# Bayesian Statistics and Uncertainty Quantification
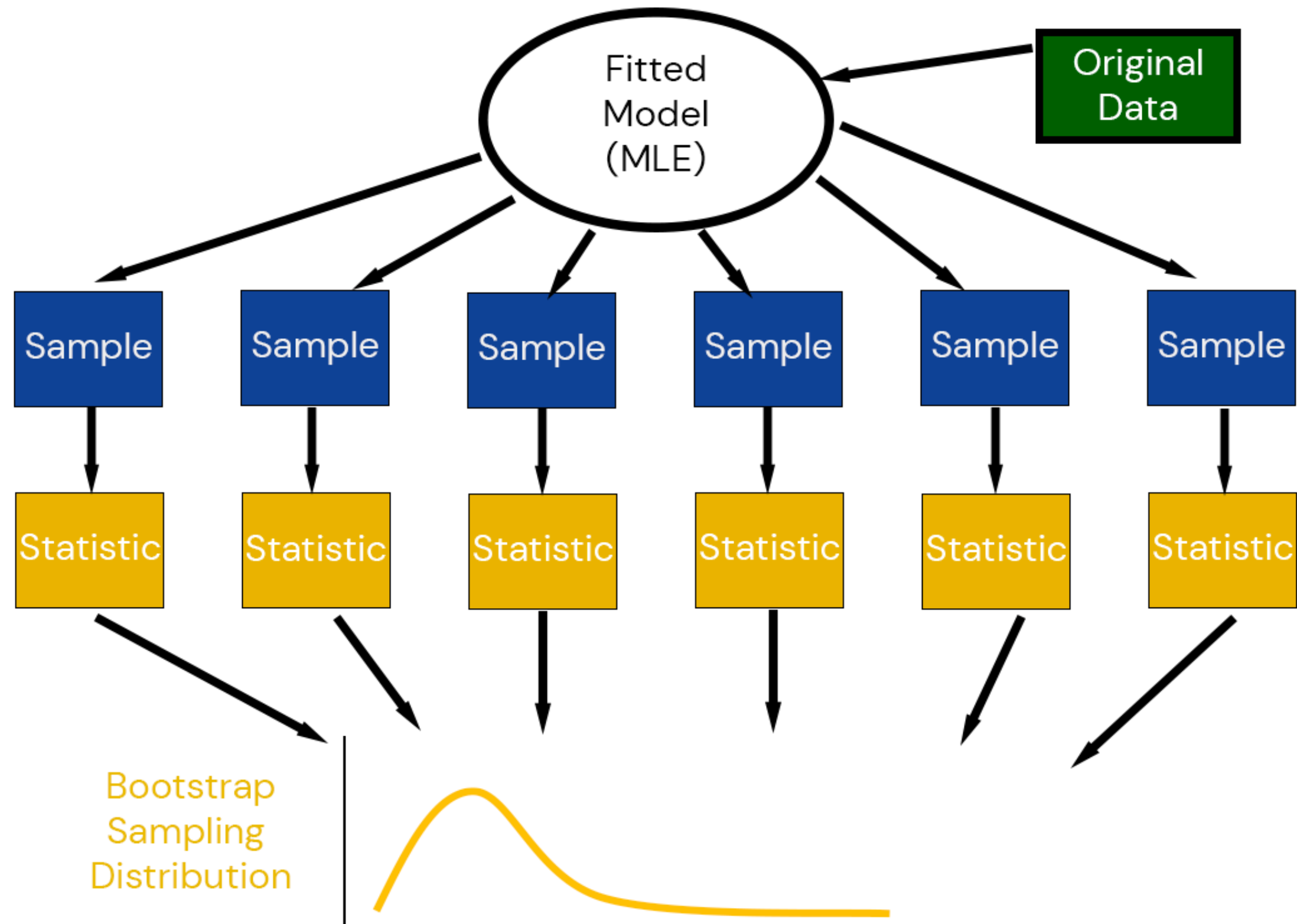
# TABLE OF CONTENTS

# REVIEW OF THE BOOTSTRAP

# LAST CLASS: FREQUENTIST STATISTICS AND SAMPLING DISTRIBUTIONS

**Frequentist UQ**: Capture the *sampling distribution* of relevant parameters.

This takes into account uncertainty in the data sample and reflects the impact of inference.

# REVIEW: THE PARAMETRIC BOOTSTRAP

# REVIEW: THE PARAMETRIC BOOTSTRAP

**Key assumptions**:

- Data is sufficiently representative of the population;

- Model structure reasonably captures data-generating process and/or residuals

# SOURCES OF BOOTSTRAP ERROR

The bootstrap has three potential sources of error:

1. **Sampling error**: error from using finitely many replications

2. **Statistical error**: error in the bootstrap sampling distribution approximation

3. **Specification error** (parametric): Error in the data-generating model

# BUT WHAT IF WE DON'T CARE ABOUT FREQUENCY?

For what might be called the "lab science paradigm," frequency properties are central to make inferences about relevant scientific laws.

But for risk analysis, do we care about them?

# But What If We Don't Care About Frequency?

For what might be called the "lab science paradigm," frequency properties are central to make inferences about relevant scientific laws.

But for risk analysis, do we care about them?

**Perhaps not!** A more relevant perspective: how much should we believe in a given level of a future or present risk? This is the perspective of *Bayesian statistics*.

# Introduction to Bayesian Statistics

# THE BAYESIAN PERSPECTIVE

From the Bayesian perspective, probability is interpreted as the degree of belief in an outcome or proposition.

# THE BAYESIAN PERSPECTIVE

From the Bayesian perspective, probability is interpreted as the degree of belief in an outcome or proposition.

There are two different two types of random quantities:

- Observable quantities, or data (also random for frequentists);

- Unobservable quantities, or parameters.

We can also speak of probabilities on model structures, rather than framing model selection as hypothesis-testing.

# CONDITIONAL PROBABILITY NOTATION

Then it makes sense to discuss the *probability* of

- model parameters $\theta$

- unobserved data $\tilde{\mathbf{y}}$

conditional on the observations $\mathbf{y}$, which we can denote:

$$p(\theta|\mathbf{y}) \text{ or } p(\tilde{\mathbf{y}}|\mathbf{y})$$

# CONDITIONING ON OBSERVATIONS

This fundamental conditioning on observations $\mathbf{y}$ is a distinguishing feature of Bayesian inference.

Compare: frequentist approaches are based on re-estimated over the distribution of possible $\mathbf{y}$ conditional on the "true" parameter value.

# Bayesian Updating

Bayesian probabilities are conditional on observations.

This means that as make new observations, we can *update* them.

