

Approximating posterior distributions using simulation

FW 891

Christopher Cahill
30 August 2023



Quantitative Fisheries Center
MICHIGAN STATE UNIVERSITY

Recap

- Last time we learned that in some special situations we could generate analytical posterior distributions (e.g., the beta-binomial model)
- However, the reality is that this only works for simple models (i.e., typically < 3 parameters)
- Today we explore a broad class of computational methods and algorithms aimed at approximating posterior distributions
- We will still use simple models today, but realize that some of these tools scale far better than others to high-dimensional space

Outline

- Computing a posterior via grid approximation
- Markov Chain Monte Carlo
- Metropolis-Hastings algorithm
- Hamiltonian Monte Carlo

Grid search approximation: the basics

- Recall that:

$$\textit{posterior} \propto \textit{likelihood} \cdot \textit{prior}$$

- For models with 1-2 parameters approximating the posterior is actually a fairly tractable problem, even when the math isn't nice
- All we do is create a “grid” of parameter values and compute the posterior
- We then search across the computed posterior values to find the posterior maximum and corresponding parameter estimates

Revisiting the sneak turtles with different priors

```
1 library(tidyverse)
2 library(ggqfc)
3 y <- c(
4   0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0,
5   0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0
6 ) # successes (turtle detections)
7 n <- length(y) # trials
8
9 # define a grid
10 grid_size <- 10
11 p_grid <- seq(from = 1e-3, to = 0.999, length.out = grid_size)
12
13 # calculate likelihood
14 lik <- dbinom(sum(y), n, p_grid)
15 lik # likelihood of having observed the data | each p_grid value
```

```
[1] 4.229830e-04 1.964028e-01 1.859940e-02 5.603740e-04 6.085257e-06
[6] 1.860800e-08 8.682232e-12 1.445413e-16 7.970008e-25 4.341304e-82
```

- see `?dunif`, `?dnorm`, `?dbinom`, etc.

Always work in log-space when dealing with likelihoods and prior probabilities!

Recognize that

$$\text{posterior} \propto \text{likelihood} \cdot \text{prior}$$

is the same as

$$\log(\text{posterior}) \propto \log(\text{likelihood}) + \log(\text{prior})$$

When we multiply likelihoods together they get smaller and may surpass machine precision, so we usually work in log-space

Back to the R side of things...

```
1 # Note that  
2 log(lik)  
  
[1] -7.768179 -1.627588 -3.984626 -7.486906 -12.009642 -17.799674  
[7] -25.469743 -36.472966 -55.488942 -187.343803  
  
1 # is the same as  
2 log_lik <- dbinom(sum(y), n, p_grid, TRUE)  
3 log_lik  
  
[1] -7.768179 -1.627588 -3.984626 -7.486906 -12.009642 -17.799674  
[7] -25.469743 -36.472966 -55.488942 -187.343803
```

- We have now defined the log-likelihood across our entire grid sequence of p
- Now we need to calculate the prior probabilities at each p_{grid} value

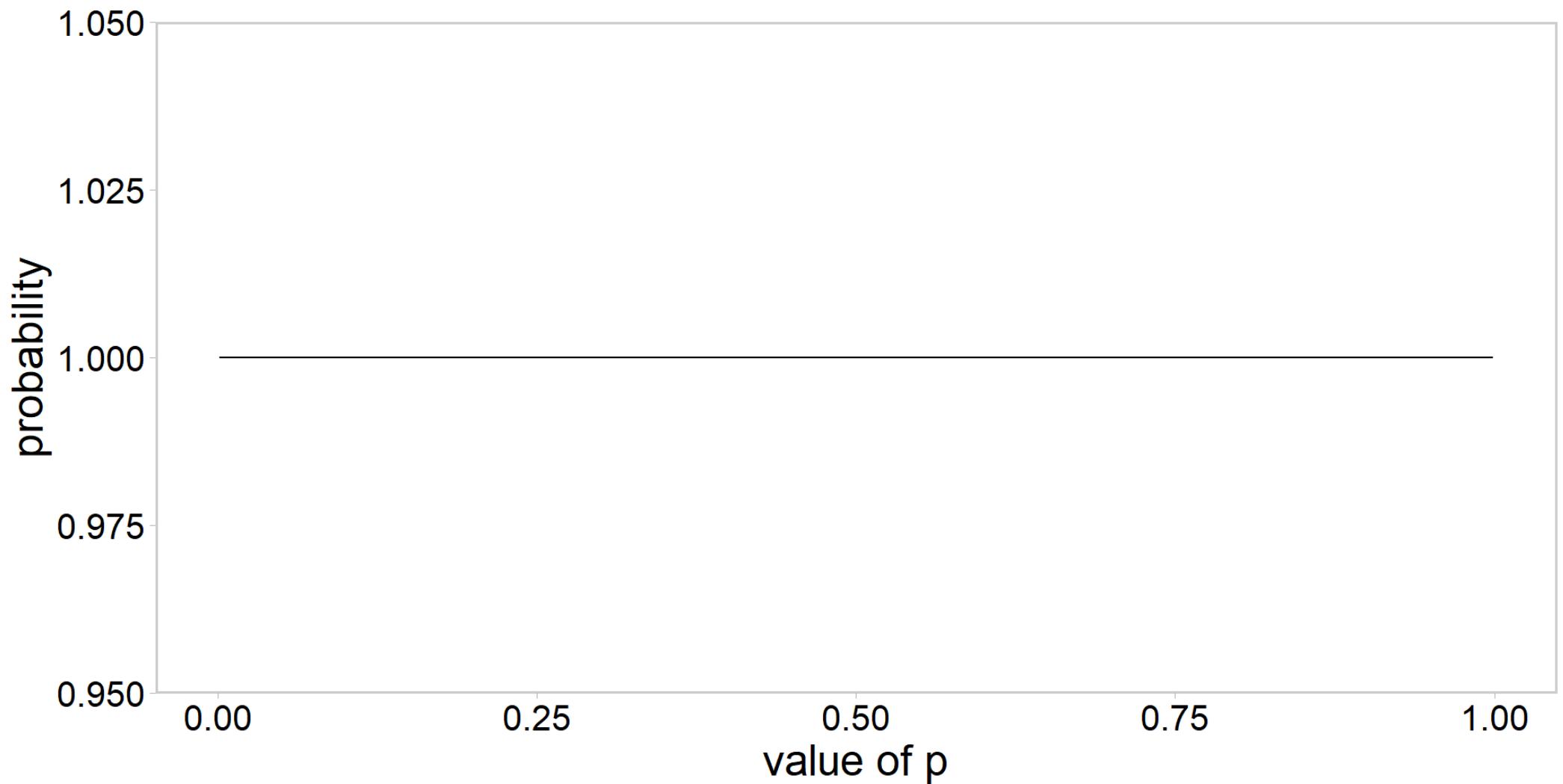
Dealing with the prior probabilities

- Assume $p \sim \text{Uniform}(0, 1)$
- We can calculate the log(prior probability) for each p in our grid given a uniform prior via:

```
1 log_prior <- dunif(p_grid, 0, 1, TRUE)
2 log_prior
[1] 0 0 0 0 0 0 0 0 0 0
```

- Is this correct?

Visualizing the Uniform prior



Calculating the posterior

- Now can calculate the posterior probabilities of each `p_grid` value given the likelihood and prior as follows:

```
1 # unstandardized log posterior
2 log_unstd_post <- log_lik + log_prior
3
4 # exponentiate the log posterior
5 unstd_post <- exp(log_unstd_post)
6
7 # standardize it (so probabilities sum to one)
8 post_prob <- unstd_post / sum(unstd_post)
9
10 # Put everything together in a data frame
11 posterior <- data.frame(
12   p = p_grid, log_lik, log_prior,
13   log_unstd_post, post_prob
14 )
```

Let's take a look at our posterior approximation

```
1 posterior
```

	p	log_lik	log_prior	log_unstd_post	post_prob
1	0.0010000	-7.768179	0	-7.768179	1.958330e-03
2	0.1118889	-1.627588	0	-1.627588	9.093073e-01
3	0.2227778	-3.984626	0	-3.984626	8.611164e-02
4	0.3336667	-7.486906	0	-7.486906	2.594424e-03
5	0.4445556	-12.009642	0	-12.009642	2.817357e-05
6	0.5554444	-17.799674	0	-17.799674	8.615147e-08
7	0.6663333	-25.469743	0	-25.469743	4.019707e-11
8	0.7772222	-36.472966	0	-36.472966	6.691985e-16
9	0.8881111	-55.488942	0	-55.488942	3.689961e-24
10	0.9990000	-187.343803	0	-187.343803	2.009941e-81

```
1 which.max(posterior$post_prob)
```

```
[1] 2
```

```
1 posterior[which.max(posterior$post_prob), ] # best estimate of p
```

	p	log_lik	log_prior	log_unstd_post	post_prob
2	0.1118889	-1.627588	0	-1.627588	0.9093073

Key points:

For each value of p in `p_grid` we get

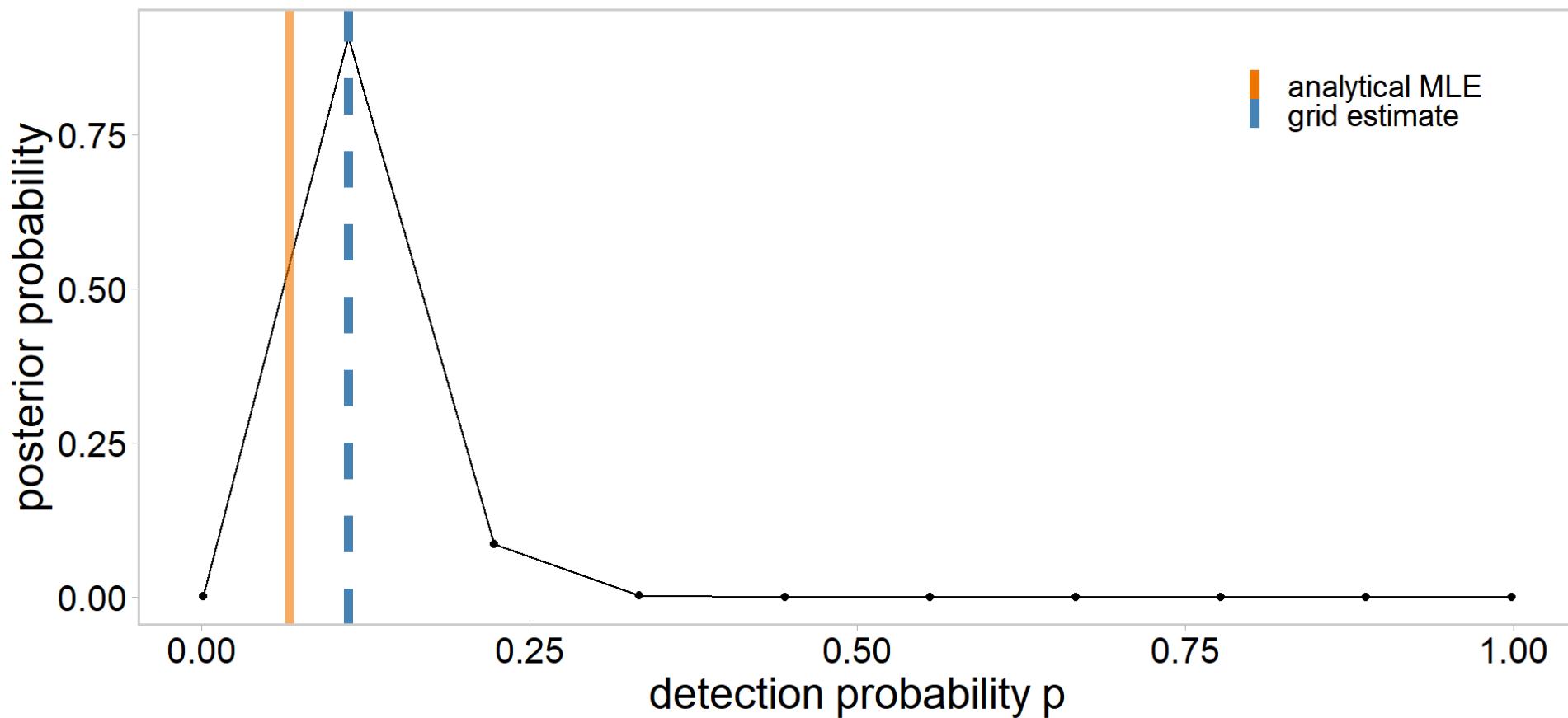
- log-likelihood of the data | that p
- log-prior probability for that p | the prior(s)
- log-posterior calculation for each p (very much not a probability)
- posterior probability for each p

