
 - Higher order polynomial will always reduce sum of residuals.
- Compare χ^2/dof in least-squares fits to use of Bayesian evidence ratios.
 - Using Laplace method here.
- Look at questions.
- How do the results depend on the data? # of data pts, range of data, size of error bars.

Quick look at Evidence for model EFT coefficients.ipynb

- Evidence with linear algebra and Gaussian integrals.
- Saturation of evidence, How to explain? (next time!)