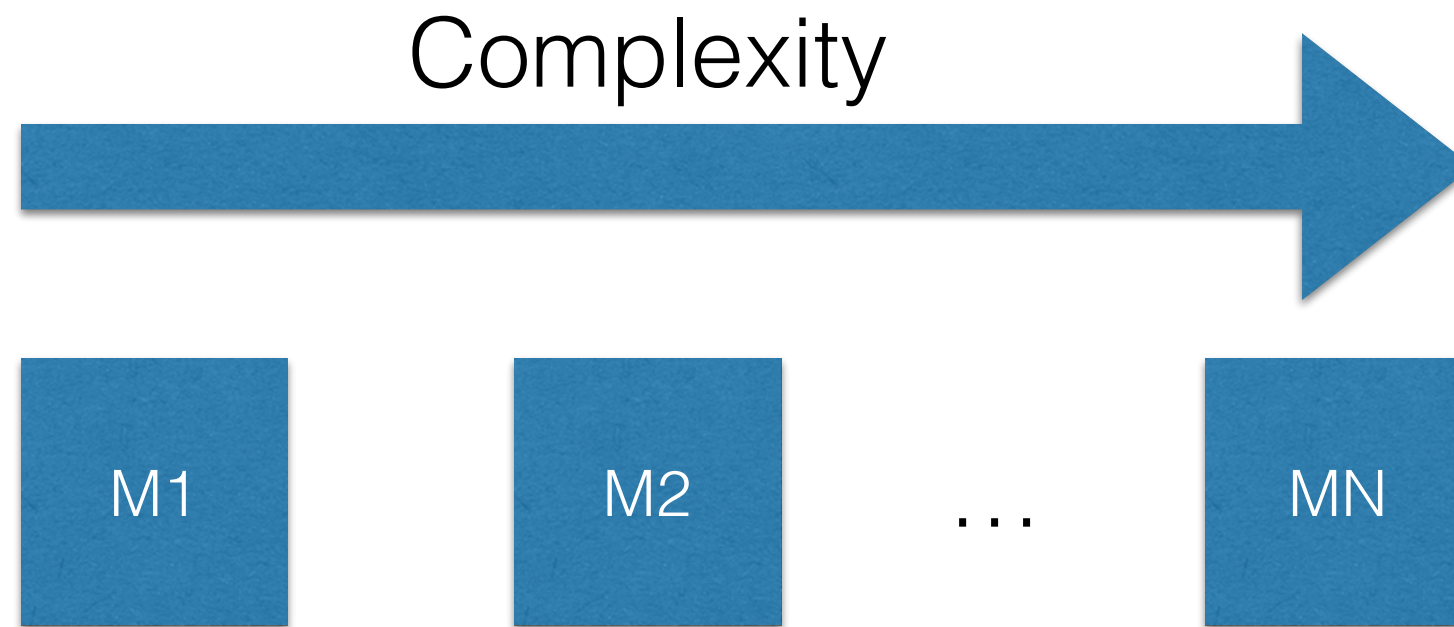


Lecture 26

Bayesian Model Comparison

The problem



What is the best model?

“We are to admit no more causes of natural things than such as are both true and sufficient to explain their appearances.”

–Isaac Newton

“Everything should be made as simple as possible, but not simpler.”

–Albert Einstein

Occam's Razor

Pick the simplest model that explains the data...

Example 1: Polynomial Regression

- M1: linear regression
- M2: quadratic regression
- M3: cubic regression
- ...

Example 2: Model Calibration

- M1: simple physical model
- M2: more complex physical model
- M3: super complex physical model
- ...

Example 3: Hypothesis Testing



AMERICAN STATISTICAL ASSOCIATION
Promoting the Practice and Profession of Statistics®

732 North Washington Street, Alexandria, VA 22314 • (703) 684-1221 • Toll Free: (888) 231-3473 • www.amstat.org • [www.twitter.com/AmstatNews](https://twitter.com/AmstatNews)

- H1

- H2

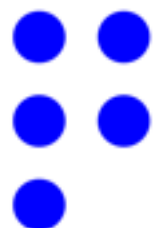
- H3

- ...