Key points continued:

- The log-likelihood refers to the joint density of all data points (i.e., the sum of the log-likelihood values for each datum)
- log-prior is the marginal distribution of the parameter(s)

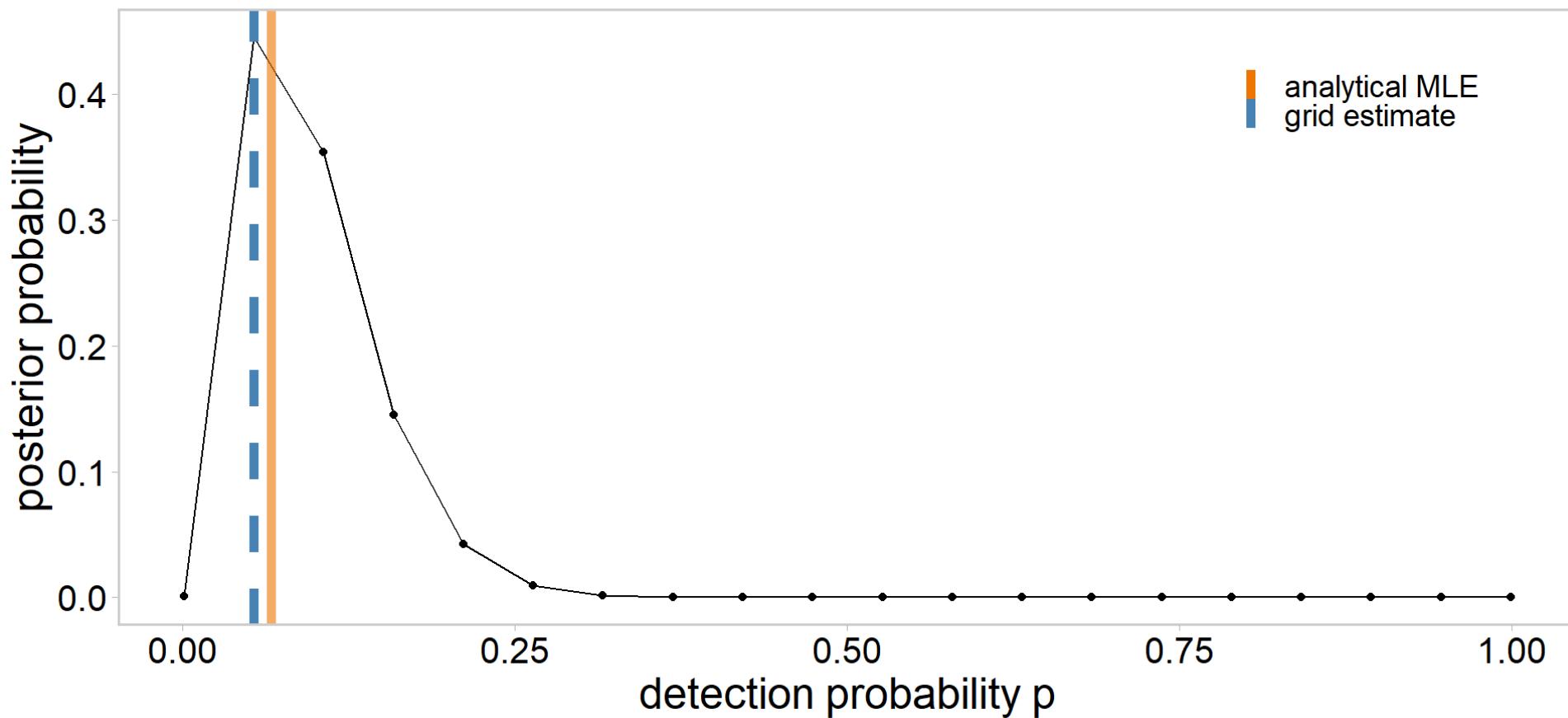
Comparing our grid approximation with the analytical MLE