We do this using **Bayes' Rule**.

# Bayes' Rule

Original version (Bayes (1763), *An Essay towards solving a Problem in the Doctrine of Chances*):

$$P(A|B) = \frac{P(B|A) \times P(A)}{P(B)} \quad \text{if} \quad P(B) \neq 0.$$
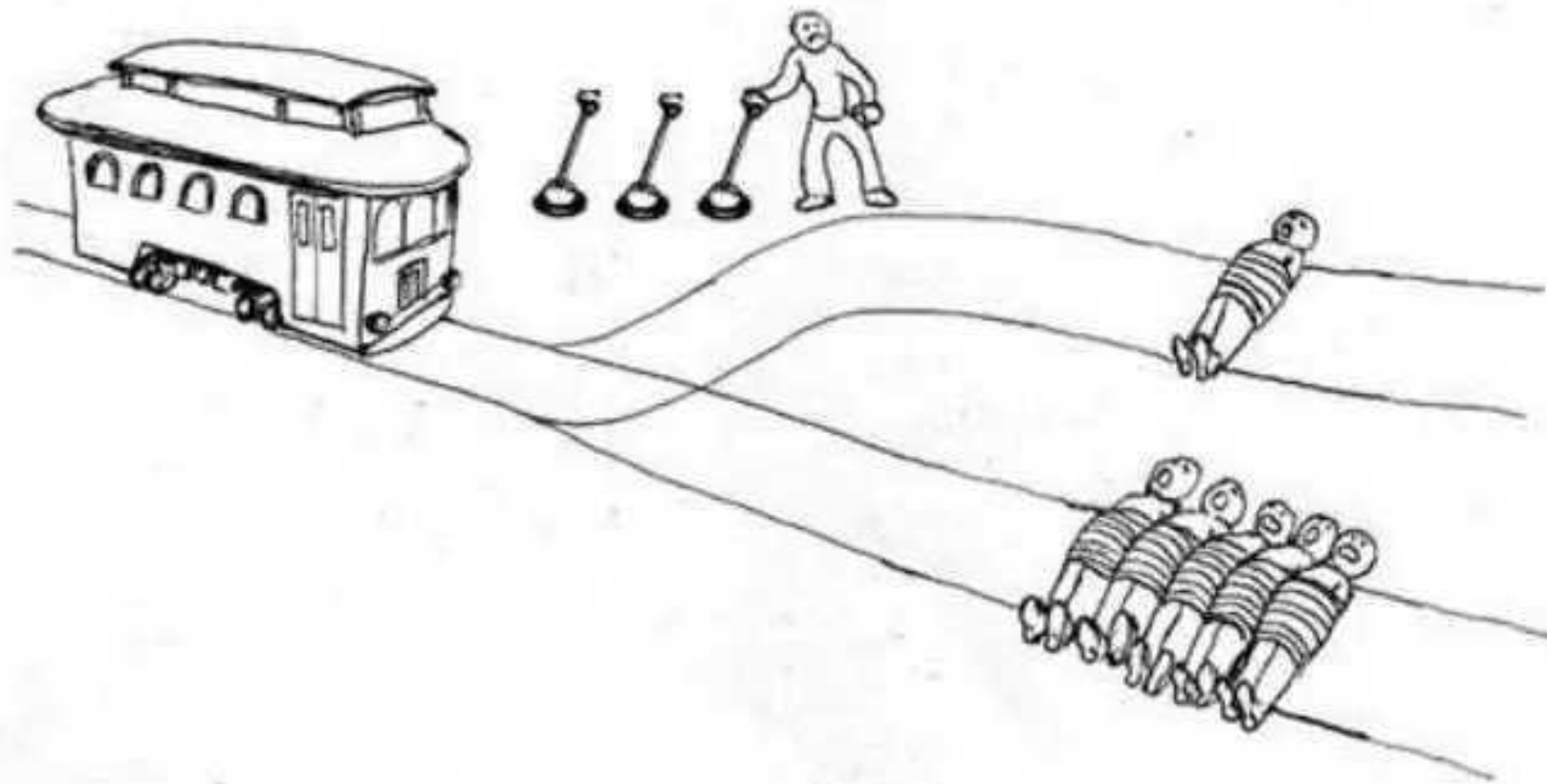
# BAYES' RULE

Original version (Bayes (1763), *An Essay towards solving a Problem in the Doctrine of Chances*):

$$P(A|B) = \frac{P(B|A) \times P(A)}{P(B)} \quad \text{if} \quad P(B) \neq 0.$$

- Standard but important result in conditional probability

- As previously seen (hopefully remembered!) in Intro Stats

  - Monty Hall Problem

# Bayes' Rule



The Trolley Hall Problem
There are three levers. Two of the levers have no
effect. You may only choose one lever to pull. You
made your choice but before you pulled your lever,
one of the other two levers is revealed as fake. Do
you switch your choice?

# Bayes' Rule

Original version (Bayes (1763), *An Essay towards solving a Problem in the Doctrine of Chances*):

$$P(A|B) = \frac{P(B|A) \times P(A)}{P(B)} \quad \text{if} \quad P(B) \neq 0.$$

- Bayes used this to estimate the distribution of a probability $p$ of a binomial outcome (think success/failure).

- Richard Price (actual writer of quite a bit of Bayes (1763); see Stigler (2018)) "rebutted" Hume by "demonstrating" we ought to believe the sun will continue to rise.

# Bayes' Rule

"Modern" version (Laplace (1774), *Mémoire sur la probabilité des causes par les événements*):

$$p(\theta|y) = \frac{p(y|\theta)}{p(y)}p(\theta)$$

# Bayes' Rule

"Modern" version (Laplace (1774), *Mémoire sur la probabilité des causes par les événements*):

$$\underbrace{p(\theta|y)}_{\text{posterior}} = \frac{\overbrace{p(y|\theta)}^{\text{likelihood}}}{\underbrace{p(y)}_{\text{normalization}}} \overbrace{p(\theta)}^{\text{prior}}$$

# ON THE NORMALIZING CONSTANT

The normalizing constant (also called the **marginal likelihood**) is the integral

$$p(y) = \int_\Theta p(y|\theta)p(\theta)d\theta.$$

Since this *generally* doesn't depend on $\theta$, it can often be ignored, as the relative probabilities don't change.

# ON THE NORMALIZING CONSTANT

The normalizing constant (also called the **marginal likelihood**) is the integral

$$p(y) = \int_\Theta p(y|\theta)p(\theta)d\theta.$$

Since this *generally* doesn't depend on $\theta$, it can often be ignored, as the relative probabilities don't change.

One big exception: model selection (will discuss later...)

