

# **Approximate computation: Monte Carlo**

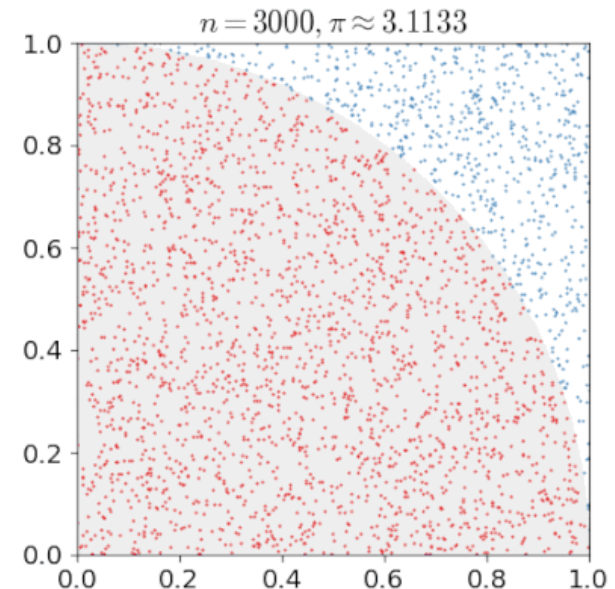
Roger Levy

9.S918: Quantitative inference in brain and cognitive sciences

24 February 2025

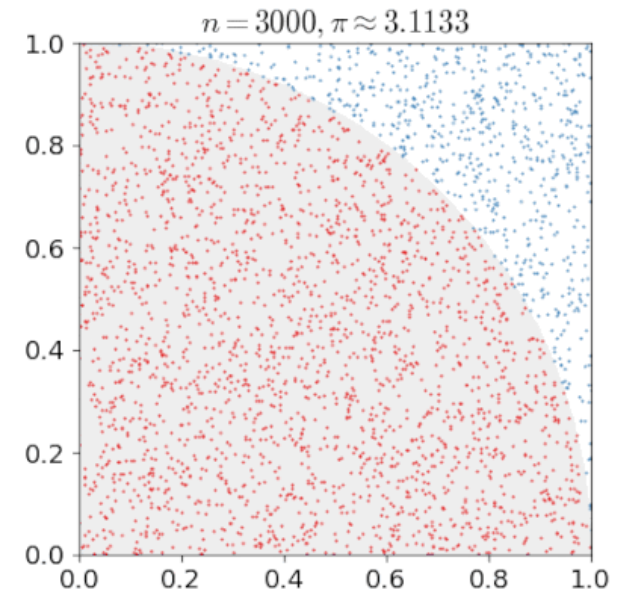
# Monte Carlo methods, or "probabilistic simulation"

- Generally speaking:
  1. Define a domain of possible inputs
  2. Generate  $n$  iid random inputs from a probability distribution on the domain
  3. Perform a deterministic computation on each randomly generated input
  4. Aggregate the results of the deterministic computation
- As  $n$  grows larger, the simulated result approaches the true value



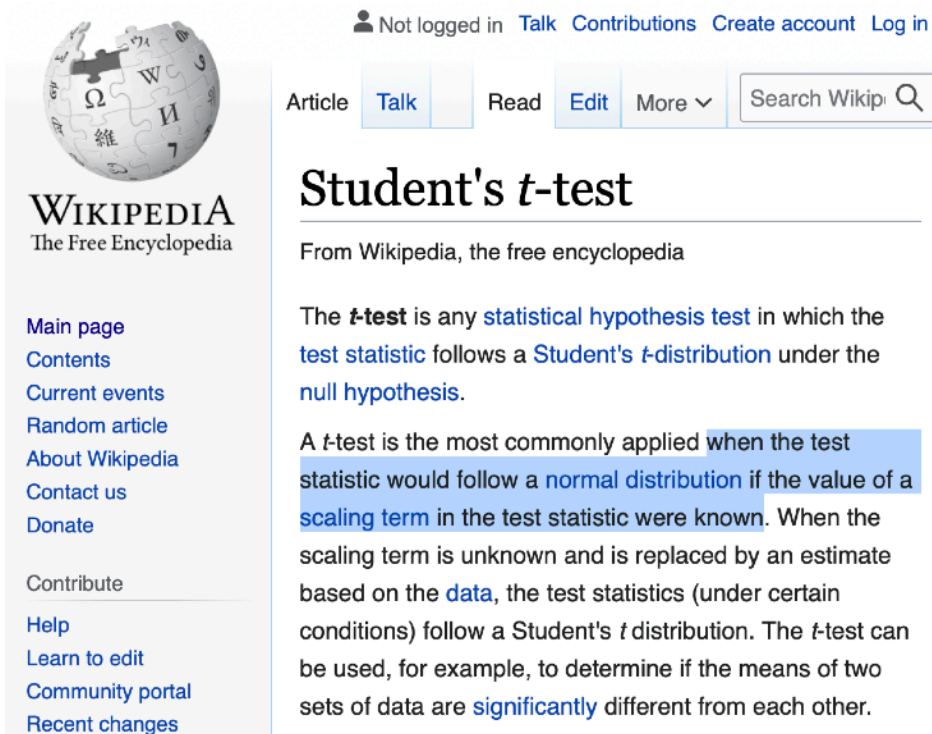
# Monte Carlo methods, or "probabilistic simulation"