Grid size: 10 points

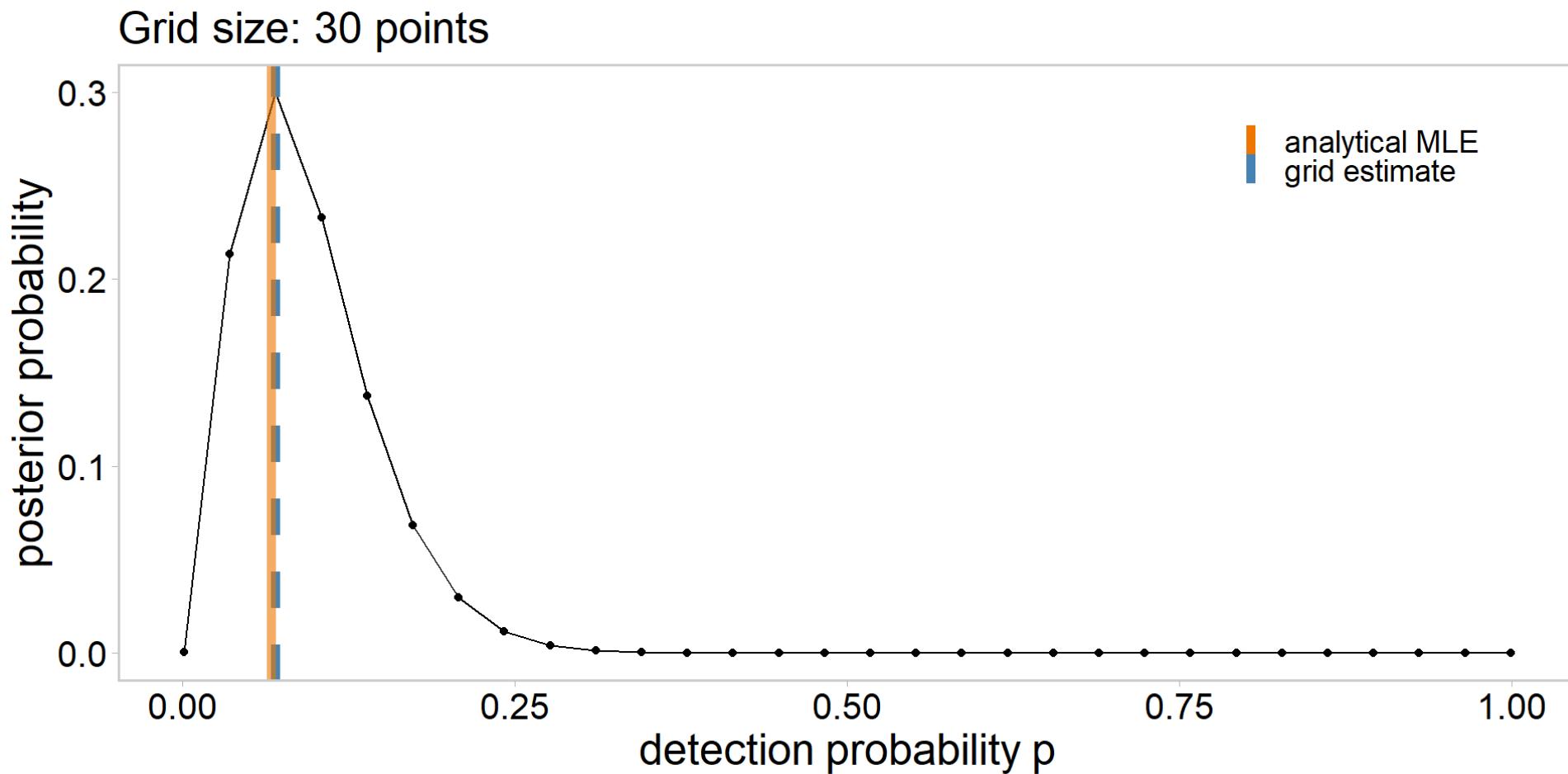


Comparing our grid approximation with the analytical MLE

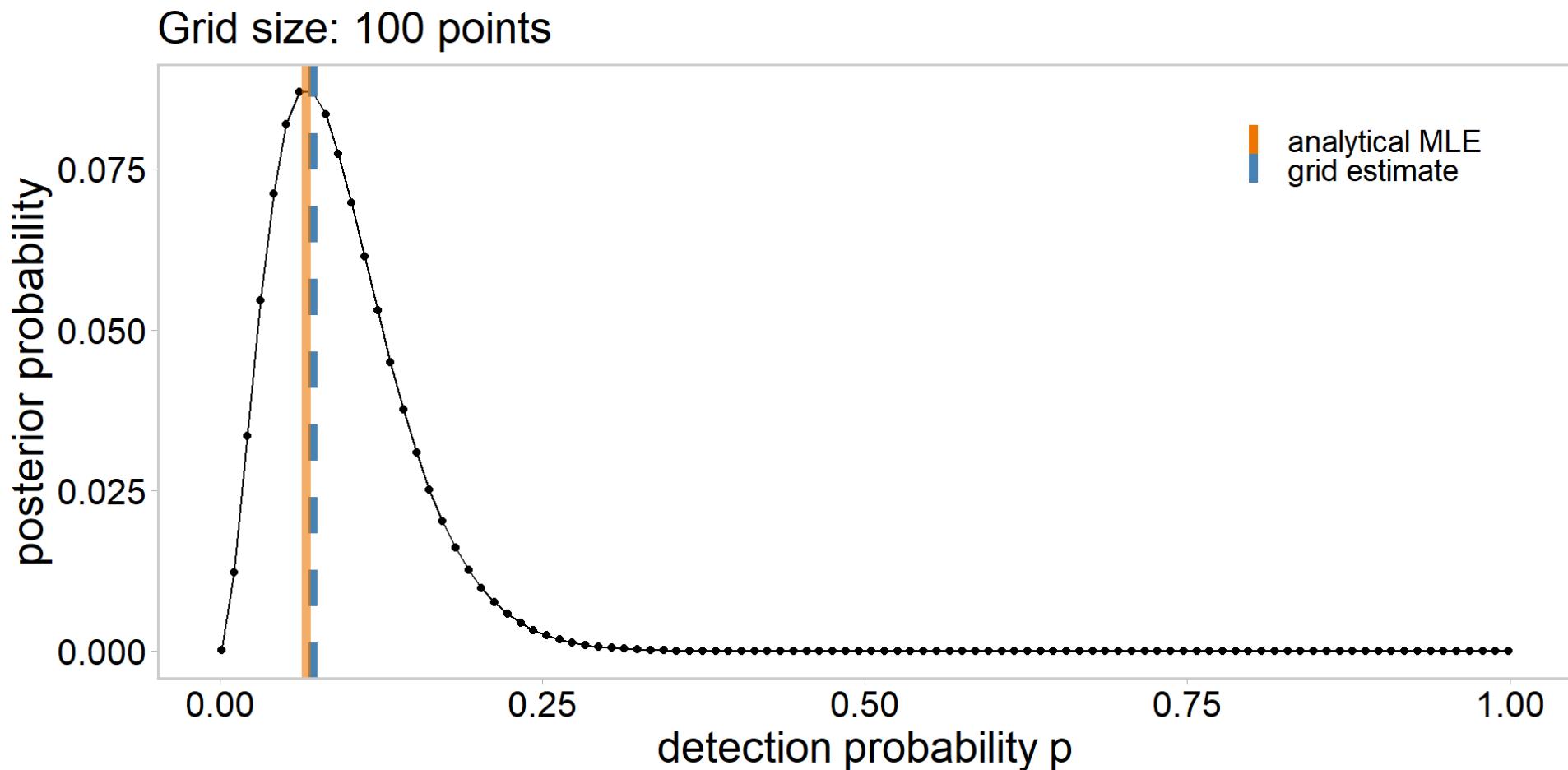
Grid size: 20 points



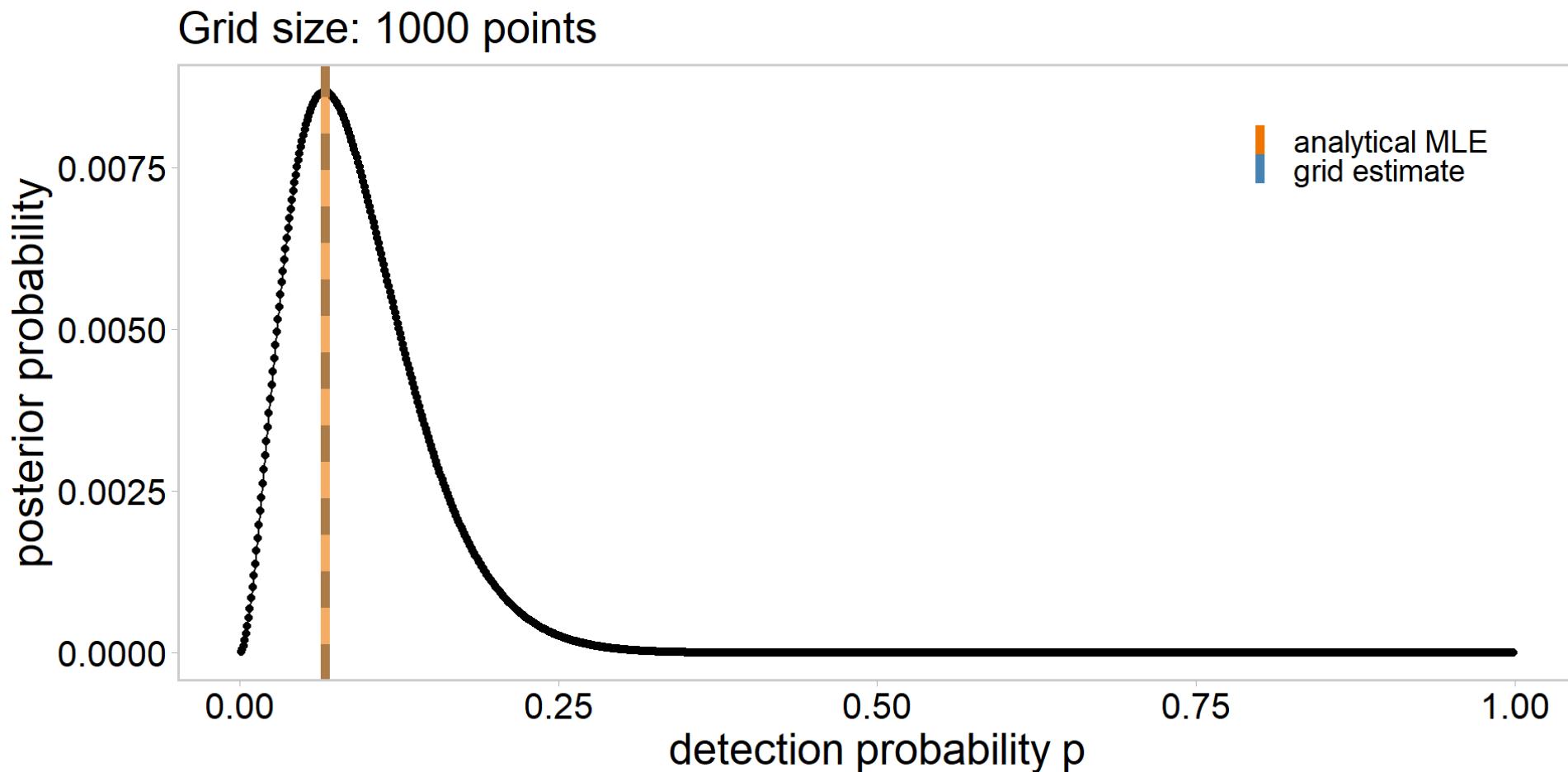
Comparing our grid approximation with the analytical MLE



Comparing our grid approximation with the analytical MLE



Comparing our grid approximation with the analytical MLE



Exercise

- A new study estimated the detection probability of sneak turtles, and you want to incorporate this information as prior information into your analysis
- Repeat the analysis above in R with a grid_size of 10000 and a $p \sim \text{Normal}(0.13, 0.12)$
- Report your best estimate of log-posterior and p given this prior
- If you get this far, try a bunch of different normal priors and think about how it influences your findings

Solution

	p	log_lik	log_prior	log_unstd_post	post_prob
747	0.07545825	-1.289803	1.098033	-0.1917702	0.0008715062

Why can't we just grid search everything?

- There are many important model types (e.g., mixed effects models) for which analytical solutions and things like grid approximation fail miserably
 - For models with many parameters and a dense grid, computational burden is intractable
 - This includes almost all models relevant to applied ecology

The Bayesian problem

$$P(\theta \mid y) = \frac{P(y \mid \theta)P(\theta)}{P(y)}$$

- $P(y)$ is called the “evidence” (i.e. the evidence that the data y was generated by this model)
- We can compute this quantity by integrating over all possible parameter values:

$$P(y) = \int_{\Theta} P(y, \theta) d\theta$$

The Bayesian problem cont'd

- This is *the* key difficulty with Bayesian statistics, and is why Bayesian statistics was not popular until computing power increased in the 1990s
 - Even for relatively simple models this thing can be intractable
- If we can't solve it, can we approximate it via **Monte Carlo**?
 - Monte Carlo = repeated random sampling
- Unfortunately this would require us to first solve Bayes' theorem because we don't know what the distribution $P(\theta | y)$ is

The Bayesian problem cont'd

- So we can't compute it using standard methods, can't sample directly from it
- Mathematicians realized this, and instead said:
 - "Let's construct an **ergodic**, reversible Markov chain that has as an equilibrium distribution which matches our posterior distribution"
- What? 🦇 💩
- Markov chain: what happens next depends only on the state of affairs now
- The remarkable fact is that this is relatively easy to do

Introduction to Markov Chain Monte Carlo (MCMC)

- What is MCMC: a broad class of algorithms that allow us to sample from an arbitrary probability distribution
- Counter-intuitive to most non-robots
- Rather than computing the $P(\theta \mid y)$ directly, MCMC draws samples of parameters from $P(\theta \mid y)$
 - This is actually **extremely** useful
- Collecting enough samples allows us to ‘map out’ or get a picture of $P(\theta \mid y)$

The Metropolis-Hastings algorithm

- Metropolis-Hastings (MH) is a classic MCMC algorithm with the following rules:
- Let $f(x)$ be a function that is *proportional* to the desired probability density function $P(x)$ (a.k.a the “target distribution”)
- Remember that we know:

$$\text{posterior} \propto \text{likelihood} \cdot \text{prior}$$

- So we just set $f(x) = P(y \mid \theta) \cdot P(\theta)$

The Metropolis-Hastings algorithm cont'd

The recipe:

- Choose an arbitrary starting value for a parameter x_t
- Choose an arbitrary probability distribution $g(x \mid y)$ to *propose* the next parameter value x given the previous value y
 - Note $g(x \mid y)$ is usually a Gaussian distribution centered at y so that points closer to y are more likely to be visited next (random walk)

The Metropolis-Hastings algorithm cont'd

- for each iteration t
 - generate a proposed value for x' by drawing a random number from $g(x' | x_t)$
 - calculate the acceptance probability (a.k.a. acceptance ratio) $\alpha = f(x') / f(x_t)$
 - accept or reject the proposed value x'
 - generate a uniform random number $u \in [0, 1]$
 - if $u > \alpha$, accept the proposal and set $x_{+1} = x'$
 - if $u < \alpha$, reject the proposal and set $x_{+1} = x_t$

Re-thinking the acceptance probability

This thing works because

$$\alpha = f(x') / f(x_t)$$

is the same as saying

$$\alpha = \frac{\frac{P(y|\theta')P(\theta')}{P(y)}}{\frac{P(y|\theta_t)P(\theta_t)}{P(y)}}$$

Note we just substituted x' and x_t for θ' and θ_t

- $P(y)$ cancels out via mathemagicks

Remember the sneak turtles !



