

Bayesian approaches to parameter estimation

CSCS 530 - Marisa Eisenberg

Bayesian approaches to parameter estimation

- Bayes' Theorem, rewritten for inference problems:

$$P(p | z) = P(\text{params} | \text{data}) = \frac{P(z | p) \cdot P(p)}{P(z)}$$

- Allows one to account for prior information about the parameters
 - E.g. previous studies in a similar population
- Update parameter information based on new data

Bayesian approaches to parameter estimation

Likelihood

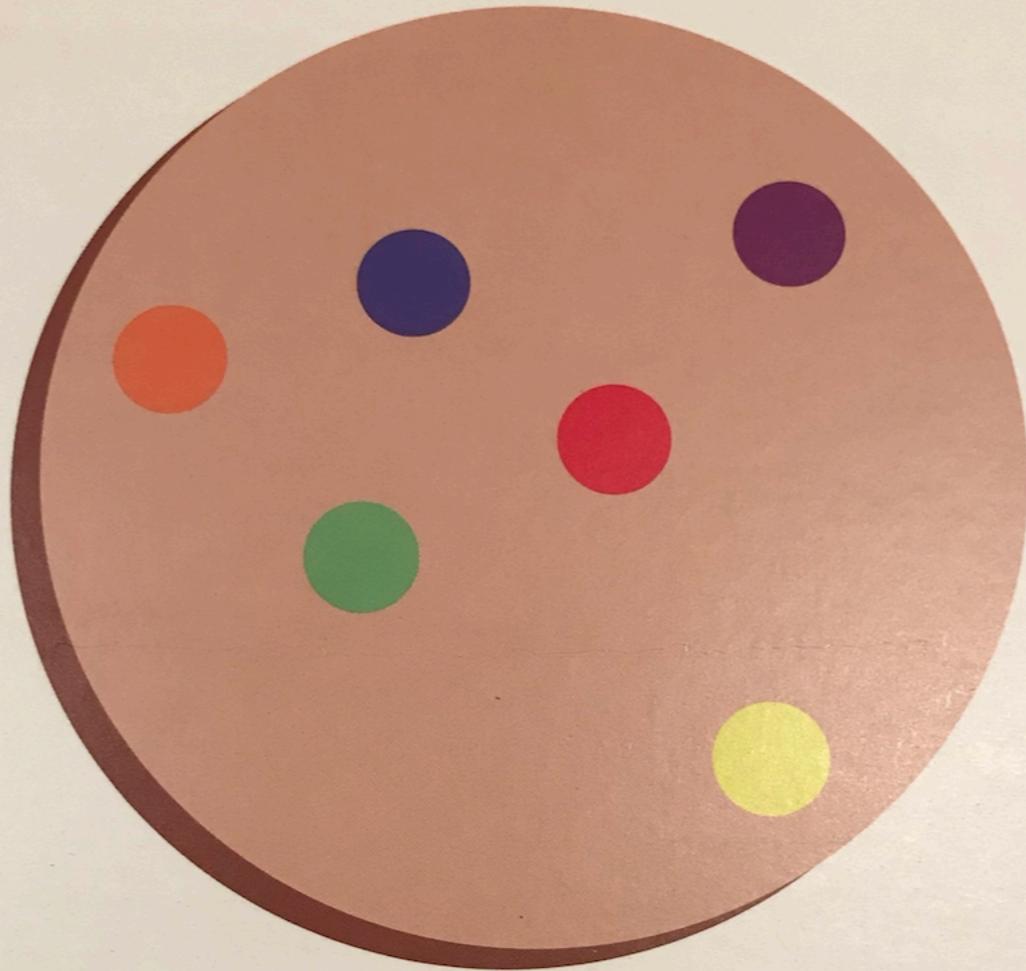
Prior distribution

$$P(p | z) = P(\text{params} | \text{data}) = \frac{P(z | p) \cdot P(p)}{P(z)}$$

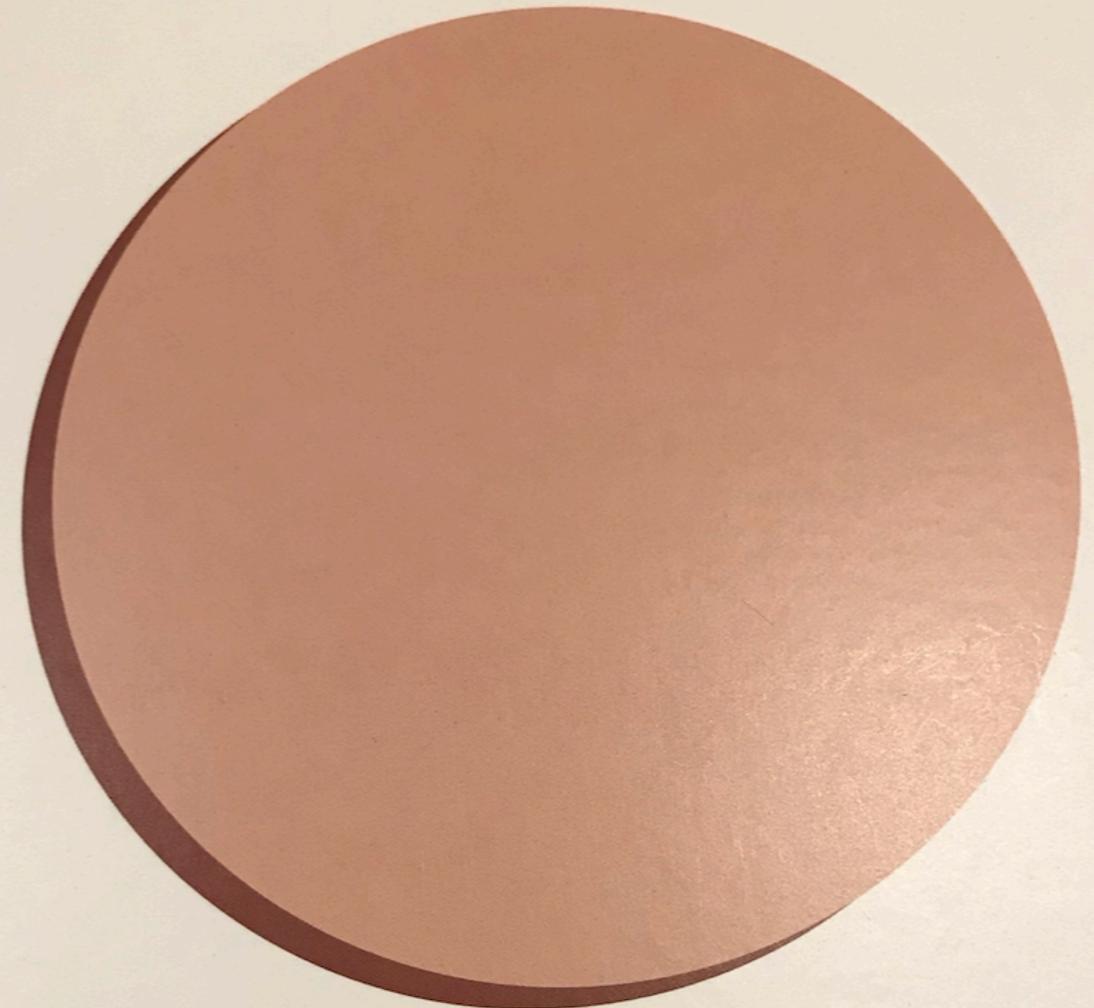
Normalizing constant
(can be difficult to calculate!)

$$P(z) = \int_p P(z, p) dp$$

Bayesian probability for babies!



Some cookies have candy.



Some don't.



Take a bite. It has no candy.



Did it come from a candy cookie?

What is the likelihood?

$$L(\text{NC cookie} \mid \text{NC bite}) = P(\text{NC bite} \mid \text{NC cookie})$$



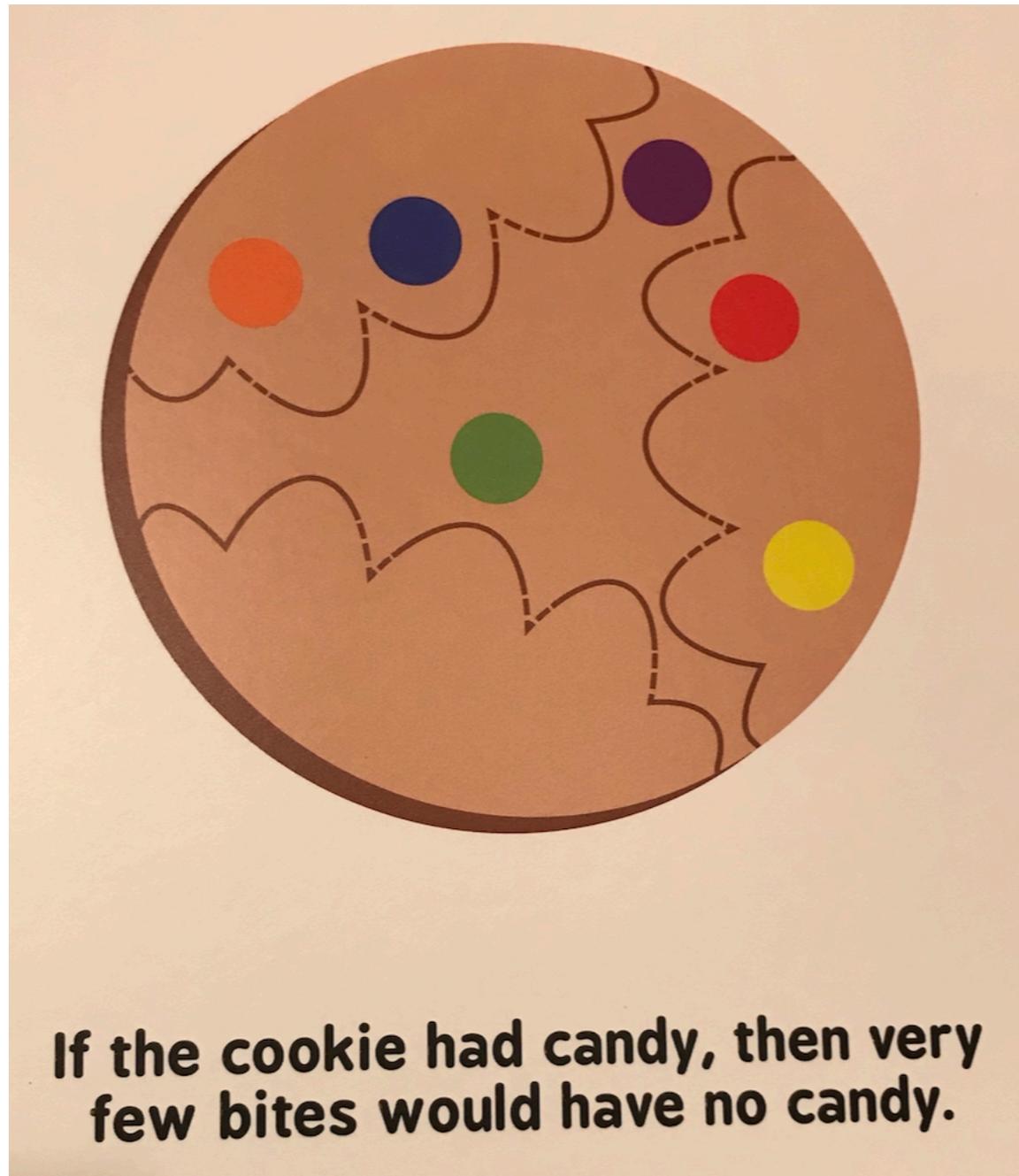
**If the cookie had no candy,
then every bite would have no candy.**

$$\text{Pr}(\text{no-candy bite} \mid \text{no-candy cookie}) = 1$$

**The probability of a no-candy bite,
given a no-candy cookie, is 1.**

What is the likelihood?

$$L(\text{C cookie} \mid \text{NC bite}) = P(\text{NC bite} \mid \text{C cookie})$$

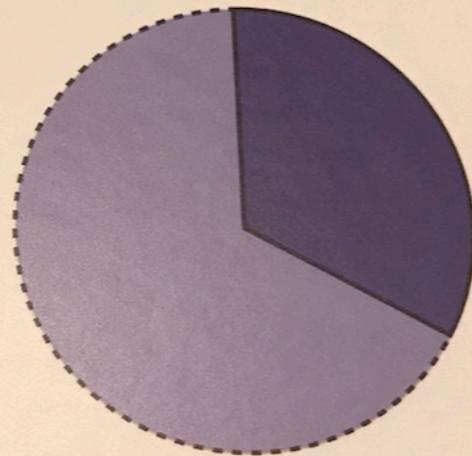
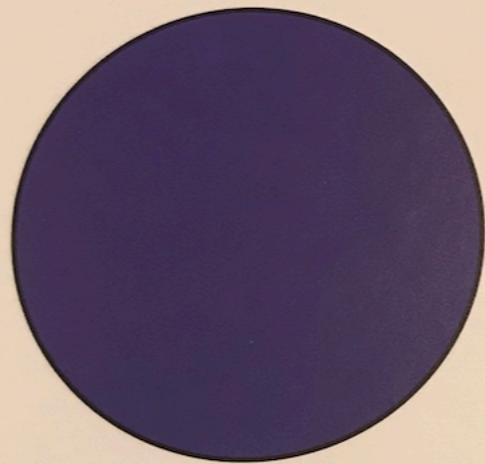


$\text{Pr}(\text{no candy bite} \mid \text{candy cookie})$
=
 $1/3$

The probability of a no-candy bite, given a candy cookie, is 1/3.

What is the maximum likelihood estimate?

$$\Pr(\text{no-candy bite} \mid \text{no-candy cookie}) > \Pr(\text{no-candy bite} \mid \text{candy cookie})$$



1 is greater than 1/3.

Maximum likelihood:



So the no-candy bite probably came from a no-candy cookie!

What about the prior distribution of cookies?



But what if we knew there were 10 cookies,



and all had candy but one?

Our data (likelihood) tells us we have a no-candy bite—how many of the bites are no candy?

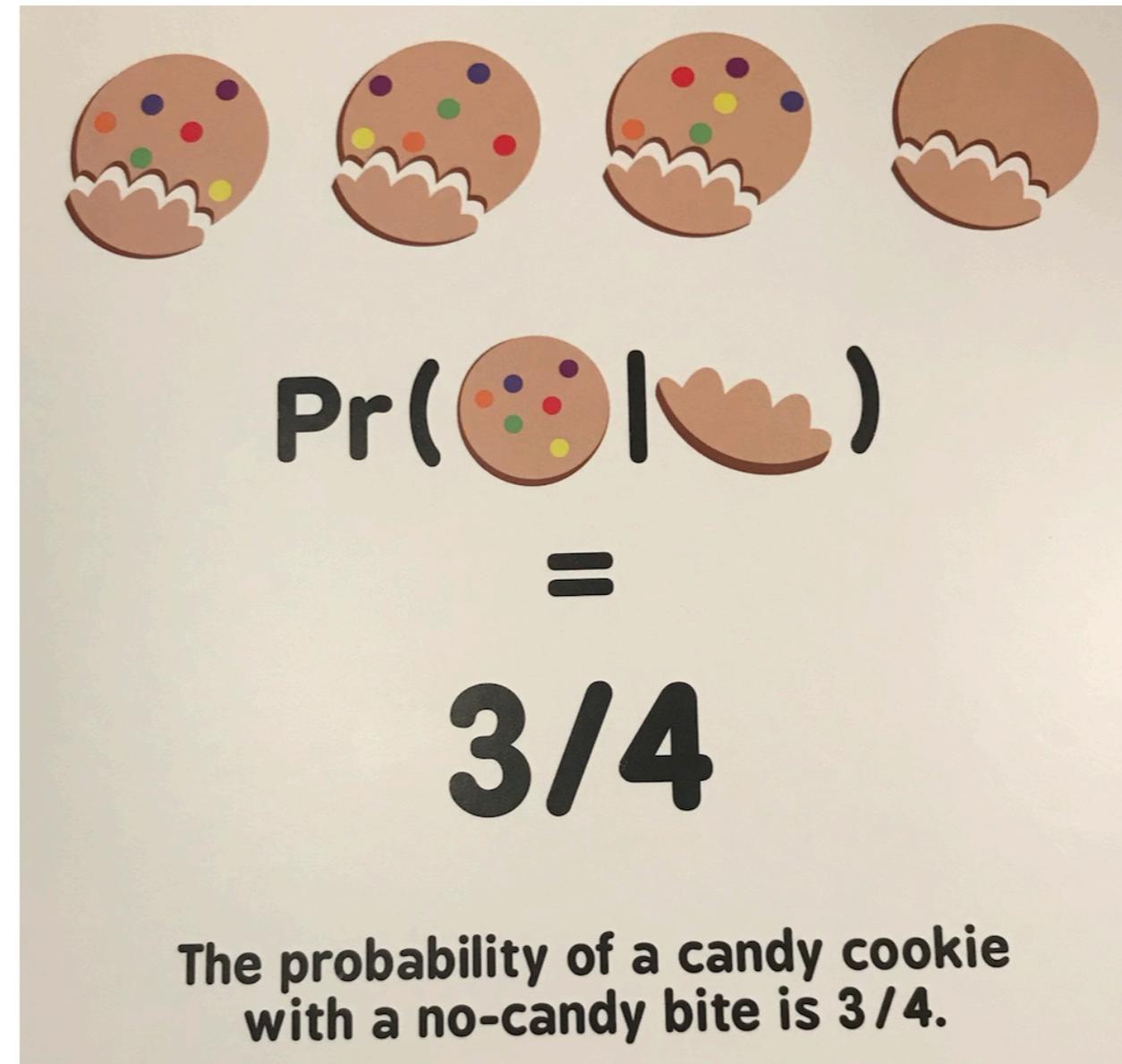
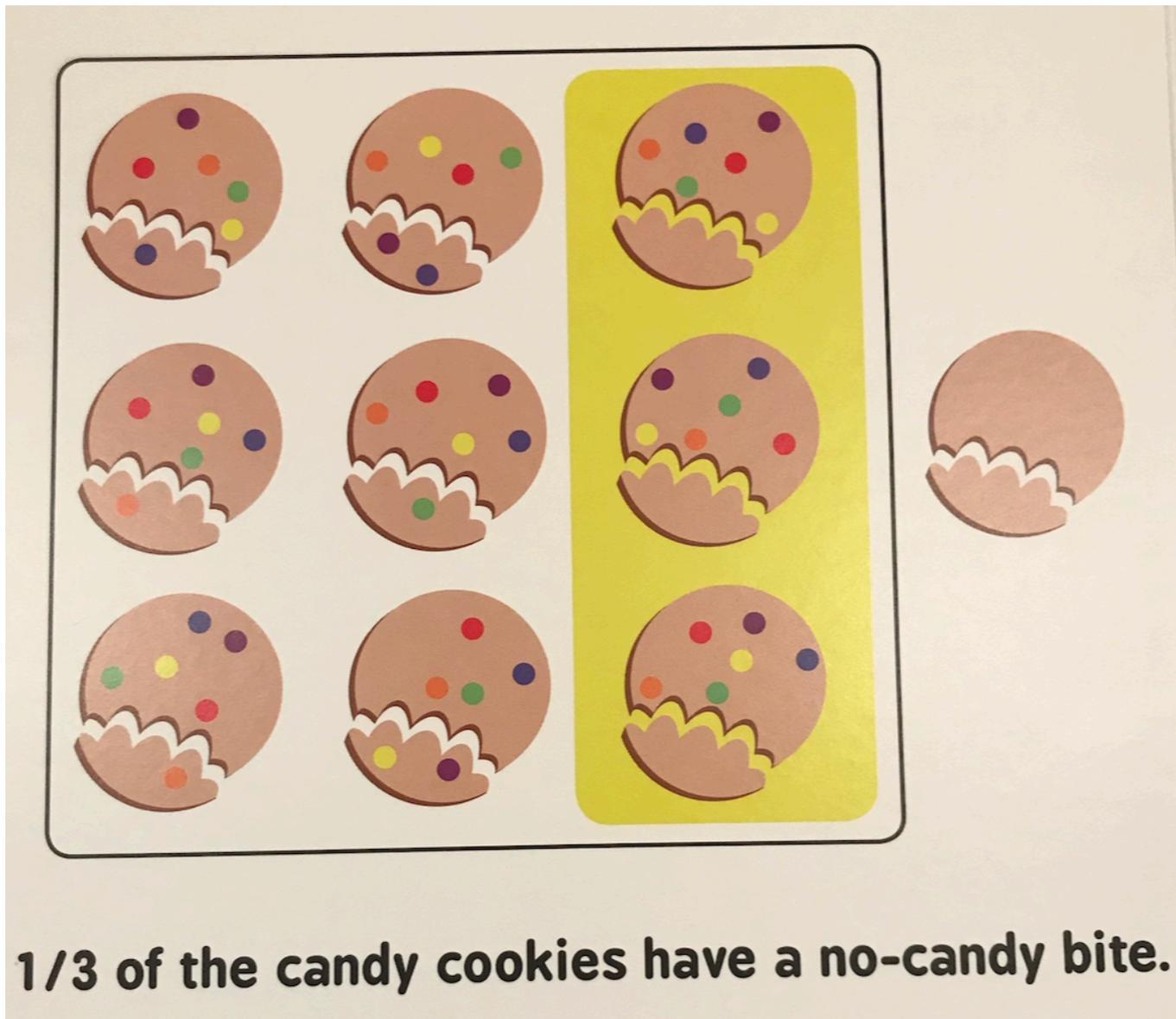


Take a bite of each.



