

Priors and approximations

ISTA 410 / INFO 510 - Bayesian Modeling and Inference

University of Arizona School of Information

August 30, 2021

Outline

Last time:

- One parameter models:
 - binomial model with a conjugate prior

coinflip
vaccine analysis

Today:

- Exploring choice of priors
- Simple approximate methods for inference
- Using SciPy to do some calculations
- Simple summaries of the posterior

Binomial model

Binomial model:

- Observed outcomes fall into one of two categories
- Model outcomes as coming from identical independent trials, with fixed probability of “success”

$$y \sim \text{Binomial}(n, \theta) \quad \left. \begin{array}{l} \text{specifies} \\ p(y | \theta) \\ = \binom{n}{y} \theta^y (1-\theta)^{n-y} \end{array} \right\}$$

θ : success probability, our unknown parameter

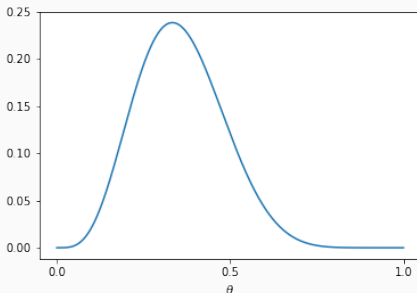
Return to the coin flip example, suppose we got 4 heads out of 12 flips

Binomial likelihood as a function of θ

obs. # heads - not variable here!

The likelihood $p(y|\theta)$ is thought of as a function of θ :

probability of heads



Bayes' theorem

Recall Bayes' theorem for densities:

$$p(\theta|y) \propto p(\theta)p(y|\theta)$$

"proportional to"

$$p(\theta|y) = p(\theta)p(y|\theta) / \underbrace{p(y)}_{\text{normalizing const.}}$$

or, in words:

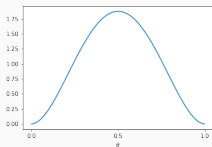
posterior \propto prior \times likelihood

un-normalized posterior = prior \times likelihood

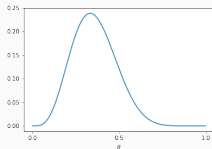
Bayes' theorem, graphically

$\text{Beta}(3, 3)$

Prior



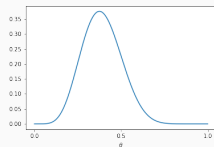
Likelihood



\times

\propto

Posterior



- Multiplicative interaction of prior and posterior
- Posterior ends up as a combination of what you knew before (prior) and what the data told you (likelihood)

Last time: using the conjugate prior

Last time we exploited the conjugate prior, a beta distribution:

$$y \sim \text{Binomial}(n, \theta)$$

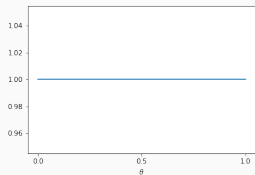
$$\theta \sim \text{Beta}(3, 3)$$

Then the posterior distribution is

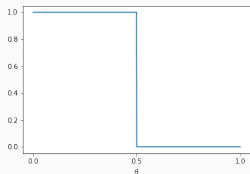
$$\theta|y \sim \text{Beta}(3 + \text{heads count}, 3 + \text{tails count})$$

$$\text{Beta}(7, 11)$$

What if we used a different prior?

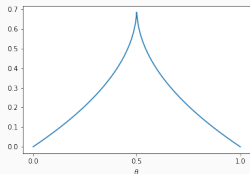


Uniform prior - encode as little info as possible - every possible θ is equally plausible



Step/cutoff prior

- Rule out $\theta > 0.5$ from the start
- this is not advisable

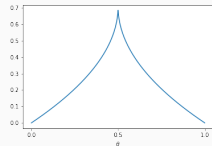
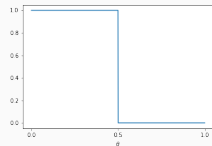
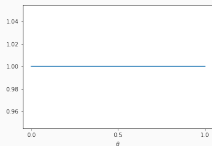