1 y

1 n

[1] 30

Building a Metropolis-Hastings algorithm for sneak turtles



First, write three helper functions:

```
1 l_prior <- function(param) { # log(prior)
2   dunif(param, min = 0, max = 1, log = T) # uniform(0,1) prior
3 }
4
5 l_lik <- function(param) { # log(likelihood)
6   dbinom(sum(y), n, param, log = T)
7 }
8
9 l_post <- function(param) { # log(posterior) = log(likelihood) + log(prior)
10   l_lik(param) + l_prior(param)
11 }
```

Building a Metropolis-Hastings algorithm for sneak turtles



Second, declare some stuff

```
1 n_iter <- 10000 # number of iterations to run the MCMC
2 chain <- rep(NA, n_iter) # place to store values of p
3 chain[1] <- 0.5 # initial guess at p
```

Building a Metropolis-Hastings algorithm for sneak turtles



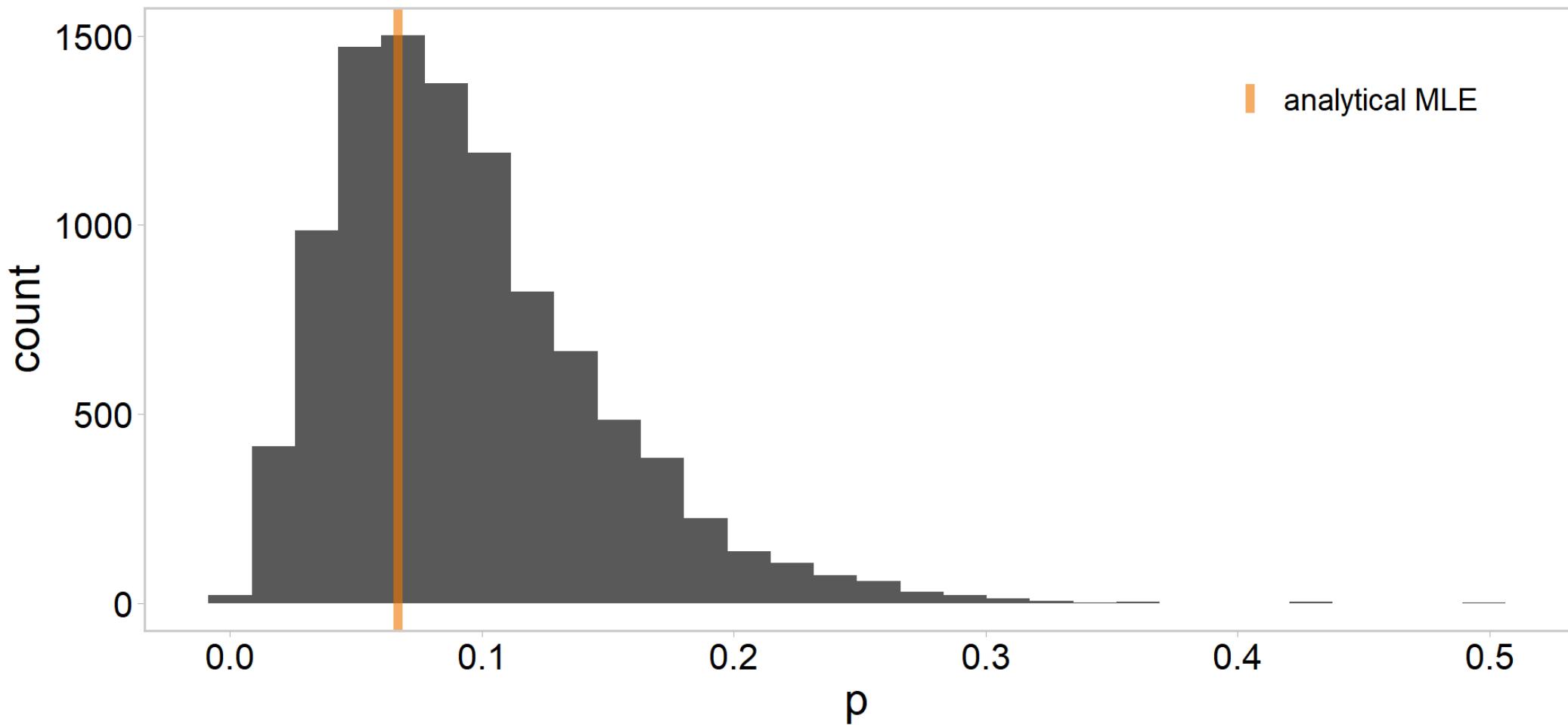
Now we run the algorithm

```
1 for (i in 2:n_iter) {  
2   p_new <- rnorm(1, chain[i - 1], 0.1) # new proposed p based on previous p  
3  
4   # next two lines ensure new p stays between [0,1] - ignore them  
5   if (p_new < 0) p_new <- abs(p_new)  
6   if (p_new > 1) p_new <- 2 - p_new  
7  
8   a_prob <- exp(l_post(p_new) - l_post(chain[i - 1])) # acceptance prob  
9  
10  if (runif(1) < a_prob) {  
11    chain[i] <- p_new # accept the proposal  
12  } else {  
13    chain[i] <- chain[i - 1] # reject the proposal  
14  }  
15 }
```

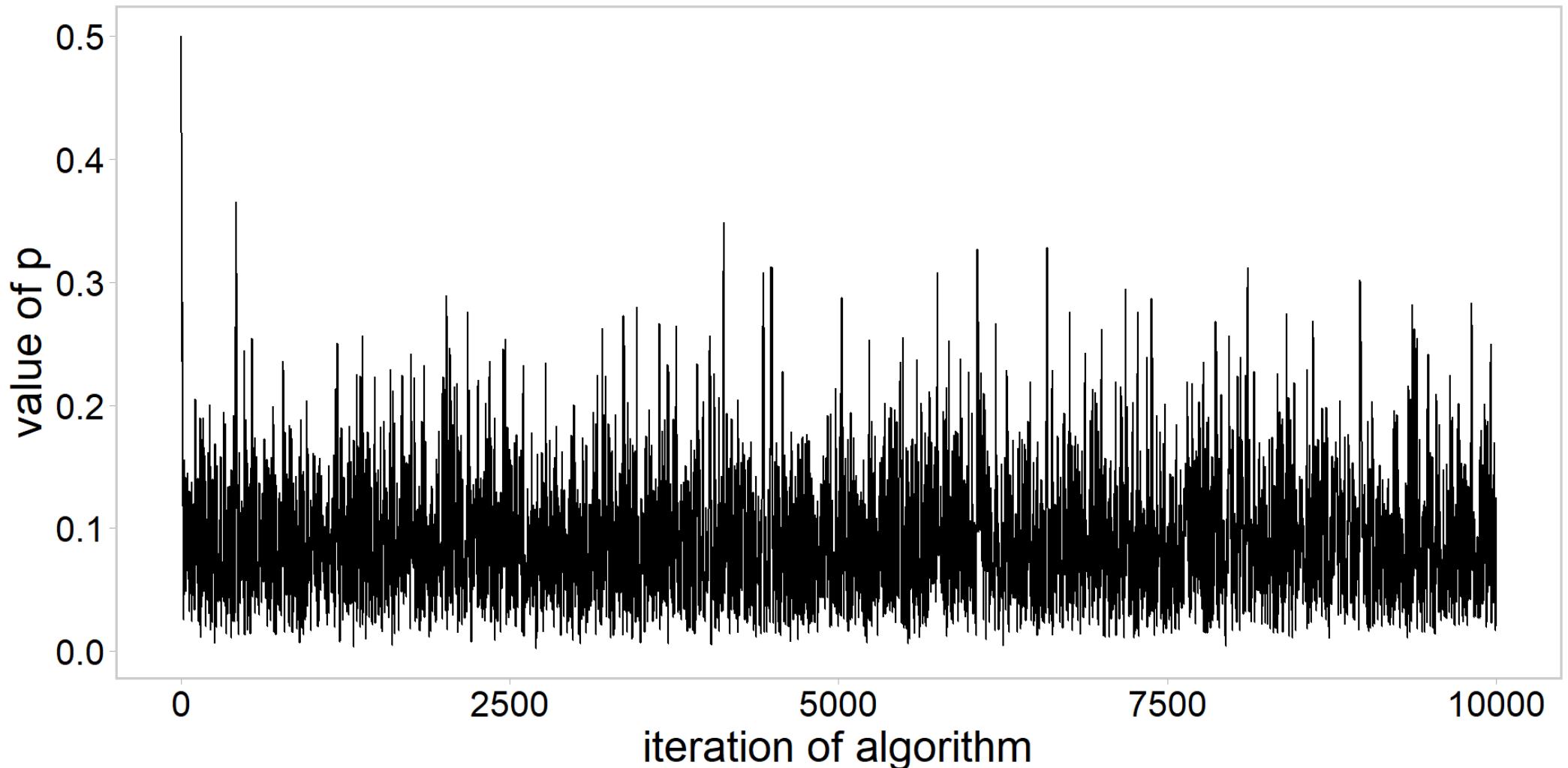
see also Gelman et al. 2003; Robert and Casella 2010

Metropolis-Hastings results

Gray histogram shows draws of p



Metropolis-Hastings results



Exercise

Get in groups of three, play with the following:

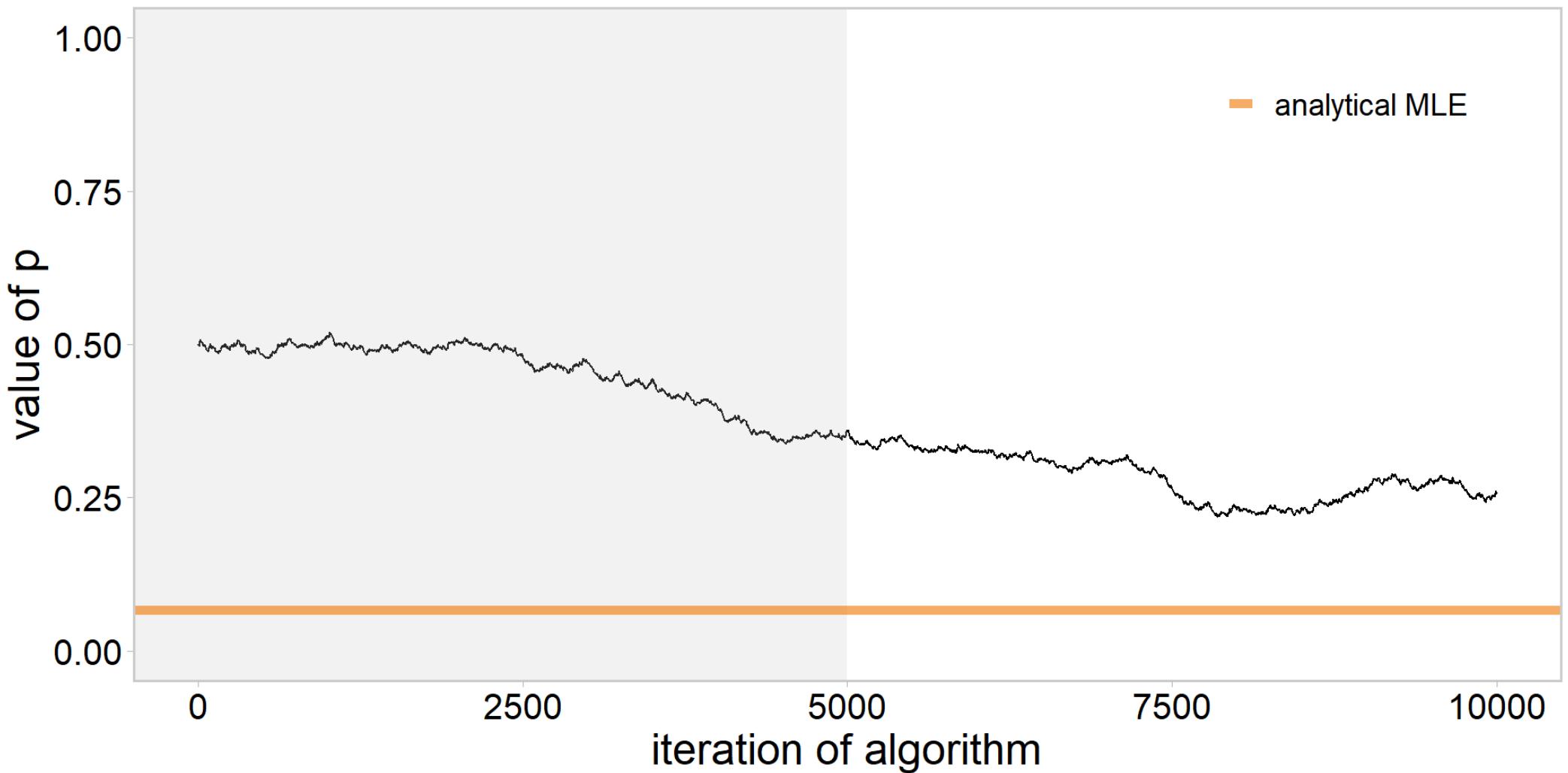
- the proposal distribution
- the initialization value
- chain length
- plot histogram of p and the chain (the fuzzy caterpillar plot)