**There are 4 no-candy bites.
3 bites are from candy cookies.
1 bite is from a no-candy cookie.**

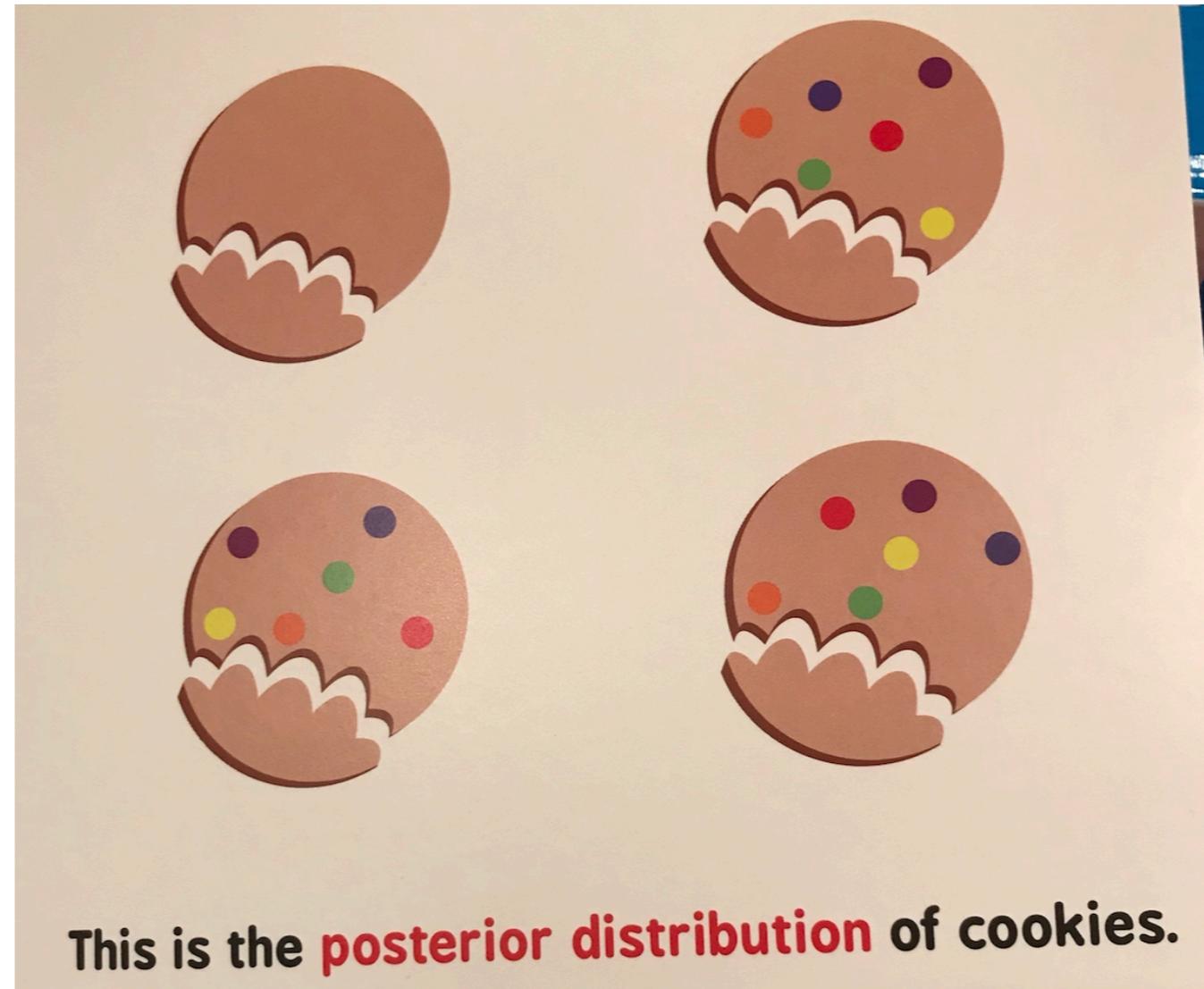
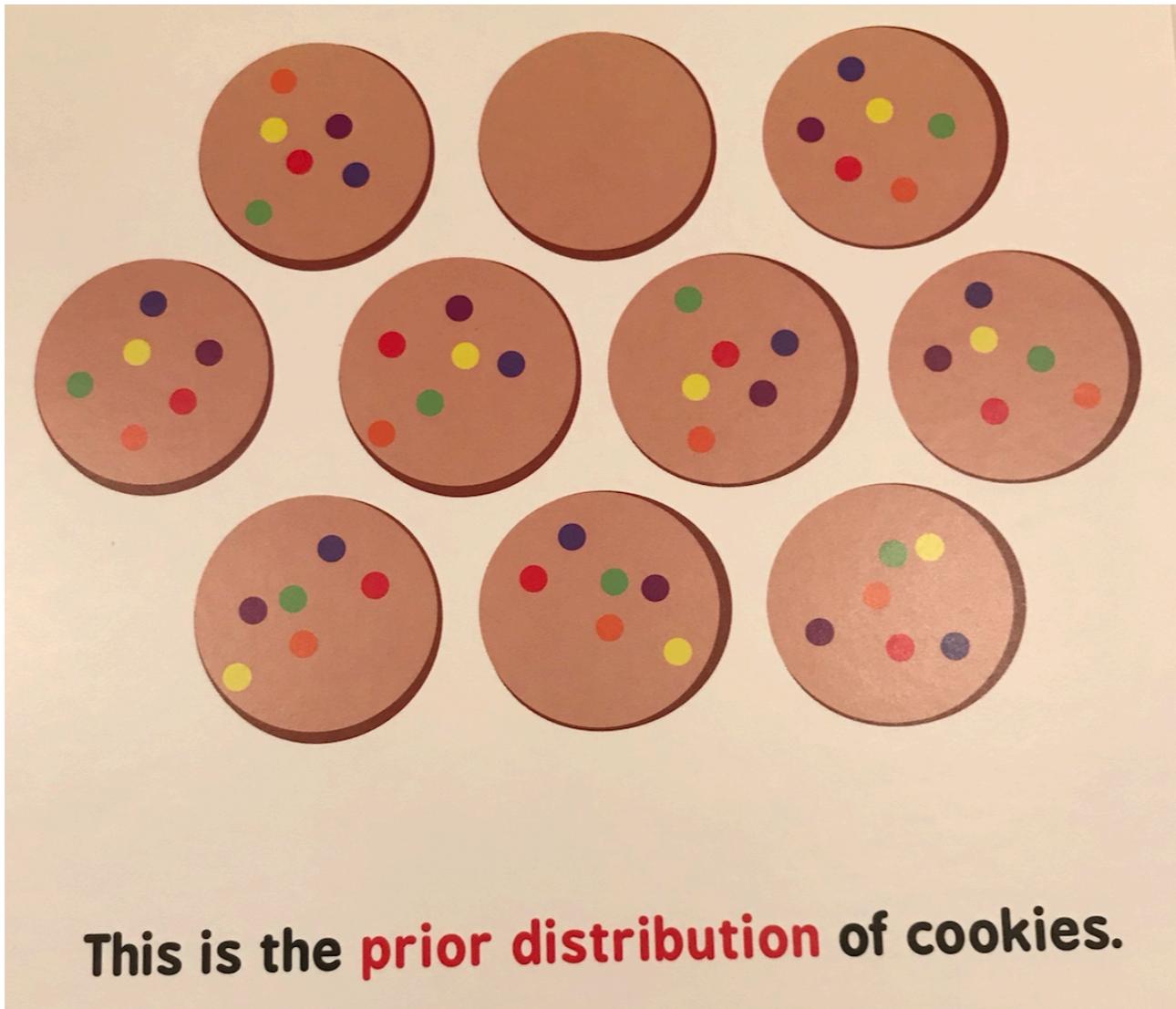
1/3 of the candy cookie bites have no candy, but there are a lot more of them



Prior x Likelihood ~ Posterior

$9 \times 1/3 = 3$ candy cookies, vs. $1 \times 1 = 1$ no-candy cookie

Bayesian estimation!



Bayesian approaches to parameter estimation

Likelihood

Prior distribution

$$P(p | z) = P(\text{params} | \text{data}) = \frac{P(z | p) \cdot P(p)}{P(z)}$$

Normalizing constant
(can be difficult to calculate!)

$$P(z) = \int_p P(z, p) dp$$

Denominator term - $P(z)$

- The denominator term:

$$P(z) = \int_p P(z, p) dp$$

- Probability of seeing the data z from the model, over all parameter space
- Often doesn't have a closed form solution—evaluating numerically can also be difficult
 - E.g. if p is a three dimensional, then if we took 1000 grid points in each direction, the grid representing the function to be integrated has $1000^3 = 10^9$ points

Maximum *a posteriori* (MAP) estimation

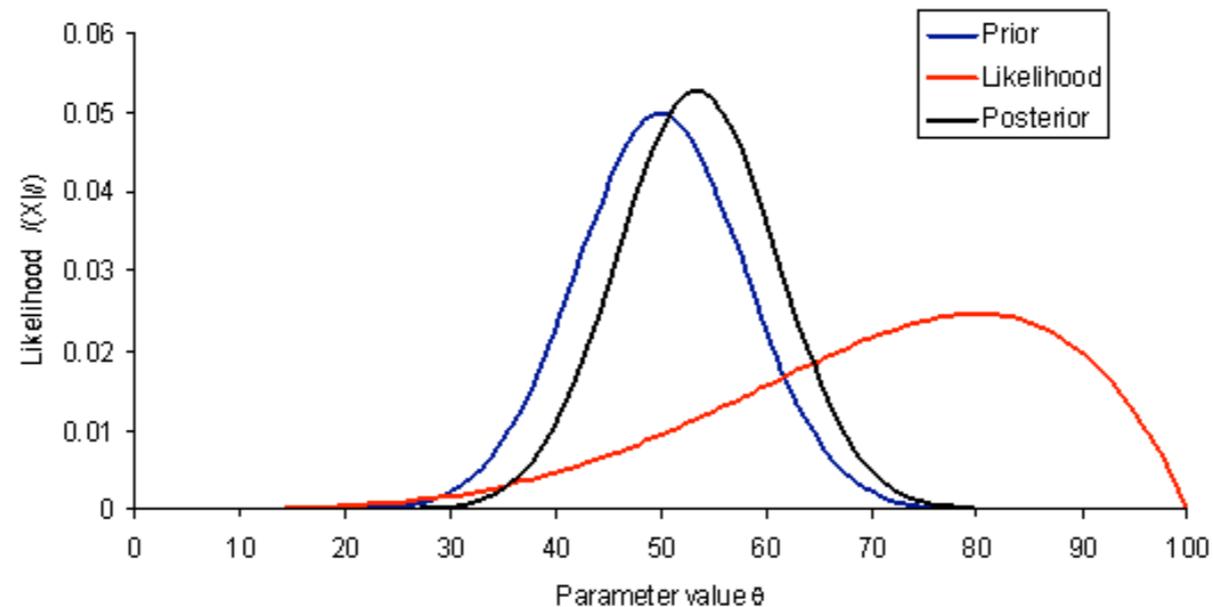
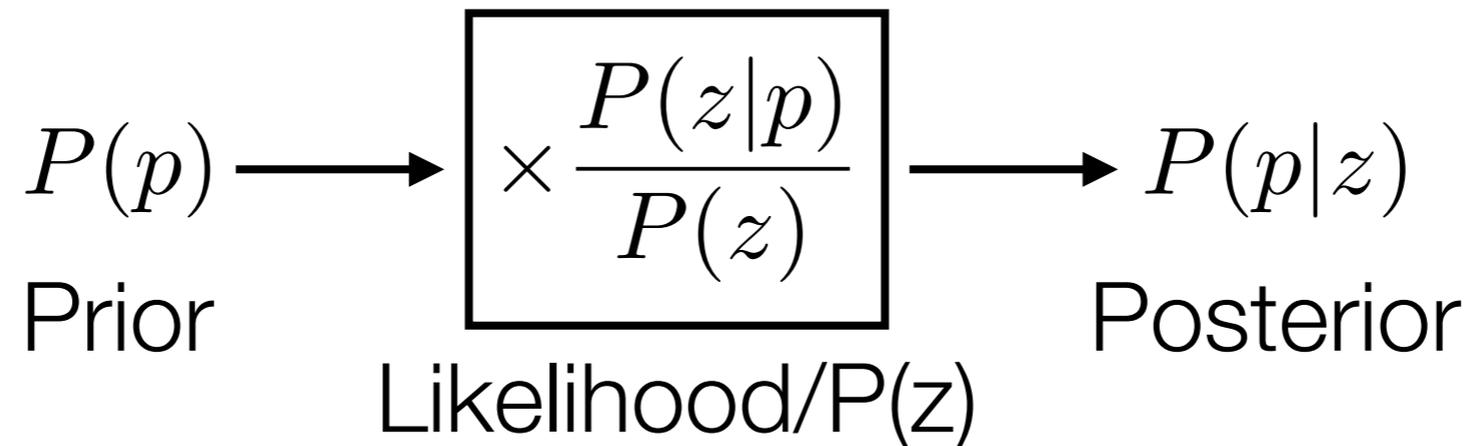
- Instead of working with the full term, just use the numerator:

$$P(p|z) = \frac{P(z|p) \cdot P(p)}{P(z)}$$

- The denominator is a constant, so the numerator is proportional to the posterior we are trying to estimate
- Then the \mathbf{p} which yields $\max(P(z|p) \cdot P(p))$ is the same \mathbf{p} that maximizes $P(p|z)$
- If we only need a point estimate, MAP gets around needing to estimate $P(z)$

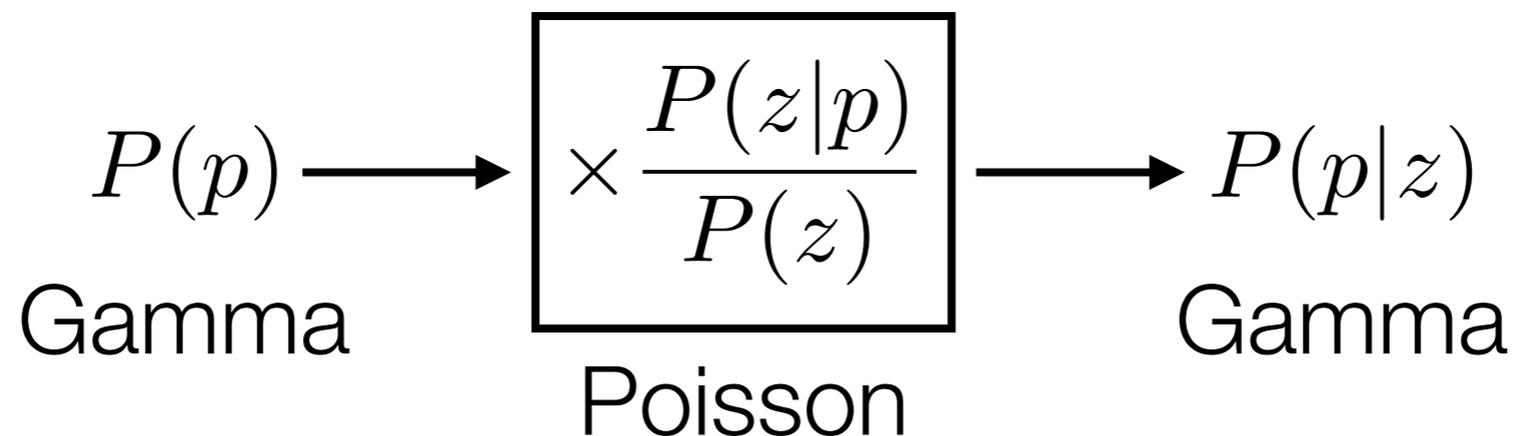
Bayesian Parameter Estimation

- Can think of Bayesian estimation as a map, where we update the prior to a new posterior based on data



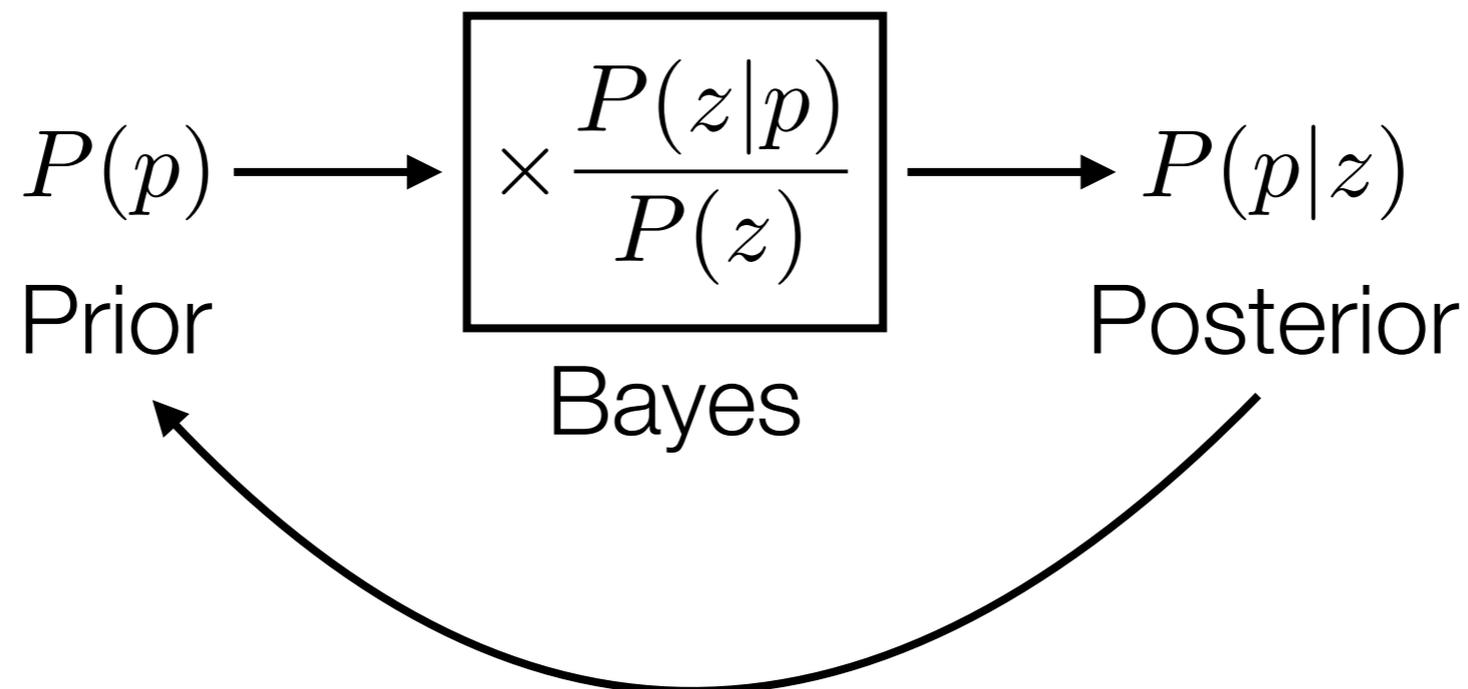
Conjugate Priors

- For a likelihood distribution, there may be a distribution family for our prior, which makes the posterior and prior come from the same type of distribution
- This is called a **conjugate prior** for that likelihood
- For example, a gamma distribution is the conjugate prior for a Poisson likelihood.



Why conjugate priors?

