

Introduction to Bayesian Inference

Felipe José Bravo Márquez

September 22, 2021

Some Critics to the Frequentist Approach

- The statistical methods that we have discussed so far are known as frequentist (or classical) methods.
- The frequentist approach requires that all probabilities be defined by connection to the frequencies of events in very large samples.
- This leads to frequentist uncertainty being premised on imaginary resampling of data.
- If we were to repeat the measurement many many times, we would end up collecting a list of values that will have some pattern to it.
- It means also that parameters and models cannot have probability distributions, only measurements can.
- The distribution of these measurements is called a sampling distribution.
- This resampling is never done, and in general it doesn't even make sense.

Bayesian Inference

There is another approach to inference called Bayesian inference [Wasserman, 2013], which is based on the following postulates:

- Probability describes **degree of belief**, not limiting frequency.
 - We can make probability statements about lots of things, not just data which are subject to random variation.
 - For example, I might say that "the probability that Albert Einstein drank a cup of tea on August 1, 1948" is .35.
 - This does not refer to any limiting frequency.
 - It reflects my strength of belief that the proposition is true.
- We can make probability statements about parameters, even though they are fixed constants.
- We make inferences about a parameter θ by producing a probability distribution for θ . Inferences, such as point estimates and interval estimates, may then be extracted from this distribution.

Statistical Rethinking

- Most of the material on Bayesian Inference in this course is based on the book “Statistical Rethinking” by Richard McElreath [McElreath, 2020].



Bayesian Inference

- In modest terms, Bayesian data analysis is no more than counting the numbers of ways the data could happen, according to our assumptions [McElreath, 2020].
- In Bayesian analysis all alternative sequences of events that could have generated our data are evaluated.
- As we learn about what did happen, some of these alternative sequences are pruned.
- In the end, what remains is only what is logically consistent with our knowledge [McElreath, 2020].
- Warning: understanding the essence of Bayesian inference can be hard.
- The following toy example taken from [McElreath, 2020] tries to explain it in a gentle way.

Counting Possibilities

- Suppose there's a bag, and it contains **four** marbles.
- These marbles come in two colors: **blue** and **white**.
- We know there are four marbles in the bag, but we don't know how many are of each color.
- We do know that there are five possibilities:
(1) [○○○○], (2) [●○○○], (3) [●●○○], (4) [●●●○], (5) [●●●●]
- These are the only possibilities consistent with what we know about the contents of the bag. Call these five possibilities the **conjectures**.
- Our goal is to figure out which of these conjectures is most **plausible**, given some **evidence** about the contents of the bag.
- Evidence: A sequence of three marbles is pulled from the bag, one at a time, replacing the marble each time and shaking the bag, in that order.
- The sequence that emerges is: ● ○ ●, which is our **data**.

Counting Possibilities

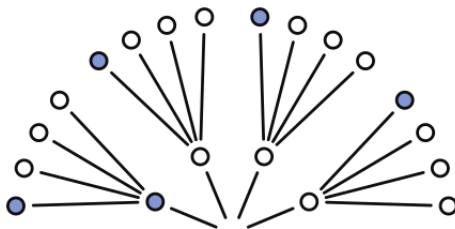
- Now, let's see how to use the data to infer what's in the bag.
- Let's begin by considering just the single conjecture, $[\bullet \circ \circ \circ]$, that the bag contains one blue and three white marbles.
- On the first draw from the bag, one of four things could happen, corresponding to one of four marbles in the bag.



- Notice that even though the three white marbles look the same from a data perspective we just record the color of the marbles, after all they are really different events.
- This is important, because it means that there are three more ways to see \circ than to see \bullet .

Counting Possibilities

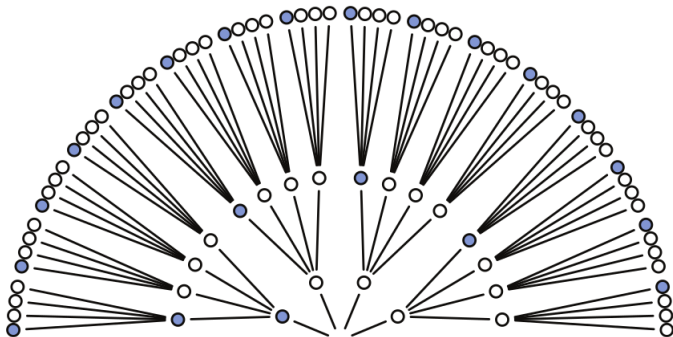
- Now consider the garden as we get another draw from the bag. It expands the garden out one layer:



- Now there are 16 possible paths through the garden, one for each pair of draws.






Counting Possibilities

- The third layer is built in the same way, and the full garden is shown below:



- There are $4^3 = 64$ possible paths in total.