- Generally speaking:
  1. Define a domain of possible inputs
  2. Generate  $n$  iid random inputs from a probability distribution on the domain
  3. Perform a deterministic computation on each randomly generated input
  4. Aggregate the results of the deterministic computation
- As  $n$  grows larger, the simulated result approaches the true value

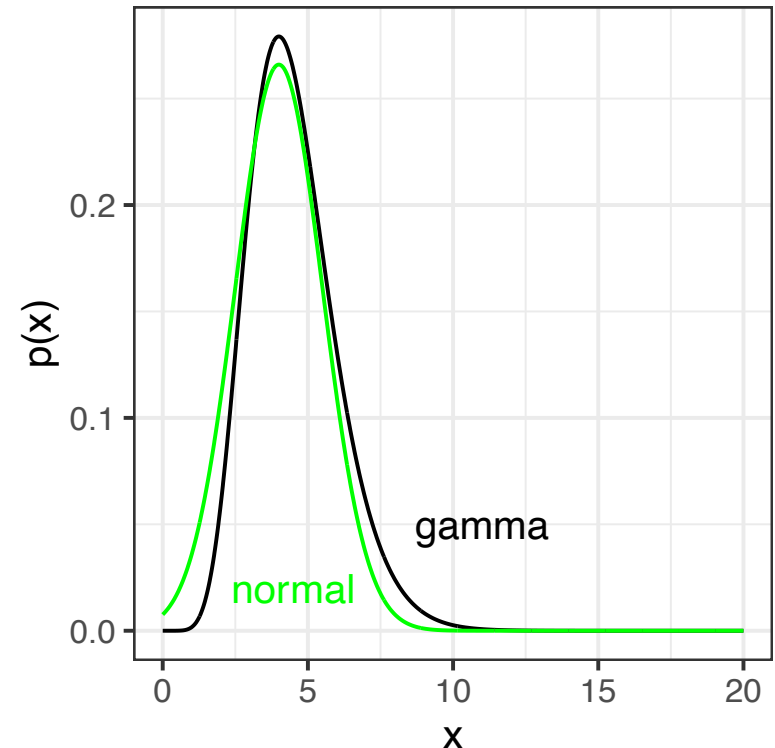


# Simple example of Monte Carlo

- Suppose I want to do a two-sample  $t$ -test but my data aren't normally distributed



The screenshot shows the Wikipedia page for "Student's  $t$ -test". The page header includes the Wikipedia logo and navigation links like "Not logged in", "Talk", "Contributions", "Create account", and "Log in". The article title is "Student's  $t$ -test". Below the title, it says "From Wikipedia, the free encyclopedia". The main text explains that the  $t$ -test is a statistical hypothesis test where the test statistic follows a Student's  $t$ -distribution under the null hypothesis. A highlighted section states: "A  $t$ -test is the most commonly applied when the test statistic would follow a normal distribution if the value of a scaling term in the test statistic were known. When the scaling term is unknown and is replaced by an estimate based on the data, the test statistics (under certain conditions) follow a Student's  $t$  distribution. The  $t$ -test can be used, for example, to determine if the means of two sets of data are significantly different from each other."



- How bad will this be for my  $t$ -test????

# Monte Carlo, in action

---

```
1 library(ggplot2)
2 library(tidyverse)
3
4 # Manually compute Student t-statistic
5 f <- function(seed,N=100,shape=9,scale=0.5) {
6   set.seed(seed)
7   y1 <- rgamma(N,shape=shape,scale=scale)
8   y2 <- rgamma(N,shape=shape,scale=scale)
9   s_p <- sqrt( (var(y1) + var(y2)) / 2 )
10  t_statistic <- ( mean(y1) - mean(y2) ) / ( s_p*sqrt(2/N) )
11  return(t_statistic)
12 }
13
14 N <- 100
15 Ts <- sapply(1:10000,f)
16
17 t_reference <- tibble(x=seq(-4,4,by=0.01),t=dt(x,df=2*(N-1)))
18
19 ggplot(data=tibble(t=Ts),aes(x=t)) +
20   geom_density() +
21   geom_line(data=t_reference,aes(x=x,y=t),color="green",linetype="dashed")
```

# Monte Carlo, in action

---

```
1 library(ggplot2)
2 library(tidyverse)
3
4 # Manually compute Student t-statistic
5 f <- function(seed,N=100,shape=9,scale=0.5) {
6   set.seed(seed) ← Reproducibility!
7   y1 <- rgamma(N,shape=shape,scale=scale)
8   y2 <- rgamma(N,shape=shape,scale=scale)
9   s_p <- sqrt( (var(y1) + var(y2)) / 2 )
10  t_statistic <- ( mean(y1) - mean(y2) ) / ( s_p*sqrt(2/N) )
11  return(t_statistic)
12 }
13
14 N <- 100
15 Ts <- sapply(1:10000,f)
16
17 t_reference <- tibble(x=seq(-4,4,by=0.01),t=dt(x,df=2*(N-1)))
18
19 ggplot(data=tibble(t=Ts),aes(x=t)) +
20   geom_density() +
21   geom_line(data=t_reference,aes(x=x,y=t),color="green",linetype="dashed")
```