# Bayes' Rule (Ignoring Normalizing Constants)

The version of Bayes' rule which matters the most for 95% (approximate) of Bayesian statistics:

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

> *"The posterior is the prior times the likelihood..."*

# Bayesian Model Components

This means that a Bayesian model specification requires two key components:

1. Probability model for the data given the parameters (the *likelihood*), $p(y|\theta)$t

2. Prior distributions over the parameters, $p(\theta)$

   ○ Can be independent or joint

# BAYESIAN MODEL COMPONENTS

This means that a Bayesian model specification requires two key components:

1. Probability model for the data given the parameters (the *likelihood*), $p(y|\theta)$t

2. Prior distributions over the parameters, $p(\theta)$

   ○ Can be independent or joint

**Bayesian updating**: Using the likelihood to "update" the prior probabilities of the parameters.

# A Coin Flipping Example

We would like to understand if a coin-flipping game is fair. We've observed the following sequence of flips:

H, H, H, T, H, H, H, H, H

# A Coin Flipping Example

We would like to understand if a coin-flipping game is fair. We've observed the following sequence of flips:

$$H, H, H, T, H, H, H, H, H$$

8/9 are heads, which might seem suspicious, but randomness can result in outliers like this.
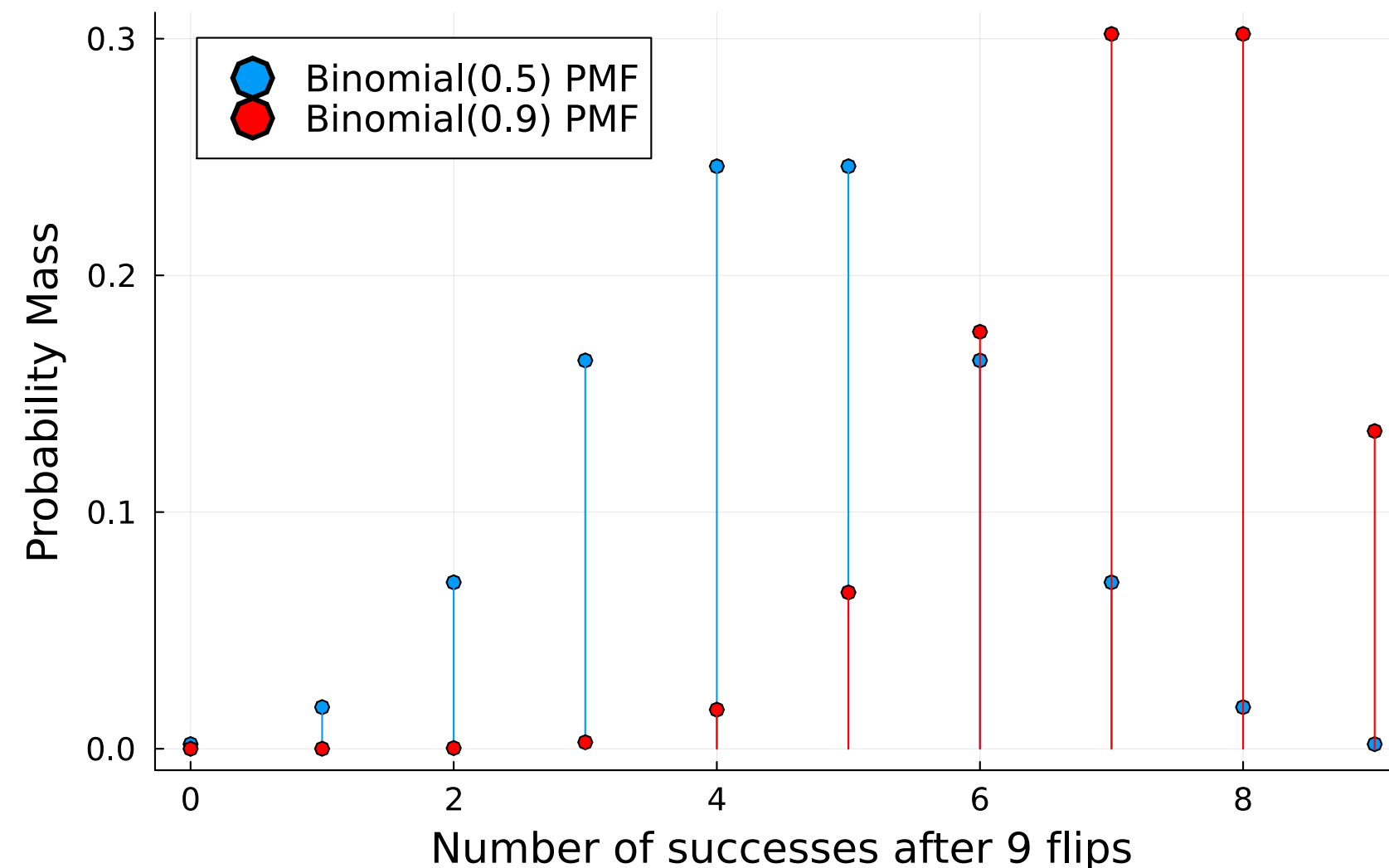
# COIN FLIPPING LIKELIHOOD

The data-generating process here is straightforward: we can represent a coin flip with a heads-probability of $\theta$ as a sample from a Binomial distribution,