Spike prior

- coin is more likely to be nearly fair
- cusp causes some weirdness in the posterior

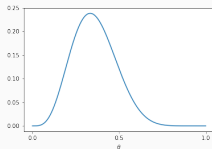
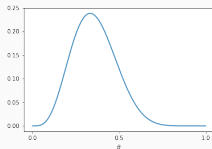
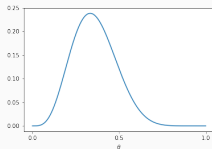
Bayes' theorem, graphically

Prior



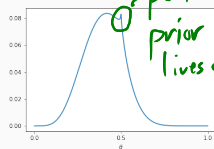
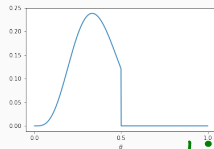
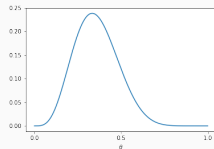
\times

Likelihood



\propto

Posterior



The action at the heart of inference in the model is *conditioning* the prior on the data.

- In principle, easy application of Bayes' theorem
- In practice, not many models can be formally conditioned; and often
- To avoid being forced into less-than-ideal choices of model or prior
 - Grid approximation
 - Quadratic approximation
 - Monte Carlo sampling

Simplest: grid approximation

1. Set in advance a finite grid of parameter values
2. Evaluate the prior and likelihood on the grid
3. Multiply the prior and likelihood vectors
component-by-component
4. Normalize the result

Grid approximation code

Some Python code:

```
grid_size = 1000
theta = np.linspace(0, 1, grid_size)
prior = sp.stats.beta.pdf(theta, 3, 3)
likelihood = sp.stats.binom.pmf(4, 12, theta)
unnorm_post = prior * likelihood
```

Quadratic approximations

Grid approximations are nicely simple, but they scale badly

- Number of grid points grows exponentially with the number of parameters
- These days, models with hundreds or thousands of parameters are common

In practice, we will usually want to use something that scales better. (Grid is nice for goofy priors, though!)

Quadratic approximations

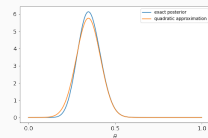
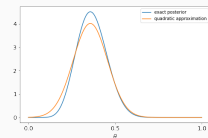
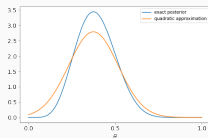
The quadratic approximation (aka normal approximation, Laplace approximation) method:

- often the posterior will closely resemble a Gaussian ("normal") distribution near its peak
- the Gaussian can be described by just two numbers: the center of the peak (mean) and spread (variance)

Quadratic because the log of the Gaussian density is a quadratic function

1. find the posterior mode (maximum density) by some kind of optimization procedure
2. estimate the curvature near the peak

Quadratic approximations



- Three plots show increasing amounts of data (12, 24, or 48 flips)
- Same sample proportion of $1/3$ heads for each
- More data \rightarrow better quadratic approximation

Let's make some posterior inferences:

- posterior mean?
- 89% interval?
- If using a conjugate prior, usually can calculate these directly
- If using the quadratic approximation, read these off from normal statistics
- If using a grid approximation, use quantiles

Normal model, known variance

We all know the normal distribution:

$$p(y|\mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(\frac{-1}{2} \frac{(y - \mu)^2}{\sigma^2}\right)$$

For now, to keep this a one-parameter model, we'll treat σ as a known constant. (We could also fix μ and use σ as the unknown parameter, in principle)

The normal model, known variance

We'll start with a normal model, and as an example case we'll use a data set for basketball scores: final scores y_i from all NCAA men's tournament games from about 1939-1995. We're interested in inferring what an “average” total score is in a game.