# Monte Carlo, in action

```
1 library(ggplot2)
2 library(tidyverse)
3
4 # Manually compute Student t-statistic
5 f <- function(seed,N=100,shape=9,scale=0.5) {
6   set.seed(seed) ← Reproducibility!
7   y1 <- rgamma(N,shape=shape,scale=scale)
8   y2 <- rgamma(N,shape=shape,scale=scale)
9   s_p <- sqrt( (var(y1) + var(y2)) / 2 )
10  t_statistic <- ( mean(y1) - mean(y2) ) / ( s_p*sqrt(2/N) )
11  return(t_statistic)
12 }
13
14 N <- 100 Monte Carlo simulation
15 Ts <- sapply(1:10000,f)
16
17 t_reference <- tibble(x=seq(-4,4,by=0.01),t=dt(x,df=2*(N-1)))
18
19 ggplot(data=tibble(t=Ts),aes(x=t)) +
20   geom_density() +
21   geom_line(data=t_reference,aes(x=x,y=t),color="green",linetype="dashed")
```

# Monte Carlo, in action

```
1 library(ggplot2)
2 library(tidyverse)
3
4 # Manually compute Student t-statistic
5 f <- function(seed,N=100,shape=9,scale=0.5) {
6   set.seed(seed)
7   y1 <- rgamma(N,shape=shape,scale=scale)
8   y2 <- rgamma(N,shape=shape,scale=scale)
9   s_p <- sqrt( (var(y1) + var(y2)) / 2 )
10  t_statistic <- ( mean(y1) - mean(y2) ) / ( s_p*sqrt(2/N) )
11  return(t_statistic)
12 }
13
14 N <- 100
15 Ts <- sapply(1:10000,f)
16
17 t_reference <- tibble(x=seq(-4,4,by=0.01),t=dt(x,df=2*(N-1)))
18
19 ggplot(data=tibble(t=Ts),aes(x=t)) +
20   geom_density() +
21   geom_line(data=t_reference,aes(x=x,y=t),color="green",linetype="dashed")
```