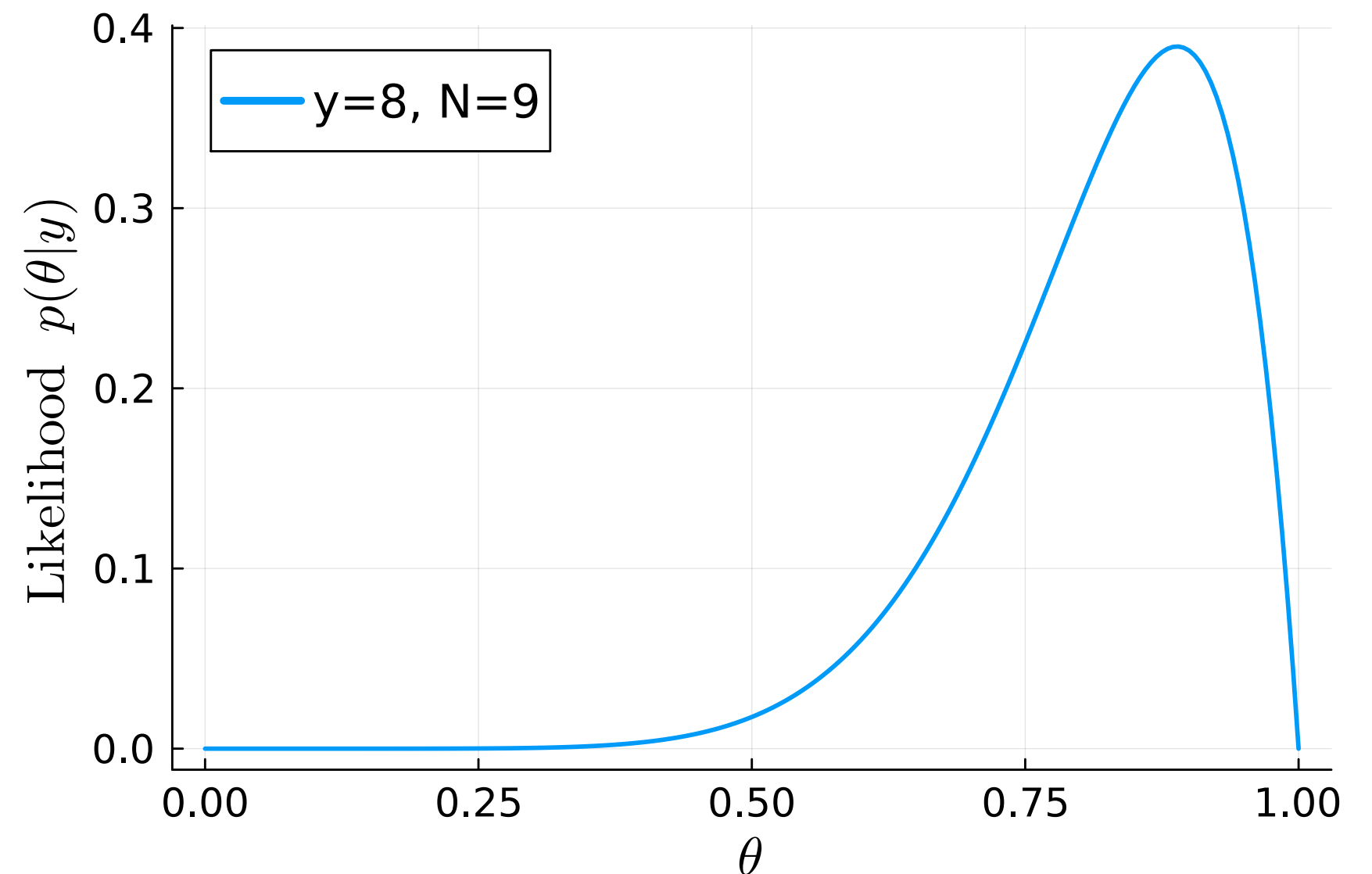$$y \sim \mathrm{Binomial}(\theta).$$

# COIN FLIPPING PROBABILITY MASS

Let's compare what the **probability mass functions** of these distributions look like for $\theta = 0.5$ and $\theta = 0.9$.

# LIKELIHOOD FUNCTION

The PMF told us what the probability of a given dataset was given a fixed parameter $\theta$. But we can view this same function from a different perspective: given the number of successes, how *likely* is a given parameter?
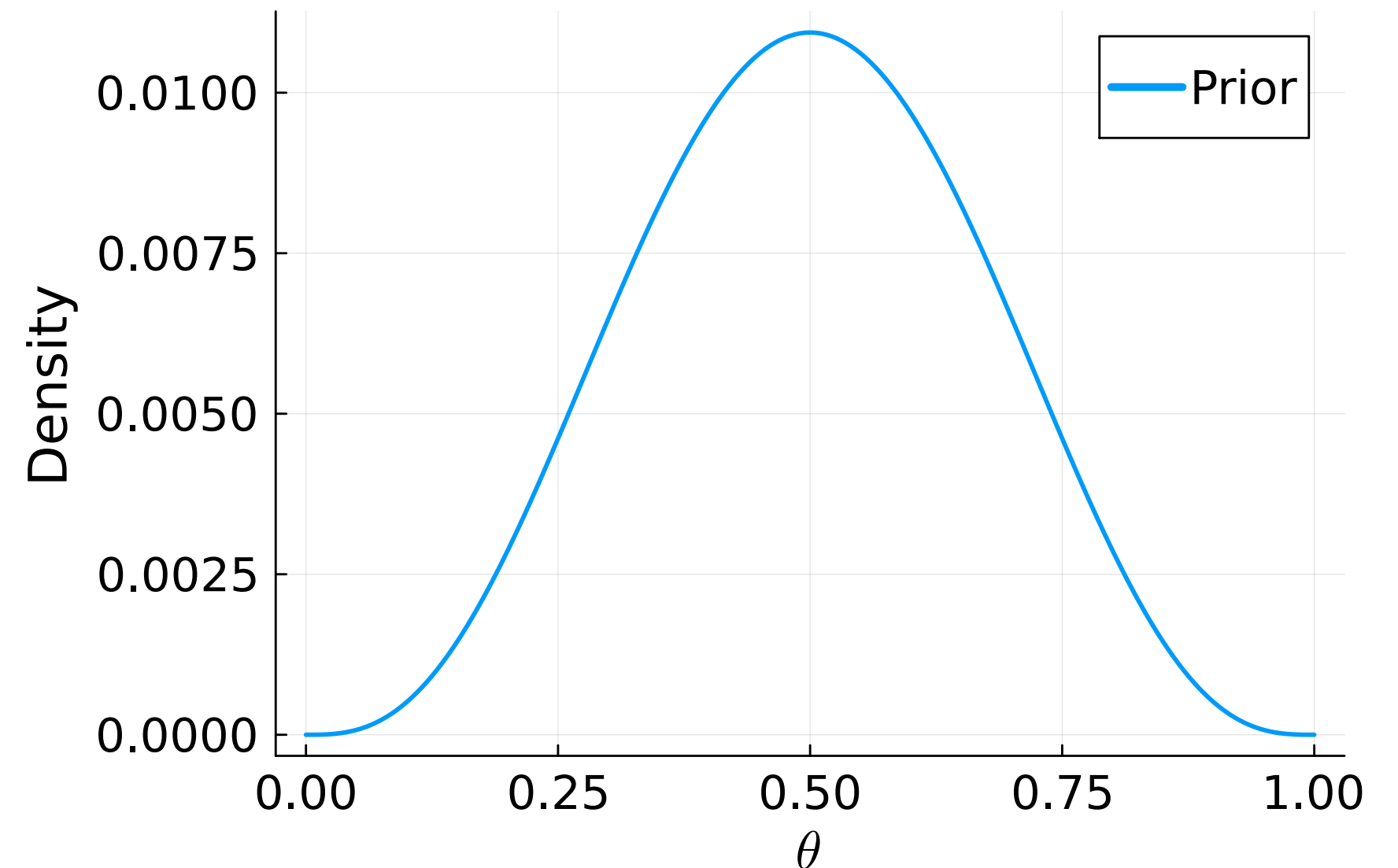
# Prior Distribution

For frequentist approaches, we could stop there and maximize the likelihood, and we'd get a maximum likelihood estimate of $\theta \approx 0.88$.