Counting Possibilities

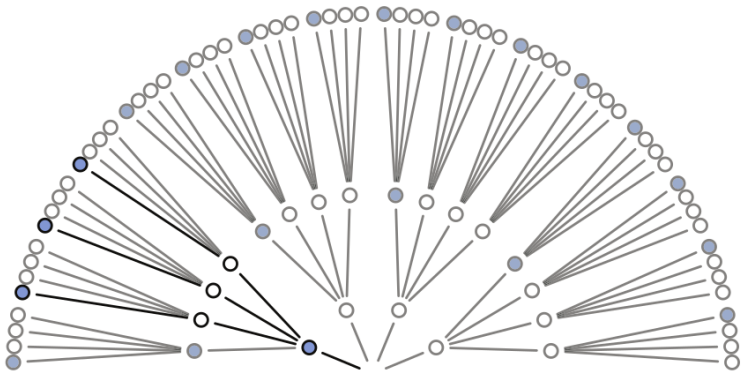
- As we consider each draw from the bag to get   , some of these paths are logically eliminated.
- The first draw turned out to be , recall, so the three white paths at the bottom are eliminated right away.
- If you imagine the real data tracing out a path, it must have passed through the one blue path near the origin.
- The second draw from the bag produces , so three of the paths forking out of the first blue marble remain.

Counting Possibilities

- As the data trace out a path, we know it must have passed through one of those three white paths (after the first blue path).
- But we don't know which one, because we recorded only the color of each marble.
- Finally, the third draw is ●.
- Each of the remaining three paths in the middle layer sustain one blue path, leaving a total of three ways for the sequence ●○● to appear, assuming the bag contains [●○○○].

Counting Possibilities

- The figure below shows the forking paths again, now with logically eliminated paths grayed out.



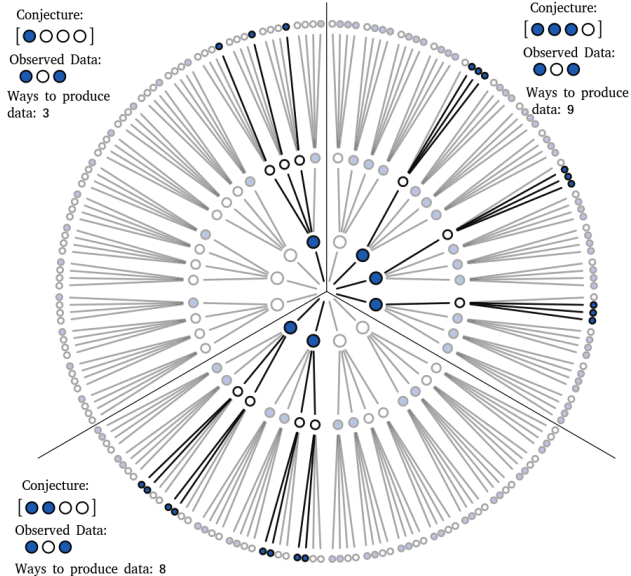
Counting Possibilities

- We can't be sure which of those three paths the actual data took.
- But as long as we're considering only the possibility that the bag contains one blue and three white marbles, we can be sure that the data took one of those three paths.
- Those are the only paths consistent with both our knowledge of the bag's contents (four marbles, white or blue) and the data (●○○●).
- This demonstrates that there are three (out of 64) ways for a bag containing [●○○○] to produce the data.
- We have no way to decide among these three ways.

Counting Possibilities

- The inferential power comes from comparing this count to the numbers of ways each of the other conjectures of the bag's contents could produce the same data.
- For example, consider the conjecture $[○○○○]$.
- There are zero ways for this conjecture to produce the observed data, because even one ● is logically incompatible with it.
- The conjecture $[●●●●]$ is likewise logically incompatible with the data.
- So we can eliminate these two conjectures, because neither provides even a single path that is consistent with the data.
- The next slide's figure displays all the paths for the remaining three conjectures:
 $[●○○○]$, $[●●○○]$, and $[●●●○]$.

Counting Possibilities



Counting Possibilities

- The number of ways to produce the data, for each conjecture, can be computed by first counting the number of paths in each “ring” of the garden and then by multiplying these counts together.

Conjecture	Ways to produce ●○○●
[○○○○]	$0 \times 4 \times 0 = 0$
[●○○○]	$1 \times 3 \times 1 = 3$
[●●○○]	$2 \times 2 \times 2 = 8$
[●●●○]	$3 \times 1 \times 3 = 9$
[●●●●]	$4 \times 0 \times 4 = 0$

- By comparing these counts, we have part a way to rate the relative **plausibility** of each conjectured bag composition.

Combining other information

- We may have additional information about the relative plausibility of each conjecture.
- This information could arise from knowledge of how the contents of the bag were generated.
- It could also arise from previous data.
- Whatever the source, it would help to have a way to combine different sources of information to update the plausibilities.
- Luckily there is a natural solution: Just multiply the counts.

Combining other information

- Suppose that each conjecture is equally plausible at the start.
- So we just compare the counts of ways in which each conjecture is compatible with the observed data: $\bullet \circ \bullet$.
- This comparison suggests that $[\bullet \bullet \bullet \circ]$ is slightly more plausible than $[\bullet \bullet \circ \circ]$, and both are about three times more plausible than $[\bullet \circ \circ \circ]$.
- Since these are our initial counts, and we are going to update them next, let's label them **prior**.
- Now suppose we draw another marble from the bag to get another observation: \bullet .
- How can we update our plausibilities about each conjecture based on this new evidence?
- There are two choices discussed next.