- Often normal models get used out of convenience or out of tradition
- When justified, usually justified by the central limit theorem: sum or average of many IID components gives rise to normal distribution

A visual inspection of the data distribution shows a normal distribution really does fit here, but it's reasonably well justified from first principles

The normal model, known variance

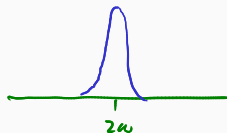
As usual, our starting point is specifying a model and priors for our parameters:

$$\begin{aligned} y_i &\sim \text{Normal}(\theta, \sigma) \\ \theta &\sim \text{Normal}(\mu_0, \tau_0) \end{aligned}$$

likelihood
prior on θ - mean total score of a game

Take $\sigma = 24$. Here, we are choosing a normal prior for convenience (it's conjugate to the normal likelihood)

How can we choose values for μ_0, τ_0 ?



$$\begin{aligned} \mu_0 &= 200 \\ \tau_0 &= 30 \end{aligned}$$

Checking our prior: prior predictive simulations

Prior predictive simulations: draw observations (i.e. values of y_i) using the prior distribution

This can be used to check the reasonableness of a prior, by making sure it doesn't produce impossible results.

- We're not looking for the prior predictions to be a perfect model for the data
- But, if our predictive draws have games with negative score, or teams scoring 500 points, maybe something is off

Calculating the posterior

Assume we start with one observation y . Since we are using a conjugate prior, the posterior is analytically expressible:

$$\begin{aligned} p(\theta|y) &\propto \exp\left(-\frac{1}{2}\frac{(y-\theta)^2}{\sigma^2}\right) \exp\left(-\frac{1}{2}\frac{(\theta-\mu_0)^2}{\tau_0^2}\right) \\ &= \exp\left(-\frac{1}{2}\left(\frac{(y-\theta)^2}{\sigma^2} + \frac{(\theta-\mu_0)^2}{\tau_0^2}\right)\right) \\ &= \exp\left(-\frac{1}{2}\left(\frac{1}{\sigma^2} + \frac{1}{\tau_0^2}\right)\theta^2 - \left(\frac{2y}{\sigma^2} + \frac{2\mu_0}{\tau_0^2}\right)\theta + \frac{y^2}{\sigma^2} + \frac{\mu_0^2}{\tau_0^2}\right) \end{aligned}$$

Then some magic happens...

Calculating the posterior

$$\theta|y \sim \text{Normal}(\mu_1, \tau_1)$$

where

$$\mu_1 = \frac{\frac{1}{\sigma^2}y + \frac{1}{\tau_0^2}\mu_0}{\frac{1}{\sigma^2} + \frac{1}{\tau_0^2}}$$

$$\frac{1}{\tau_1} = \frac{1}{\sigma^2} + \frac{1}{\tau_0^2}$$

The inverse variances $1/\sigma^2, 1/\tau^2$ are called the *precisions* of these distributions

(Where's the magic? Complete the square (exercise 2.14(a)) in BDA3)

The posterior as a compromise

Three ways of writing the posterior mean of θ :

$$\mu_1 = \frac{\frac{1}{\sigma^2}y + \frac{1}{\tau_0^2}\mu_0}{\frac{1}{\sigma^2} + \frac{1}{\tau_0^2}}$$

$$\mu_1 = \mu_0 + (y - \mu_0)\frac{\tau_0^2}{\sigma^2 + \tau_0^2}$$

$$\mu_1 = y - (y - \mu_0)\frac{\sigma^2}{\sigma^2 + \tau_0^2}$$

- Weighted average of μ_0 and y
- Prior mean μ_0 adjusted toward the data
- Data “shrunk” toward the prior mean