But suppose that we spoke to a friend who knows something about coins, and she tells us that it is extremely difficult to make a passable weighted coin which comes up heads more than 75% of the time. Since we have a relatively small amount of data, this seems like valuable information to include, and we can do this through our prior.
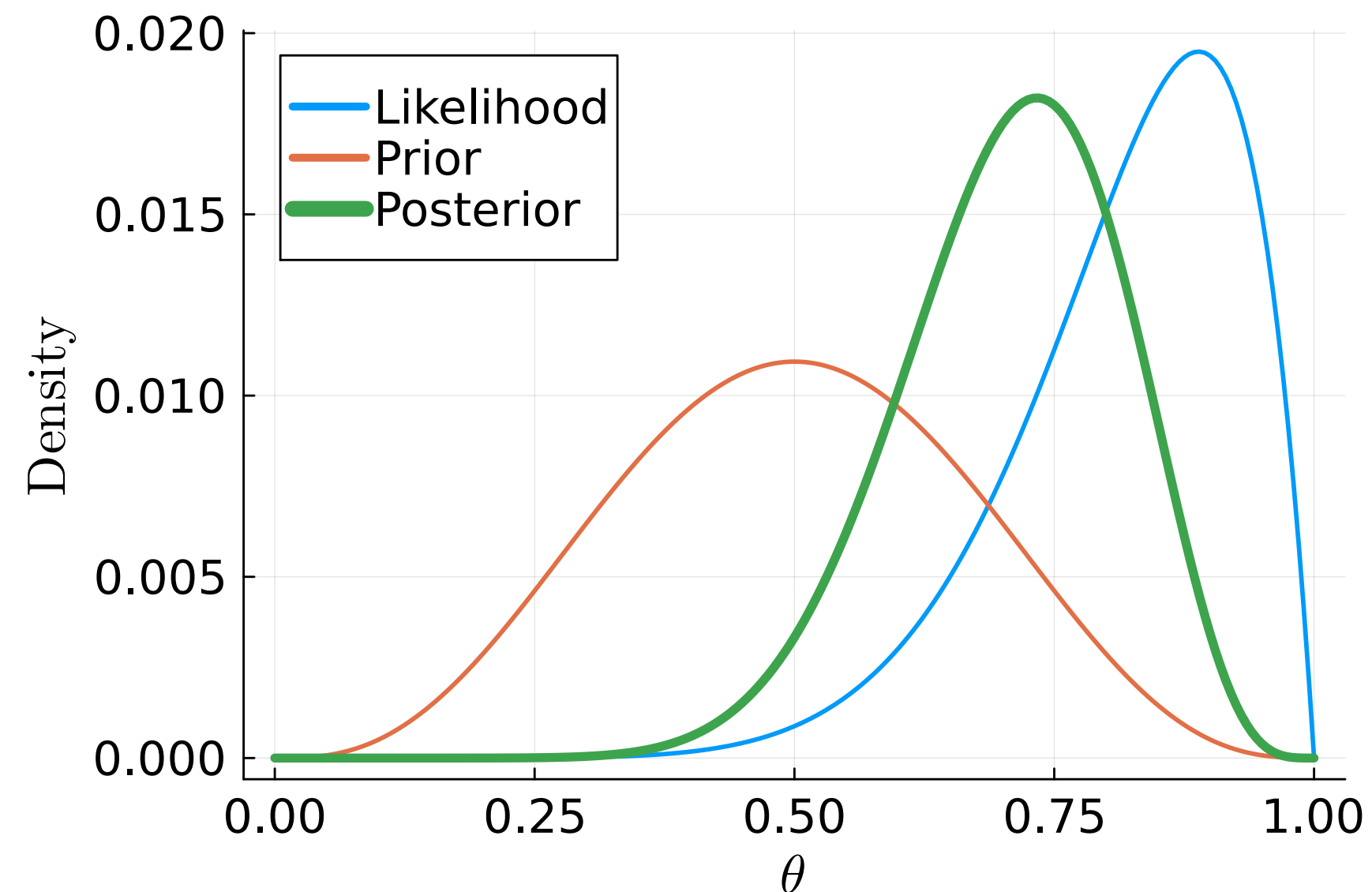
# PRIOR DISTRIBUTION

Since $\theta$ is bounded between 0 and 1, we'll settle on a Beta distribution for our prior, specifically $\mathrm{Beta}(4, 4)$, which covers a reasonable spread of possibilities while maintaining symmetry.

# POSTERIOR DISTRIBUTION

Combining using Bayes' rule gives us a **maximum** *a posteriori* **(MAP)** estimate of $\theta \approx 0.74$.

# Bayesian Updating As An Information Filter

- The posterior is a "compromise" between the prior and the data.

- The posterior mean is a weighted combination of the data and the prior mean

- The weights depend on the prior and the likelihood variances

- More data makes the posterior more confident (lower variance)

# Representing Uncertainty

As with frequentist approaches, can reflect posterior inferences through a point estimate (mean, median, or some other *Bayes estimator*).

But more often, we want to capture the degree of uncertainty associated with a particular value.

# CREDIBLE INTERVALS

Bayesian **credible intervals** are straightforward to interpret: $\theta$ is in $I$ with probability $\alpha$.