- If we have a conjugate prior, we can calculate the posterior directly from the likelihood and the prior— handles the issue with calculating the denominator $P(z)$
- Also makes it easier to repeat Bayesian estimation— making the posterior the prior and updating as new data comes in



Conjugate prior example: coin flip

- Let z be the data—i.e. the coin flip outcome, $z = 1$ if it's heads, $z = 0$ if it's tails
- Let θ be the probability the coin shows heads
- Likelihood: Bernoulli distribution

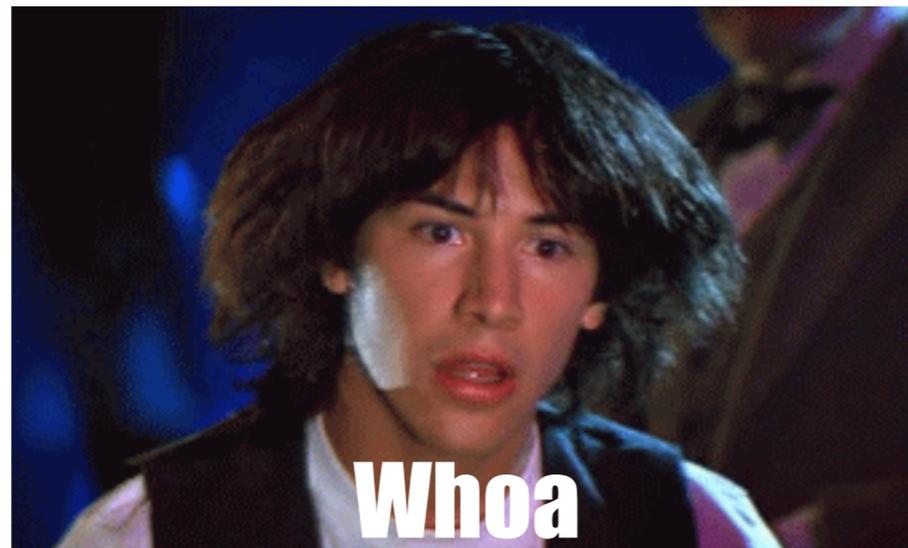
$$P(z|\theta) = \theta^z (1 - \theta)^{1-z}$$

Conjugate prior example: coin flip

- **Conjugate prior:** beta distribution

$$P(\theta|\alpha, \beta) = \frac{\theta^{\alpha-1} (1 - \theta)^{\beta-1}}{\int_0^1 \theta^{\alpha-1} (1 - \theta)^{\beta-1} d\theta}$$

- α and β are **hyperparameters** - shape parameters that describe the distribution of the model parameters



How does the posterior work out to be a beta distribution as well?

$$\begin{aligned} P(\theta|z) &= \frac{P(z|\theta)P(\theta|\alpha, \beta)}{P(z)} \\ &= \frac{\theta^z (1 - \theta)^{1-z} \frac{\theta^{\alpha-1} (1-\theta)^{\beta-1}}{\int_0^1 \theta^{\alpha-1} (1-\theta)^{\beta-1} d\theta}}{P(z)} \\ &= \frac{\theta^z (1 - \theta)^{1-z} \frac{\theta^{\alpha-1} (1-\theta)^{\beta-1}}{\int_0^1 \theta^{\alpha-1} (1-\theta)^{\beta-1} d\theta}}{\int_0^1 P(z, \theta) d\theta} \\ &= \frac{\theta^z (1 - \theta)^{1-z} \frac{\theta^{\alpha-1} (1-\theta)^{\beta-1}}{\int_0^1 \theta^{\alpha-1} (1-\theta)^{\beta-1} d\theta}}{\int_0^1 \theta^z (1 - \theta)^{1-z} d\theta} \end{aligned}$$

Etc.—but you can see it will work out to be beta distributed

Coin flip example - Posterior

- **Beta distributed** with posterior hyperparameters:

$$\alpha_{post} = \alpha + z \qquad \beta_{post} = \beta + 1 - z$$

- If we take multiple data points, this works out to be:

$$\alpha_{post} = \alpha + \sum_{i=1}^n z_i \qquad \beta_{post} = \beta + n - \sum_{i=1}^n z_i$$

Sampling methods: approximating a distribution

- What if we want priors that aren't conjugate? Or what if our likelihood is more complicated and it isn't clear what the conjugate prior is?
- Now we need some way to get the posterior, even though the denominator term is annoying
- How to approximate the distribution?

Markov Chain Monte Carlo (MCMC)

- Sampling-based methods—in particular, **Markov chain Monte Carlo (MCMC)**
- Also used for many other things! Can approximate distributions more generally—used in cryptography, calculating neutron diffusion, all sorts of things

Markov Chain Monte Carlo (MCMC)

- **MCMC** is a method for sampling from a distribution
- **Markov chain:** a type of (discrete) Markov process
 - Markov: memoryless, i.e. what happens at the next step only depends on the current step
- **Monte Carlo methods** are a class of algorithms that use sampling/randomness—often used to solve deterministic problems (such as approximating an integral)

Markov Chain Monte Carlo (MCMC)

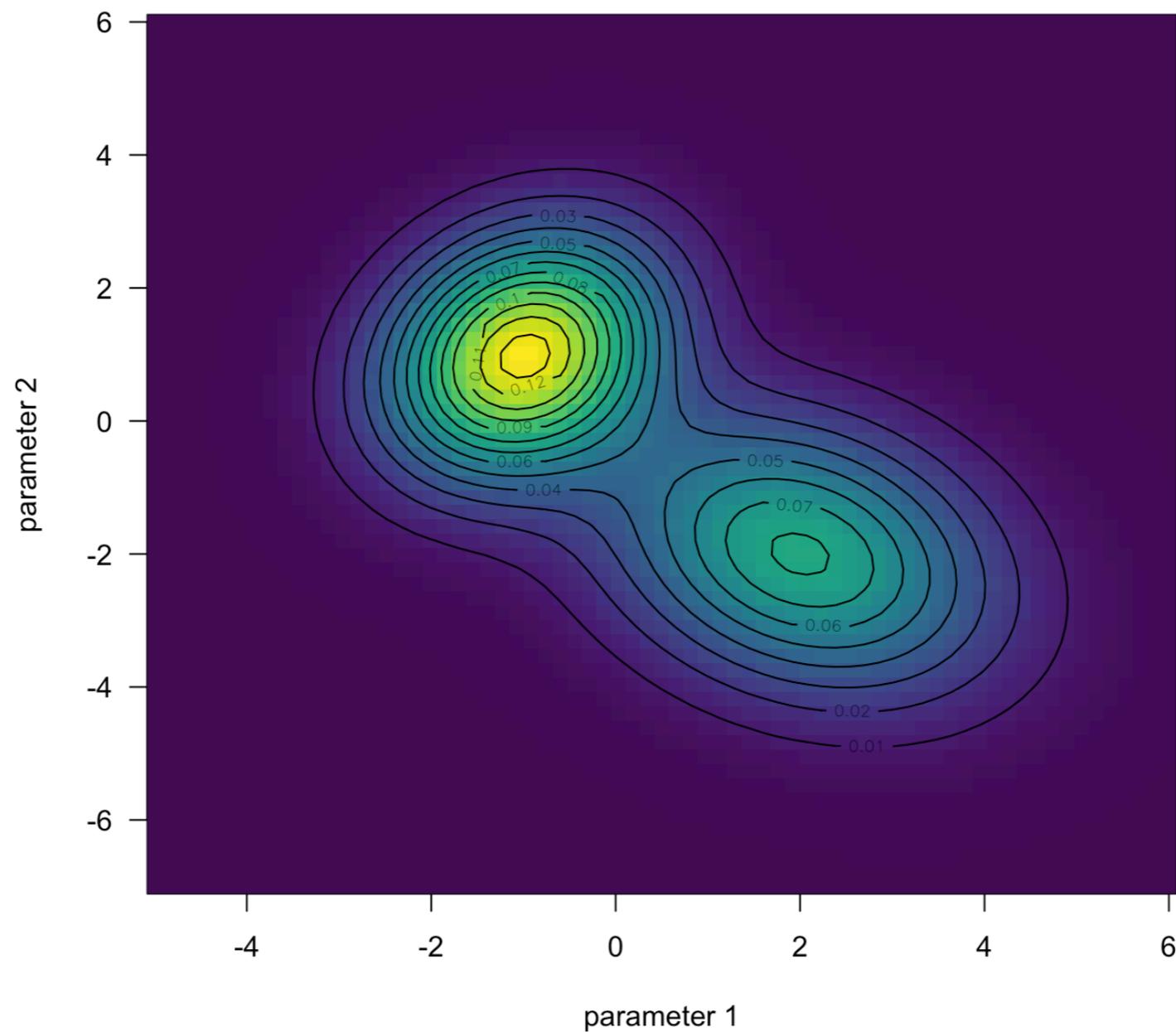
- **Main idea:** make a Markov chain that converges to the distribution we're trying to sample from—in this case, the posterior distribution!
- The Markov chain will have some transient dynamics (burn-in), and then reach an equilibrium distribution which is the one we're trying to approximate

Markov Chain Monte Carlo (MCMC)

- Many MCMC methods are based on random walks
 - Set up walk to spend more time in higher probability regions
- Typically don't need the actual distribution for this, just something proportional—so we can get the relative probability density at two points
 - So we don't need to calculate $P(z)$! We can just use the numerator

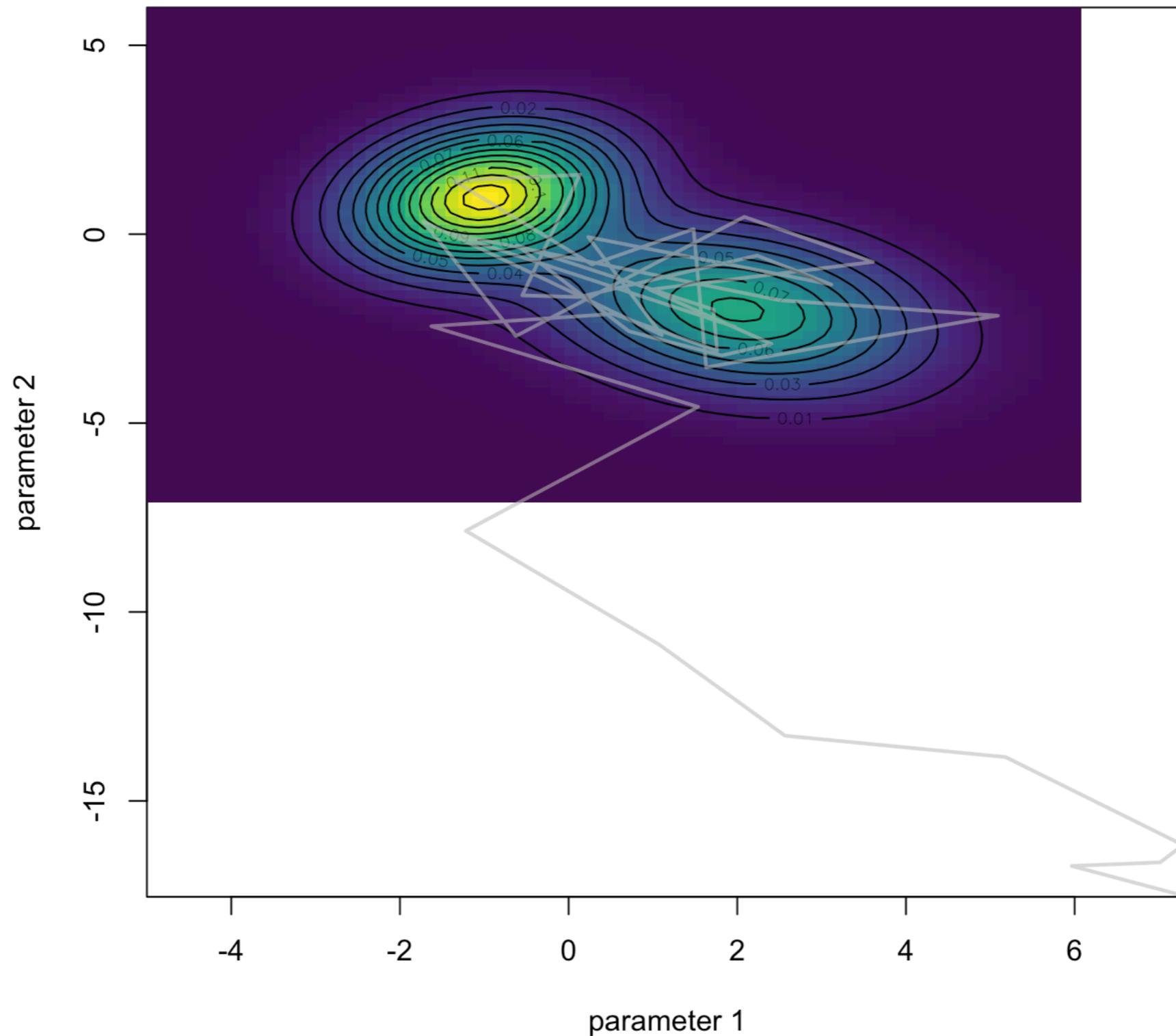
Example

- Suppose two parameters, with likelihood \times prior:

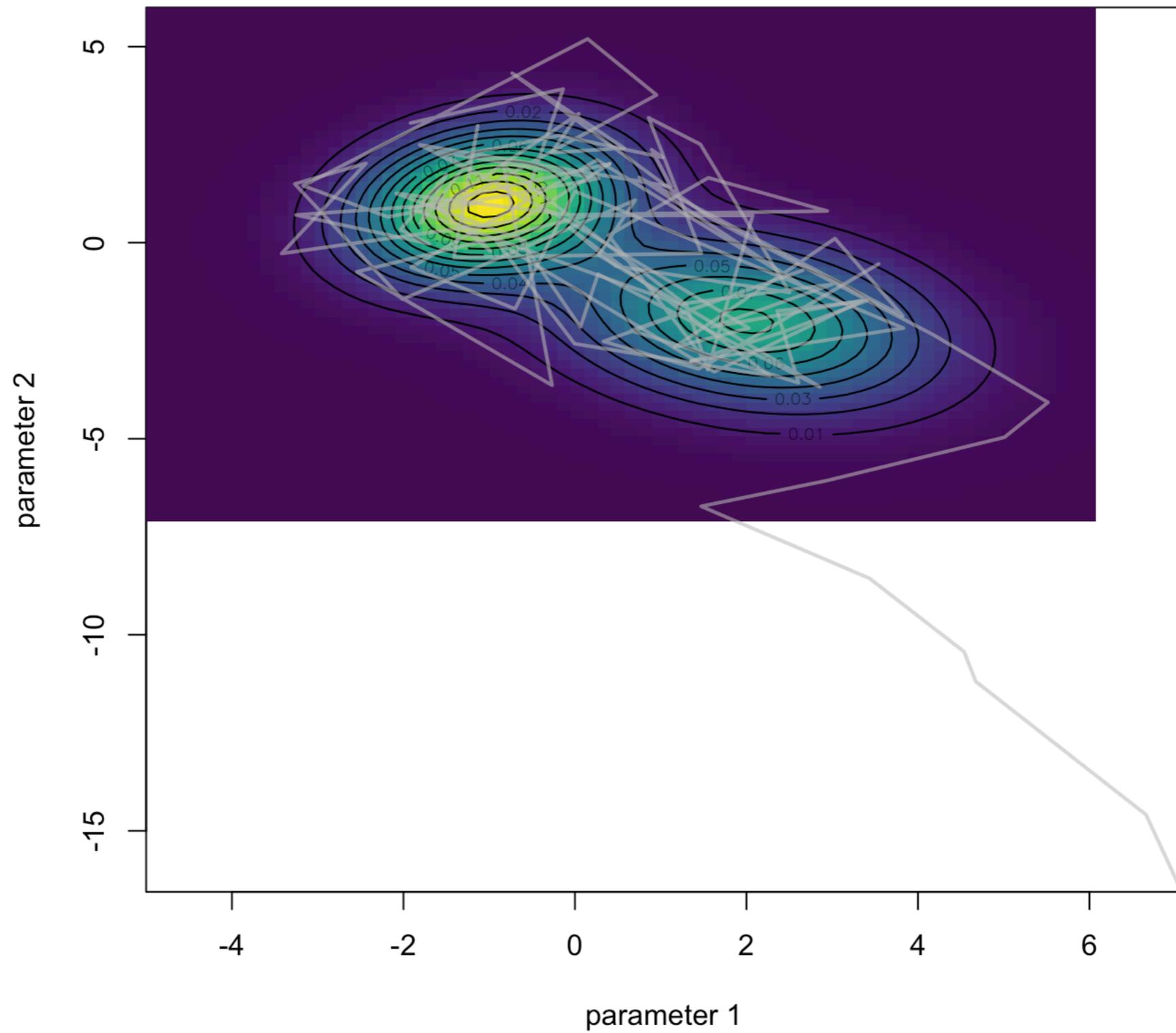


Sample path

We start our parameter values at a random guess
The random walk the MCMC traverses is shown as
the grey line

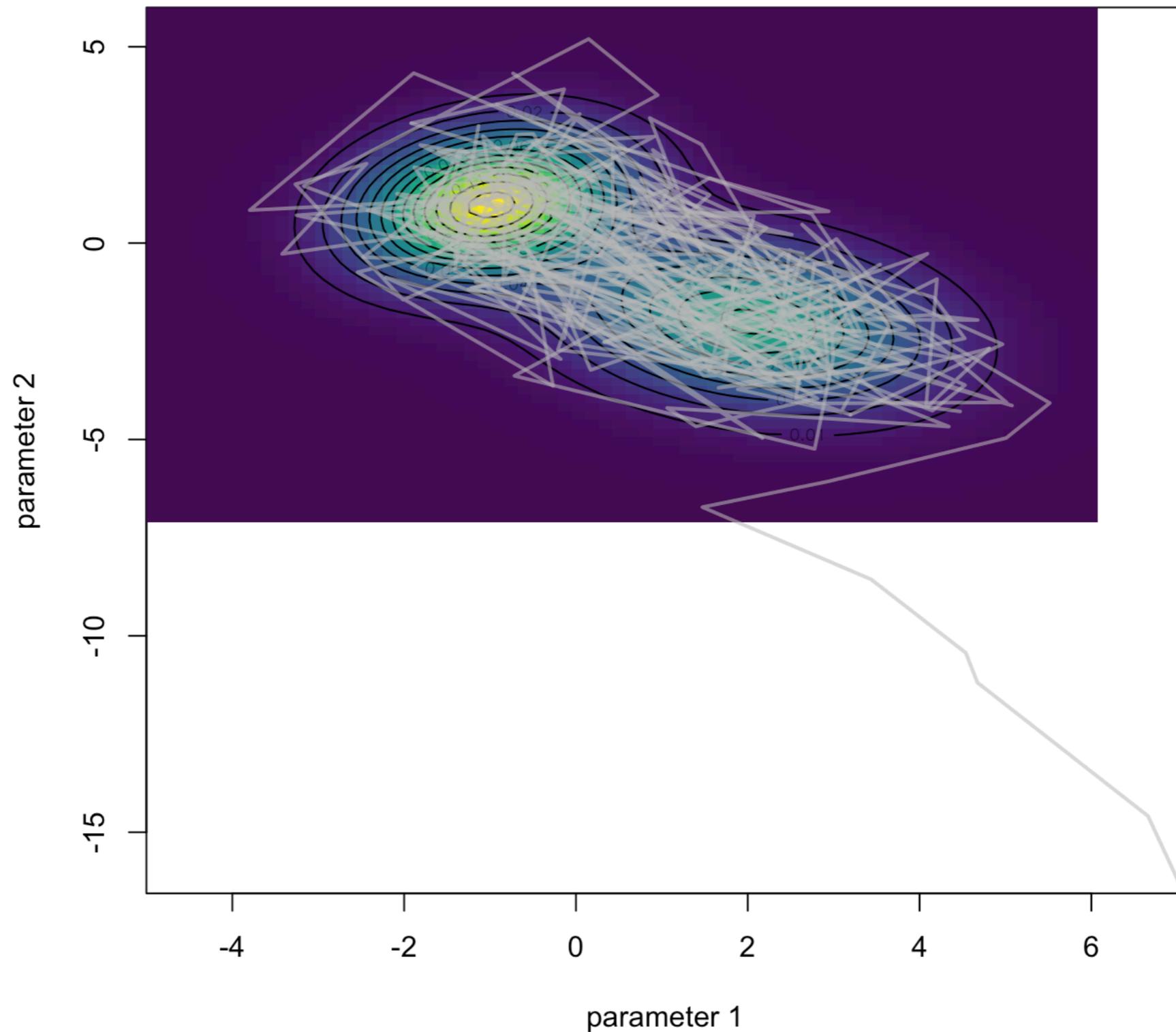


Sample path



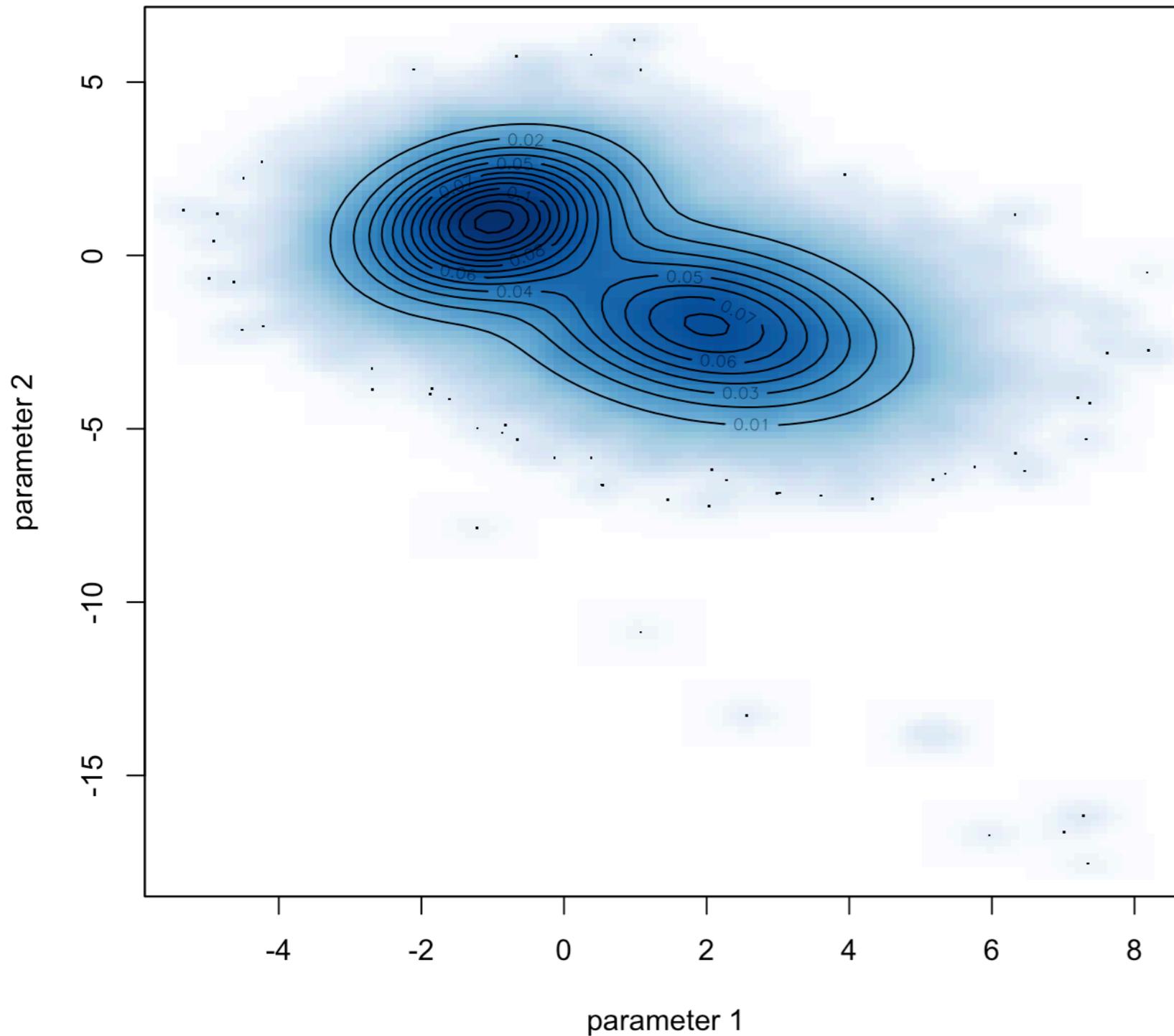
Sample path

Over time, the random walk accrues samples of the posterior distribution, proportional to the probability of those values. In other words, we get more samples from higher probability regions