AMERICAN STATISTICAL ASSOCIATION RELEASES STATEMENT ON STATISTICAL SIGNIFICANCE AND P-VALUES

*Provides Principles to Improve the Conduct and Interpretation of Quantitative
Science*

March 7, 2016

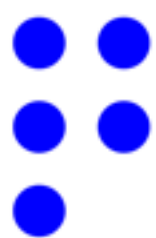


PREDICTIVE
SCIENCE LABORATORY

Example 3: Hypothesis Testing

The statement's six principles, many of which address misconceptions and misuse of the p -value, are the following:

1. *P -values can indicate how incompatible the data are with a specified statistical model.*
2. *P -values do not measure the probability that the studied hypothesis is true, or the probability that the data were produced by random chance alone.*
3. *Scientific conclusions and business or policy decisions should not be based only on whether a p -value passes a specific threshold.*
4. *Proper inference requires full reporting and transparency.*
5. *A p -value, or statistical significance, does not measure the size of an effect or the importance of a result.*
6. *By itself, a p -value does not provide a good measure of evidence regarding a model or hypothesis.*



Bayesian Model Selection

Prior over models:

$$M_i \sim p(M_i)$$

Prior
state of knowledge

Prior over model parameters:

$$\theta_i | M_i \sim p(\theta_i | M_i)$$

Likelihood of the data:

$$D | \theta_i, M_i \sim p(D | \theta_i, M_i)$$

Measurement
process

Bayesian Model Selection

Prior

state of knowledge

$$M_i \sim p(M_i)$$

$$\theta_i | M_i \sim p(\theta_i | M_i)$$

Measurement

process

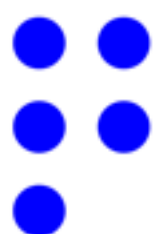
$$D | \theta_i, M_i \sim p(D | \theta_i, M_i)$$

Posterior over the parameters:

$$p(\theta_i | D, M_i) = \frac{p(D | \theta_i, M_i) p(\theta_i | M_i)}{p(D | M_i)}$$

where the **evidence** is the log of the normalizing constant:

$$p(D | M_i) = \int p(D | \theta_i, M_i) p(\theta_i | M_i) d\theta_i$$



Bayesian Model Selection

Prior

state of knowledge

$$M_i \sim p(M_i)$$

$$\theta_i | M_i \sim p(\theta_i | M_i)$$

$$p(\theta_i | D, M_i) = \frac{p(D | \theta_i, M_i) p(\theta_i | M_i)}{p(D | M_i)}$$

Measurement

process

$$D | \theta_i, M_i \sim p(D | \theta_i, M_i)$$

$$p(D | M_i) = \int p(D | \theta_i, M_i) p(\theta_i | M_i) d\theta_i$$

Posterior over models:

$$p(M_i | D) \propto p(D | M_i) p(M_i)$$

How do we calculate the model evidence is the big problem!

Model Averaging

Assume you have this:

$$p(M_i | D) \propto p(D | M_i) p(M_i)$$

How do you make predictions?

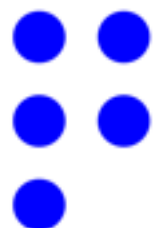
Predictive distribution:

$$p(y | x) = \sum_i p(y | x, D, M_i) p(M_i, D)$$



$$p(y | x, D_i, M_i) = \int p(y | x, \theta_i, M_i) p(\theta_i | D, M_i) d\theta_i$$

Model predictive distribution.



Bayesian Model Selection

Assume you have this:

$$M_{i^*} = \operatorname{argmax}_i p(M_i \mid D) = \operatorname{argmax}_i p(D \mid M_i) p(M_i)$$

How do you make predictions?

$$p(y \mid x, D, M_{i^*}) = \int p(y \mid x, \theta_{i^*}, M_{i^*}) p(\theta_{i^*} \mid D, M_{i^*}) d\theta_{i^*}$$

Model Evidence Calculation

$$p(D | M_i) = \int p(D | \theta_i, M_i) p(\theta_i | M_i) d\theta_i$$

Analytically...
if you are lucky enough

The Laplace Approximation

Why don't we just approximate this:

$$p(\theta_i | D, M_i) = \frac{p(D | \theta_i, M_i) p(\theta_i | M_i)}{p(D | M_i)}$$

by a Gaussian!