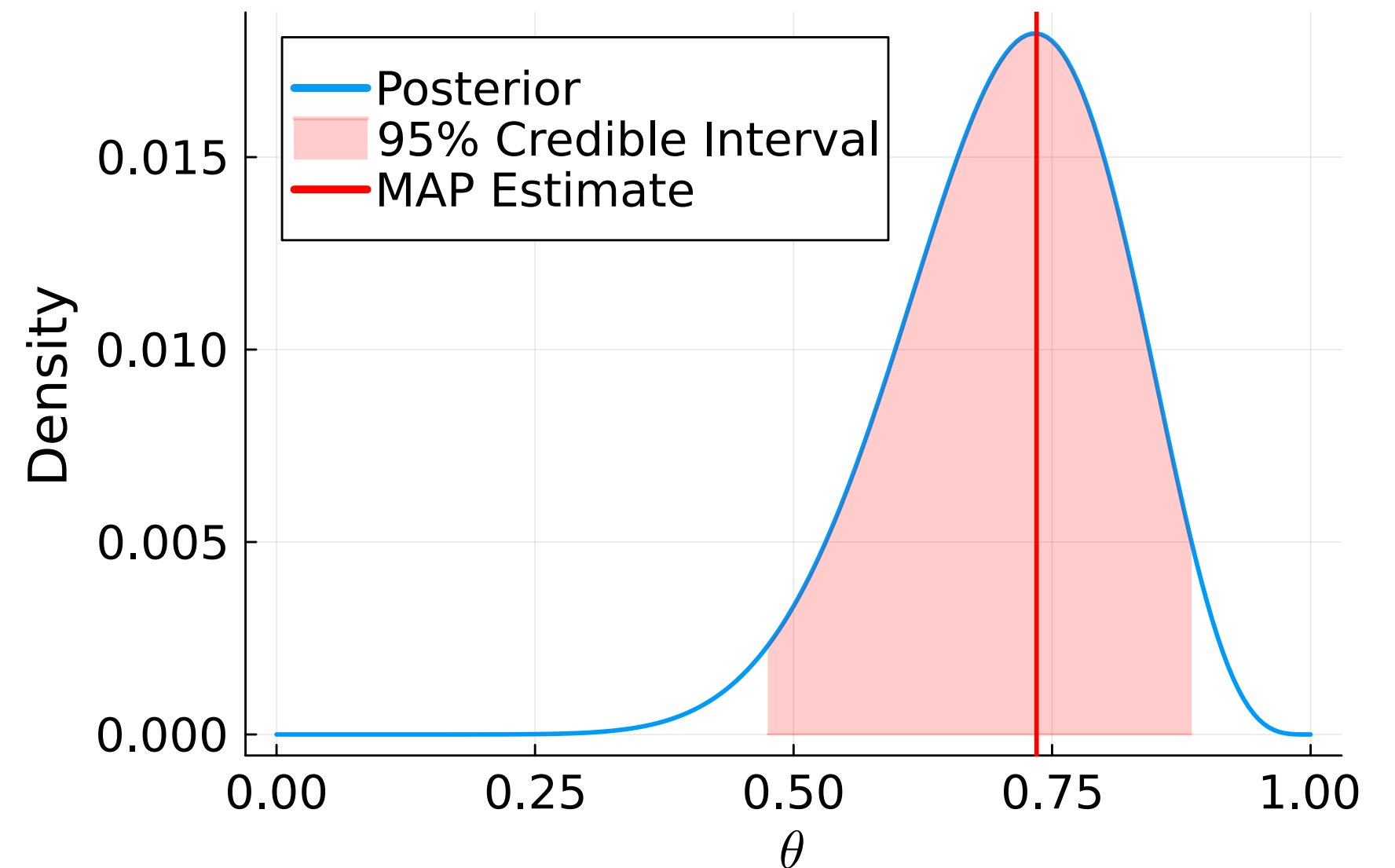In other words, choose $I$ such that

$$p(\theta \in I | \mathbf{y}) = \alpha.$$

This is not usually a unique choice, but the "equal-tailed interval" between the $(1 - \alpha)/2$ and $(1 + \alpha)/2$ quantiles is a common choice.

# CREDIBLE INTERVALS: COIN FLIPPING EXAMPLE

In the case of our coin flipping example, how much uncertainty is there? Let's say we want to capture the 95% credible interval, which is $(0.48, 0.89)$.

# CREDIBLE INTERVALS: COIN FLIPPING EXAMPLE

But that was for a simple example where it was easy to compute the posterior for a large number of values.

The easiest way to do this in general is through Monte Carlo: draw a lot of samples from the posterior and compute the empirical quantiles. We'll discuss later today/next week.

# MORE COMPLEX MODELS

One advantage of the Bayesian framework is it can be extended to more complex problems:

- Models with heteroskedastic residual structures:

$$y(t) = f(\theta, t) + R(\phi, t)$$
$$R(\phi, t) = \zeta(\phi, t) + \varepsilon t$$

# More Complex Models

One advantage of the Bayesian framework is it can be extended to more complex problems:

- Hierarchical models:

$$y_j | \theta_j, \phi \sim P(y_j | \theta_j, \phi)$$
$$\theta_j | \phi \sim P(\theta_j | \phi)$$
$$\phi \sim P(\phi)$$

# Generative Modeling

Bayesian models can also be used to generate new data $\tilde{y}$ through the *posterior predictive distribution*:

$$p(\tilde{y}|\mathbf{y}) = \int_{\Theta} p(\tilde{y}|\theta)p(\theta|\mathbf{y})d\theta$$

This allows us to test the model through simulation (*e.g.* hindcasting) and generate probabilistic predictions.

# Role of the Prior

There are two views on how to select prior distributions:

1. Priors as capturing by constraints/prior knowledge;

# ROLE OF THE PRIOR

There are two views on how to select prior distributions:

1. Priors as capturing by constraints/prior knowledge;

2. Priors as part of the data-generating process.

# ON UNIFORM PRIORS

What about uniform priors?

- Unbounded uniform priors: often chosen to reflect ignorance, but nonsensical for data-generation;

- Bounded uniform priors: sudden transition from positive to zero probability is rarely justified.

# CONSIDERATIONS WHEN SELECTING PRIORS

- **Informativeness**: How much information does the prior encode?

- **Structure**: Does the prior encode modeling features (*e.g.* symmetry)?

- **Regularization**: Does the prior yield more "stable" inferences (*e.g.* penalizing extreme parameter values)?

Also: what values are assigned zero prior probability? These values are ruled out from the posterior.

# BAYESIAN COMPUTATION: A PREVIEW

# GOALS OF BAYESIAN COMPUTATION

1. Sampling from the *posterior* distribution

$$p(\theta | \mathbf{y})$$

2. Sampling from the *posterior predictive* distribution

$$p(\tilde{y} | \mathbf{y})$$

by generating data.

# BAYESIAN COMPUTATION AND MONTE CARLO

In other words, Bayesian computation involves Monte Carlo simulation from the posterior (predictive) distribution.