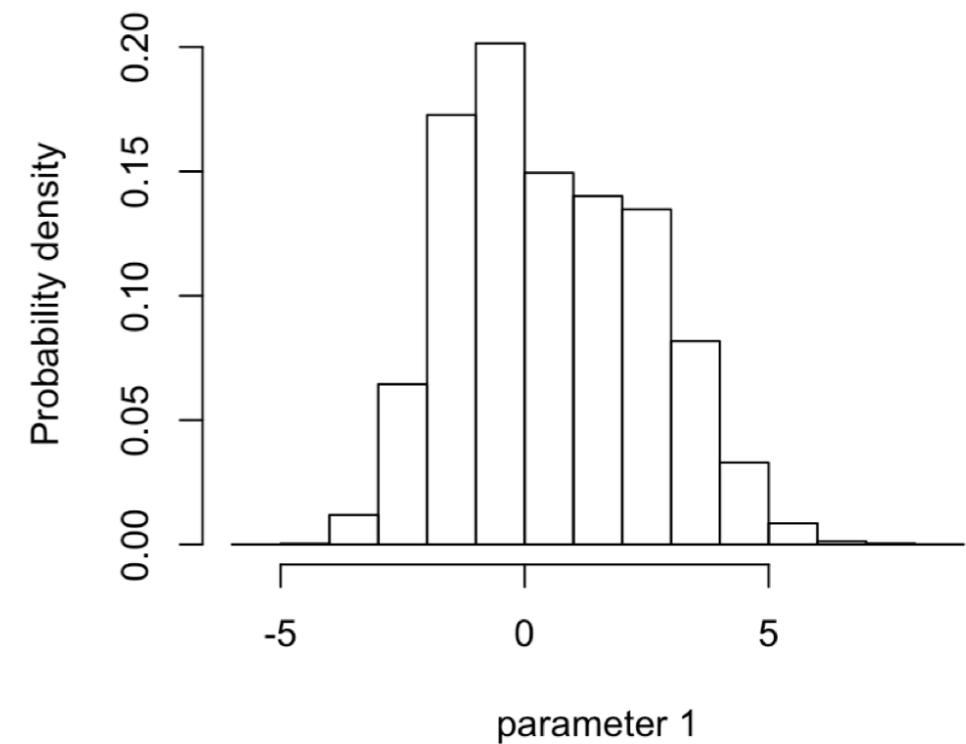


Sampled density

Eventually, the sampled values recreate the posterior distribution! And we didn't need the denominator term, only the numerator term, so these are relatively easy to calculate



Can also get marginals:



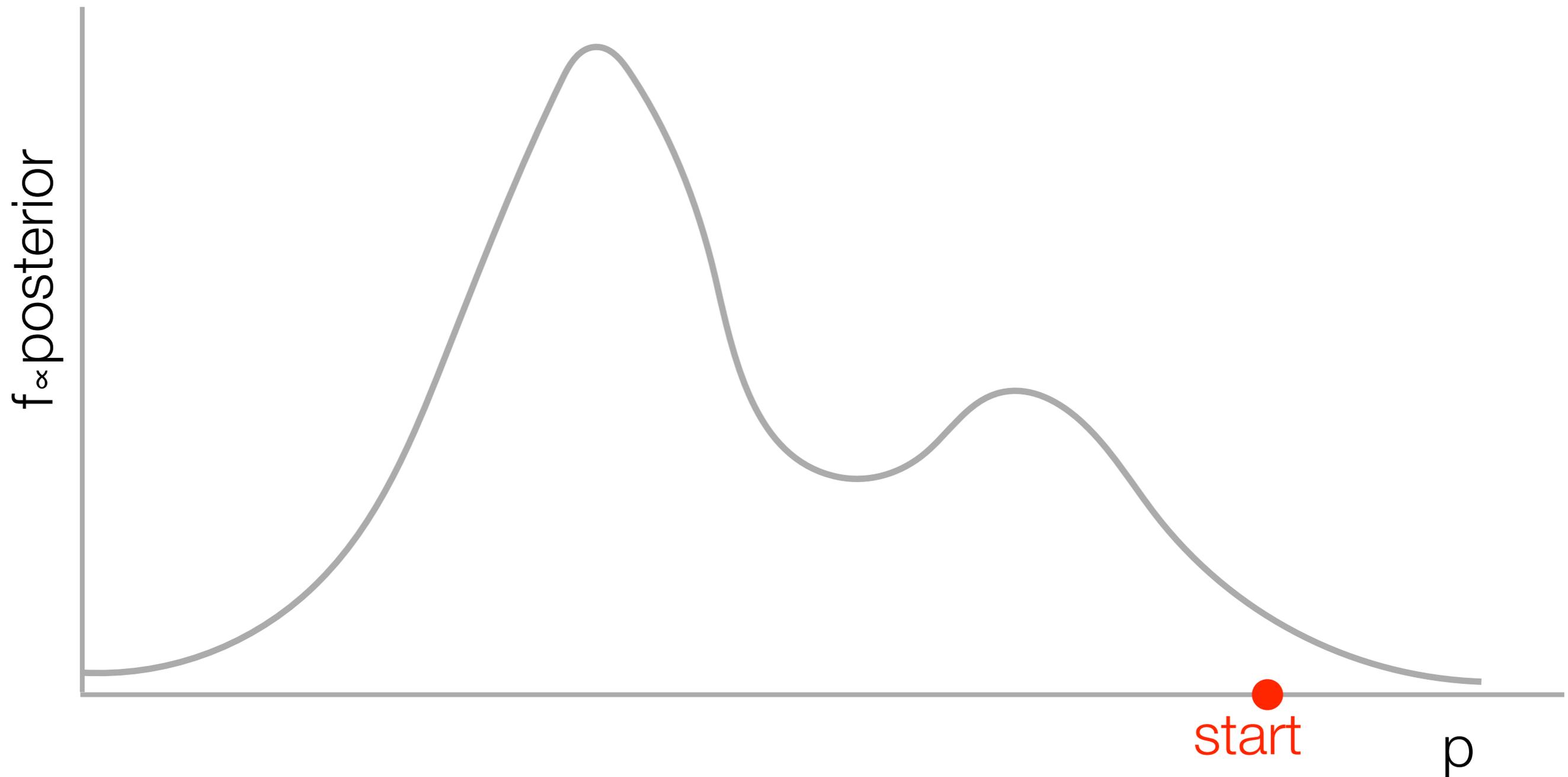
Example: Metropolis Algorithm

- Idea is to ‘walk’ randomly through parameter space, spending more time in places that are higher probability— that way, the overall distribution draws more from higher probability spots
- Setup— we need
 - A function $f(p)$ proportional to the distribution we want to sample, in our case $f(p) = P(z|p) \cdot P(p)$
 - A proposal distribution (how we choose the next point from the current one) - more on this in a minute

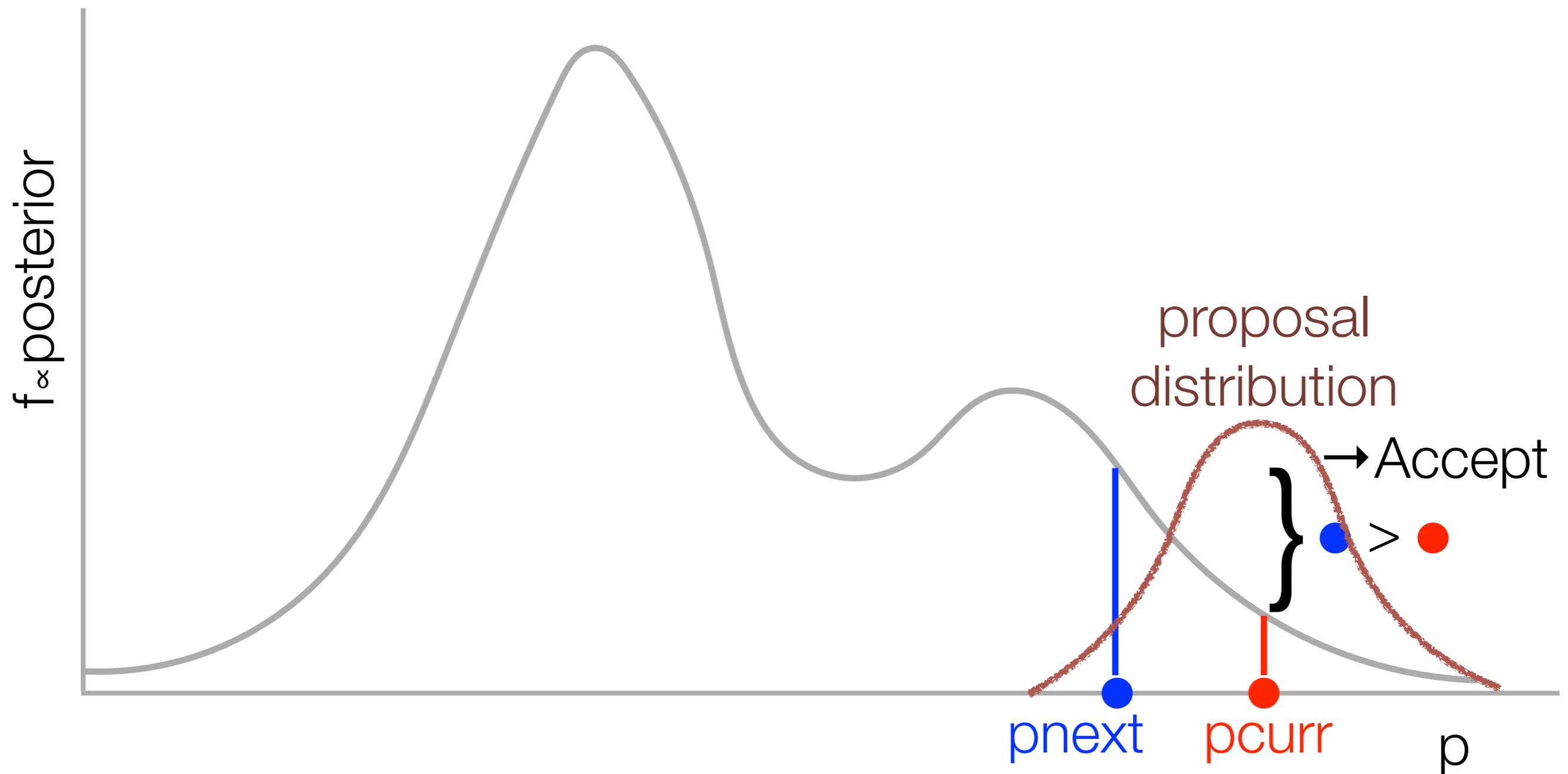
Metropolis Algorithm

- Start at some point in parameter space
- For each iteration
 - Propose a new random point p_{next} based on the current point p_{curr} (using the proposal distribution)
 - Calculate the **acceptance ratio**, $\alpha = f(p_{next})/f(p_{curr})$
 - If $\alpha \geq 1$, the new point is as good or better—accept
 - If $\alpha < 1$, accept with probability α

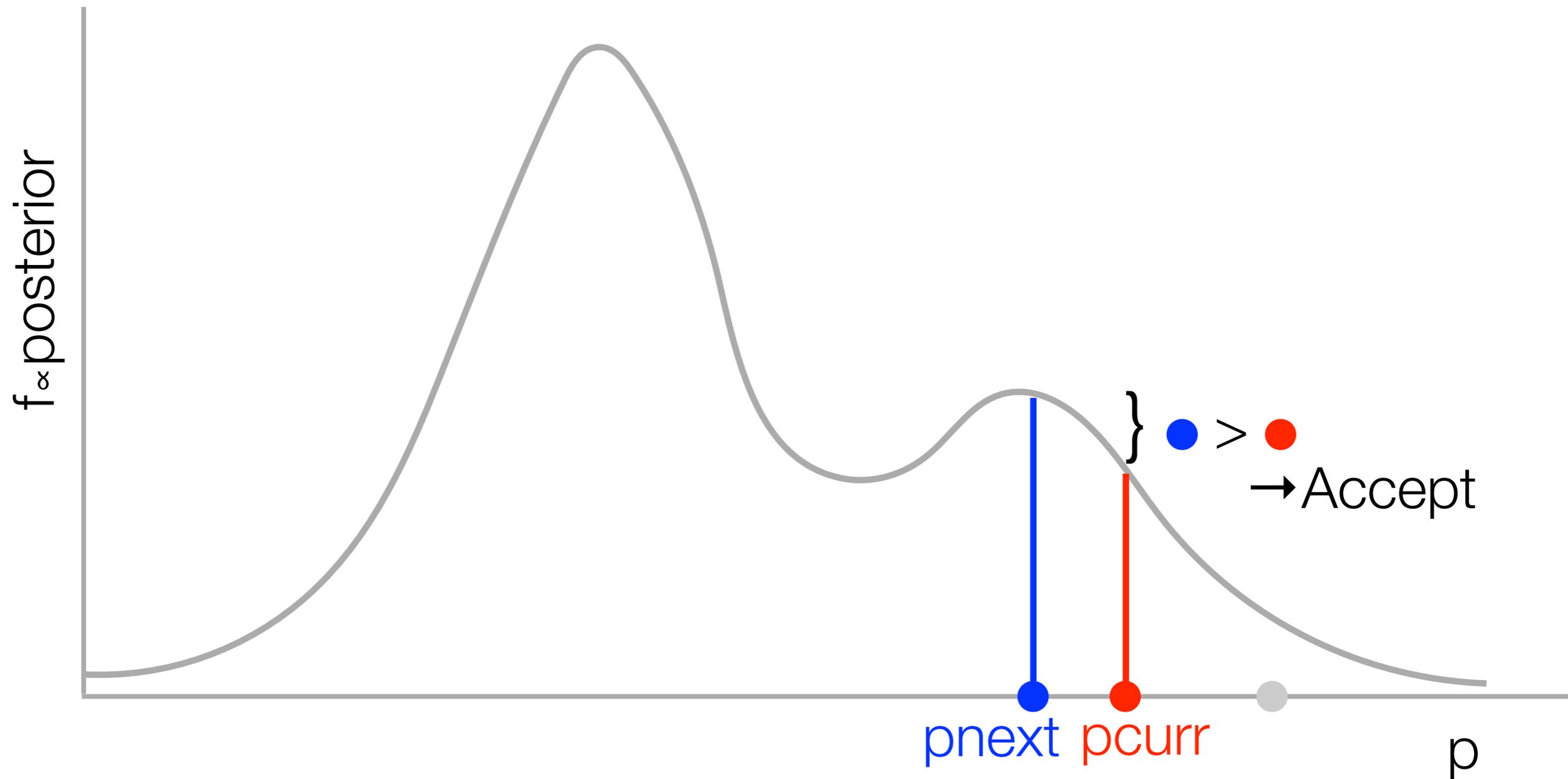
What does the metropolis algorithm do?



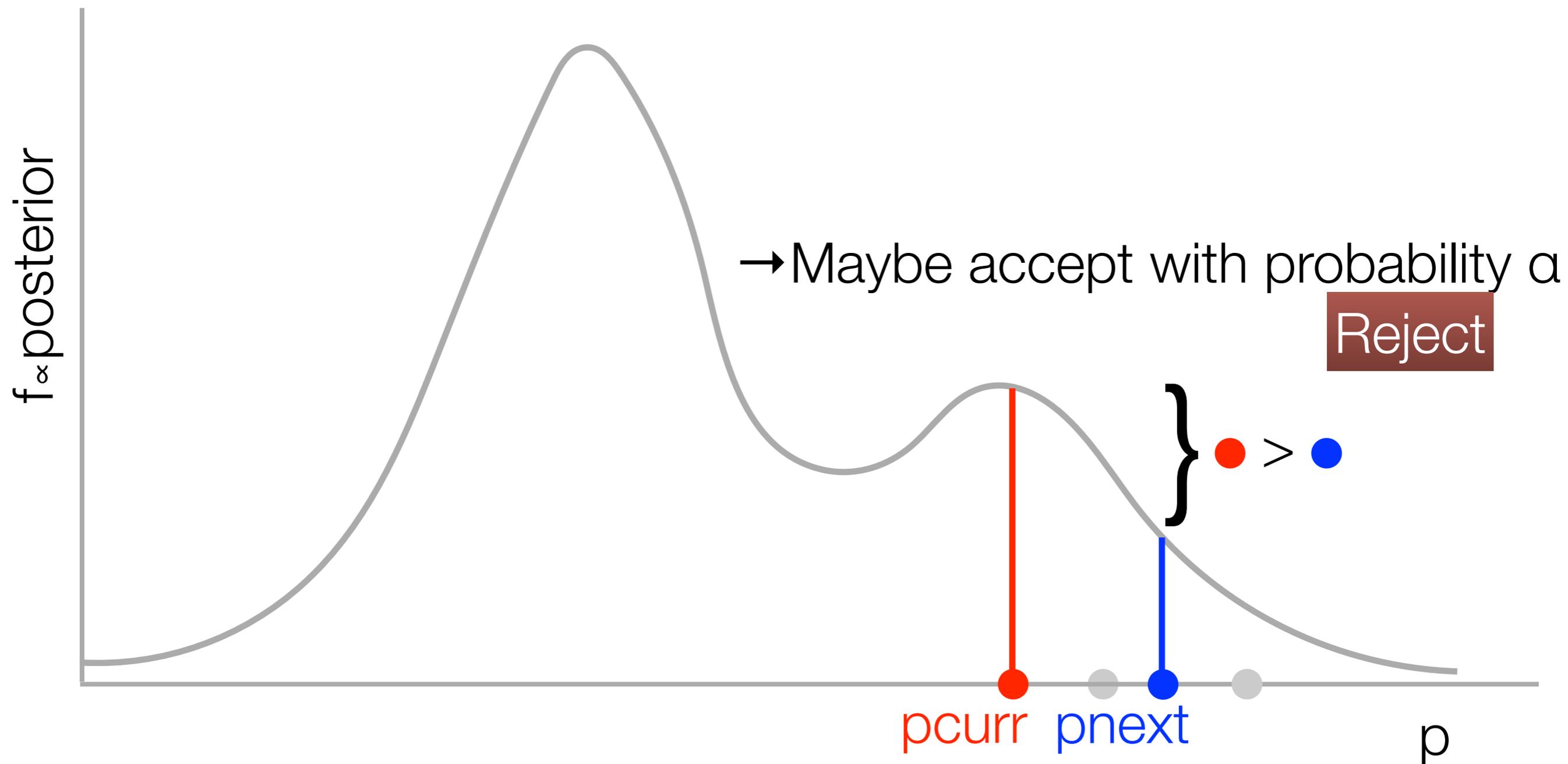
What does the metropolis algorithm do?