Report your findings back to the group

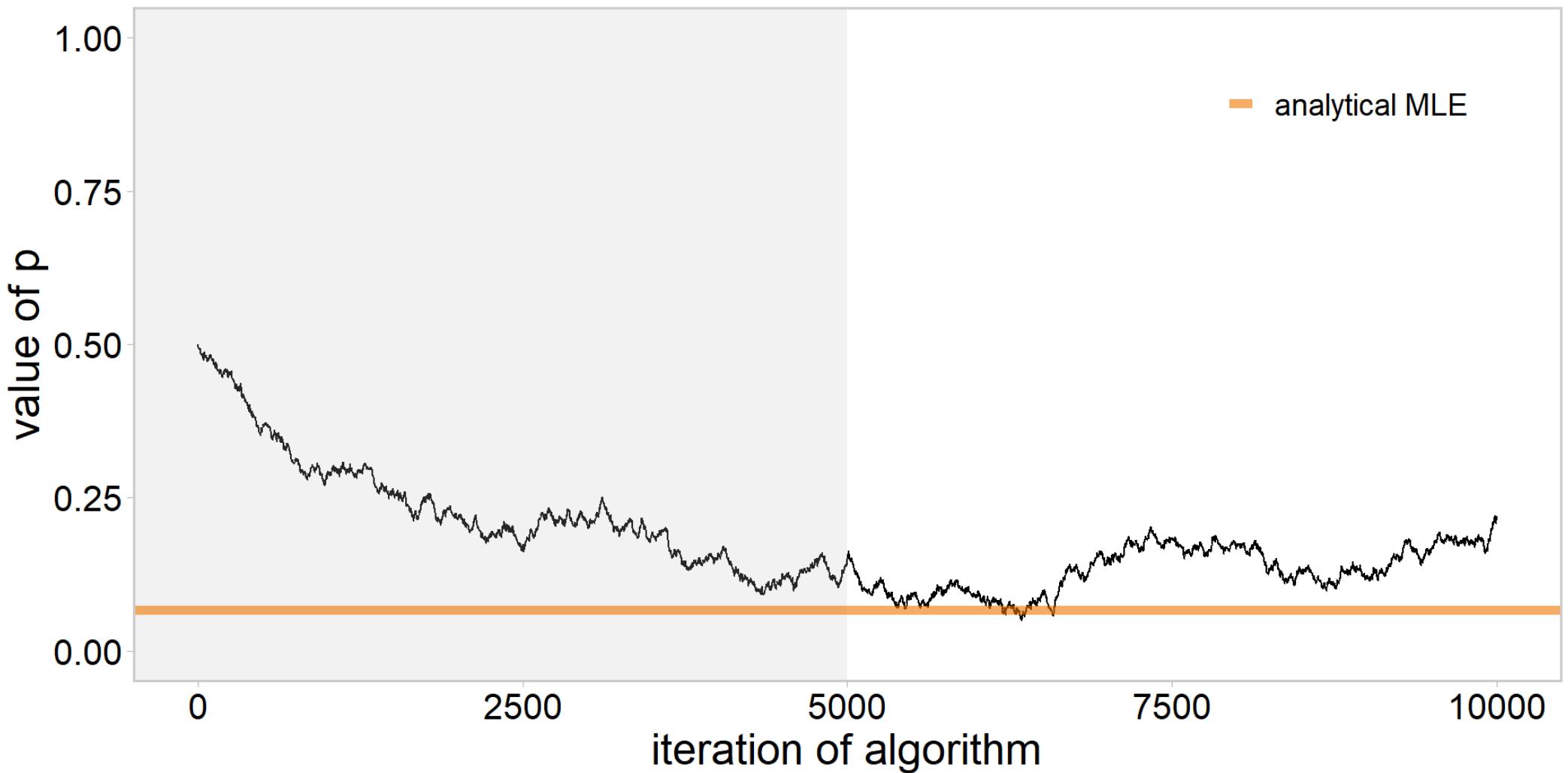
MCMC algorithm diagnostics

- Many of the parameters we were playing with affect our approximation of the posterior!
- Remember, this is a dead-simple example. Eventually we will use MCMC to approximate high-dimensional integrals which is much more difficult
- Clearly takes some time to “converge” on a stable posterior distribution
 - With MH MCMC this period between initial values and convergence to some distribution is known as the “burn-in” period, and we usually discard it