These samples can then be analyzed to identify estimators, credible intervals, etc.
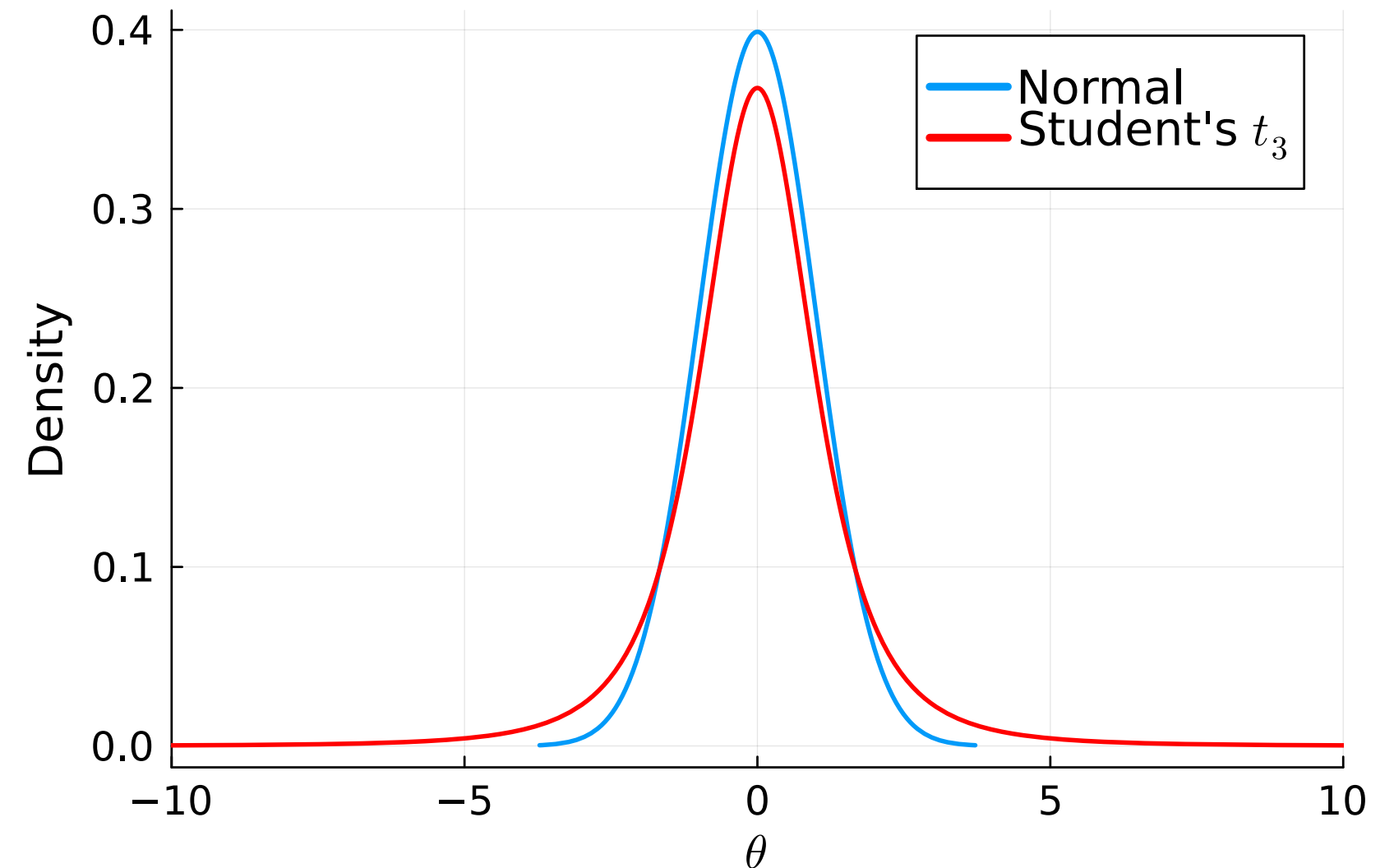
# SAMPLING FROM THE POSTERIOR

Trivial for extremely simple problems: low-dimensional and with "conjugate" priors (which make the posterior a closed-form distribution).

What to do when problems are more complex and/or we don't want to rely on priors for computational convenience?

# A First Algorithm: Rejection Sampling

Idea:

1. Generate proposed samples from another distribution $g(\theta)$ which covers the target $p(\theta|\mathbf{y})$;

2. Accept those proposals based on the ratio of the two distributions.

# REJECTION SAMPLING ALGORITHM

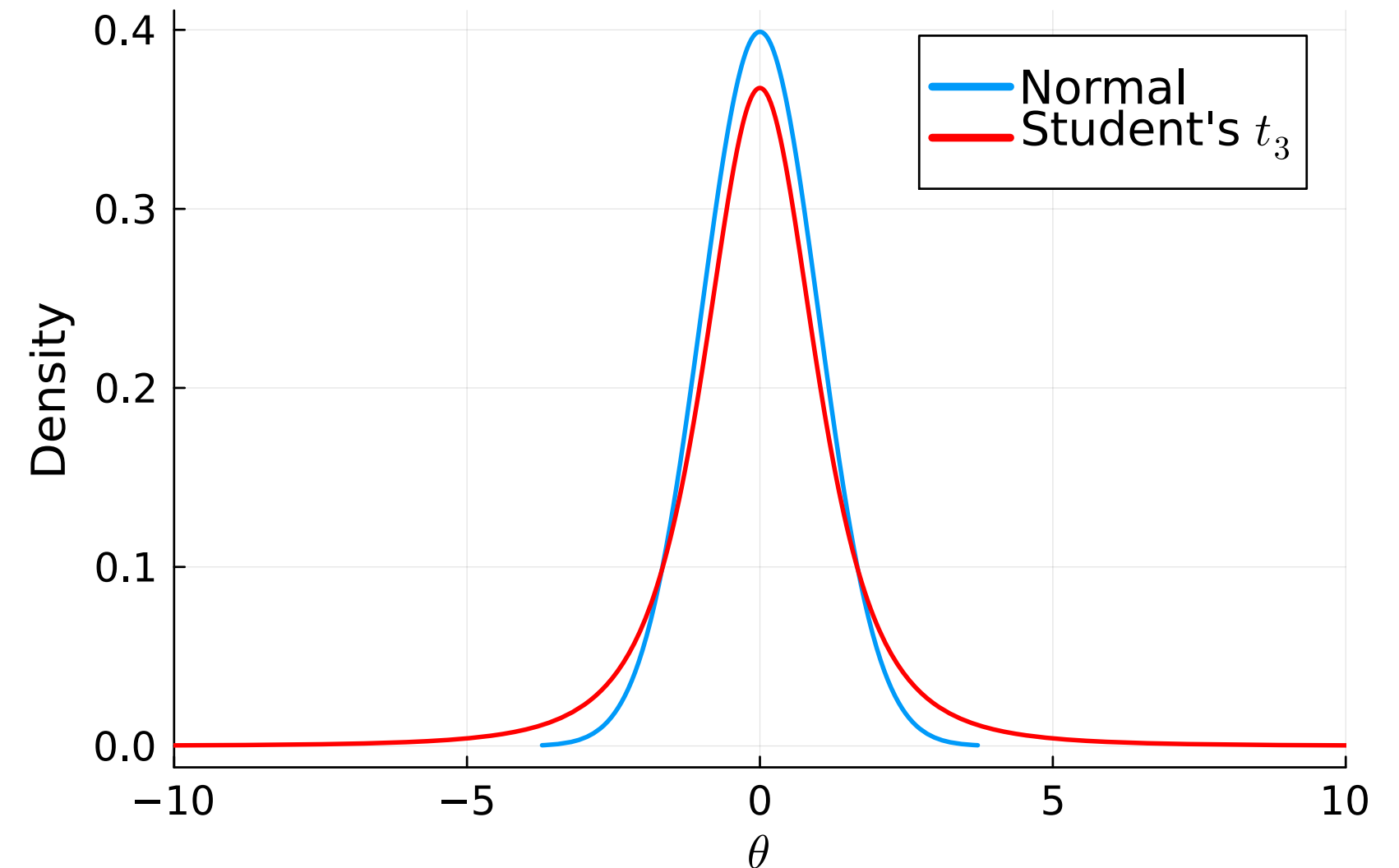Suppose $p(\theta|\mathbf{y}) \leq Mg(\theta)$ for some $1 < M < \infty$.