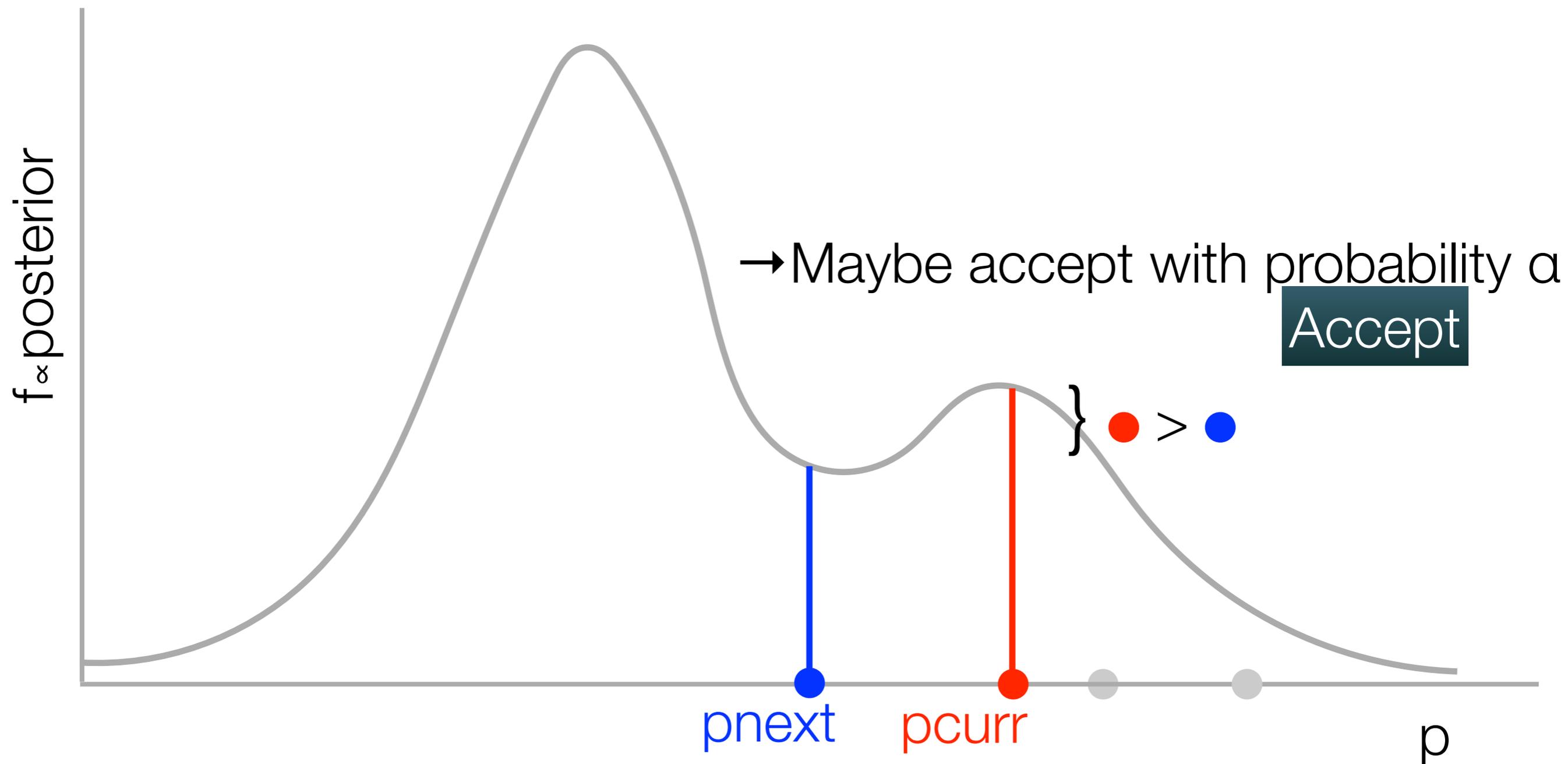
What does the metropolis algorithm do?



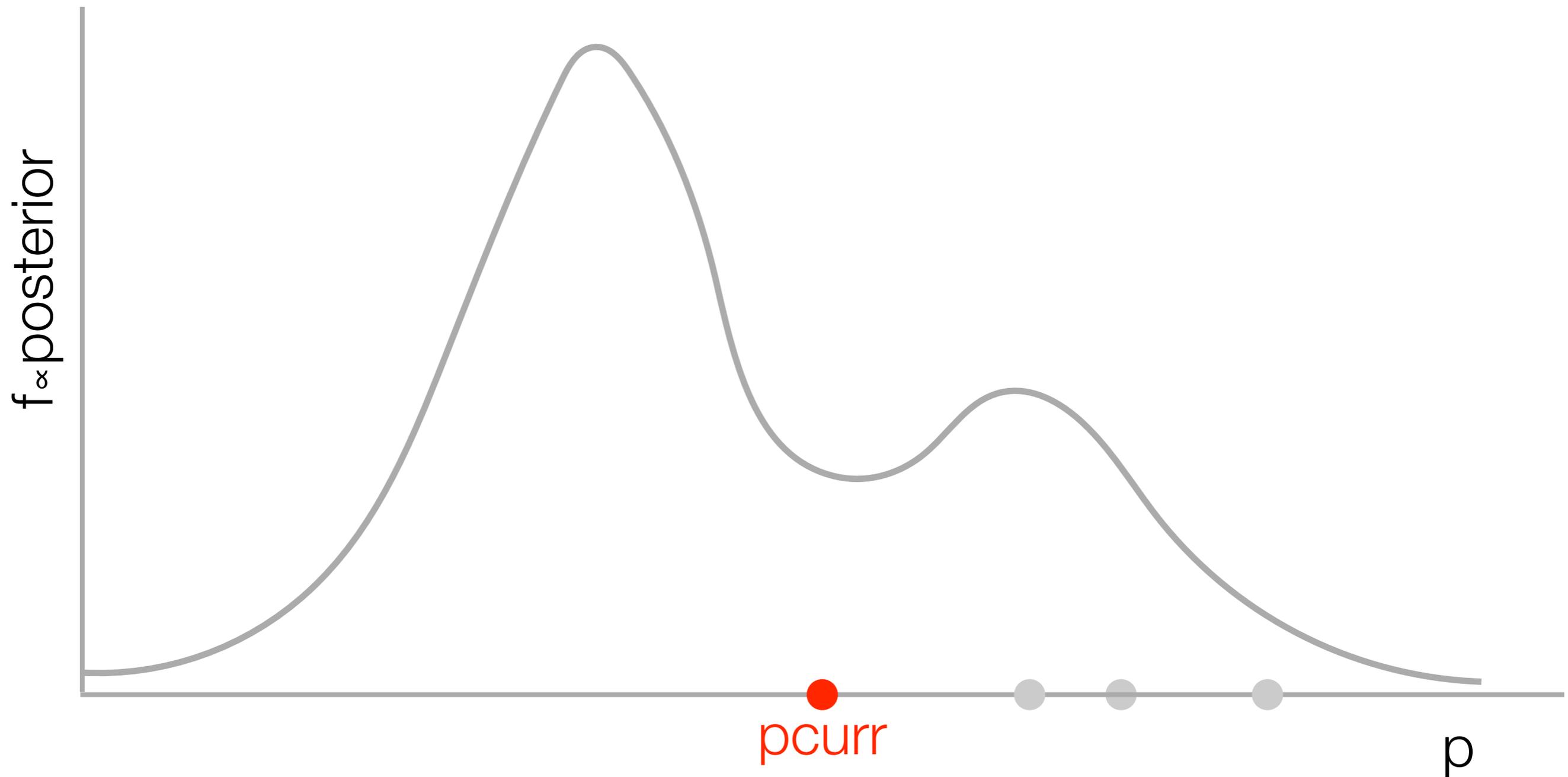
What does the metropolis algorithm do?



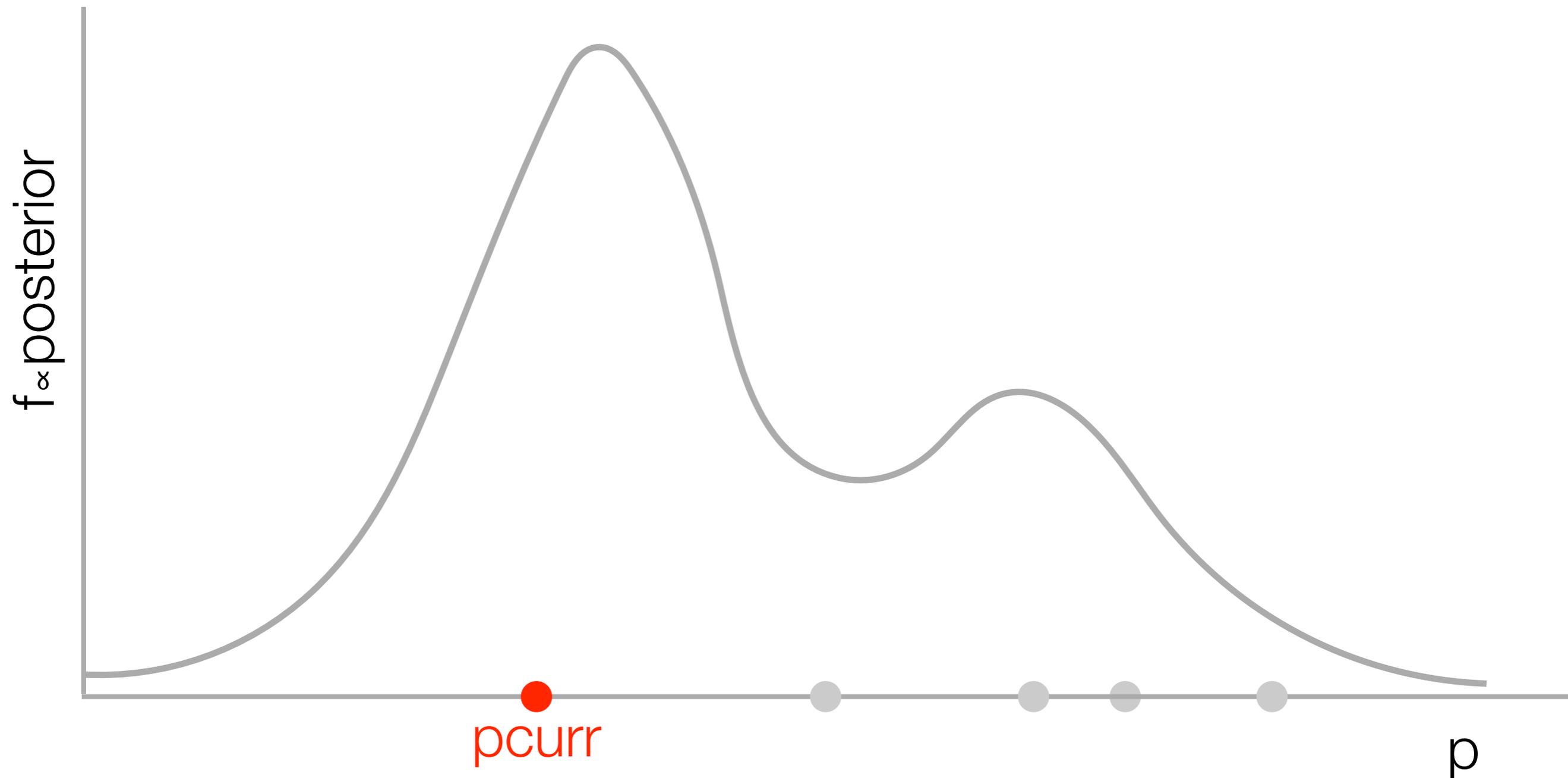
What does the metropolis algorithm do?



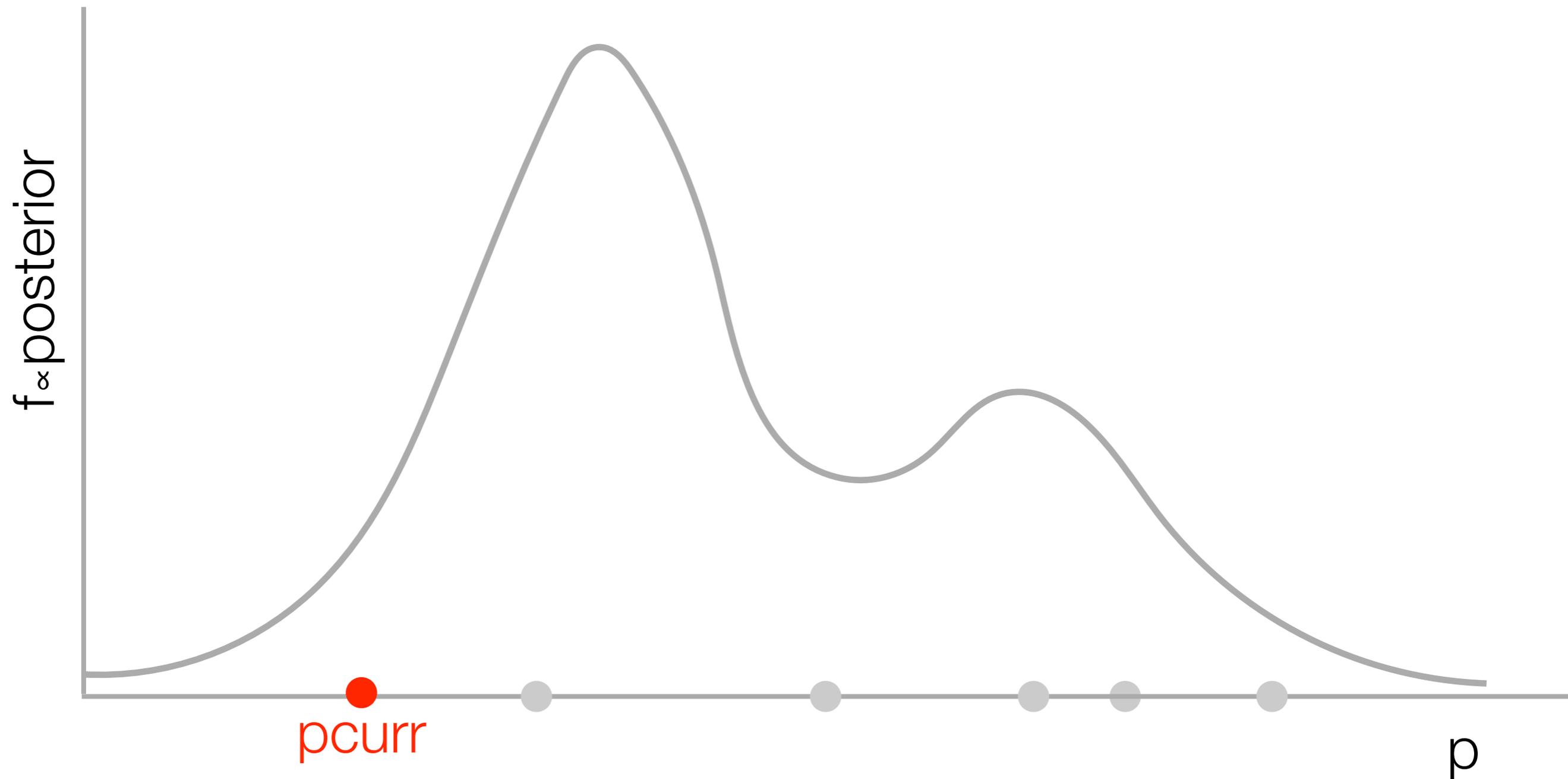
What does the metropolis algorithm do?



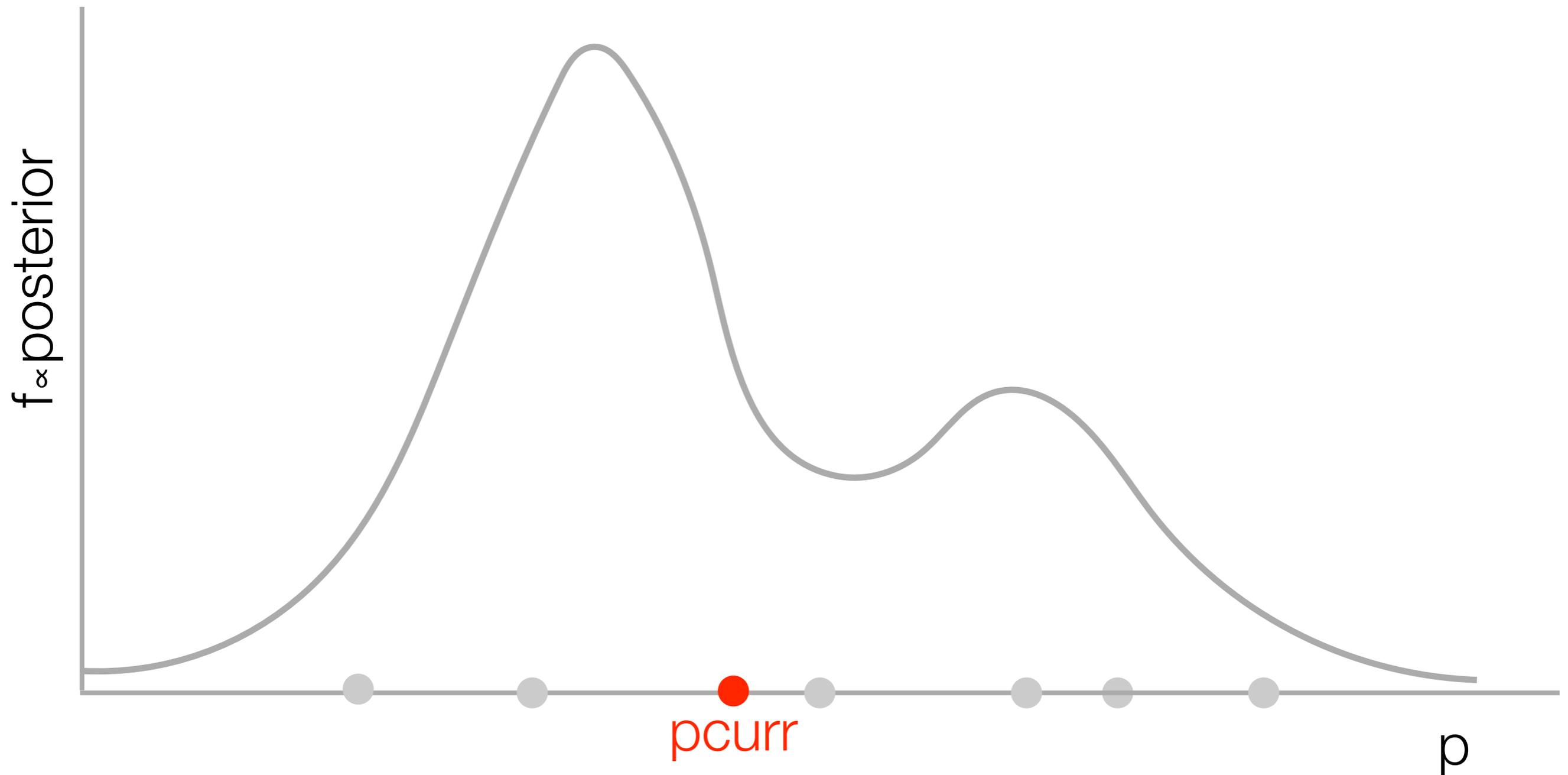
What does the metropolis algorithm do?



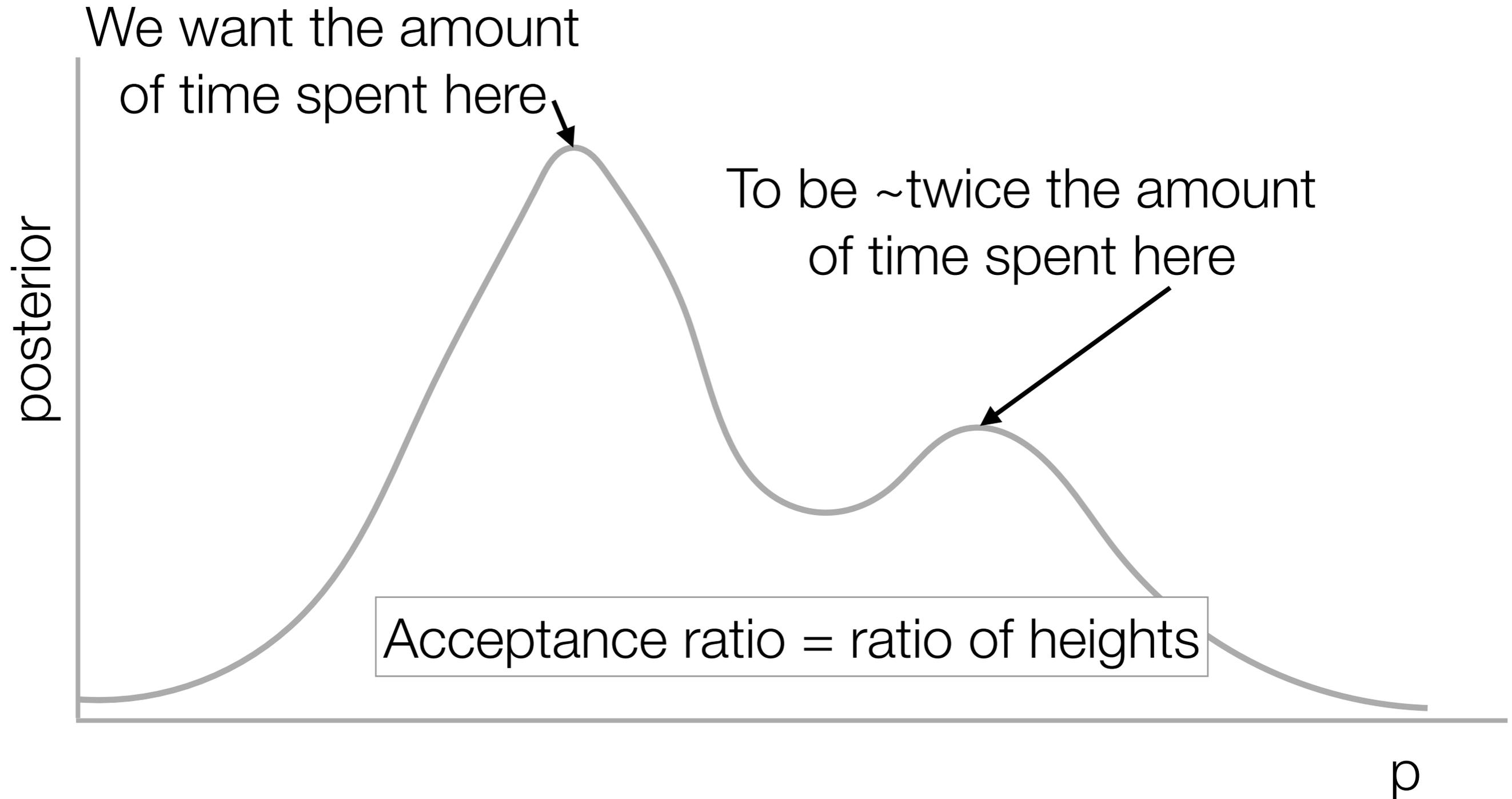
What does the metropolis algorithm do?