**Reproducibility!**

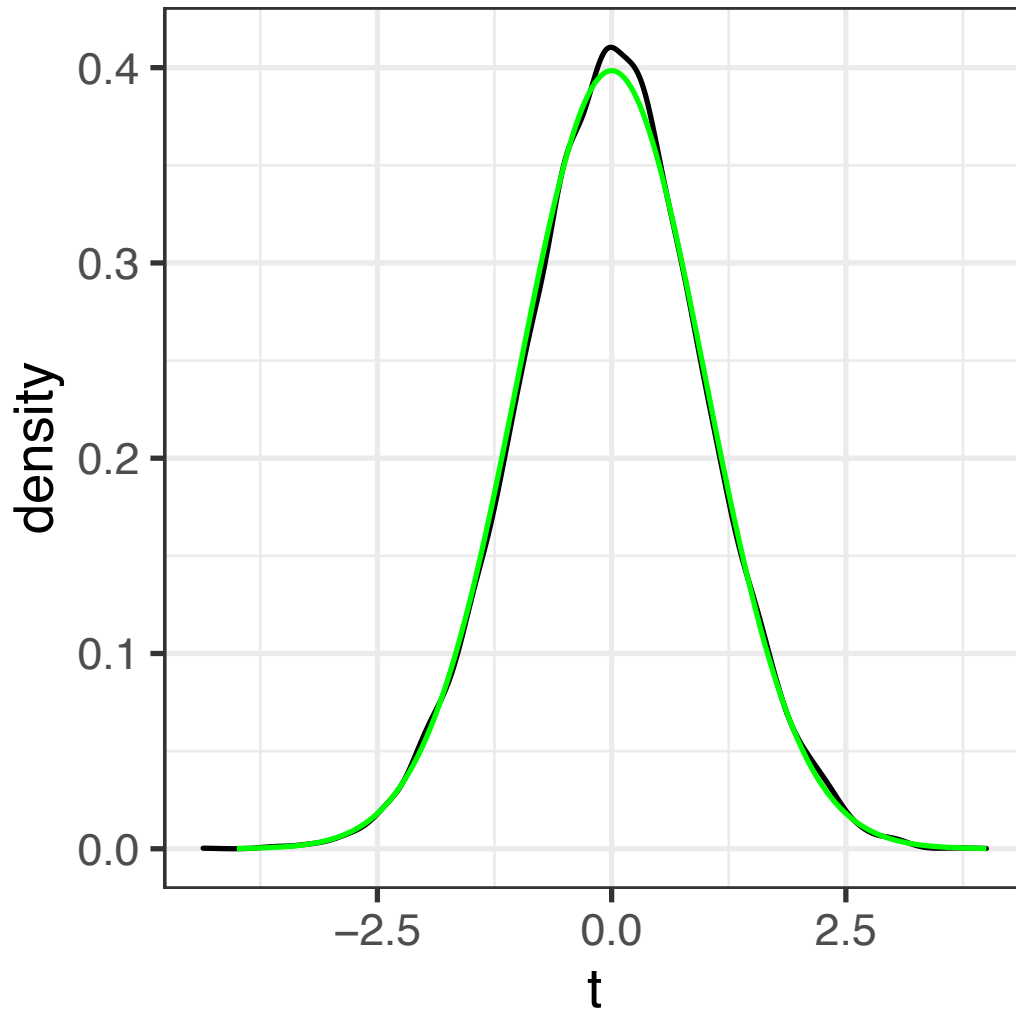
**Monte Carlo simulation**

**Compare against Student's  $t$  distribution**



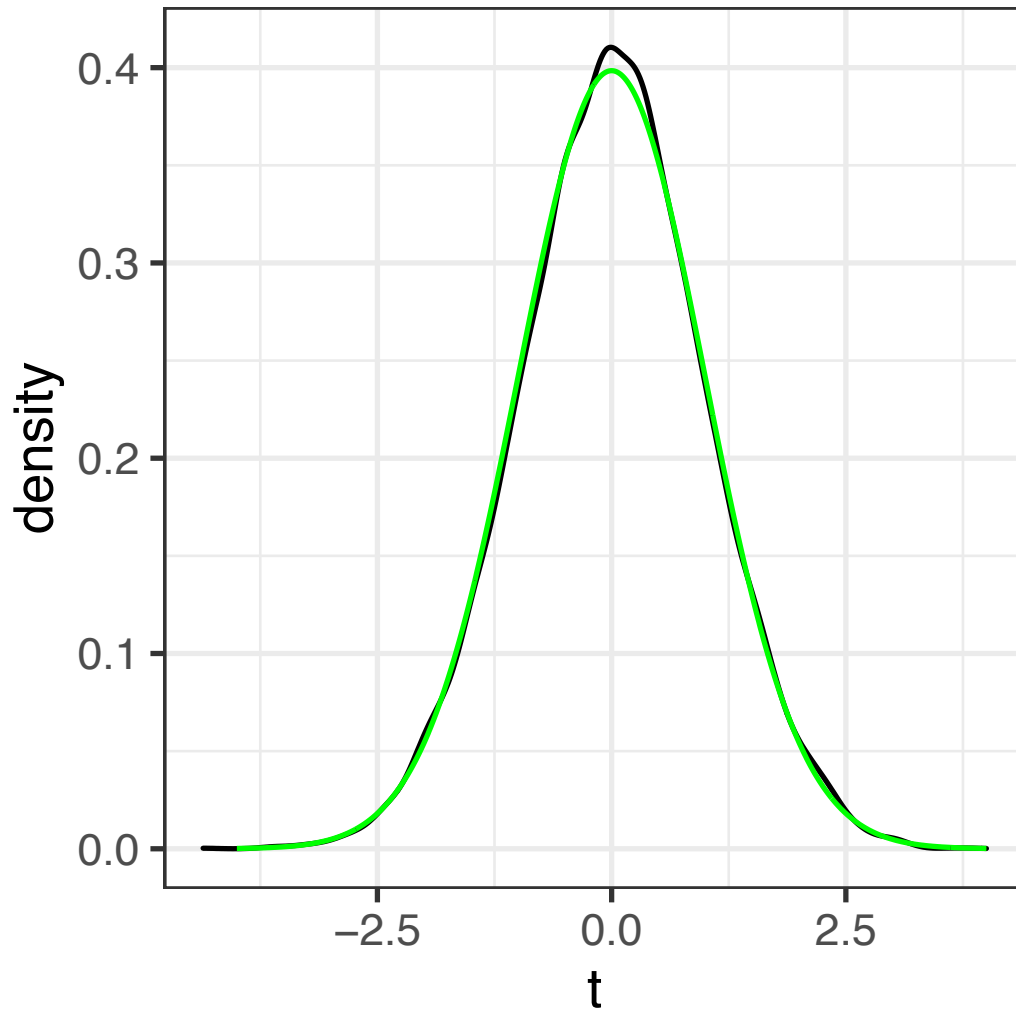
# Monte Carlo, in action

---



# Monte Carlo, in action

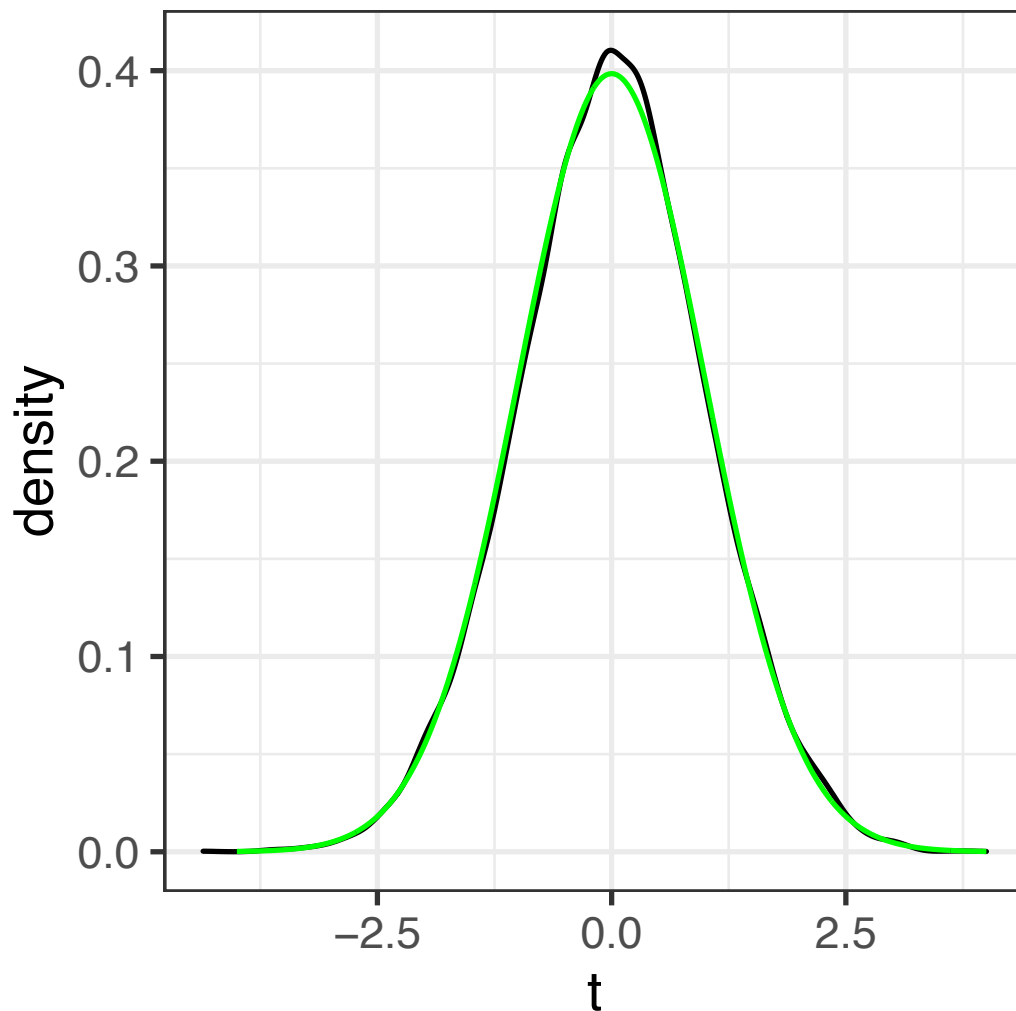
---



- The  $t$  distribution is still a pretty good approximation of the distribution of the  $t$  statistic, even when the underlying distribution is gamma!

# Monte Carlo, in action

---



- The  $t$  distribution is still a pretty good approximation of the distribution of the  $t$  statistic, even when the underlying distribution is gamma!
- This exemplifies what is meant when people say that the  $t$  test is **robust to deviations from normality**

# Unnormalizable posteriors

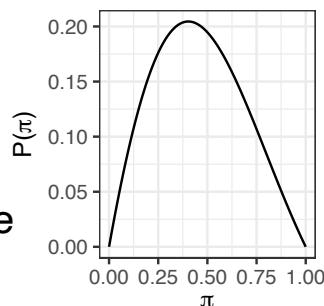