1. Simulate $u \sim \mathbf{Unif}(0,1)$.

2. Simulate a proposal $\hat{\theta} \sim g(\theta)$.

3. If

$$u < \frac{p(\hat{\theta}|\mathbf{y})}{Mg(\hat{\theta})},$$

accept $\hat{\theta}$. Otherwise reject.

# Rejection Sampling Intuition

We want to keep more samples from the areas where $g(\theta) < p(\theta|\mathbf{y})$ and reject where $g$ is heavier (in this case, the tails).
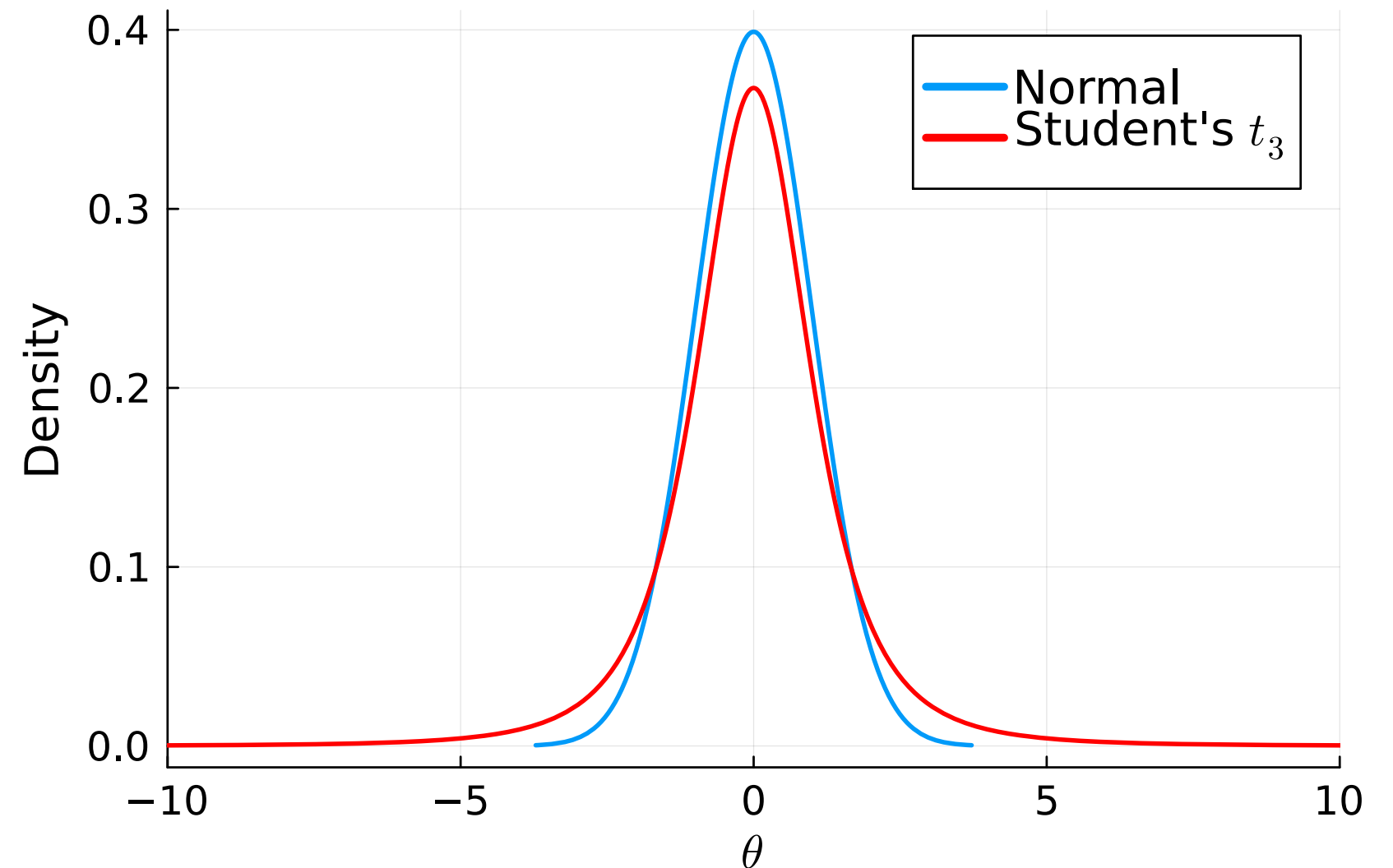
# REJECTION SAMPLING INTUITION

In the bulk: $Mg(\theta)$ is closer to $p(\theta|\mathbf{y})$, so the acceptance ratio is higher.

In the tails:

$$Mg(\theta) \gg p(\theta|\mathbf{y}),$$

so the acceptance ratio is much lower.

# REJECTION SAMPLING CONSIDERATIONS

1. Probability of accepting a sample is $1/M$, so the "tighter" the proposal distribution coverage the more efficient the sampler.

2. Need to be able to compute $M$.

Finding a good proposal and computing $M$ may not be easy for complex posteriors!

**How can we do better?**

# KEY TAKEAWAYS

# KEY TAKEAWAYS: BAYESIAN STATISTICS

- Probability as degree of belief

- Emphasis on explicit conditioning on the data

- Bayes' Rule as the fundamental theorem of conditional probability

- Bayesian updating as an information filter

- Prior selection important: lots to consider!

- Rejection sampling as a first Monte Carlo algorithm for sampling from "arbitrary" distributions.

# Upcoming Schedule

# Upcoming Schedule

**Next Monday**: Markov chains and Markov chain Monte Carlo