What does the metropolis algorithm do?



Why does this recover the posterior distribution?
Key is the acceptance ratio α



Why does this recover the posterior distribution?

- The acceptance ratio $\alpha = f(p_{next}) / f(p_{curr})$
- Note it is equal to $P(p_{next}|z) / P(p_{curr}|z)$ since the denominators cancel
- Suppose we're at the peak
 - If $f(p_{curr}) = 2 f(p_{next})$, then $\alpha = 1/2$, i.e. we accept with 1/2 probability
- Overall, will mean the number of samples we take from a region will be proportional to the height of the distribution

Proposal Distribution

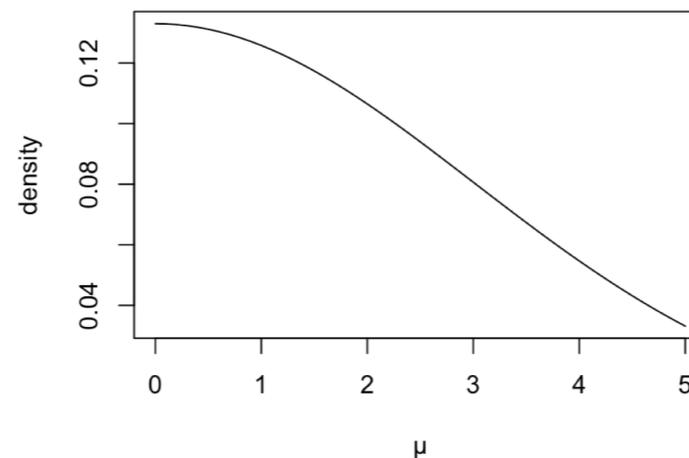
- A distribution that lets us choose our next point randomly from our current one
- For Metropolis algorithm, must be symmetric
- Common to choose a normal distribution centered on current point
- Width (SD) of normal = proposal width
 - Choice of proposal width can strongly affect how the Markov chain behaves, how well it converges, mixes, etc.

Example

- Model: normal distribution $\mathcal{N}(\mu, \sigma)$
 - Suppose σ is known, μ to be estimated

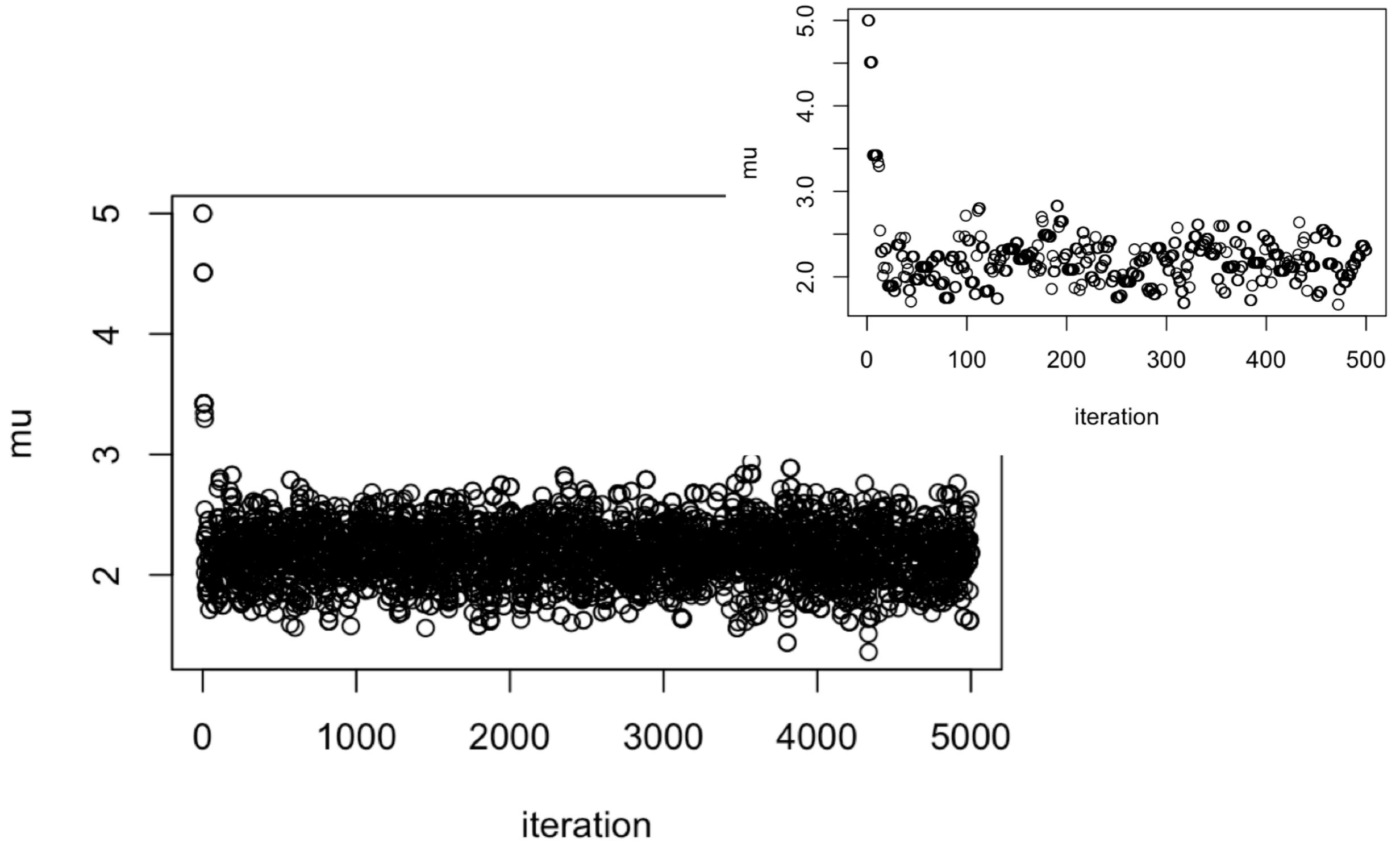
- Likelihood: $P(z_i | \mu, 1) = f(z_i | \mu, 1) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(z_i - \mu)^2}{2}}$ $P(z | \mu) = \prod_{i=1}^n f(z_i | \mu, 1)$

- Prior: $\mu \sim \mathcal{N}(0, 3)$



- Suppose we have 20 data points

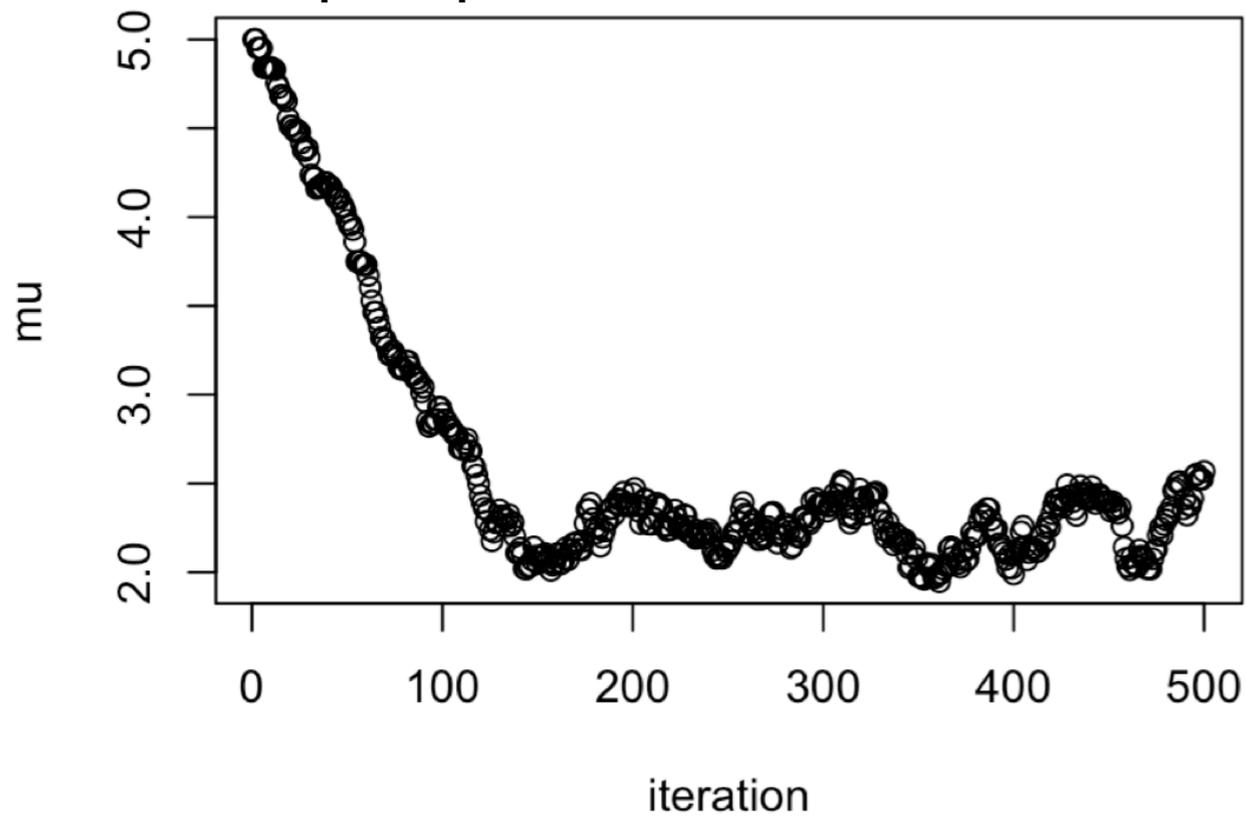
Example - proposal width: $SD = 0.5$



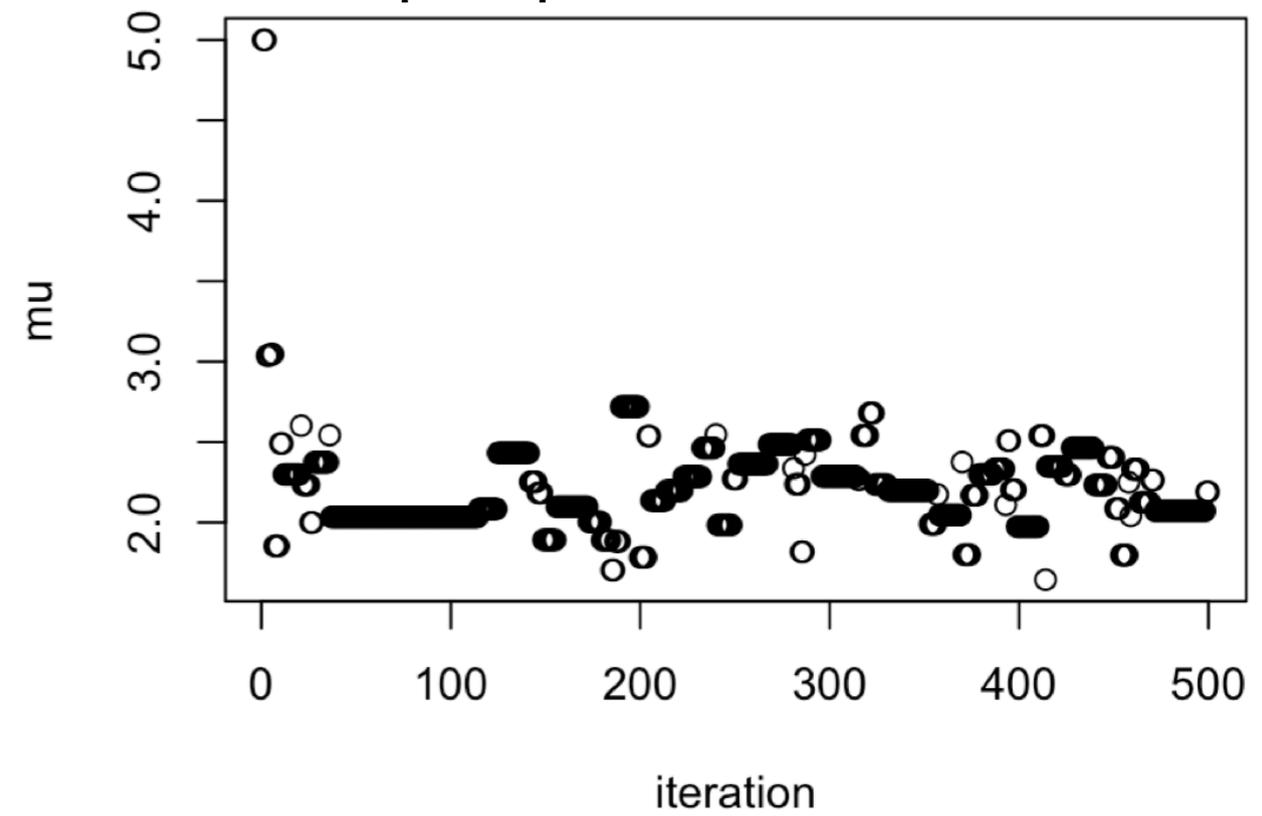
Goldilocks problem:

What happens if we change the proposal width?

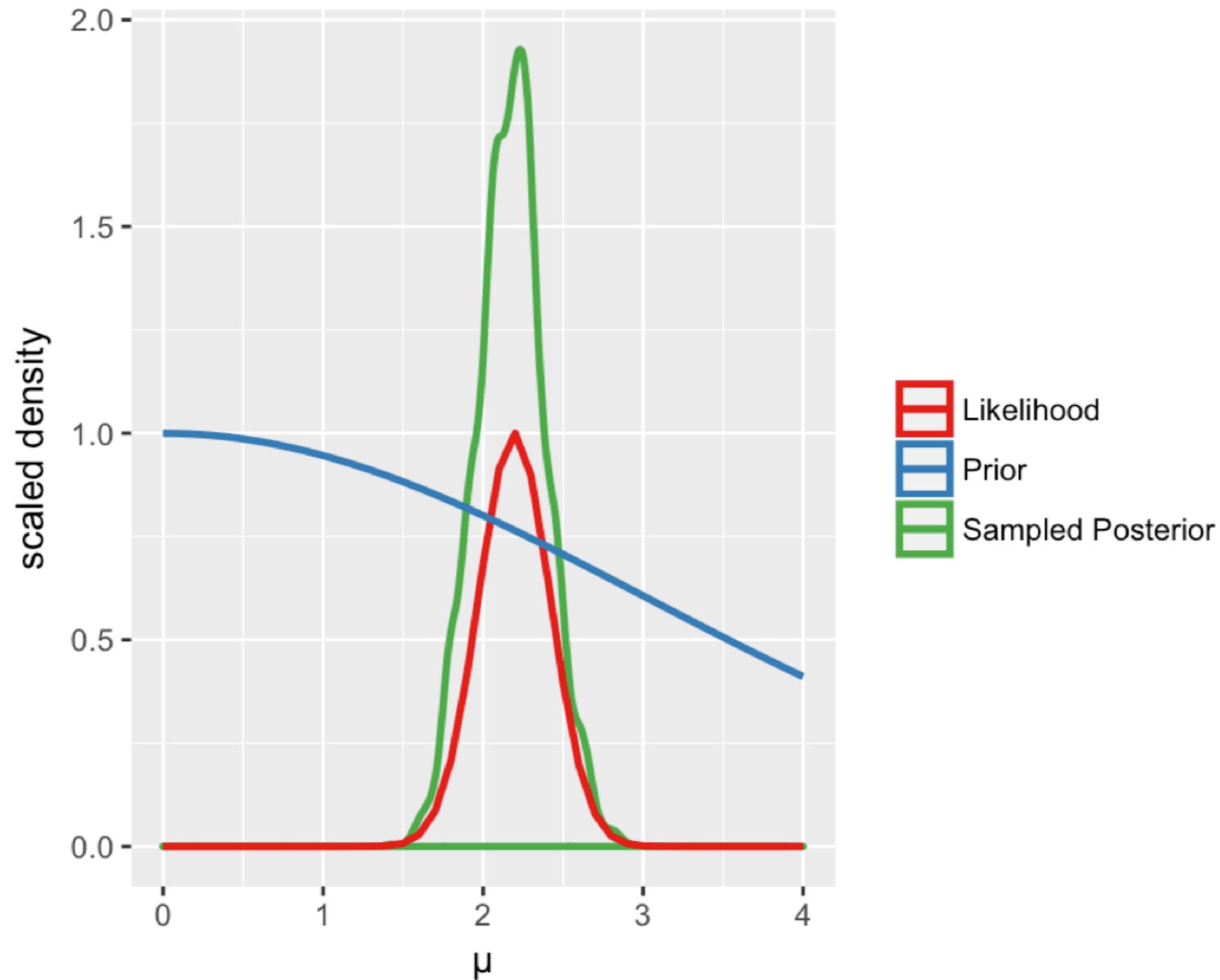
proposal SD = 0.05



proposal SD = 2



Example: prior, likelihood, and posterior (all scaled)



MCMC

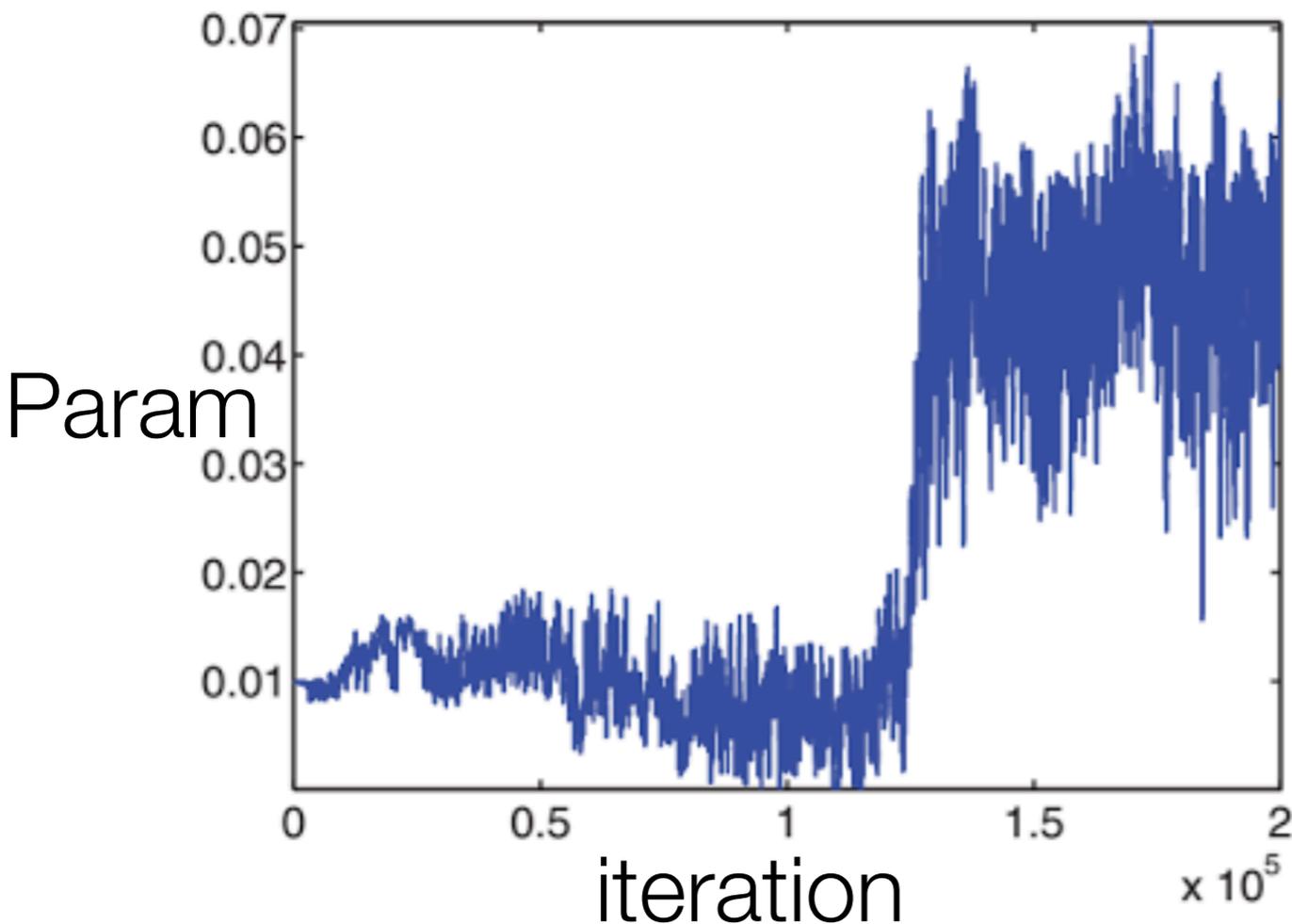
- MCMC improves many of the problems that other optimization methods face (getting trapped in local minima, etc.)
- However, those issues can still cause problems for MCMC too
- How to know when you've run the MCMC long enough and collected enough samples to reflect the distribution?
- How to know if you have explored the space sufficiently?