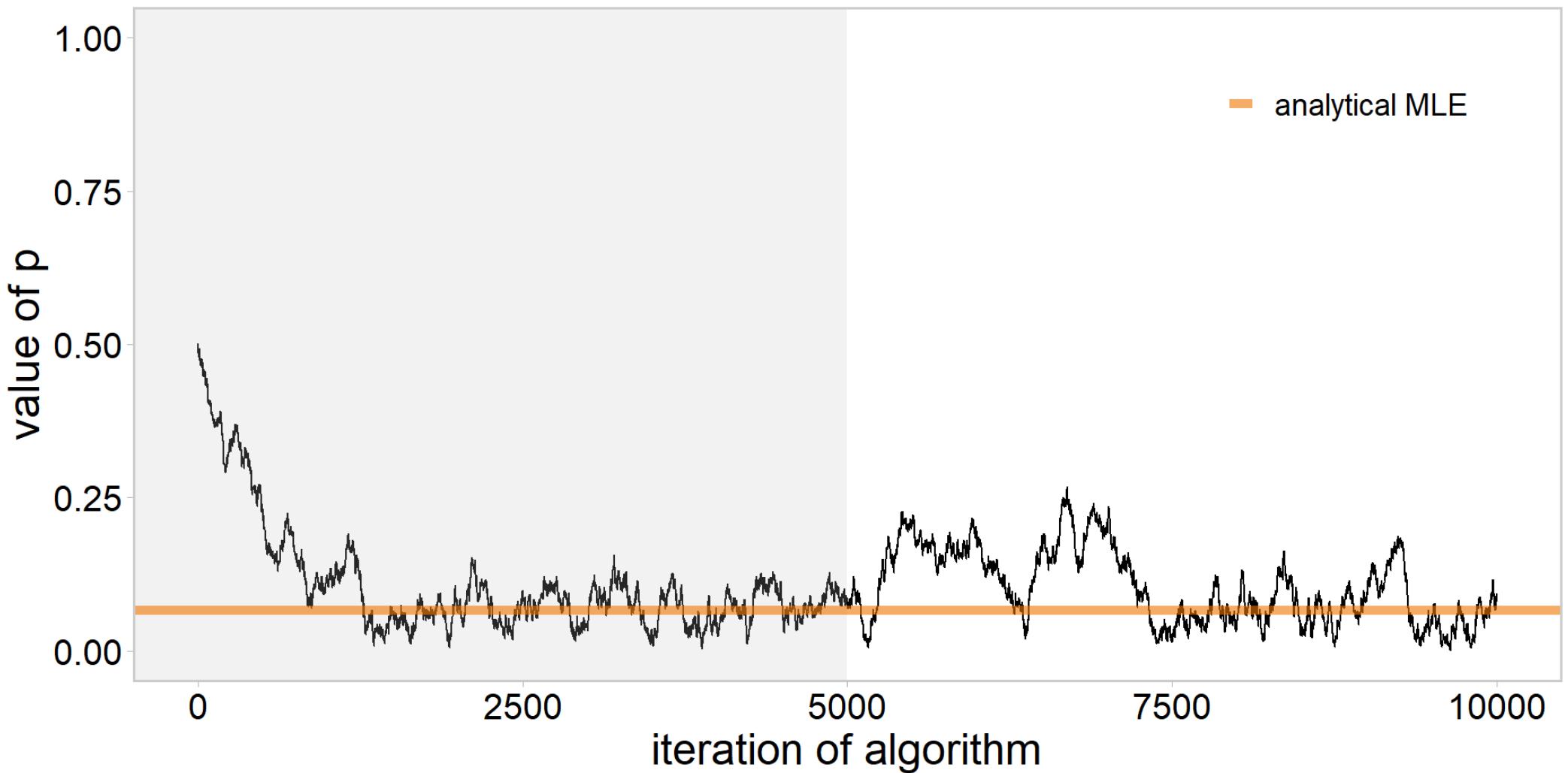
Proposal distribution sd = 0.001



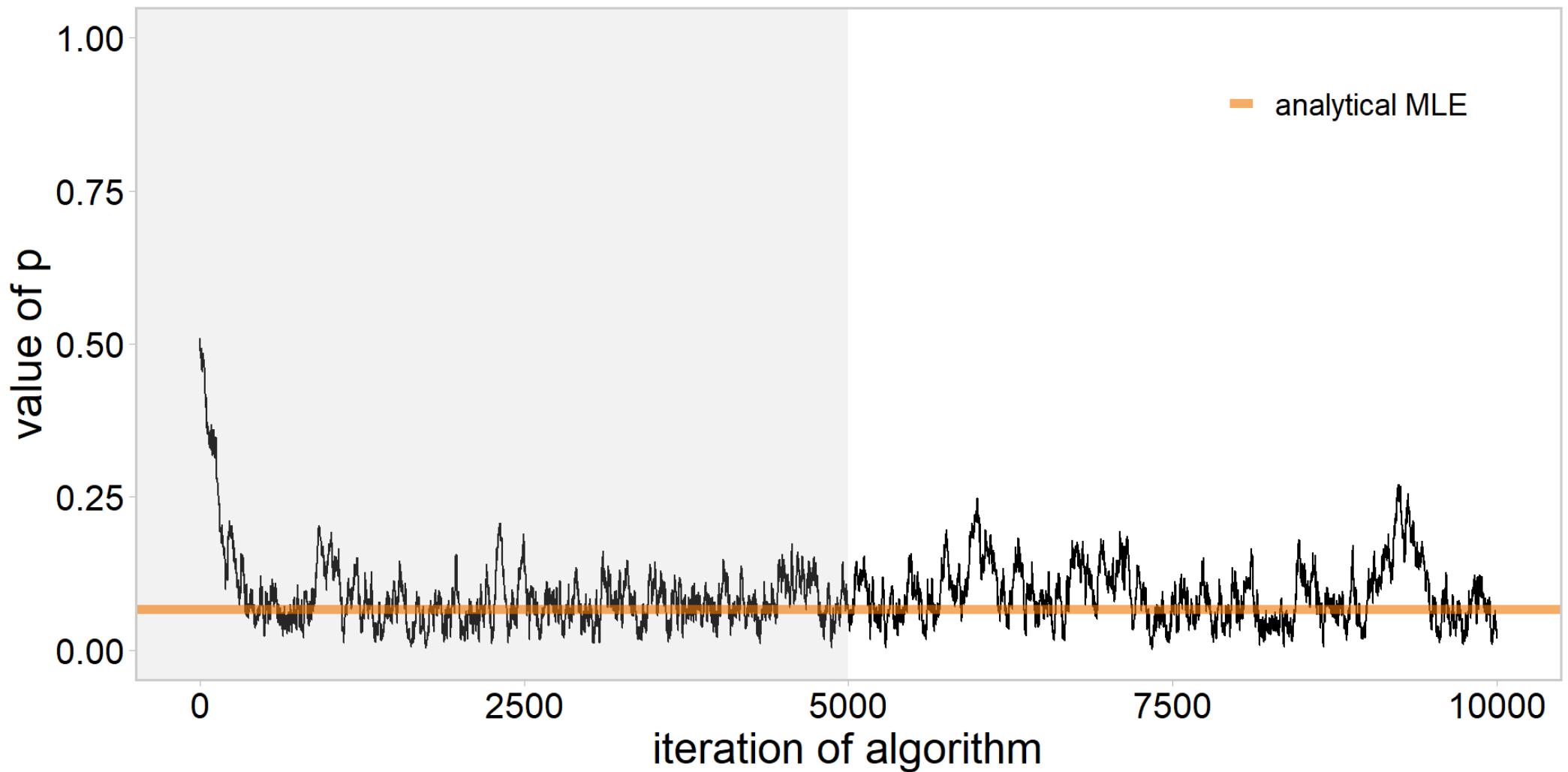
Proposal distribution sd = 0.002



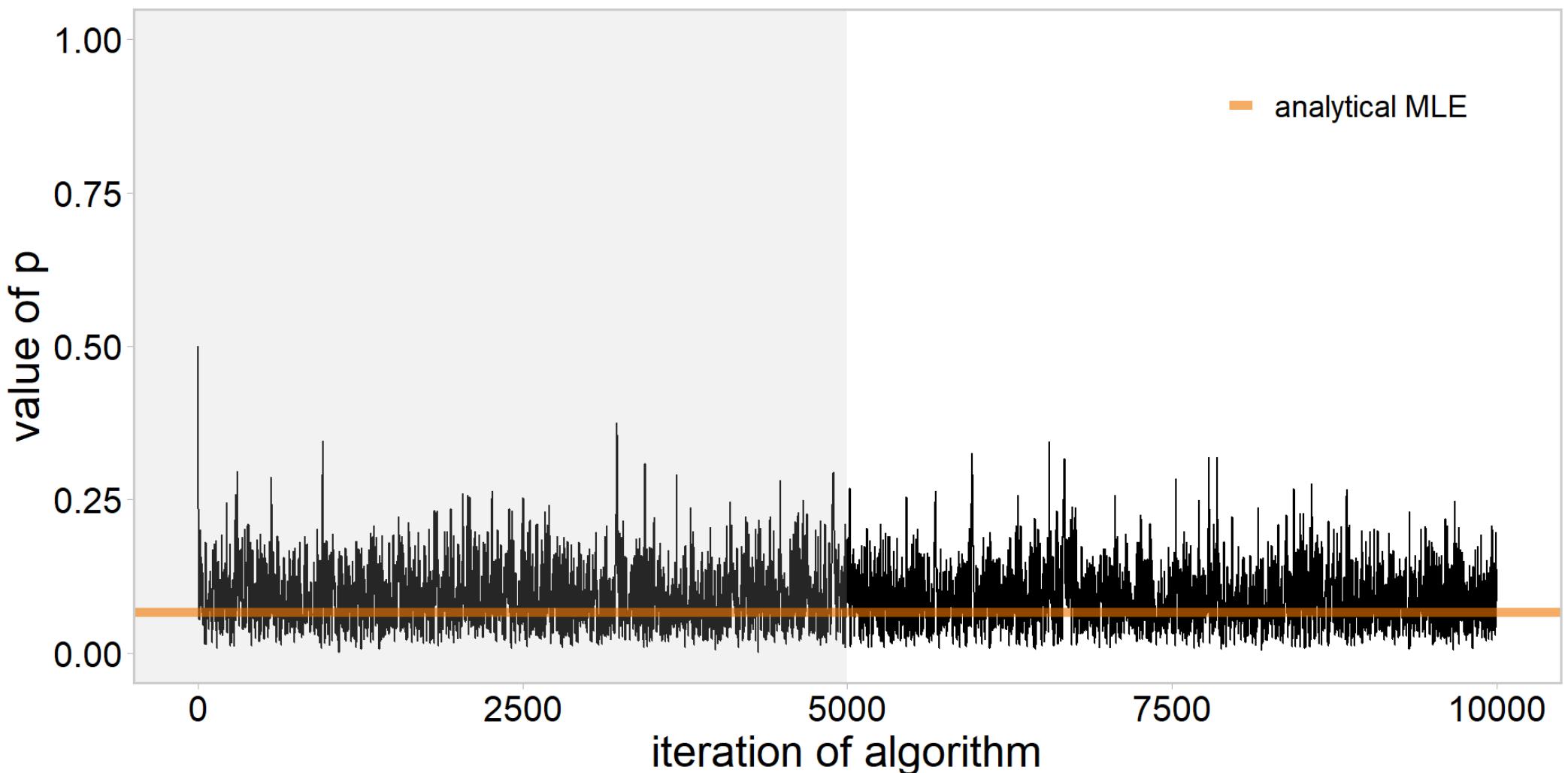
Proposal distribution sd = 0.005



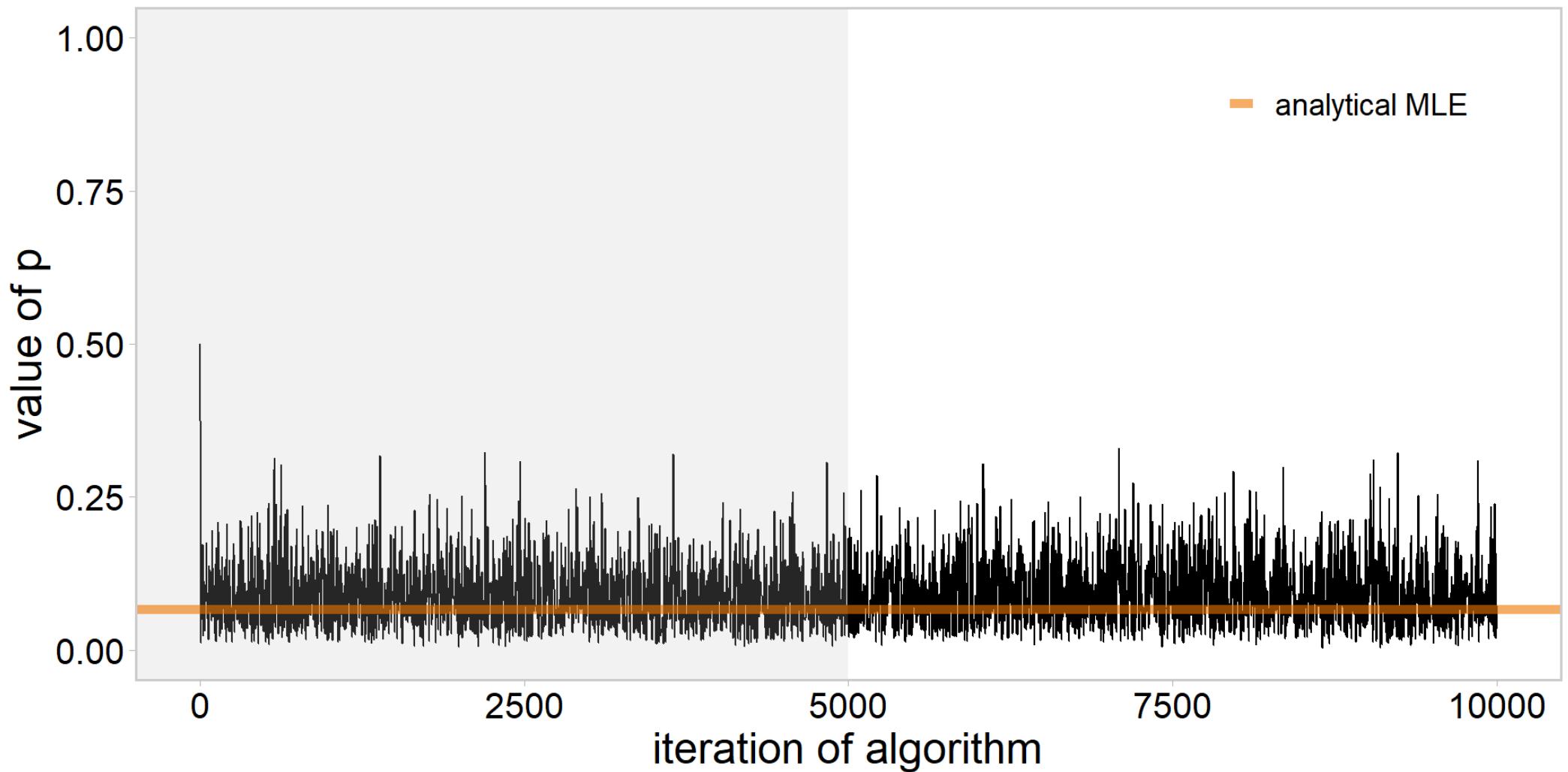
Proposal distribution sd = 0.01



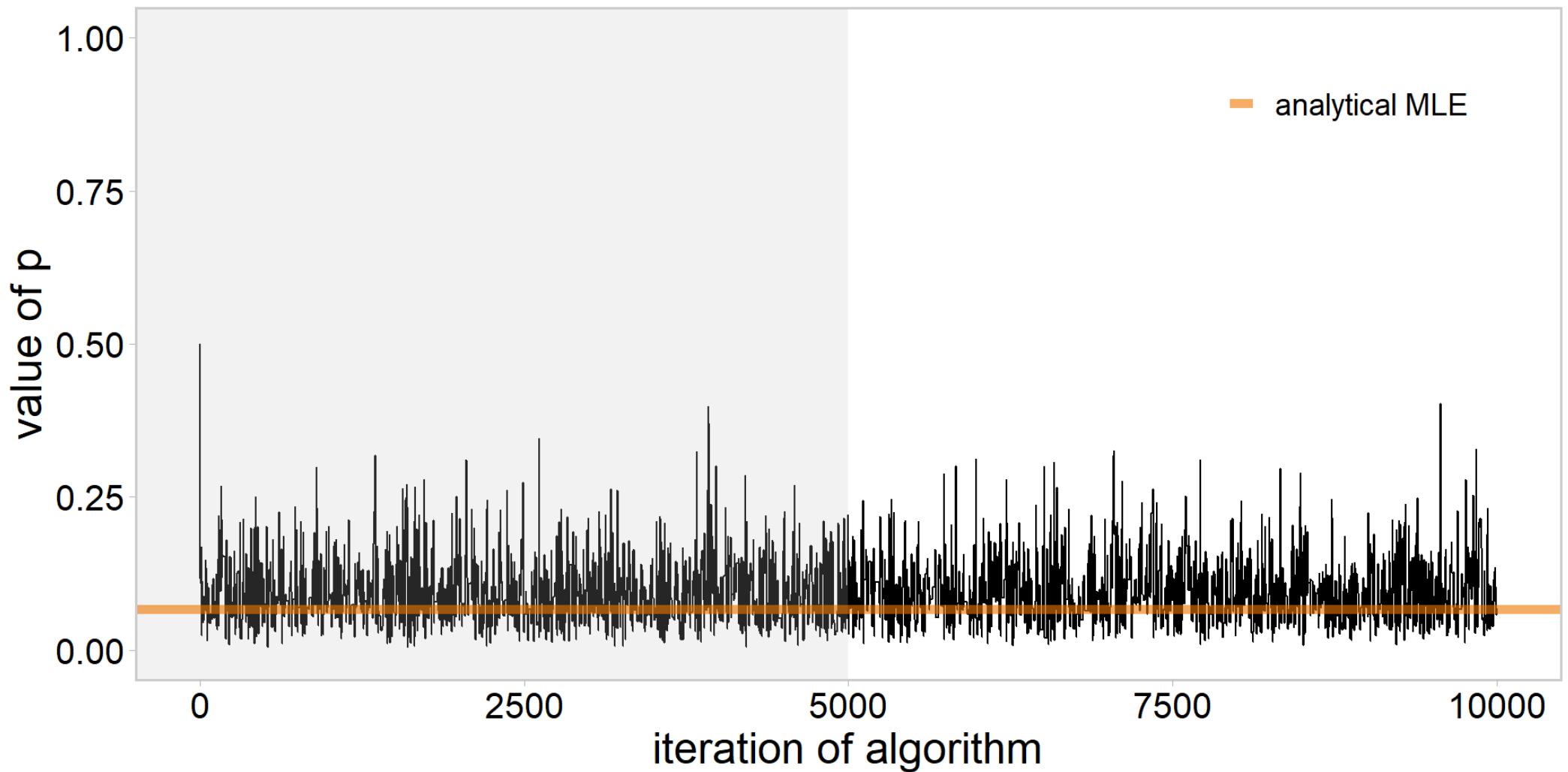
Proposal distribution sd = 0.1



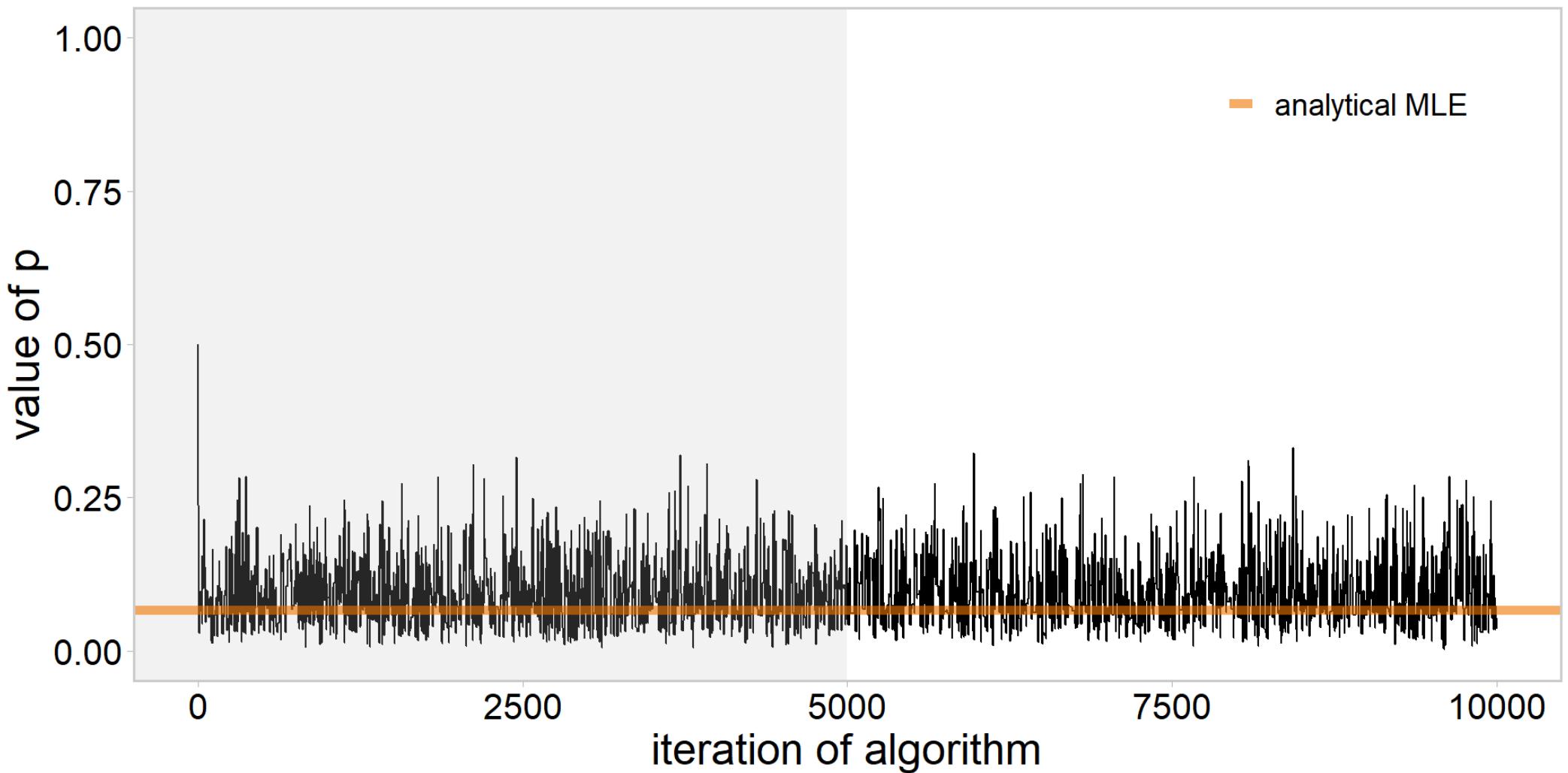
Proposal distribution sd = 0.2



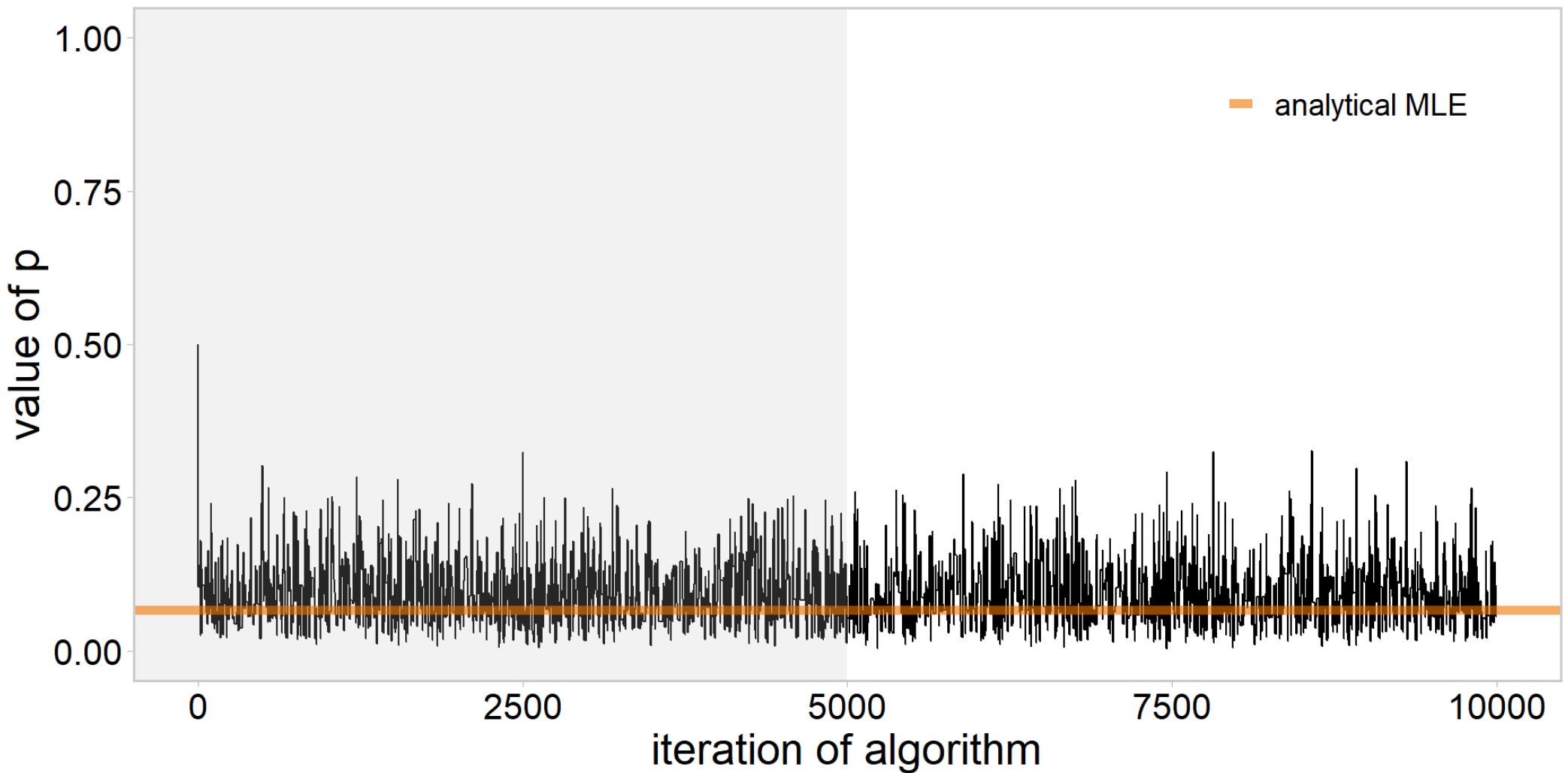
Proposal distribution sd = 0.4



Proposal distribution sd = 0.42



Proposal distribution sd = 0.44



Proposal distribution $sd > 0.44$

Breaks the following lines and dies because gives p outside [0,1]

```
1 if (p_new < 0) p_new <- abs(p_new)
2 if (p_new > 1) p_new <- 2 - p_new
```

- There are smarter ways to keep a value between [0,1]
- Key points:
 - tuning the algorithm impacts its ability to sample the posterior
 - correlation in the Markov chain means you don't get independent samples of the parameters of interest