The Laplace Approximation

Let's do it for an arbitrary distribution:

$$p(z) = \frac{f(z)}{Z}$$

Use some Taylor expansion stuff:

$$\log f(z) \approx \log f(z_0) - \frac{1}{2}(z - z_0)^T \nabla^2 \log f(z_0)(z - z_0).$$

where:

$$z_0 = \operatorname{argmax}_z \log f(z).$$

Taking the exponential:

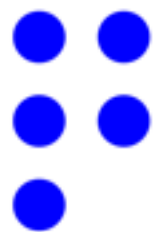
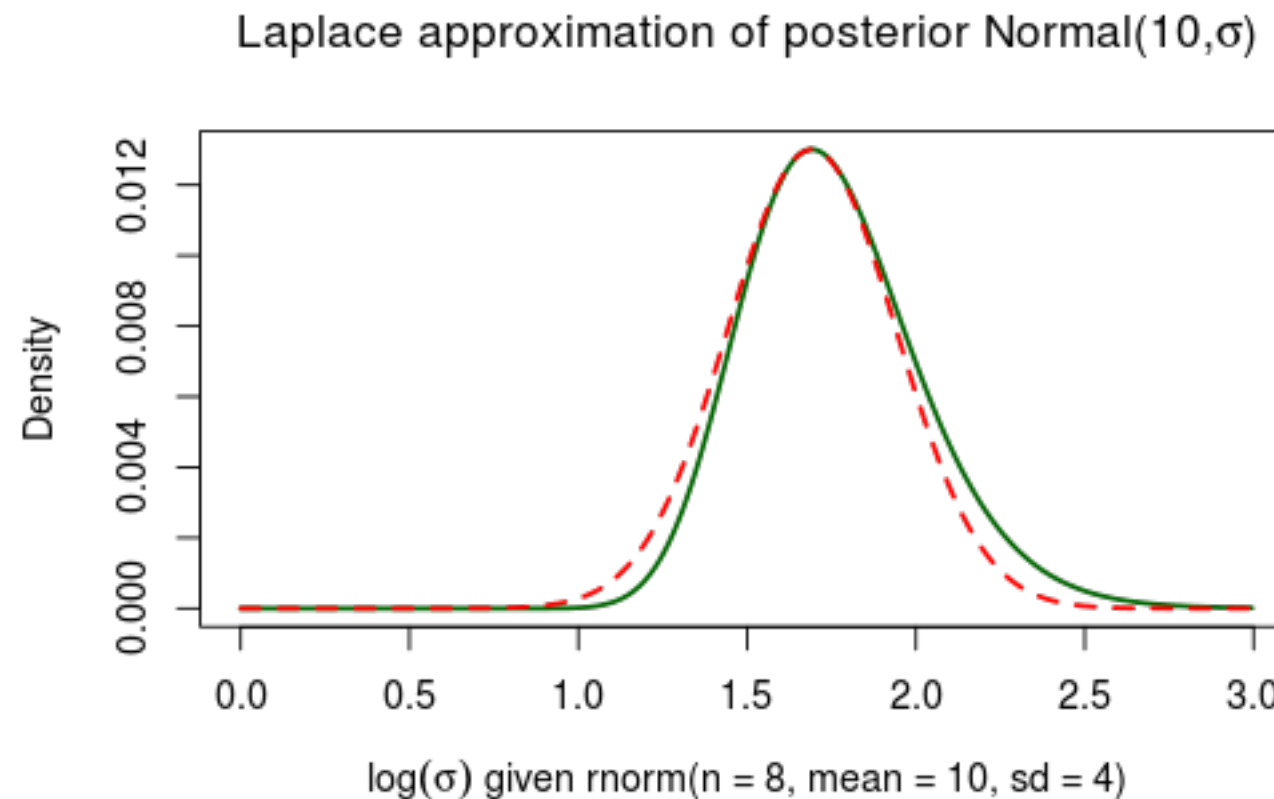
$$f(z) \approx f(z_0) \exp \left\{ -\frac{1}{2}(z - z_0)^T \nabla^2 \log f(z_0)(z - z_0) \right\} = \mathcal{N}(z | f(z_0), \nabla^2 \log f(z_0)^{-1})$$

Thus:

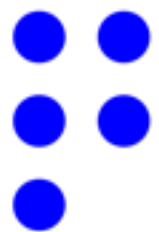
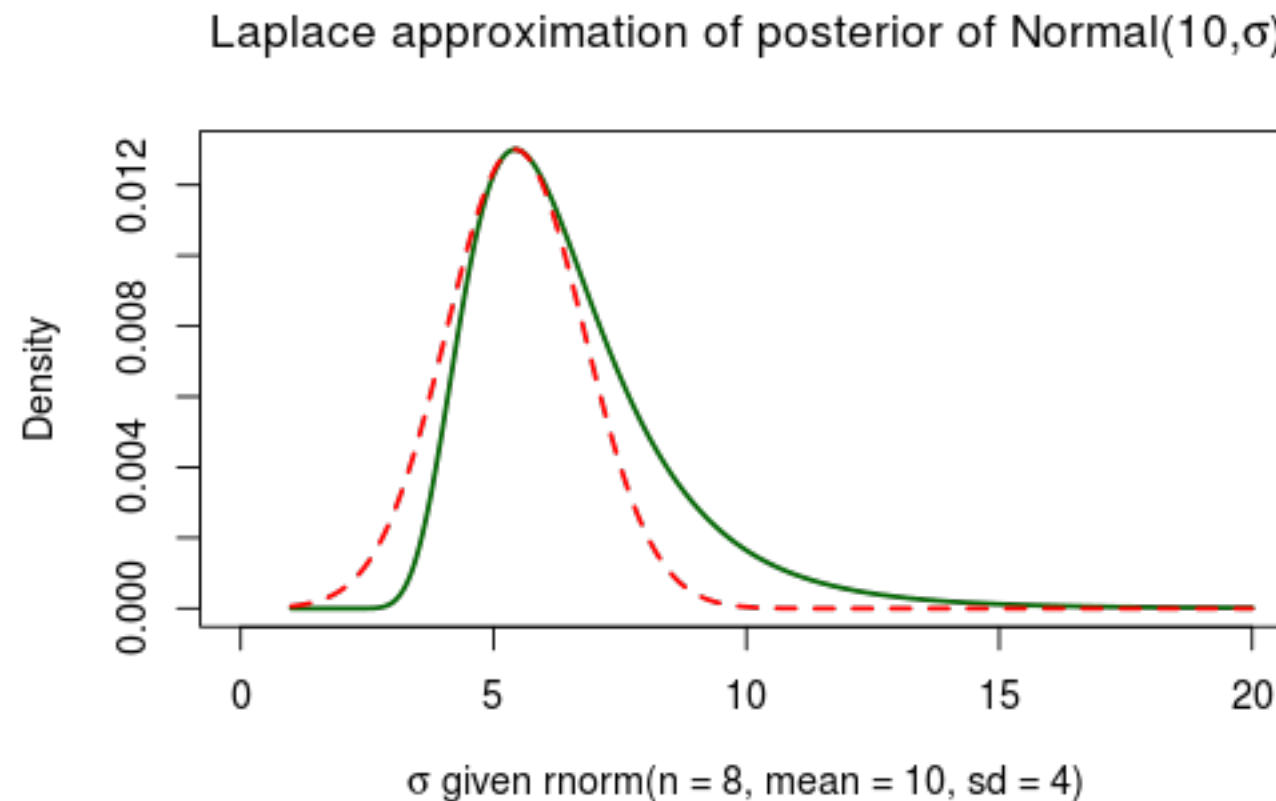
$$Z \approx \int f(z) dz = f(z_0) \frac{(2\pi)^{d/2}}{|\nabla^2 \log f(z_0)|^{1/2}}.$$

→ dimension of z

The Laplace Approximation



The Laplace Approximation



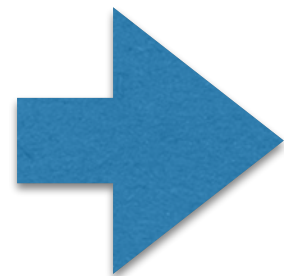
The Laplace Approximation

Why don't we just approximate this:

$$p(\theta_i | D, M_i) = \frac{p(D | \theta_i, M_i) p(\theta_i | M_i)}{p(D | M_i)}$$

by a Gaussian!

$$p(z) = \frac{f(z)}{Z}$$

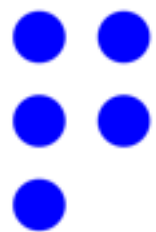


Laplace

$$\log Z \approx \log f(z_0) + \frac{1}{2} \log |\nabla^2 \log f(z_0)| + \frac{d}{2} \log 2\pi.$$

$$\log p(D | M_i) \approx \log(p(D | \theta_i^*, M_i) p(\theta_i^* | M_i)) + \frac{1}{2} \left| \nabla^2 \log(p(D | \theta_i^*, M_i) p(\theta_i^* | M_i)) \right| + \frac{d}{2} \log 2\pi$$