- Our motivation: Bayesian posterior inference

$$P(\theta | \mathbf{y}, I) = \frac{P(\mathbf{y} | \theta, I)P(\theta | I)}{P(\mathbf{y} | I)}$$

- Sometimes  $P(\mathbf{y} | I)$  can't be calculated exactly. Example

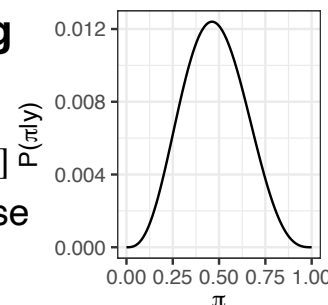
**Bernoulli data with non-conjugate prior:**

$$P(\pi) \propto \begin{cases} \pi(1 - \pi)e^{-\pi^2} & \pi \in [0,1] \\ 0 & \text{otherwise} \end{cases}$$

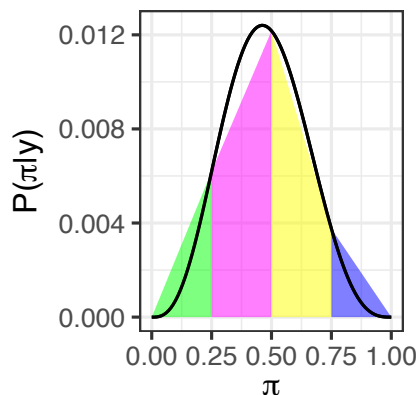


**Posterior after observing 2 heads, 2 tails:**

$$P(\pi) \propto \begin{cases} \pi^3(1 - \pi)^3e^{-\pi^2} & \pi \in [0,1] \\ 0 & \text{otherwise} \end{cases}$$



- In simple cases like this, we can numerically approximate the integral:



- But in high dimension and/or unbounded ranges, difficult or even impossible!