What do we want the chains to do

- We want chains to look like fuzzy caterpillars

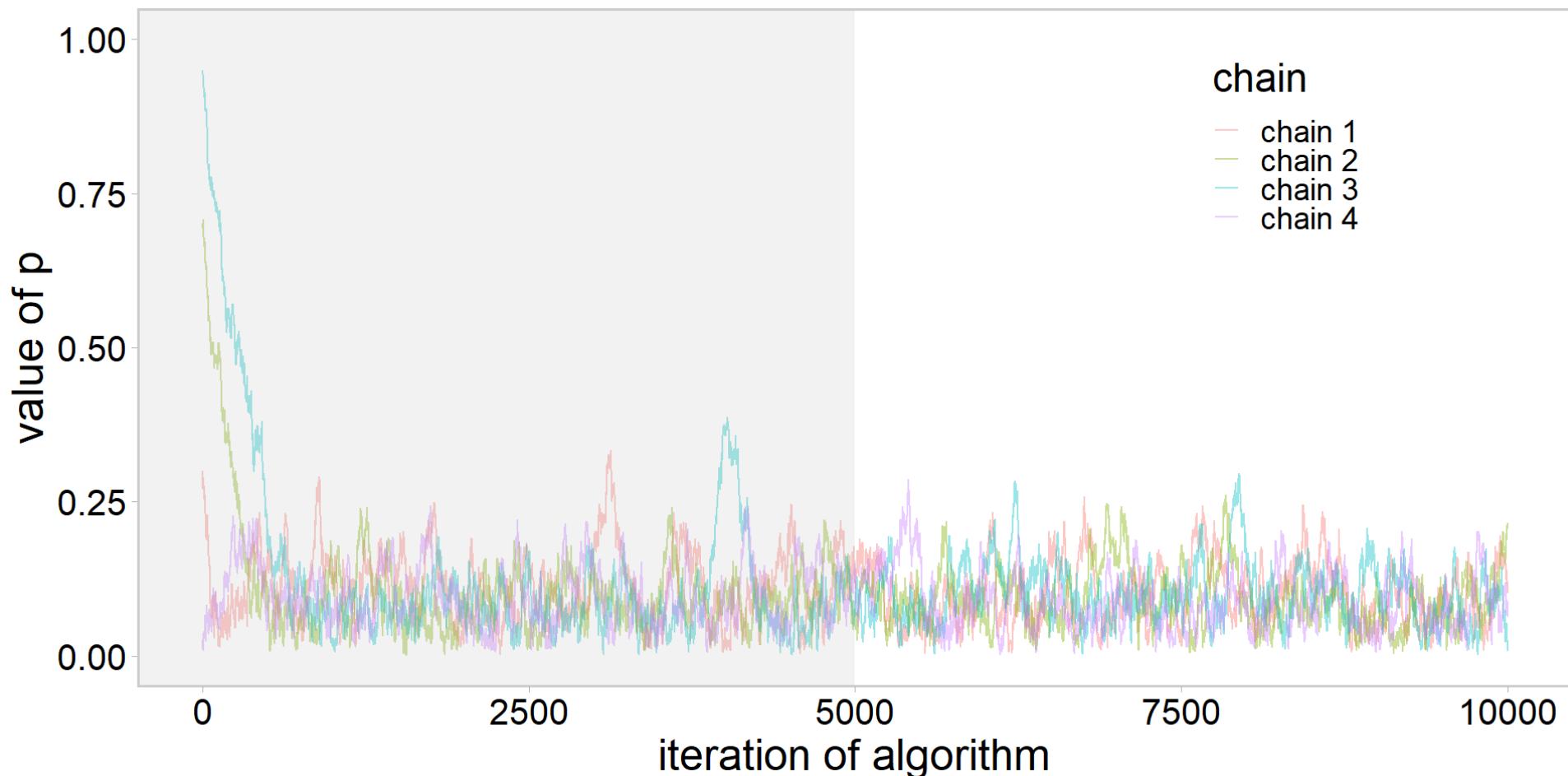


- Good indication that the MCMC is **efficiently** sampling from a maximum in the underlying distribution

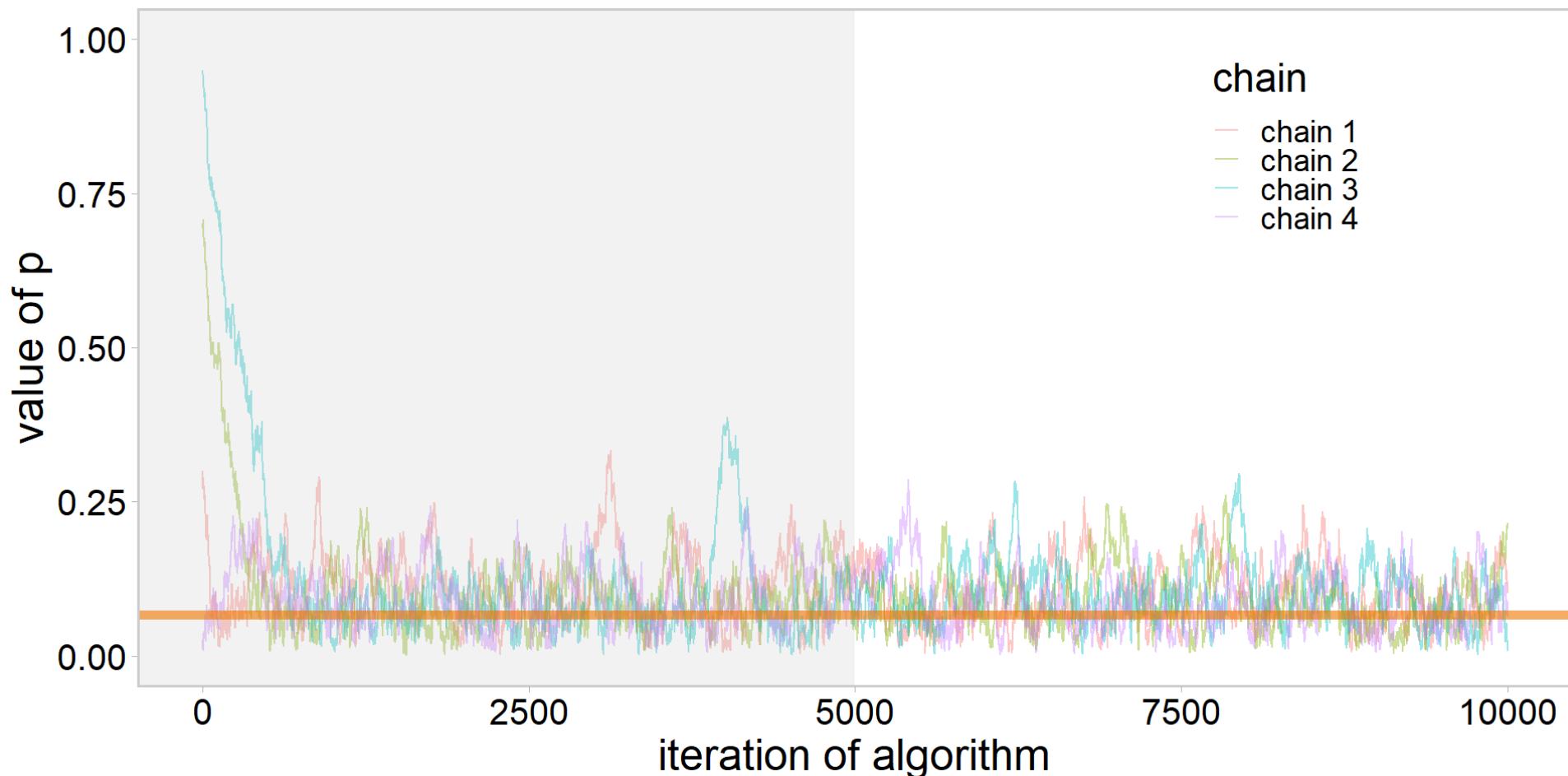
Fuzzy caterpillars are good indication that a chain is behaving, but...

- We could be stuck in a local maximum
 - If we ran it longer we might find a better solution
 - One way to deal: run multiple, independent chains initialized from different starting locations
- Should get overlapping fuzzy caterpillars if everything goes well
- Because chains are independent, we often run them in parallel

Running multiple chains each with different initialization values



Running multiple chains each with different initialization values



Software that can run MCMC for you

- BUGS/JAGS
- NIMBLE (de Valpine et al. 2017)
- MCMCpack
- LaplacesDemon
- Stan (what we will use in this class)
- Many others

These tools implement many different MCMC routines
Ecologists can (mostly) treat these tools as black boxes

Hamiltonian Monte Carlo (HMC)

The following draws heavily from Chapter 9 in McElreath 2023

It appears to be a quite general principle that, wherever there is a randomized way of doing something, then there is a nonrandomized way that delivers better performance but requires more thought-E.T. Jaynes

- No free lunch



Hamiltonian Monte Carlo (HMC)

- HMC is much more computationally demanding than MH, but its proposals are much more efficient
- HMC requires fewer samples to map out the posterior distribution

The bottom line:

You need less computer time in total, even though each sample requires more time than simpler algorithms

- Really outshines other algorithms for models with 10s to 1000s of parameters
- Don't need to know everything (lots of math in references)

HMC: understanding some concepts

- Leverages tools from **Hamiltonian dynamics** and **differential geometry** to generate better proposals
- Uses information on the gradient (curvature) of the log-posterior
- HMC creates a θ' that is uncorrelated with θ_t and which has a high probability of being accepted
 - Does this using a physics simulation and pretends parameters are tiny, frictionless particles
- See also Neal (2011) **MCMC using Hamiltonian Dynamics**

HMC hockey puck analogy

- Imagine you have a standard two parameter model
 - This creates a posterior shaped like a bowl
- In this case, HMC is like a hockey puck launched in random directions around an ice rink shaped like a bowl (the posterior distribution)

[Go to MCMC interactive gallery](#)

Stan does the HMC sneakery for you

- In the Stan language, stochastic models are described by specifying stochastic or deterministic relationships between quantities such as parameters and data
- Stan (software) constructs the log-posterior for you if you specify prior(s) and likelihood(s)
- Applies HMC to your problem, reports chains, log-posterior, and diagnostics



Summary and outlook

Posteriors are difficult to approximate:

1. Analytical solutions (beta-binomial)
2. Brute force (grid search approximations)
3. MCMC (draw samples from the posterior)

We have laid the ground work, and now we get to play and fit Bayesian models in Stan



References

1. Gelman et al. 2003. Bayesian Data Analysis. Appendix C.
2. Green et al. 2020. Introduction to Bayesian Methods in Ecology and Natural Resources. Appendix B.
3. Kery and Schaub. 2012. Bayesian Population Analysis using WinBUGS.
4. McElreath 2023. Statistical Rethinking. Second Edition, Chapters 2 and 9.
5. Neal 2011. MCMC using Hamiltonian Dynamics. In: Handbook of Markov Chain Monte Carlo.
6. Robert and Casella 2010. Introduction to Monte Carlo methods with R.

Useful web link:

[MCMC visualization app](#)

Michael Betancourt [seminar](#) on using Hamiltonian Monte Carlo for Bayesian inference