Combining other information

- Option 1: draw a forking path with four layers and do the counting again.
- Option 2: Update previous counts (0, 3, 8, 9, 0) with the new information by multiplying the new count by the old count.
- Both approach are mathematically identical as long as the new observation is logically independent of the previous observations.

Conjecture	Ways to produce ●	Prior counts	New count
[○○○○]	0	0	$0 \times 0 = 0$
[●○○○]	1	3	$3 \times 1 = 3$
[●●○○]	2	8	$8 \times 2 = 16$
[●●●○]	3	9	$9 \times 3 = 27$
[●●●●]	4	0	$0 \times 4 = 0$

Combining other information

- In the previous example, the prior data and new data are of the same type: marbles drawn from the bag.
- But in general, the prior data and new data can be of different types.
- Suppose for example that someone from the marble factory tells you that blue marbles are rare.
- So for every bag containing $[\bullet\bullet\bullet\circ]$, they made two bags containing $[\bullet\bullet\circ\circ]$ and three bags containing $[\bullet\circ\circ\circ]$.
- They also ensured that every bag contained at least one blue and one white marble.

Combining other information

- We can update our counts again:

Conjecture	Prior count	Factory count	New count
[○○○○]	0	0	$0 \times 0 = 0$
[●○○○]	3	3	$3 \times 3 = 9$
[●●○○]	16	2	$16 \times 2 = 32$
[●●●○]	27	1	$27 \times 1 = 27$
[●●●●]	0	0	$0 \times 0 = 0$

- Now the conjecture [●●○○] is most plausible, but barely better than [●●●○].
- Is there a threshold difference in these counts at which we can safely decide that one of the conjectures is the correct one?
- We will explore this question next.

From counts to probability

- So far, we have defined the updated plausibility of each possible composition of the bag, after seeing the data, as:

$$\begin{aligned} &\text{plausibility of } [\bullet \circ \circ \circ] \text{ after seeing } \bullet \circ \bullet \\ &\quad \propto \\ &\quad \text{ways } [\bullet \circ \circ \circ] \text{ can produce } \bullet \circ \bullet \\ &\quad \times \\ &\quad \text{prior plausibility } [\bullet \circ \circ \circ] \end{aligned}$$

- The problem of representing plausibilities as counts is that these numbers grow very quickly as the amount of data grows.
- It is better to standardize them to turn them into probabilities.

From counts to probability

- Now we will formalize the Bayesian framework using probabilities.
- Let index our conjecture with a parameter θ defined as the fractions of marbles from the bag that are blue:

$\theta = 0 \rightarrow [\text{O O O O}], \theta = 0.25 \rightarrow [\text{● O O O}], \theta = 0.5 \rightarrow [\text{● ● O O}], \theta = 0.75 \rightarrow [\text{● ● ● O}], \theta = 1 \rightarrow [\text{● ● ● ●}].$

- Let's call our data $\text{● O ● } d$.
- We construct probabilities by standardizing the plausibility so that the sum of the plausibilities for all possible conjectures will be one.

$$\text{plausibility of } \theta \text{ after } d = \frac{\text{ways } \theta \text{ can produce } d \times \text{prior plausibility } \theta}{\text{sum of products}} \quad (1)$$

- This is essentially the Bayes theorem:

$$\mathbb{P}(\theta|d) = \frac{\mathbb{P}(d|\theta) \times \mathbb{P}(\theta)}{\mathbb{P}(d)} \quad (2)$$

From counts to probability

- The denominator $\mathbb{P}(d)$ (that standardizes values to sum one) can be expressed by the law of total probabilities as:

$$\mathbb{P}(d) = \sum_{\theta} \mathbb{P}(d|\theta) \times \mathbb{P}(\theta) \quad (3)$$

- Let's consider the prior assumptions that all conjectures are equally plausible at the start, then $\mathbb{P}(\theta)$ is uniformly distributed.

θ	$\mathbb{P}(\theta)$	Ways to Produce Data	$\mathbb{P}(d \theta)$	$\mathbb{P}(\theta d) = \mathbb{P}(d \theta) * \mathbb{P}(\theta) / \mathbb{P}(d)$
0	1/5	0	0/64	$\frac{0/64 * 1/5}{0.0625} = 0$
0.25	1/5	3	3/64	$\frac{3/64 * 1/5}{0.0625} = 0.15$
0.5	1/5	8	8/64	$\frac{8/64 * 1/5}{0.0625} = 0.4$
0.75	1/5	9	9/64	$\frac{9/64 * 1/5}{0.0625} = 0.45$
1	1/5	0	0/64	$\frac{0/64 * 1/5}{0.0625} = 0$

- where $\mathbb{P}(d) = 1/5 * 0/64 + 1/5 * 3/64 + 1/5 * 8/64 + 1/5 * 9/64 + 1/5 * 0/64 = 0.0625$