# Unnormalizable posteriors

- Our motivation: Bayesian posterior inference

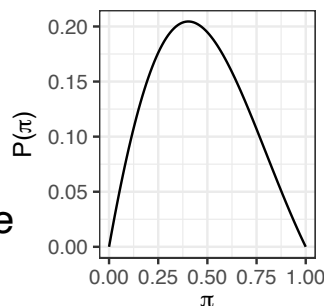
$$P(\theta | \mathbf{y}, I) = \frac{P(\mathbf{y} | \theta, I) P(\theta | I)}{P(\mathbf{y} | I)}$$

*Model parameters* →  $\theta$

- Sometimes  $P(\mathbf{y} | I)$  can't be calculated exactly. Example

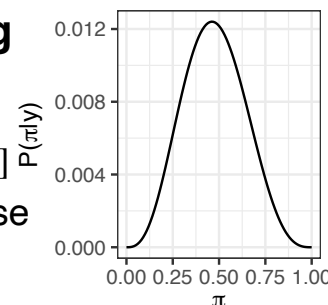
**Bernoulli data with non-conjugate prior:**

$$P(\pi) \propto \begin{cases} \pi(1-\pi)e^{-\pi^2} & \pi \in [0,1] \\ 0 & \text{otherwise} \end{cases}$$

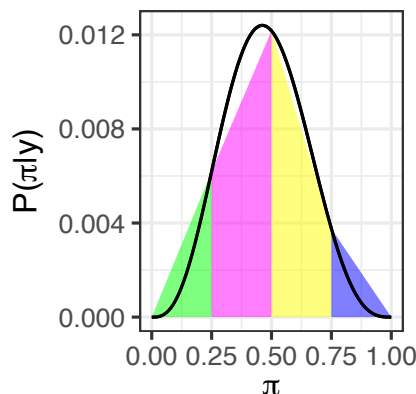


**Posterior after observing 2 heads, 2 tails:**

$$P(\pi) \propto \begin{cases} \pi^3(1-\pi)^3e^{-\pi^2} & \pi \in [0,1] \\ 0 & \text{otherwise} \end{cases}$$



- In simple cases like this, we can numerically approximate the integral:



- But in high dimension and/or unbounded ranges, difficult or even impossible!

# Unnormalizable posteriors

- Our motivation: Bayesian posterior inference

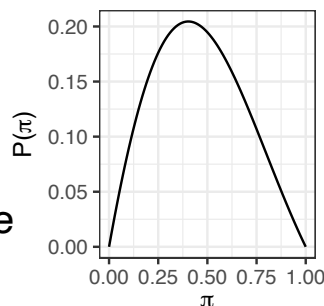
$$P(\theta | \mathbf{y}, I) = \frac{P(\mathbf{y} | \theta, I) P(\theta | I)}{P(\mathbf{y} | I)}$$

*Model parameters* →  $\theta$        $I$  ← *Background knowledge*

- Sometimes  $P(\mathbf{y} | I)$  can't be calculated exactly. Example

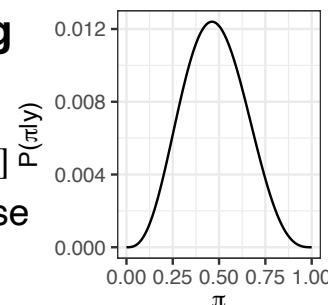
**Bernoulli data with non-conjugate prior:**

$$P(\pi) \propto \begin{cases} \pi(1-\pi)e^{-\pi^2} & \pi \in [0,1] \\ 0 & \text{otherwise} \end{cases}$$

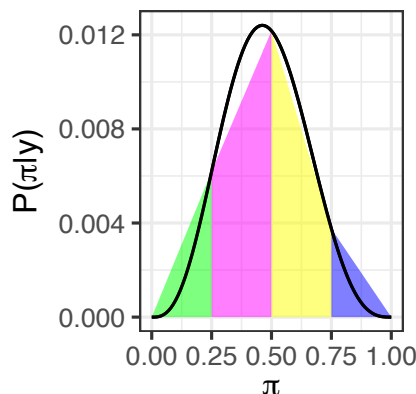


**Posterior after observing 2 heads, 2 tails:**

$$P(\pi) \propto \begin{cases} \pi^3(1-\pi)^3e^{-\pi^2} & \pi \in [0,1] \\ 0 & \text{otherwise} \end{cases}$$



- In simple cases like this, we can numerically approximate the integral:



- But in high dimension and/or unbounded ranges, difficult or even impossible!

# Unnormalizable posteriors

- Our motivation: Bayesian posterior inference

$$P(\theta | \mathbf{y}, I) = \frac{P(\mathbf{y} | \theta, I) P(\theta | I)}{P(\mathbf{y} | I)}$$

*Observed data* points to  $\mathbf{y}$

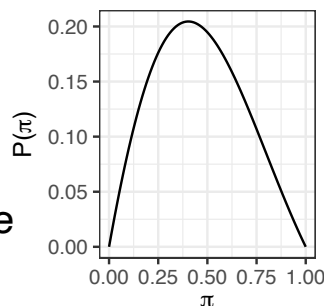
*Model parameters* points to  $\theta$

*Background knowledge* points to  $I$

- Sometimes  $P(\mathbf{y} | I)$  can't be calculated exactly. Example

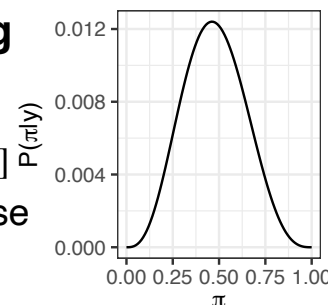
**Bernoulli data with non-conjugate prior:**

$$P(\pi) \propto \begin{cases} \pi(1-\pi)e^{-\pi^2} & \pi \in [0,1] \\ 0 & \text{otherwise} \end{cases}$$