Assessing convergence

- MCMC methods will let us sample the posterior once they've converged to their equilibrium distribution
- How to know once we've reached equilibrium?
 - Visual evaluation of **burn-in**
 - Autocorrelation of elements in chain k iterations apart
- Also approaches to use in combination with/instead of burn-in: start with MAP estimation, multiple chains, etc.

Assessing convergence

- Often done visually
- Although, this can be misleading:



Chain shifts after 130,000 iterations due to a local min in sum of squares
(Example from R. Smith, *Uncertainty Quantification*)

Metropolis & Metropolis-Hastings Caveats

- Assessing convergence—how long is burn-in?
 - What about when you have unidentifiability or multiple minima?
- Correlated samples
- How to choose a proposal width? (~size of next jump)

Wide range of methods

- Metropolis–Hastings
- Gibbs sampling
- Variations of the above: prior optimization, multi-start, adaptive methods, delayed rejection
 - DRAM (Delayed Rejection Adaptive Metropolis-Hastings)
- Many more!

Examples



American Journal of Epidemiology

© The Author(s) 2017. Published by Oxford University Press on behalf of the Johns Hopkins Bloomberg School of Public Health. All rights reserved. For permissions, please e-mail: journals.permissions@oup.com.

Vol. 186, No. 12

DOI: 10.1093/aje/kwx217

Advance Access publication:

June 9, 2017

Practice of Epidemiology

Application of an Individual-Based Transmission Hazard Model for Estimation of Influenza Vaccine Effectiveness in a Household Cohort

Joshua G. Petrie*, Marisa C. Eisenberg, Sophia Ng, Ryan E. Malosh, Kyu Han Lee, Suzanne E. Ohmit, and Arnold S. Monto

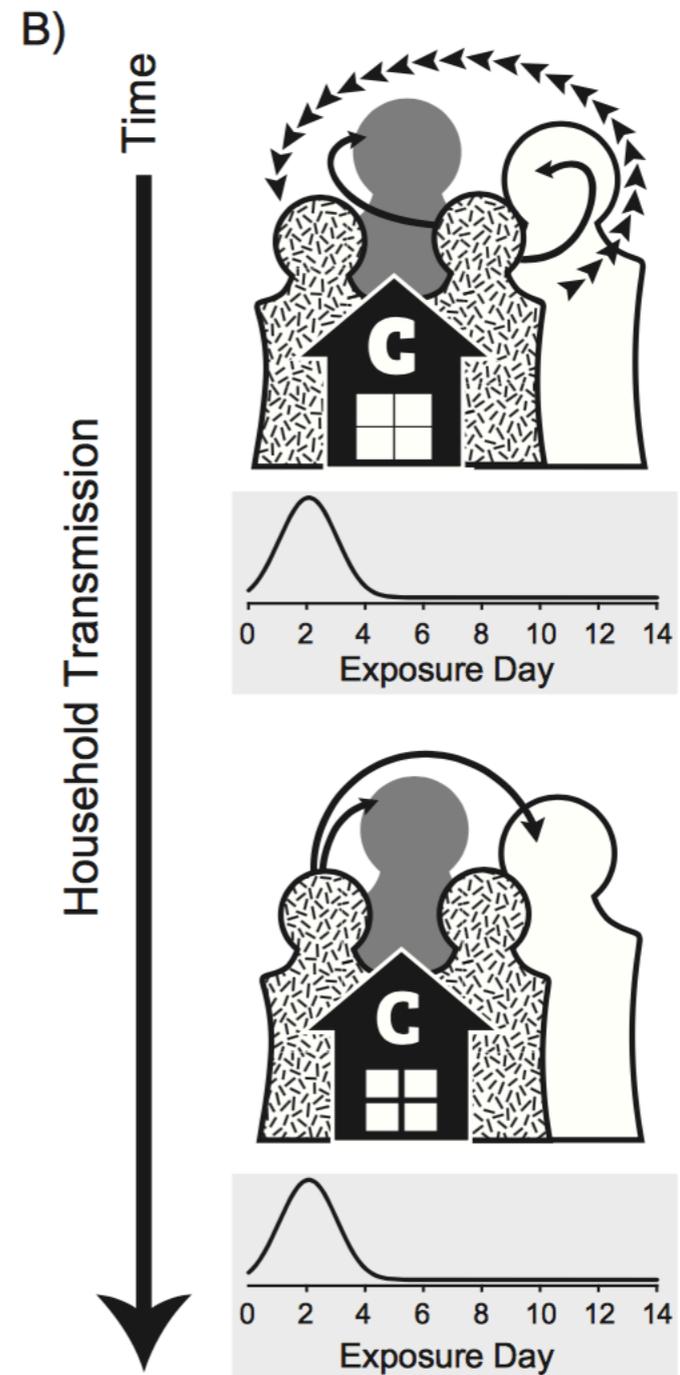
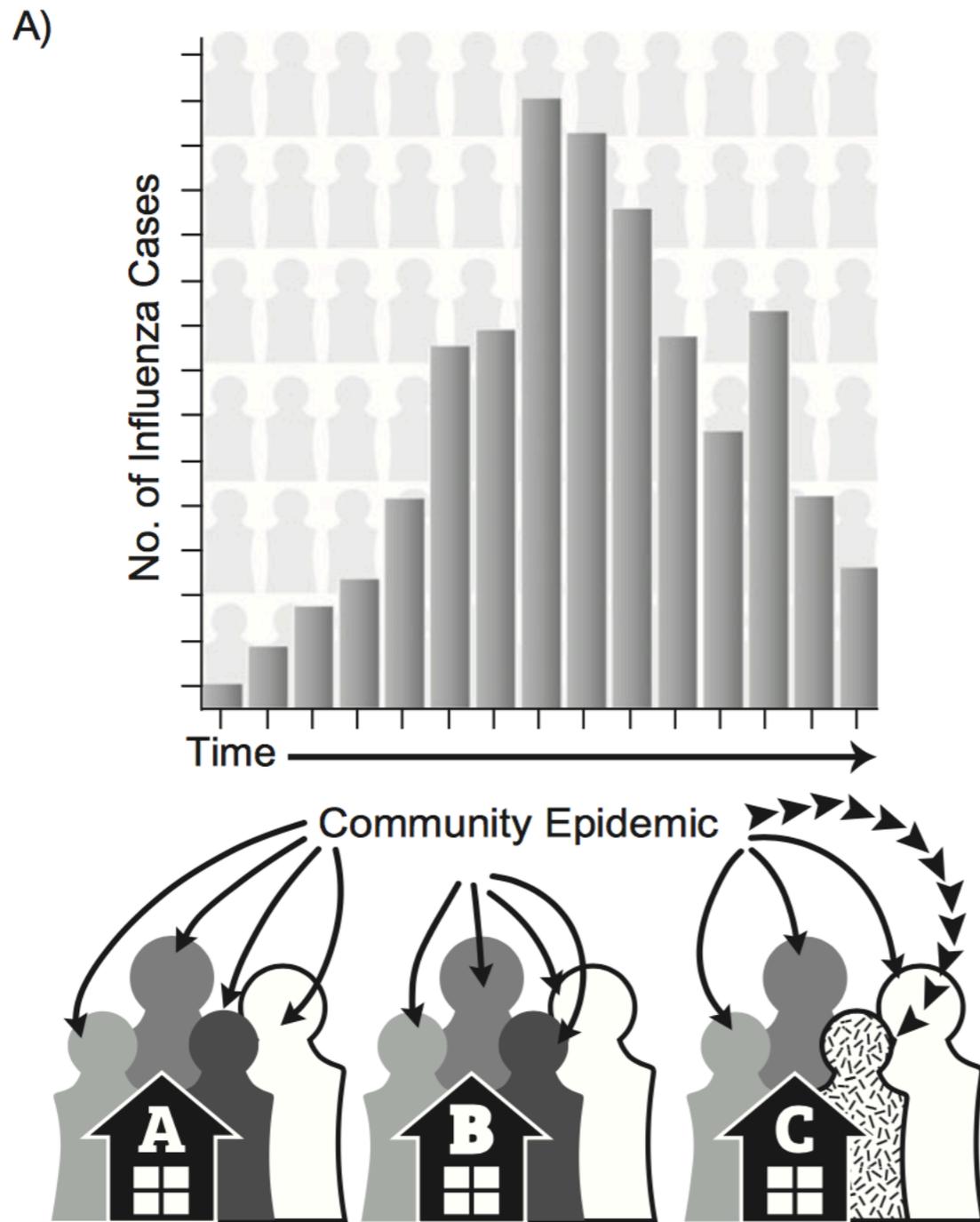


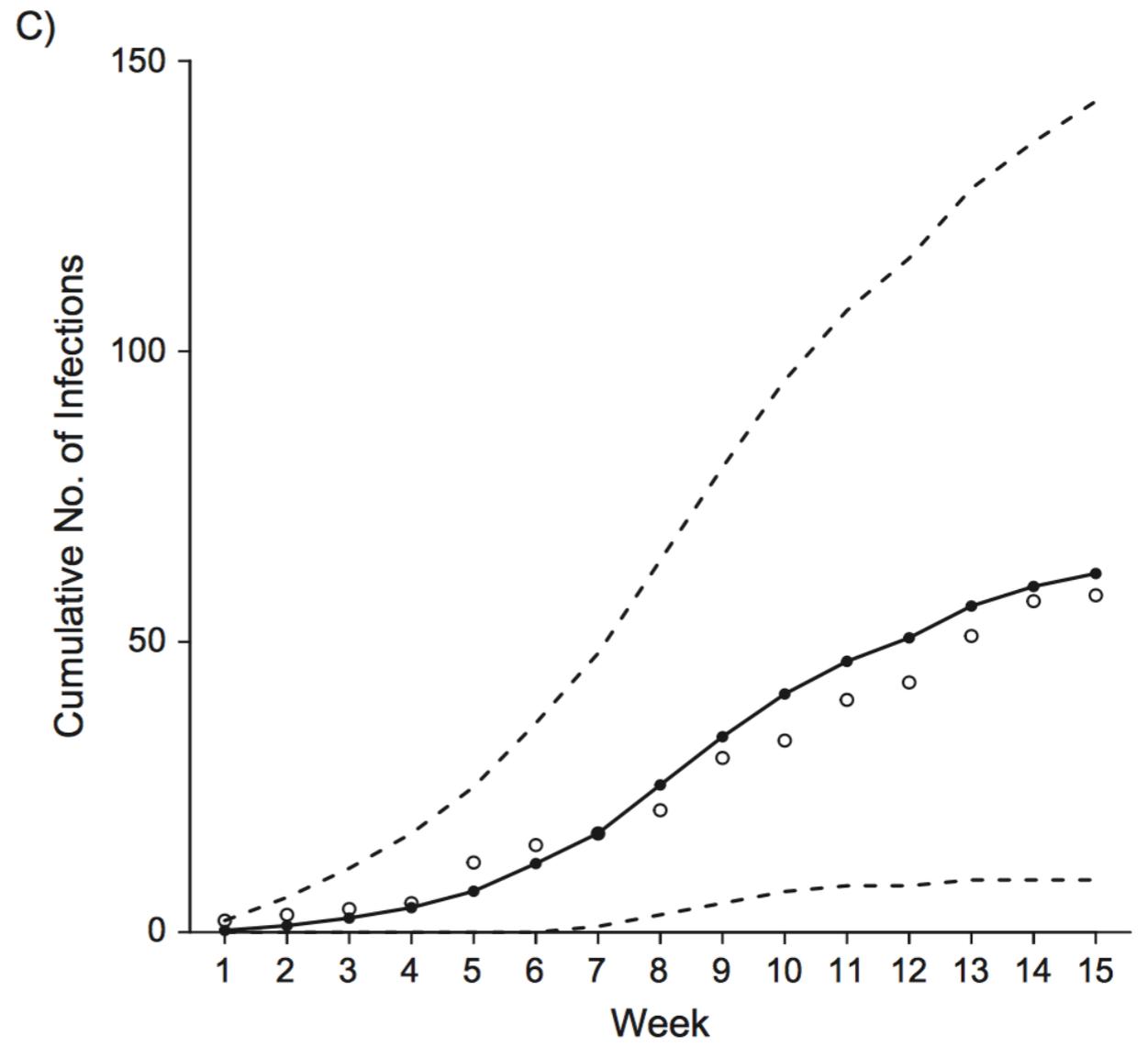
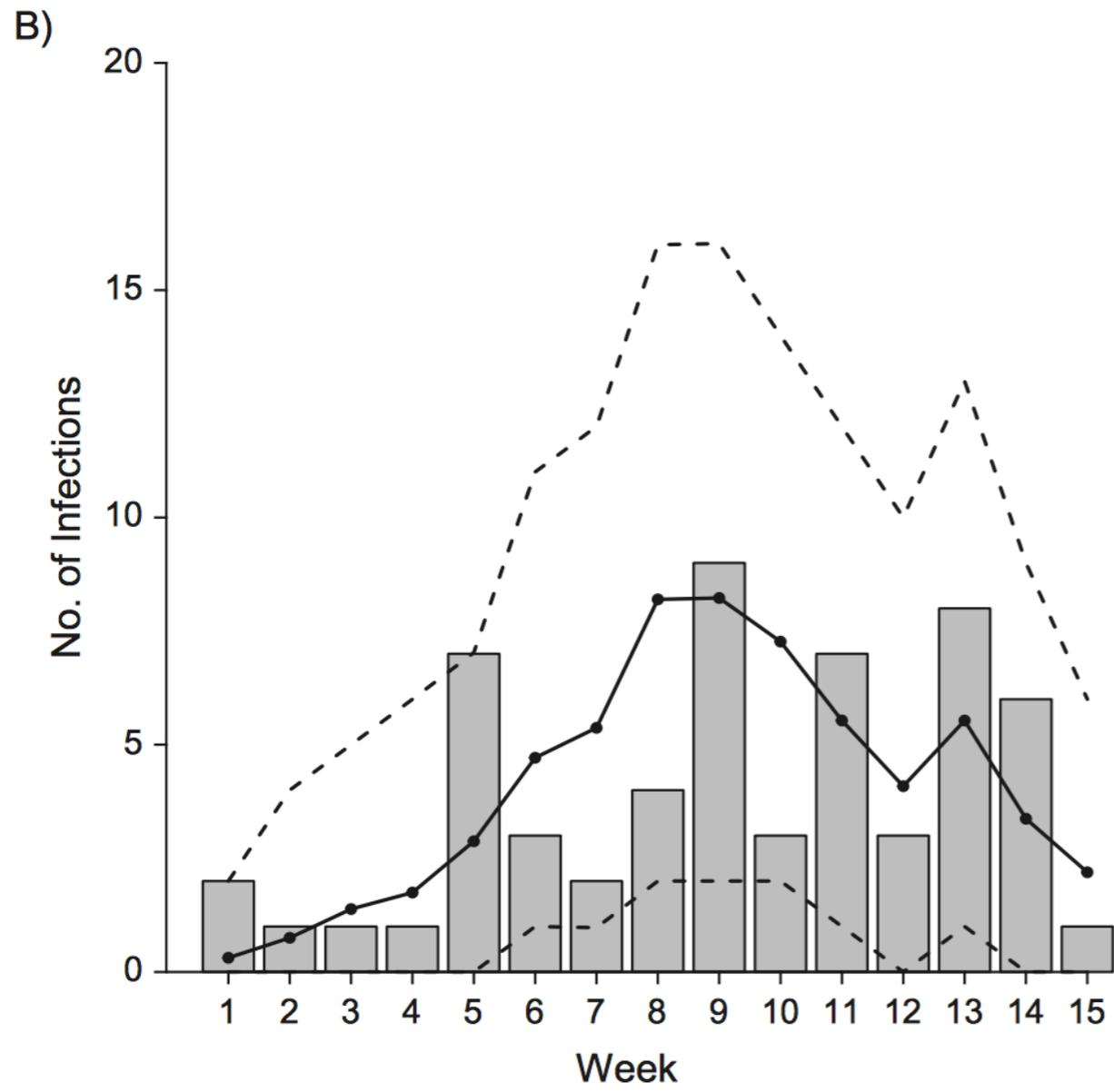
Table 2. Observed and Individual-Based Transmission Hazard Model–Predicted Influenza A(H3N2) Infections According to Infection Source, Age, Presence of High-Risk Health Condition, and Influenza Vaccination Status, Household Influenza Vaccine Effectiveness Study, Ann Arbor, Michigan, 2010–2011