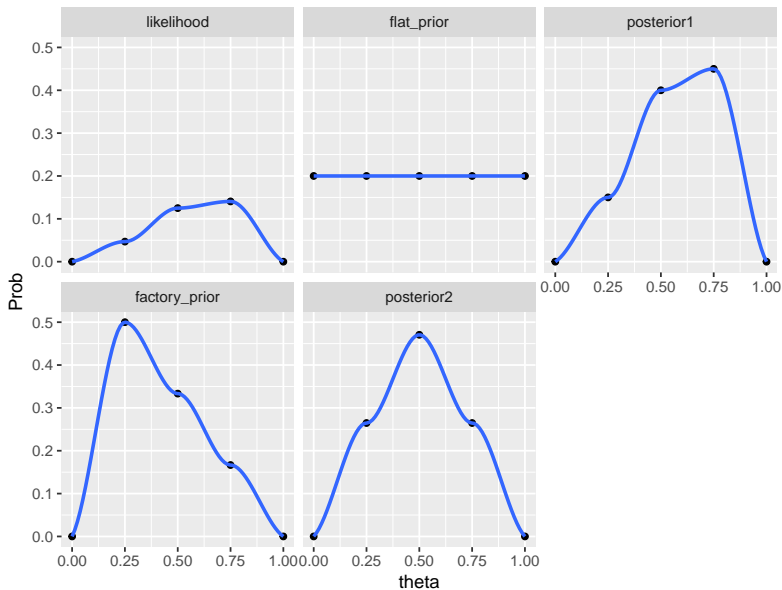
From counts to probability

- Let's use the factory counts information (blue marbles are rare) now in our prior assumptions of $\mathbb{P}(\theta)$.
- This can be done by normalizing the factory counts.
- Notice that this new prior assumption doesn't affect the ways each conjecture can generate the data and $\mathbb{P}(d|\theta)$ remains unchanged.

θ	Factory count	$\mathbb{P}(\theta)$	$\mathbb{P}(d \theta)$	$\mathbb{P}(\theta d) = \mathbb{P}(d \theta) * \mathbb{P}(\theta) / \mathbb{P}(d)$
0	0	0/6	0/64	$\frac{0/64 * 0/6}{0.08854167} = 0$
0.25	3	3/6	3/64	$\frac{3/64 * 3/6}{0.08854167} = 0.2647059$
0.5	2	2/6	8/64	$\frac{8/64 * 2/6}{0.08854167} = 0.4705882$
0.75	1	1/6	9/64	$\frac{9/64 * 1/6}{0.08854167} = 0.2647059$
1	0	0/6	0/64	$\frac{0/64 * 0/6}{0.08854167} = 0$

- where $\mathbb{P}(d) = 0/6 * 0/64 + 3/6 * 3/64 + 2/6 * 8/64 + 1/6 * 9/64 + 0/6 * 0/64 = 0.08854167$
- Two different prior assumptions led us to different values of $\mathbb{P}(\theta|d)$.

From counts to probability



Bayesian Components

Now we will introduce the names of the components of our Bayesian model.

Density and Mass functions

Because the Bayesian framework applies to both discrete and continuous random variables, we will use function f (instead of \mathbb{P}) to refer to both probability mass and density functions.

- **Parameter** θ : A way of indexing possible explanations of the data. In our example θ is a conjectured proportion of blue marbles.
- **Likelihood** $f(d|\theta)$: The relative number of ways that a value θ can produce the data. It is derived by enumerating all the possible data sequences that could have happened and then eliminating those sequences inconsistent with the data.
- **Prior probability** $f(\theta)$: The prior plausibility of any specific value of θ .
- **Posterior probability** $f(\theta|d)$: The new, updated plausibility of any specific θ .
- **Evidence or Average Likelihood** $f(d)$: the average probability of the data averaged over the prior. It's job is just to standardize the posterior, to ensure it sums (integrates) to one.

Bayesian Components

- It is important to remark that in the Bayesian setting a parameter θ is random a variable, so we can make probability statements about it.
- Whereas in the frequentist approach parameters are considered unknown quantities.
- This is an important property of Bayesian inference: despite θ is an **unobserved variable** we can treat it as a random variable and calculate $f(\theta)$ or $f(\theta|d)$.
- The likelihood function $f(d|\theta)$ is very similar to the likelihood function in the frequentist approach $f(d; \theta)$ but now we can condition on θ instead of just using it as function parameter.
- All the probability functions of a Bayesian model can correspond to either 1) a probability mass or 2) a density functions depending if the variable (observed or unobserved) is discrete or continuous.

Bayesian Components

- The general equation that relates all Bayesian components (for both density and mass functions) is the following:

$$f(\theta|d) = \frac{f(d|\theta) \times f(\theta)}{f(d)} \quad (4)$$

- This equation is essentially the Bayes theorem (for both density and mass functions).
- It says that the probability of any particular value of θ considering the data d , is proportional to the product of the relative plausibility of the data, conditional on θ , and the prior plausibility of θ .
- This product is then divided by the average probability of the data to produce a valid probability distribution for the posterior (to sum or integrate to one).
- We must bear in mind that Bayesian statistics is not only about using Bayes theorem.
- There are many non-Bayesian techniques that use this theorem.
- Bayesian inference uses the Bayes theorem more generally, to quantify uncertainty about unobserved variables such as parameters.