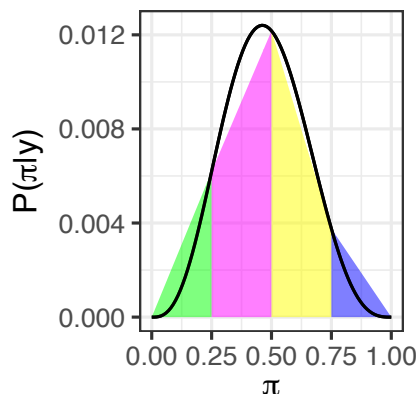


**Posterior after observing 2 heads, 2 tails:**

$$P(\pi) \propto \begin{cases} \pi^3(1-\pi)^3e^{-\pi^2} & \pi \in [0,1] \\ 0 & \text{otherwise} \end{cases}$$



- In simple cases like this, we can numerically approximate the integral:



- But in high dimension and/or unbounded ranges, difficult or even impossible!

# Unnormalizable posteriors

- Our motivation: Bayesian posterior inference

*Observed data*  $\rightarrow$

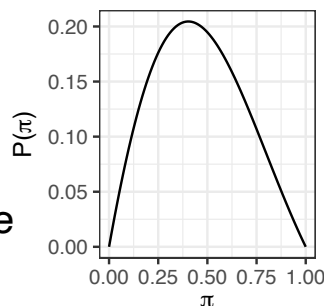
$$P(\theta | \mathbf{y}, I) = \frac{P(\mathbf{y} | \theta, I) P(\theta | I)}{P(\mathbf{y} | I)}$$

*Model parameters*  $\rightarrow$   $\theta$   $\rightarrow$   $\mathbf{y}$   $\rightarrow$   $I$   $\rightarrow$  *Background knowledge*

- Sometimes  $P(\mathbf{y} | I)$  can't be calculated exactly. Example

**Bernoulli data with non-conjugate prior:**

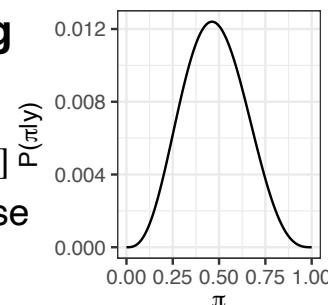
$$P(\pi) \propto \begin{cases} \pi(1-\pi)e^{-\pi^2} & \pi \in [0,1] \\ 0 & \text{otherwise} \end{cases}$$



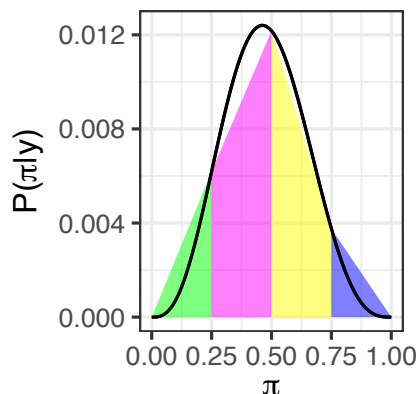
**Posterior after observing 2 heads, 2 tails:**

$$P(\pi) \propto \begin{cases} \pi^3(1-\pi)^3e^{-\pi^2} & \pi \in [0,1] \\ 0 & \text{otherwise} \end{cases}$$

*No closed form!*



- In simple cases like this, we can numerically approximate the integral:

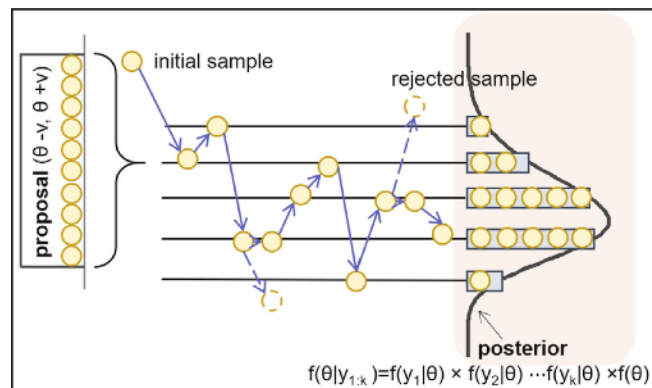


- But in high dimension and/or unbounded ranges, difficult or even impossible!



# Markov chain Monte Carlo

- However, we can often **take samples from the posterior** even when we can't compute normalized probabilities
- One general and widely used approach: **Markov chain Monte Carlo (MCMC)**
- MCMC is a mathematically principled random walk on a non-negative function, directed toward regions where the function takes on a larger value