MAP estimate



Bayesian Information Criterion

$$\log p(D | M_i) \approx \log(p(D | \theta_i^*, M_i) p(\theta_i^* | M_i)) + \frac{1}{2} \left| \nabla^2 \log(p(D | \theta_i^*, M_i) p(\theta_i^* | M_i)) \right| + \frac{d}{2} \log 2\pi$$

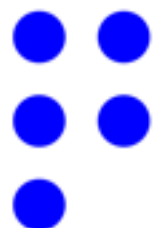


Crudest possible approximation

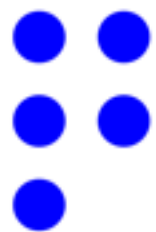
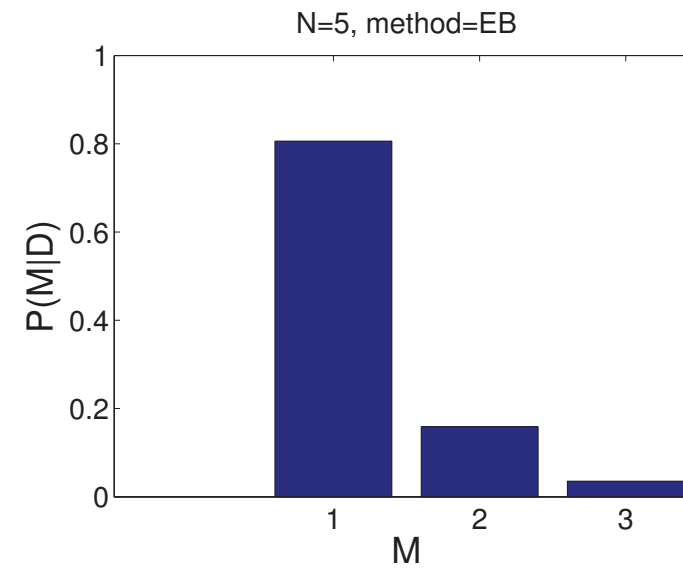
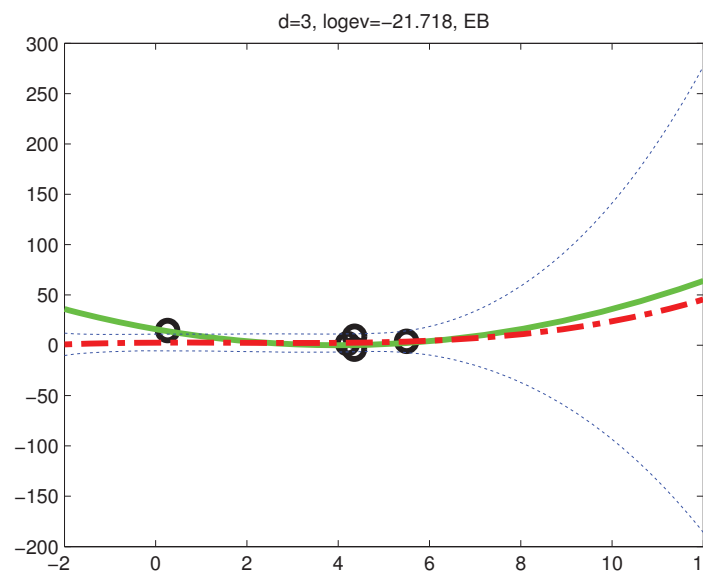
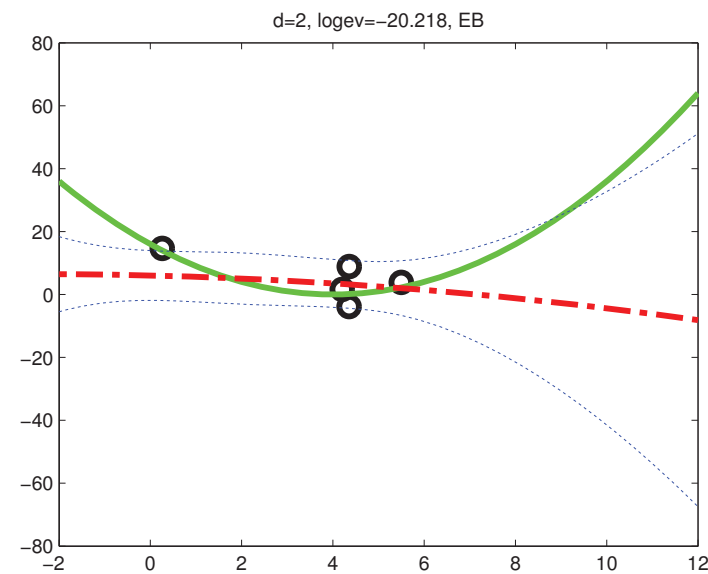
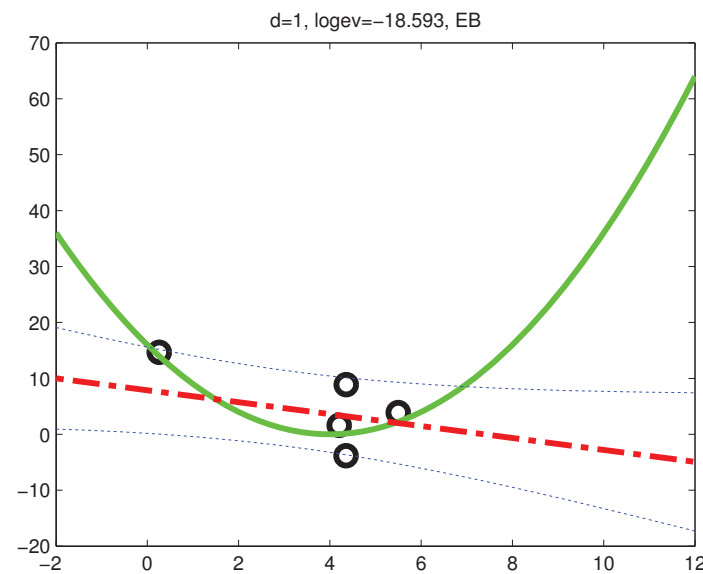
(flat prior, full rank Hessian, etc.)