Bayesian Components

- In the marble example θ is discrete so the prior and the posterior are probability mass functions.
- When θ is continuous, the prior and the posterior are density functions, and the **evidence** or **average likelihood** is calculated with an integral called **marginal**

$$f(d) = \int_{\theta} f(d|\theta)f(\theta)d\theta \quad (5)$$

- In most cases this integral doesn't have a closed solution.
- However, there are nice computational methods available that can efficiently approximate the posterior even when the evidence cannot be calculated (e.g., MCMC, Variational Inference).
- Next, we will go deeper into these concepts by building another Bayesian toy model.

A Globe Model

- We have a globe representing our planet.
- We want to estimate much of the surface is covered in water.
- We adopt the following strategy: we toss the globe up in the air, we catch it, record whether or not the surface under your right index finger is water or land.
- Then we toss the globe up in the air again and repeat the procedure.
- The first nine samples are: W L W W W L W L W where W indicates water and L indicates land.
- We observed 6 W and 3 L. This is our data.



Designing a simple Bayesian model benefits from a design loop with three steps.

- 1 Data story: Motivate the model by narrating how the data might arise.
- 2 Update: Educate your model by feeding it the data.
- 3 Evaluate: All statistical models require supervision, leading to model revision. ¹

¹We won't elaborate on this part until later in the course..

- You can motivate your data story by trying to explain how each piece of data is born.
- This usually means describing aspects of the underlying reality as well as the sampling process.
- The data story in this case is simply a restatement of the sampling process:
 - 1 The true proportion of water covering the globe is p .
 - 2 A single toss of the globe has a probability p of producing W and $1 - p$ of producing L.
 - 3 Each toss of the globe is independent of the others.
- The data story is then translated into a formal probability model where we assign distributions to our Bayesian components.
- Keep in mind that distribution functions are essentially shortcuts to the process of counting forking paths of the previous example.

Let's define the variables of our model:

- The first variable is the unobserved parameter p , the proportion of water on the globe which is our target of inference.
- The other variables are observed in our data: the count of water W and the count of land L .
- The sum of these two variables is the number of globe tosses: $N = W + L$

Now, we can assign a **likelihood** function to our observed variables given the parameter that respects the two assumptions of our data story:

- 1 Every toss is independent of the other tosses.
- 2 The probability of W is the same on every toss.

A Globe Model

- The binomial distribution is the de facto discrete distribution for this kind of “coin tossing” problem:

$$f(W, L|p) = \frac{(W + L)!}{W!L!} p^W (1 - p)^L$$

- This can be also written as $W \sim \text{Binomial}(W + L, p)$.
- Next, we need to assign initial probability values (our beliefs before observing data) for each possible value of p using a **prior** distribution.
- Recall that p (the proportion of water) can take any real value between 0 and 1.
- We will assume that all possible values of p are equally likely, which implies that p follows a **continuous Uniform distribution** between 0 and 1, $p \sim \text{Uniform}(0, 1)$:

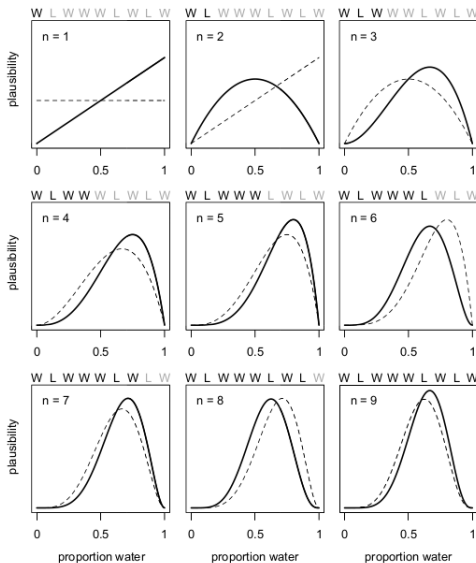
$$f(p) = \text{Uniform}(0, 1) = 1/(b - a) \text{ where } b=1, \text{ and } a=0 = 1$$

- This flat prior assumes that $p = 0$, $p = 0.5$ and $p = 1$ are all equally plausible.
- This is not the best prior information we can declare, considering that we already know that the earth cannot be completely covered by land ($p = 0$) or by water ($p = 1$).

Bayesian Updating

- Now that we have defined our model: variables, likelihood and prior, we can feed it with our data to obtain the **posterior distribution** of p .
- The process of going from the prior $f(\theta)$ to the posterior $f(\theta|d)$ is called **Bayesian Updating**.
- We can view the prior as our initial belief of the possible values that θ can take.
- Then we collect some data d and update our prior using the likelihood to obtain the posterior.
- This process can be repeated iteratively: the posterior becomes a new prior, we collect more data and update our posterior.
- In practice we feed the data only once to our statistical model, but it is important to think that Bayesian updating is an iterated learning process.
- The next slide shows the Bayesian updating process for the Globe tossing example.
- The dashed line shows the prior (or the previous posterior) and the solid line shows the current posterior after seeing each example.

Bayesian Updating



Calculating the Posterior

- The posterior distribution encodes updated plausabilities (or beliefs) for all parameter values conditioned on the data.
- As we have already seen, it can be obtained using the Bayes formula:

$$f(p|W, L) = \frac{f(W, L|p) * f(p)}{f(W, L)}$$

- Where the denominator (evidence) makes sure that the posterior is a valid density function that integrates to one.
- It is not always possible to compute the posterior analytically unless we constrain our prior to special forms that are easy to do mathematics with.
- But bear in mind that in many of the interesting models in contemporary science we will need to approximate the posterior using computational techniques such as Markov Chain Montecarlo.
- This example is one the cases where the posterior can be found analytically as shown next.