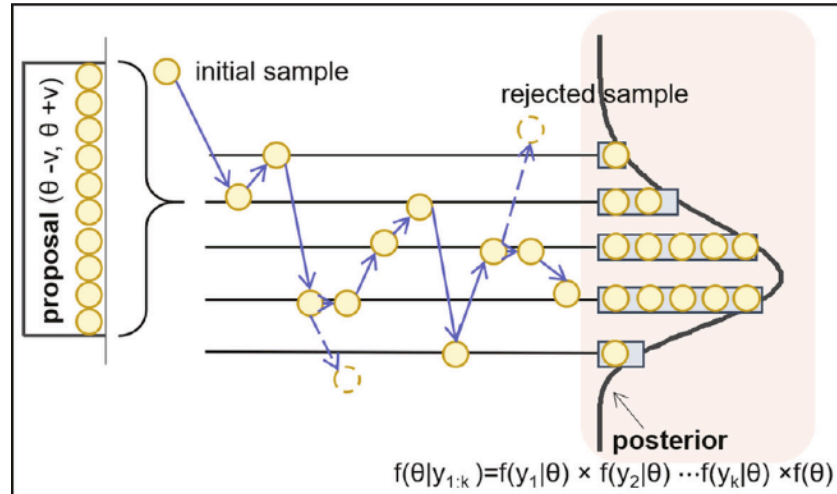


- Asymptotically, the random walk gives us samples from in proportion to the height of the function

# MCMC for posterior sampling

- We use the *unnormalized* form of the posterior:

$$P(\theta | \mathbf{y}, I) \propto P(\mathbf{y} | \theta, I)P(\theta | I)$$

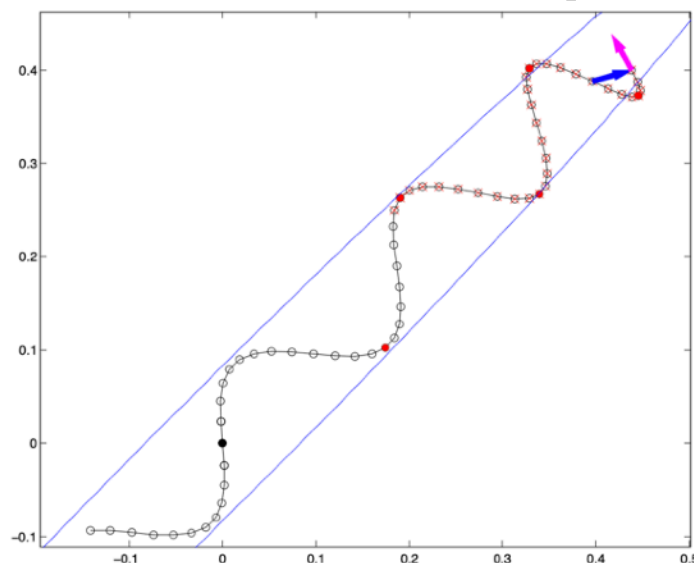


- We run MCMC and then treat the chain of values as samples from the posterior
- The full set of samples is not iid (nearby values on the chain are correlated), but methods exist for estimating "effectively" how many independent samples we have

# Stan, HMC, and NUTS

---

- There are many different MCMC algorithms (e.g., Metropolis, Gibbs Sampling)
- We will use the probabilistic programming language **Stan** for Bayesian inference about model parameters
- Stan uses an algorithm called **Hamiltonian Monte Carlo (HMC)** with the **No U-Turn Sampler (NUTS)**



*(Hoffman & Gelman, 2014)*

- This algorithm tends to be particularly efficient for many problems we'll face

# Bayesian posterior inference with Stan

# Bayesian posterior inference with Stan

---

1. Define the generative model you assume underlies the data you want to analyze

# Bayesian posterior inference with Stan

---

1. Define the generative model you assume underlies the data you want to analyze
2. Choose a prior distribution for your model parameters  $\theta$

# Bayesian posterior inference with Stan

---

1. Define the generative model you assume underlies the data you want to analyze
2. Choose a prior distribution for your model parameters  $\theta$
3. Encode the model structure and prior in a Stan program

# Bayesian posterior inference with Stan

---

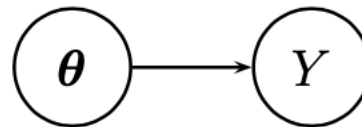
1. Define the generative model you assume underlies the data you want to analyze
2. Choose a prior distribution for your model parameters  $\theta$
3. Encode the model structure and prior in a Stan program
4. Provide the data  $Y$  you want to analyze, and ask Stan to sample from the posterior  $P(\theta | Y, I) \propto P(Y | \theta)P(\theta | I)$  (often written as  $P(\theta | Y) \propto P(Y | \theta)P(\theta)$ , i.e. eliding  $I$ )



# Bayesian posterior inference with Stan

---

1. Define the generative model you assume underlies the data you want to analyze
2. Choose a prior distribution for your model parameters  $\theta$
3. Encode the model structure and prior in a Stan program
4. Provide the data  $Y$  you want to analyze, and ask Stan to sample from the posterior  $P(\theta | Y, I) \propto P(Y | \theta)P(\theta | I)$  (often written as  $P(\theta | Y) \propto P(Y | \theta)P(\theta)$ , i.e. eliding  $I$ )



# Bayesian posterior inference with Stan

---

1. Define the generative model you assume underlies the data you want to analyze
2. Choose a prior distribution for your model parameters  $\theta$
3. Encode the model structure and prior in a Stan program
4. Provide the data  $Y$  you want to analyze, and ask Stan to sample from the posterior  $P(\theta | Y, I) \propto P(Y | \theta)P(\theta | I)$  (often written as  $P(\theta | Y) \propto P(Y | \theta)P(\theta)$ , i.e. eliding  $I$ )