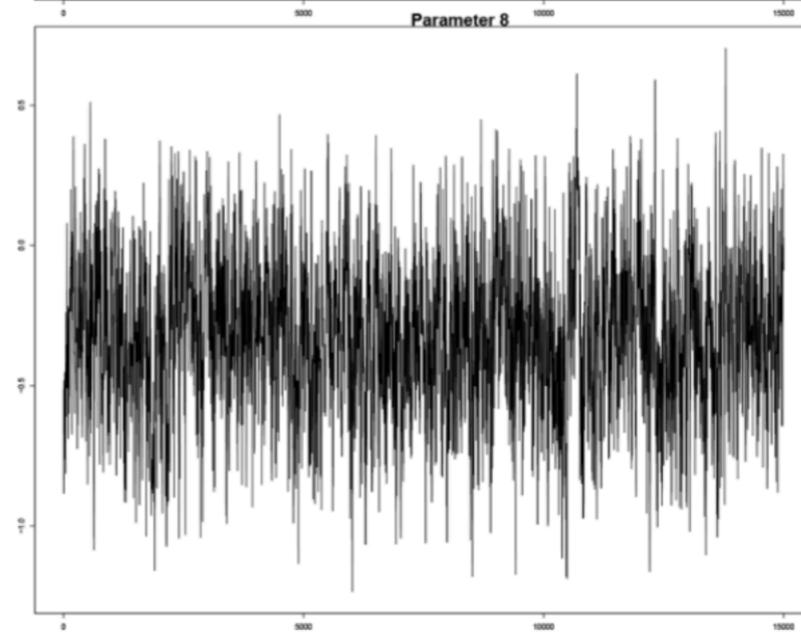
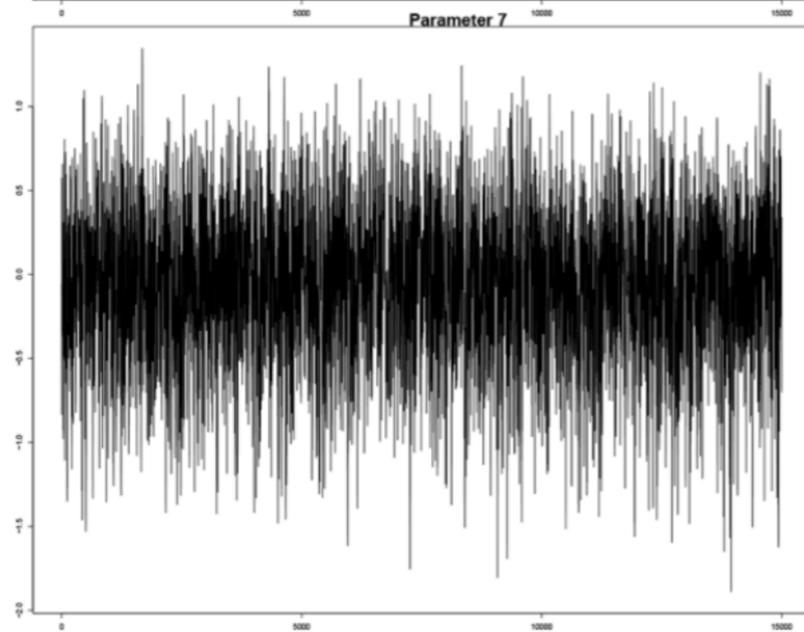
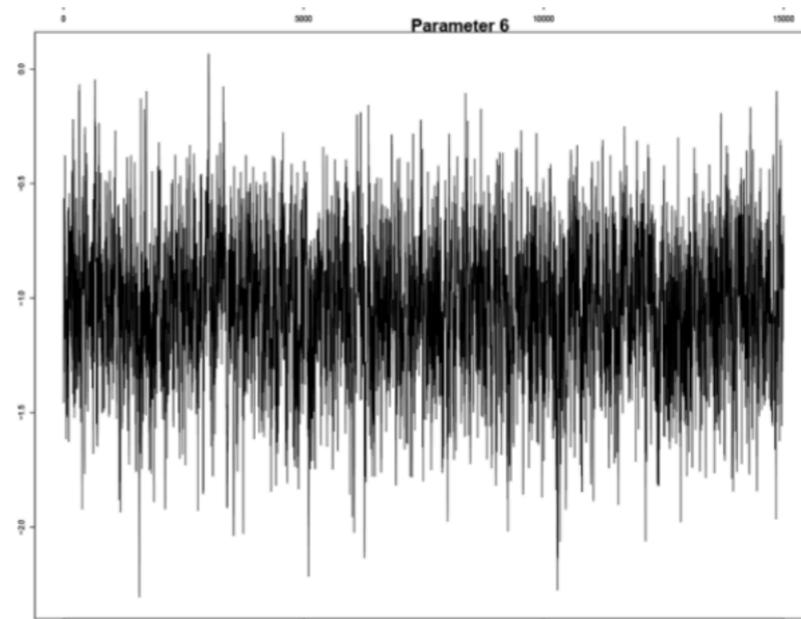
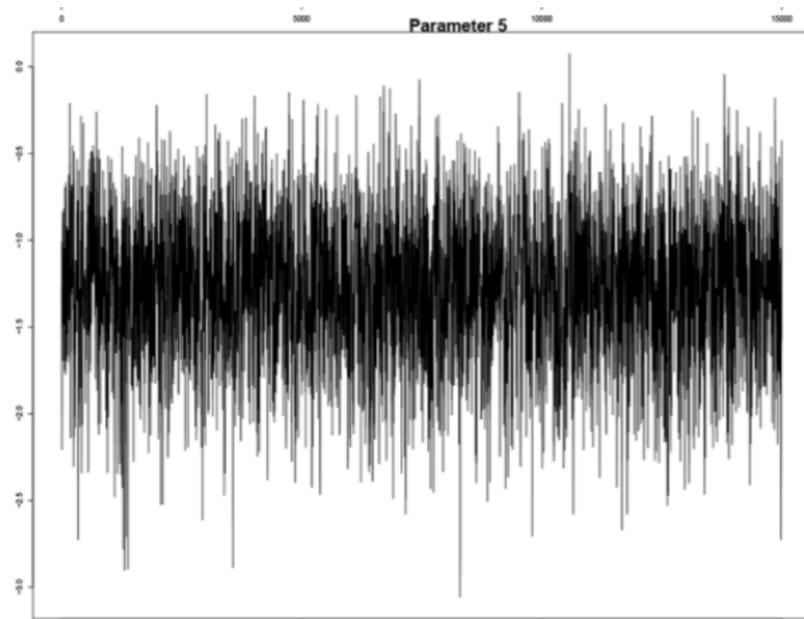
Characteristic	Observed Data			TH Model Predictions				P Value ^a
	No. of Cases (n = 58)	Total No. Exposed (n = 1,441)	% Positive	Median No. of Cases	95% CrI	% Positive	95% CrI	
Community-acquired	41	1,441	2.8	43	31, 55	3.0	2.2, 3.8	0.70
Household-acquired	17	111	15.3	18	9, 30	13.2	6.6, 20.5	
Secondary	N/O	N/O		15	7, 24			
Tertiary	N/O	N/O		3	0, 9			
Quaternary	N/O	N/O		0	0, 0			
Age category, years								0.80
<9	32	468	6.8	36	22, 50	7.7	4.7, 10.7	
9–17	8	371	2.2	8	3, 14	2.2	0.8, 3.8	
≥18	18	602	3.0	18	9, 27	3.0	1.5, 4.5	
Documented high-risk health condition								0.49
Any	6	162	3.7	5	1, 11	3.1	0.6, 6.8	
None	52	1,279	4.1	56	38, 76	4.4	3.0, 5.9	
Documented influenza vaccination ^b								0.45
Yes	33	864	3.8	32	19, 48	3.7	2.2, 5.6	
No	25	577	4.3	29	16, 44	5.0	2.8, 7.6	
Overall model predictions				62	42, 82	4.3	2.9, 5.7	

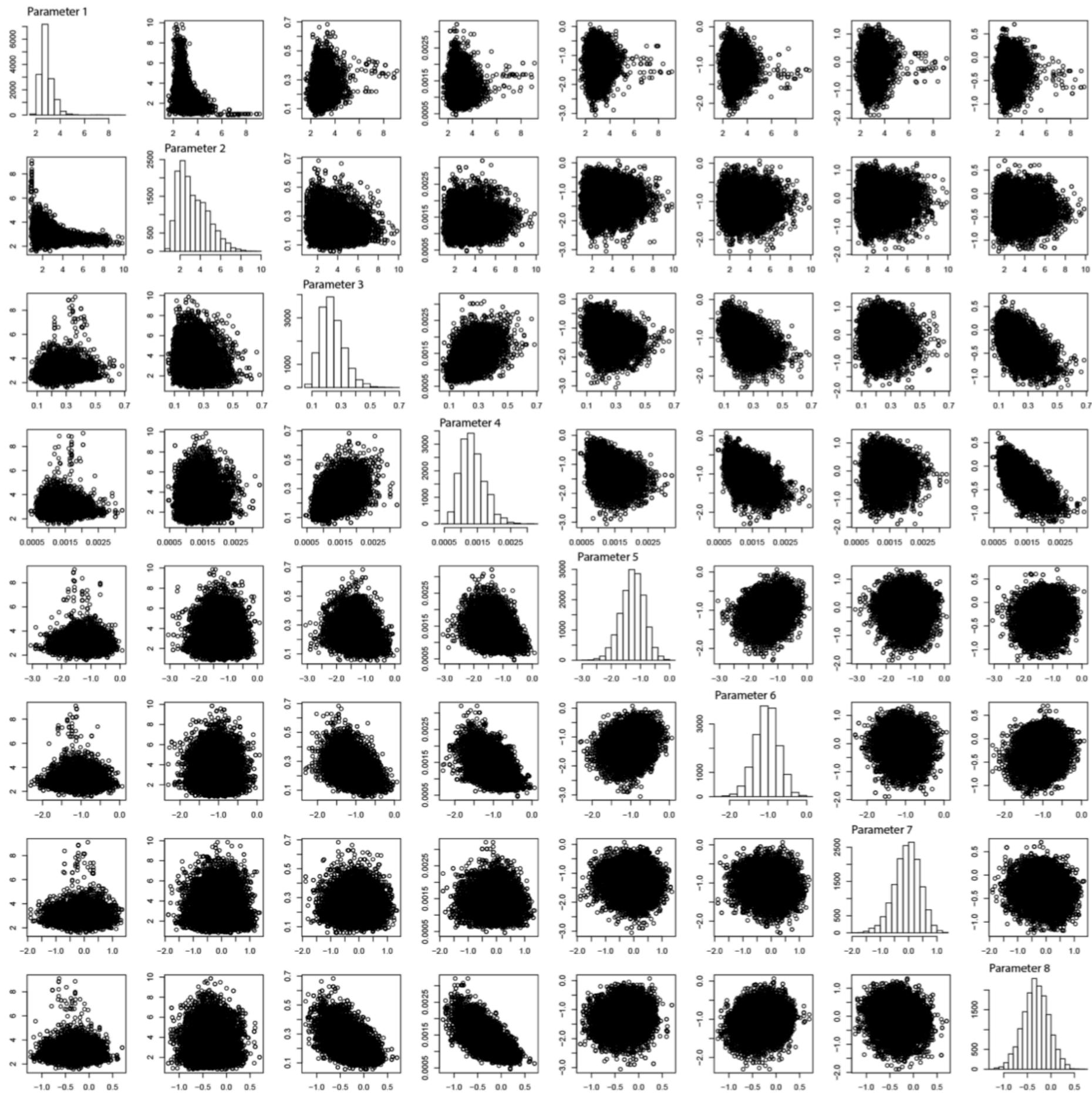
Abbreviations: CrI, credible interval; N/O, not observed; TH, transmission hazard.

^a Simulation-based χ^2 test.

^b At least 1 dose of 2010–2011 influenza vaccine documented in the electronic medical record or state registry; vaccination must have occurred ≥14 days prior to illness onset for influenza A(H3N2) infected subjects.





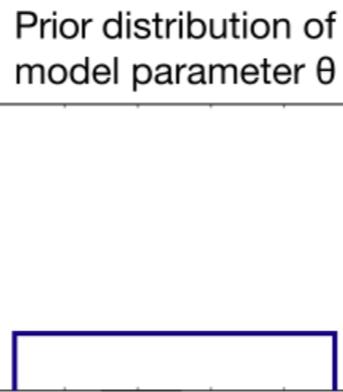
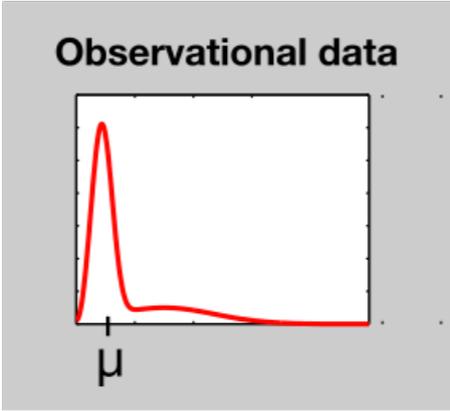


Sample Importance Resampling and Approximate Bayesian Computation

- MCMC can be slow — another approach to getting a rough sample of parameter space that matches the data is sample importance resampling
 - Can be used with the true likelihood
 - Or with an approximating function (approximate Bayesian computation)
 - E.g. may take a threshold based on distance between the model and observed data
- One of a bunch of related approaches in importance sampling/approximate Bayesian computation/etc)

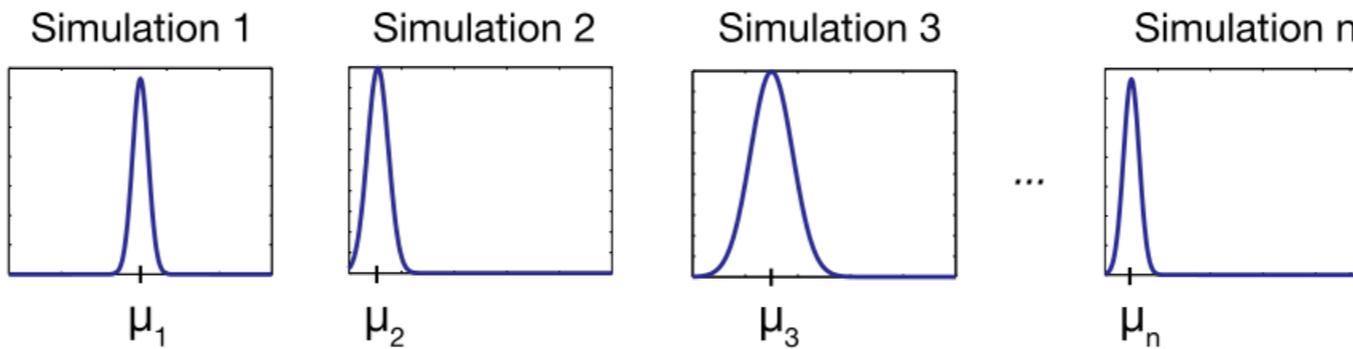
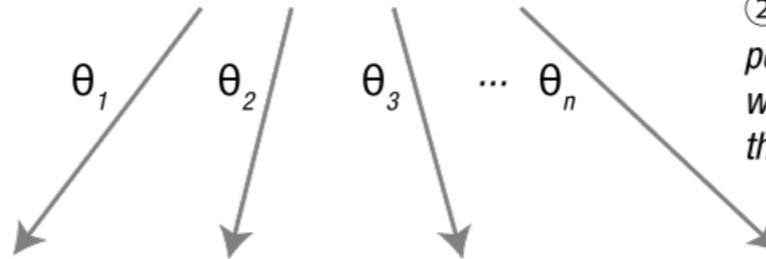
Basic idea

- Draw a sample of parameters from your prior (either drawing at random or with LHS/sobol/etc. sampling)
- Run the model for each sample
- Calculate the likelihood value (or approximation of it) for each sample
- Weight the samples based on the likelihood
- Resample to get the final set of samples



① Compute summary statistic μ from observational data

② Given a certain model, perform n simulations, each with a parameter drawn from the prior distribution



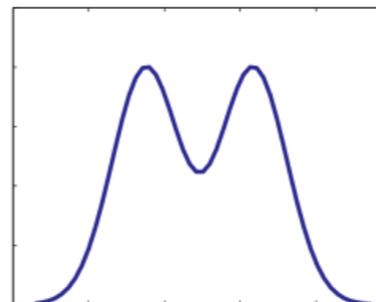
③ Compute summary statistic μ_i for each simulation

$$\rho(\mu_i, \mu) \stackrel{?}{\leq} \varepsilon$$



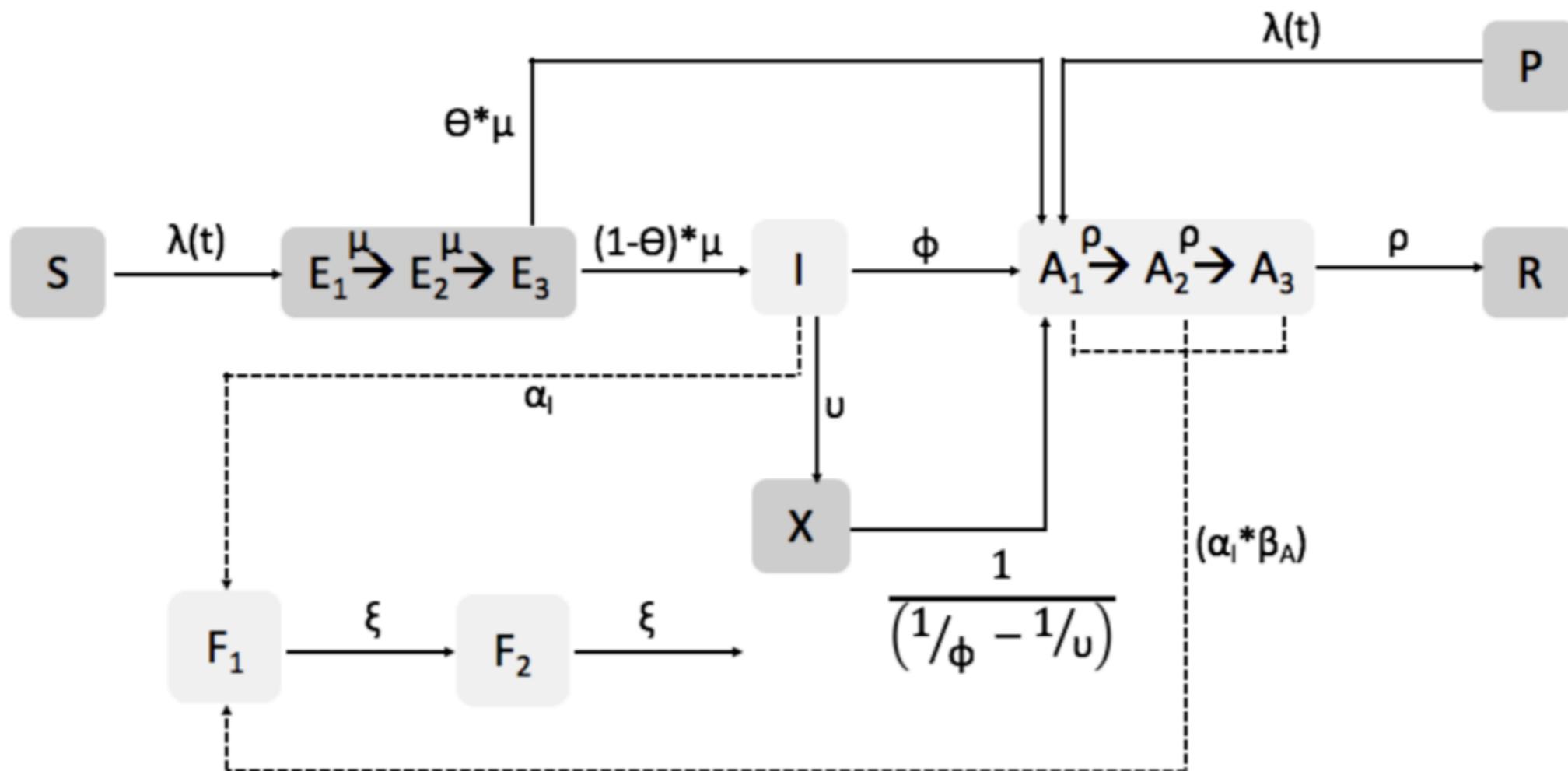
④ Based on a distance $\rho(\cdot, \cdot)$ and a tolerance ε , decide for each simulation whether its summary statistic is sufficiently close to that of the observed data.

Posterior distribution of model parameter θ

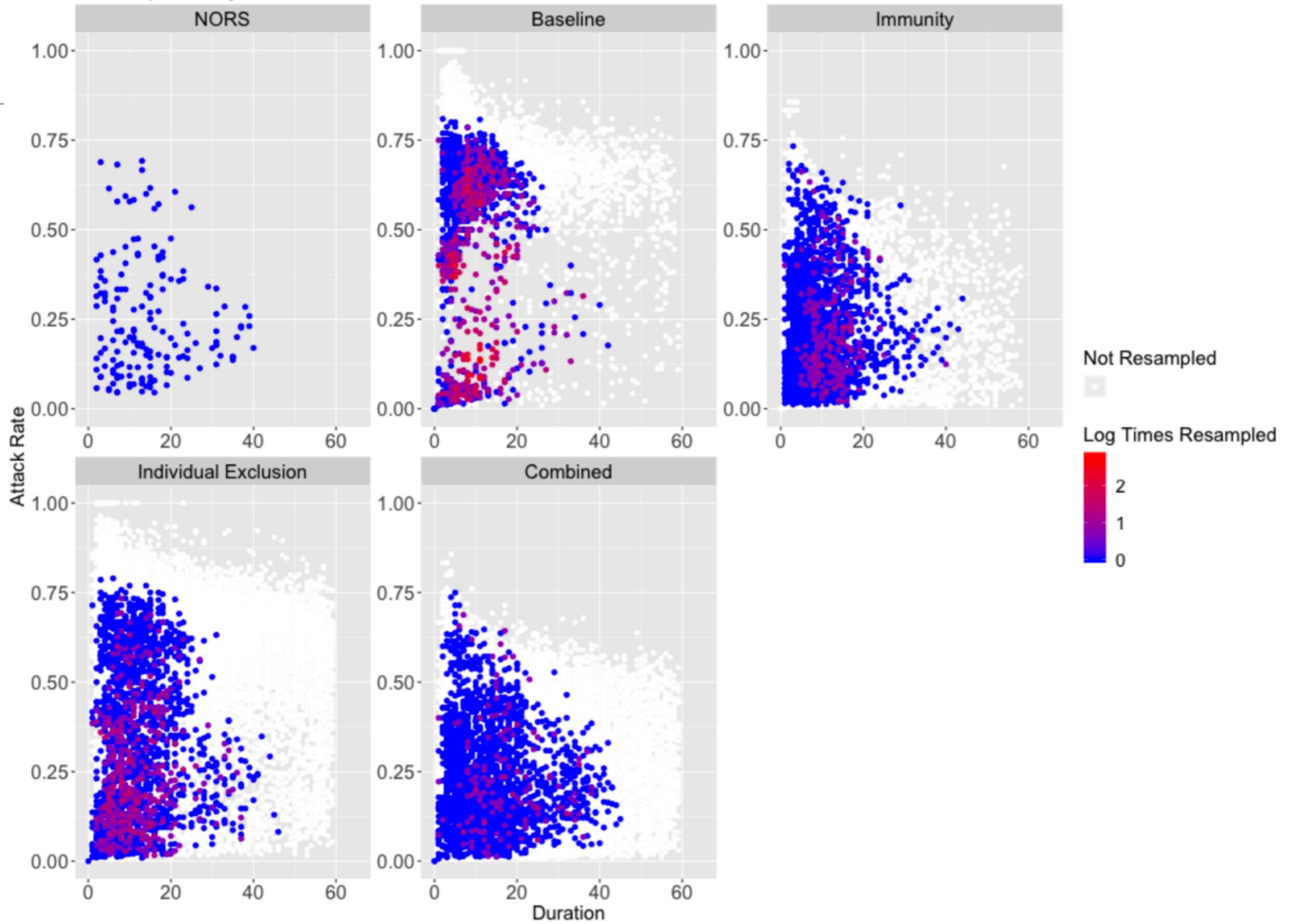


⑤ Approximate the posterior distribution of θ from the distribution of parameter values θ_i associated with accepted simulations.

Example: Norovirus model



Resampled Daycare Attack Rates vs. Outbreak Durations



Readings

- Menzies NA, Soeteman DI, Pandya A, Kim JJ. Bayesian methods for calibrating health policy models: a tutorial. *Pharmacoeconomics*. 2017 Jun 1;35(6):613-24.