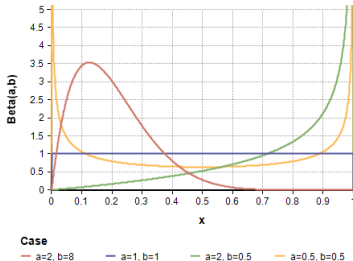
Calculating the Posterior

- The posterior of our globe model with binomial likelihood and uniform prior has a closed form which is a Beta distribution.

$$\mathbb{P}(p|W, L) = \text{Beta}(W + 1, L + 1)$$

- This distribution is defined on the interval $[0, 1]$ and is parameterized by two positive shape parameters, denoted by α and β .
- The Beta distribution is a continuous distribution on probabilities.

$$\text{Beta}(\alpha, \beta) = \frac{p^{\alpha-1}(1-p)^{\beta-1}}{\frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}}$$



Calculating the Posterior

- For any positive integer (such as W and L) the gamma function $\Gamma(n) = (n - 1)!$
- Hence,

$$\text{Beta}(W + 1, L + 1) = \frac{p^W(1 - p)^L}{\frac{\Gamma(W+1)\Gamma(L+1)}{\Gamma(W+1+L+1)}} = \frac{p^W(1 - p)^L}{\frac{W!L!}{(W+L)!}} = \frac{(W + L)!}{W!L!} p^W(1 - p)^L$$

- This surprisingly looks identical to the binomial distribution.
- This is because both distributions are very similar. The binomial distribution models the number of successes (W) and the beta distribution models the probability p of success.
- Let's build our posterior from the likelihood and the prior:

$$f(p|W, L) = \frac{f(W, L|p) * f(p)}{f(W, L)} = \frac{f(W, L|p) * f(p)}{\int_0^1 f(W, L|p) * f(p) dp}$$

Calculating the Posterior

- Since $f(p) = 1$ (uniform prior) we have that

$$f(p|W, L) = \frac{f(W, L|p)}{\int_0^1 f(W, L|p) dp}$$

- The integral of the denominator is equal to 1 (essentially we are integrating a Beta distribution over its complete space of p):



integrate $\Gamma(W+1+L+1)/(\Gamma(W+1)\Gamma(L+1))p^W(1-p)^L$ dp 0 to 1

Extended Keyboard Upload Examples Random

Definite integral:

$$\int_0^1 \frac{\Gamma(W+1+L+1) p^W (1-p)^L}{\Gamma(W+1) \Gamma(L+1)} dp = 1 \text{ for } \operatorname{Re}(L) > -1 \wedge \operatorname{Re}(W) > -1$$

$\Gamma(x)$ is the gamma function
 $\operatorname{Re}(z)$ is the real part of z
 $e_1 \wedge e_2 \wedge \dots$ is the logical AND function

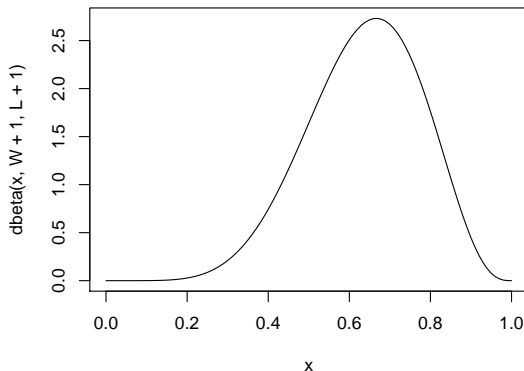
- So, we get

$$f(p|W, L) = \frac{(W+L)!}{W!L!} p^W (1-p)^L = \text{Beta}(W+1, L+1)$$

Calculating the Posterior

- So, in our globe tossing model ($W = 6, L = 3$) we can calculate the posterior distribution analytically

$$f(p|W = 6, L = 3) = \text{Beta}(7, 4)$$



- We could calculate the posterior analytically because a property called **conjugate priors**.
- First of all, we must understand the Beta distribution is very flexible and can model a uniform distributions by setting α and β to 1, $\text{Beta}(1,1)=\text{Uniform}$
- We can change now our prior to a more general one using a Beta distribution $f(p) = \text{Beta}(\alpha, \beta)$.
- Now we can consider values of α, β that are more in line with our prior beliefs.
- The nice thing here is that when the prior follows a Beta distribution and the likelihood $f(W, L|p)$ a Binomial one, the posterior takes the form of another Beta distribution with parameters $(\alpha + W, \beta + L)$.
- The hyper-parameters of our Beta prior α and β can be seen as “pseudo-counts” of successes and failures before collecting the data.
- This shows that our previous result for the uniform prior was a special case of this property.