Generalizing to many observations

We don't have to iterate this process a thousand times to incorporate our thousand games (although the ability to incorporate observations one by one can be considered a feature of the Bayesian approach); the posterior depends on y_1, y_2, \dots only through the sample mean \bar{y} ¹

$$\theta | (y_1, y_2, \dots, y_n) \sim \text{Normal}(\mu_n, \tau_n)$$

$$\mu_n = \frac{\frac{1}{\tau_0^2} \mu_0 + \frac{n}{\sigma^2} \bar{y}}{\frac{1}{\tau_0^2} + \frac{n}{\sigma^2}}$$

$$\frac{1}{\tau_n^2} = \frac{1}{\tau_0^2} + \frac{n}{\sigma^2}$$

¹ \bar{y} is called a *sufficient statistic* in this model

Some issues about uninformative priors:

- uninformative doesn't always mean “flat” / uniform
 - a prior that is flat in one parameterization may be non-flat if you change variables
 - flat priors can be improper
 - flat priors can be practically nonsensical

Weakly informative priors

A compromise between the informative and uninformative priors is so-called “weakly informative” priors, which generally attempt to include enough outside knowledge to ensure that the prior is proper and sensible, but the information in the prior is intentionally weaker than the available outside informations.

- Our basketball prior example: I asked you to come up with rough bounds, but we used them loosely
- Example: in the coin spinning problem, take $\text{Beta}(3, 3)$ in place of uniform or $\text{Beta}(1, 1)$.
- Example (from the book): in estimating the proportion of female births, choose a prior with the probability mass concentrated between, say, 0.4 and 0.6 (e.g. $\text{Normal}(0.5, 0.1)$)

What to do?

We'll return to problems of prior choice frequently, but for now:

- Weakly informative often a good choice
- Flat priors can be problematic
- Avoid assigning probability 0 to anything unless you are really sure
- Prior predictive checks to avoid truly nonsensical values

Summary

Today:

- Normal model, known variance
- Informative vs. uninformative priors

Next week:

- Some multi-parameter models
- More on priors

Proper and improper prior distributions

The prior precision $1/\tau_0^2$ is the weight given to the prior mean in the posterior distribution; if τ_0^2 is very large, then this weight is very small and the posterior is dominated by the data.

Specifically, if $n/\sigma^2 \gg 1/\tau_0^2$, then the posterior distribution is approximately

$$p(\theta|y) \sim \text{Normal}(\bar{y}, \sigma/\sqrt{n})$$

and so the prior has essentially no input in the posterior.

To eliminate the influence of the prior, we could assign a completely flat/uniform prior on θ . Problem: has an infinite integral, so it's not really a probability distribution.

Improper prior distributions can produce proper posteriors

This example shows that even with an improper uniform prior on θ , the posterior distribution is proper – i.e. $p(\theta|y)$ has a finite integral for any possible data y (as long as there is at least one observation).

- This must be checked any time you use an improper prior
- Most reasonable interpretation of the posterior: as an approximation, valid when the likelihood dominates the prior density
- This is generally dependent on both sufficient amount of data and sufficiently localized likelihood

Posterior predictive distribution

One feature of the posterior we haven't met yet: the *posterior predictive distribution*

- probability distribution of future observations
- will the next case observed in the vaccine trial come from the vaccine or placebo group?
- what will the combined score be of the next game?

In some cases, we can calculate the posterior predictive distribution explicitly; but if not, we can sample in stages (like in the dice problem)

Posterior predictions

The posterior predictive distribution is (unsurprisingly) also normal (details in section 2.5 of BDA) with

$$E(y|y_{\text{obs}}) = \mu_n$$

$$\text{var}(y|y_{\text{obs}}) = \sigma^2 + \tau_n^2$$

Intuitively:

- mean prediction is posterior mean of θ
- uncertainty of prediction is the uncertainty in θ (epistemic uncertainty, τ_n^2) plus the uncertainty of individual observations (aleatoric uncertainty, σ^2)

Informative vs. uninformative priors

Most often, priors are categorized as *informative* or *uninformative priors* depending on whether they incorporate outside scientific information

- informative priors: bring in knowledge about the application domain, or results of previous study, as a starting point for estimation and inference
- uninformative priors: avoid using external knowledge, “let the data speak for itself”