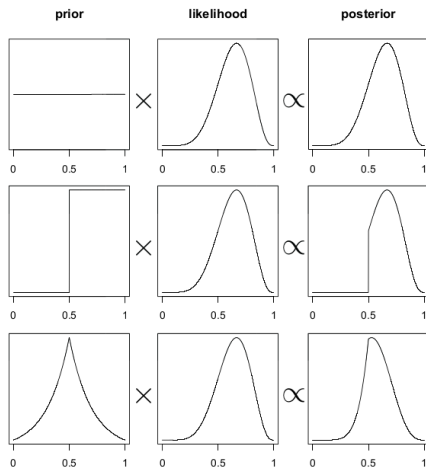
Conjugate Priors

- The Beta distribution is conjugate distribution to binomial distribution, which means that the posterior distribution in the same probability distribution family as the prior.
- In simple words there are some families of conjugate distributions that can be used to calculate posterior distributions analytically.
- A very complete table of conjugate distributions is given in https://en.wikipedia.org/wiki/Conjugate_prior.
- In essence, conjugate priors constrain our choice of prior to special forms that are easy to do mathematics with.
- However, there are numerical techniques that allow us to accommodate any prior that is most useful for our inference problem, such as Markov Chain Monte Carlo (MCMC).

Subjective Bayesian approach

- Priors are engineering assumptions, chosen to help the machine learn.
- They are also scientific assumptions, chosen to reflect what we know about a phenomenon.
- **Subjective Bayesian approach:** a school of Bayesian inference that emphasizes choosing priors based upon the personal beliefs of the analyst.
- This subjective Bayesian approach is rare in the sciences.
- The prior is considered to be just part of the model.
- It should be chosen, evaluated, and revised just like all of the other components of the model.
- The following diagram shows the effect of changing the prior in the globe tossing example.

Different Priors



Top: A flat prior constructs a posterior that is simply proportional to the likelihood.

Middle: A step prior, assigning zero probability to all values less than 0.5, results in a truncated posterior.

Bottom: A peaked prior that shifts and skews the posterior, relative to the likelihood.

- In many interesting models in contemporary science the posterior cannot be calculated analytically, no matter your skill in mathematics.
- Below are some numerical techniques for approximating the mathematics that follows from the definition of the posterior.
 - 1 Grid approximation
 - 2 Laplace approximation
 - 3 Markov chain Monte Carlo (MCMC)
 - 4 Variational Inference²

²Beyond the scope of this course.

Grid Approximation

- Idea: consider only a finite grid of parameter values, and evaluate the posterior for all these points.
- This approach scales very poorly, as the number of parameters increases.
- It is mainly useful as a pedagogical tool, since learning it forces you to really understand the nature of Bayesian updating.

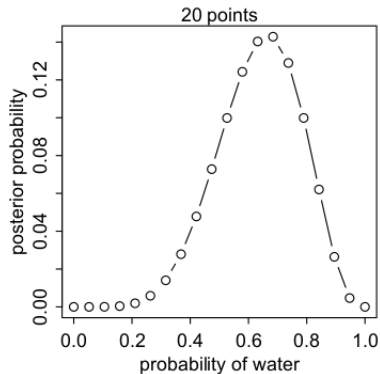
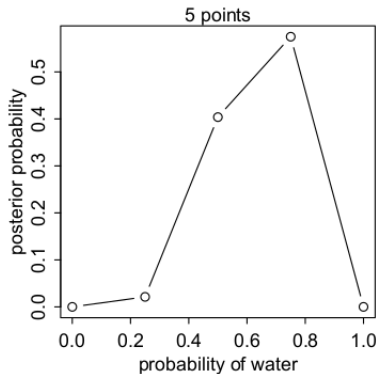
Process

- 1 Define the grid. This means you decide how many points to use in estimating the posterior, and then you make a list of the parameter values on the grid.
- 2 Compute the value of the prior at each parameter value on the grid.
- 3 Compute the likelihood at each parameter value.
- 4 Compute the unstandardized posterior at each parameter value, by multiplying the prior by the likelihood.
- 5 Finally, standardize the posterior, by dividing each value by the sum of all values.

Grid Approximation

```
# define grid
p_grid <- seq( from=0 , to=1 , length.out=20 )
# define prior
prior <- rep( 1 , 20 )
# compute likelihood at each value in grid
likelihood <- dbinom( 6 , size=9 , prob=p_grid )
# compute product of likelihood and prior
unstd.posterior <- likelihood * prior
# standardize the posterior, so it sums to 1
posterior <- unstd.posterior / sum(unstd.posterior)
plot( p_grid , posterior,type="b",
xlab="probability of water",ylab="posterior probability")
mtext("20 points")
```

Grid Approximation



Laplace Approximation

- Under quite general conditions, the region near the peak of the posterior distribution will be nearly Gaussian in shape.
- Idea: approximate the posterior distribution by a Gaussian distribution.
- Gaussians are convenient because can be completely described by two parameters: μ and σ .
- The Laplace approximation is also called “quadratic approximation” because the logarithm of a Gaussian distribution forms a parabola.
- A parabola is a quadratic function.
- Laplace approximation essentially represents any log-posterior with a parabola.

Process

- 1 Maximum a Posteriori (MAP): Find the posterior mode using some optimization algorithm, a procedure that virtually “climbs” the posterior distribution.
- 2 Once you find the peak of the posterior, you must estimate the curvature near the peak.

Laplace Approximation

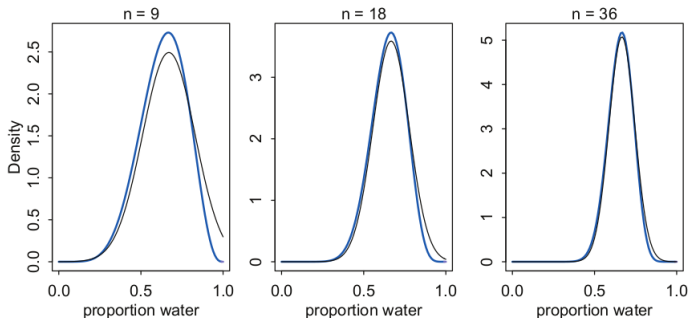
- Laplace approximation is implemented in function **quap** from the **rethinking** R package.
- This package will be used extensively in the remainder of this course.
- Installation instructions: <https://github.com/rmcelreath/rethinking>.

```
library(rethinking)
globe.qa <- quap(
  alist(
    W ~ dbinom( W+L ,p) , # binomial likelihood
    p ~ dunif(0,1)      # uniform prior
  ) ,
  data=list(W=6,L=3) )

# display summary of quadratic approximation
> precis( globe.qa )
  mean    sd 5.5% 94.5%
p 0.67 0.16 0.42  0.92
```

Assuming the posterior is Gaussian, it is maximized at 0.67, and its standard deviation is 0.16.

Laplace Approximation



- Blue Curve: the exact posterior distribution.
- Black Curve: the Laplace approximation.
- Left: The globe tossing data with $n = 9$ tosses and $w = 6$ waters.
- Middle: Double the amount of data, with the same fraction of water, $n = 18$ and $w = 12$.
- Right: Four times as much data, $n = 36$ and $w = 24$.

Maximum a Posteriori and Maximum Likelihood

- The Laplace approximation, either with a uniform prior or with a lot of data, is often equivalent to a maximum likelihood estimate (MLE) and its standard error.
- This is because maximum a posteriori with a uniform prior is equivalent to maximum likelihood estimation.
- By using a uniform prior, our posterior is essentially the likelihood multiplied by a constant, so it makes sense that the maximum value of the posterior and the likelihood are the same.
- More info:
<https://wiseodd.github.io/techblog/2017/01/01/mle-vs-map/>.
- This helps re-interpret many non-Bayesian models in Bayesian terms.

Markov Chain Monte Carlo

- There are lots of important model types, like multilevel (mixed-effects) models, where Laplace approximation doesn't work.
- Such models may have hundreds or thousands or tens-of-thousands of parameters so we can't use Grid approximation either.
- Multilevel models do not always allow us to write down a single, unified function for the posterior distribution.
- This means that the function to maximize (when finding the MAP) is not known, but must be computed in pieces.
- As a result, various counterintuitive model fitting techniques have arisen.
- The most popular of these is Markov chain Monte Carlo (MCMC).

Markov Chain Monte Carlo

- Instead of attempting to compute or approximate the posterior distribution directly, MCMC techniques merely draw samples from the posterior.
- You end up with a collection of parameter values, and the frequencies of these values correspond to the posterior plausibilities.
- You can then build a picture of the posterior from the histogram of these samples. We nearly always work directly with these samples, rather than first constructing some mathematical estimate from them.
- It is fair to say that MCMC is largely responsible for the resurgence of Bayesian data analysis that began in the 1990s.
- While MCMC is older than the 1990s, affordable computer power is not, so we must also thank the engineers.
- Probabilistic programming: environments for designing Bayesian models and performing inference using numerical techniques (e.g., STAN, Pyro).

History

- Bayesian inference was initially developed by English reverend Thomas Bayes and French scientist Pierre Simon Laplace in the 18th century.



- Both scientist worked on the problem of determining the posterior distribution of a binomial likelihood with a uniform prior.
- In those times, Bayesian inference was known as “Inverse Probability”.

- In the 20th century, the frequentist approach became mainstream.
- The founding figures of the frequentist approach (e.g., Fisher, Pearson) did not embrace the Bayesian school.
- Ronald Fisher, made the following statement about inverse probability: “The theory of inverse probability is founded upon an error, and must be wholly rejected” [Fisher, 1925]
- Perhaps this is the reason why “Bayesian inference” was neglected for most of the first half of the 20th century.
- Nowadays, Bayesian methods are widely used by many practicing scientists, for example, in computer science and machine learning.

Conclusions

- The target of inference in Bayesian inference is a posterior probability distribution.
- Posterior probabilities state the relative numbers of ways each conjectured cause of the data could have produced the data.
- These relative numbers indicate plausibilities of the different conjectures.
- These plausibilities are updated in light of observations through Bayesian updating.
- Bayesian models are fit to data using numerical techniques.
- Each method imposes different trade-offs.

References I



Fisher, R. A. (1925).
Statistical methods for research workers.
Oliver & Boyd.



McElreath, R. (2020).
Statistical rethinking: A Bayesian course with examples in R and Stan.
CRC press.



Wasserman, L. (2013).
All of statistics: a concise course in statistical inference.
Springer Science & Business Media.