$$\log p(D | M_i) \approx \log(p(D | \theta_i^*, M_i)) - \frac{d}{2} \log n$$

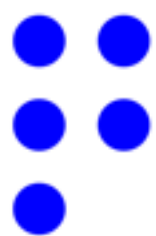
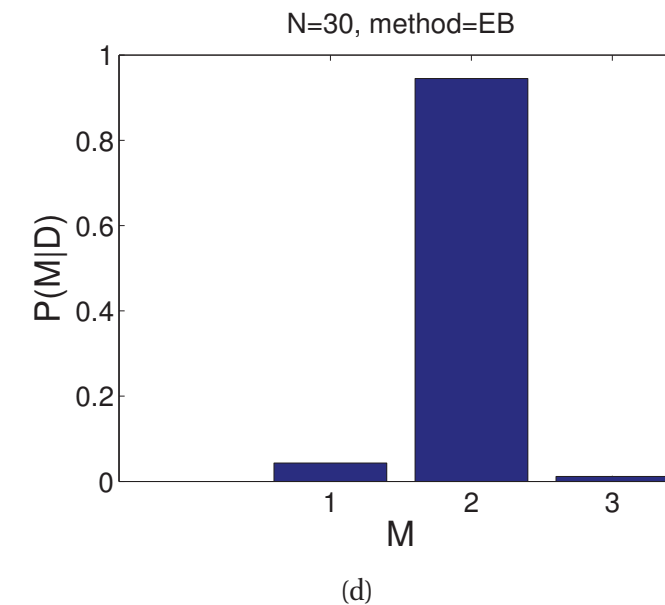
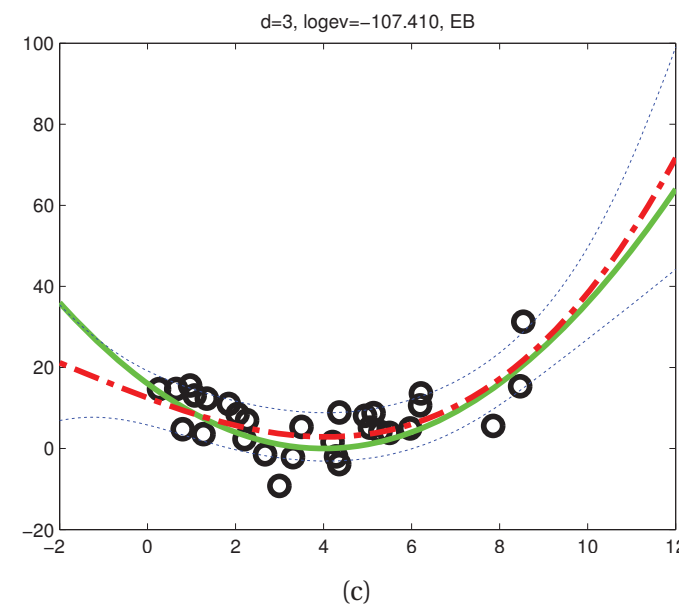
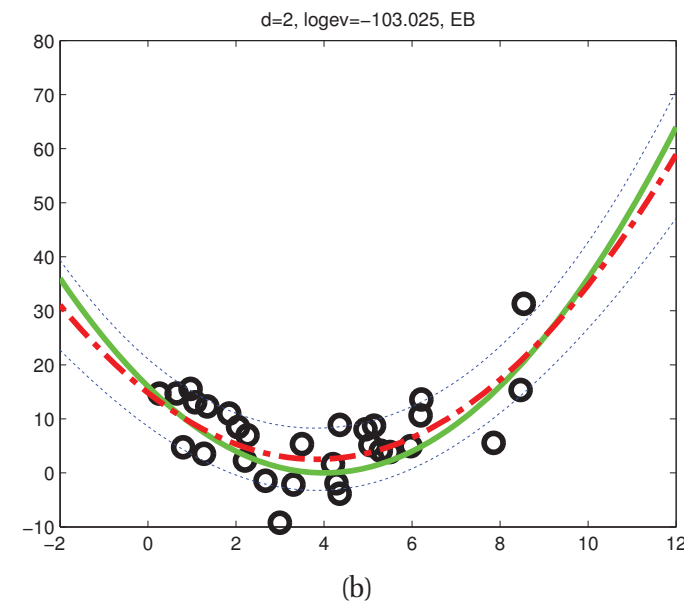
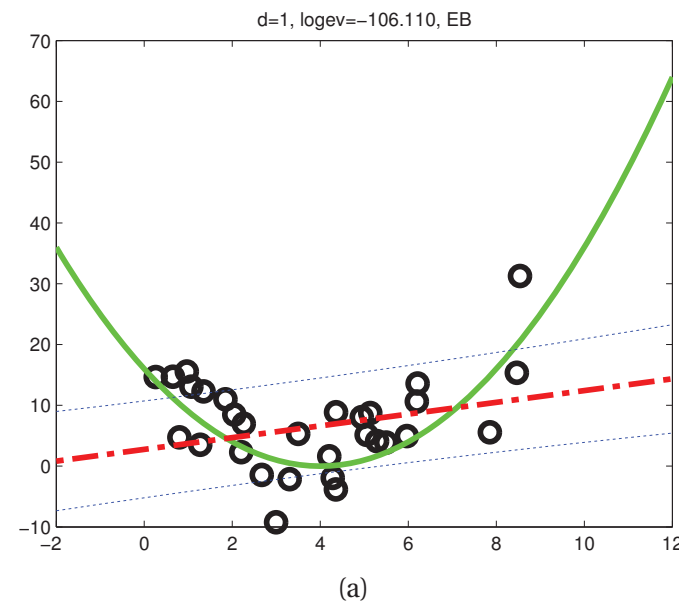
number of samples



Example: Polynomial Regression



Example: Polynomial Regression



Something more accurate?

What about MCMC?

- MCMC (on its own) cannot be used to compute model evidence.
- There is an estimator called the **harmonic mean** estimator to the evidence, but it is garbage... (it should never be used).
- So, what do we do?

Sequential Monte Carlo Sampling Techniques

The problem

Sample from this:

$$p(\theta | D) = \frac{p(D | \theta)p(\theta)}{Z}$$

Do it efficiently.

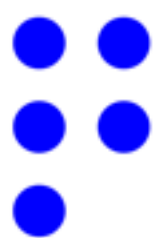
Compute the evidence (which can be used for model selection):

$$Z = \int p(D | \theta)p(\theta)d\theta$$

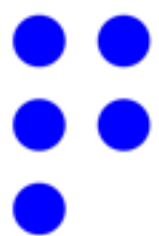
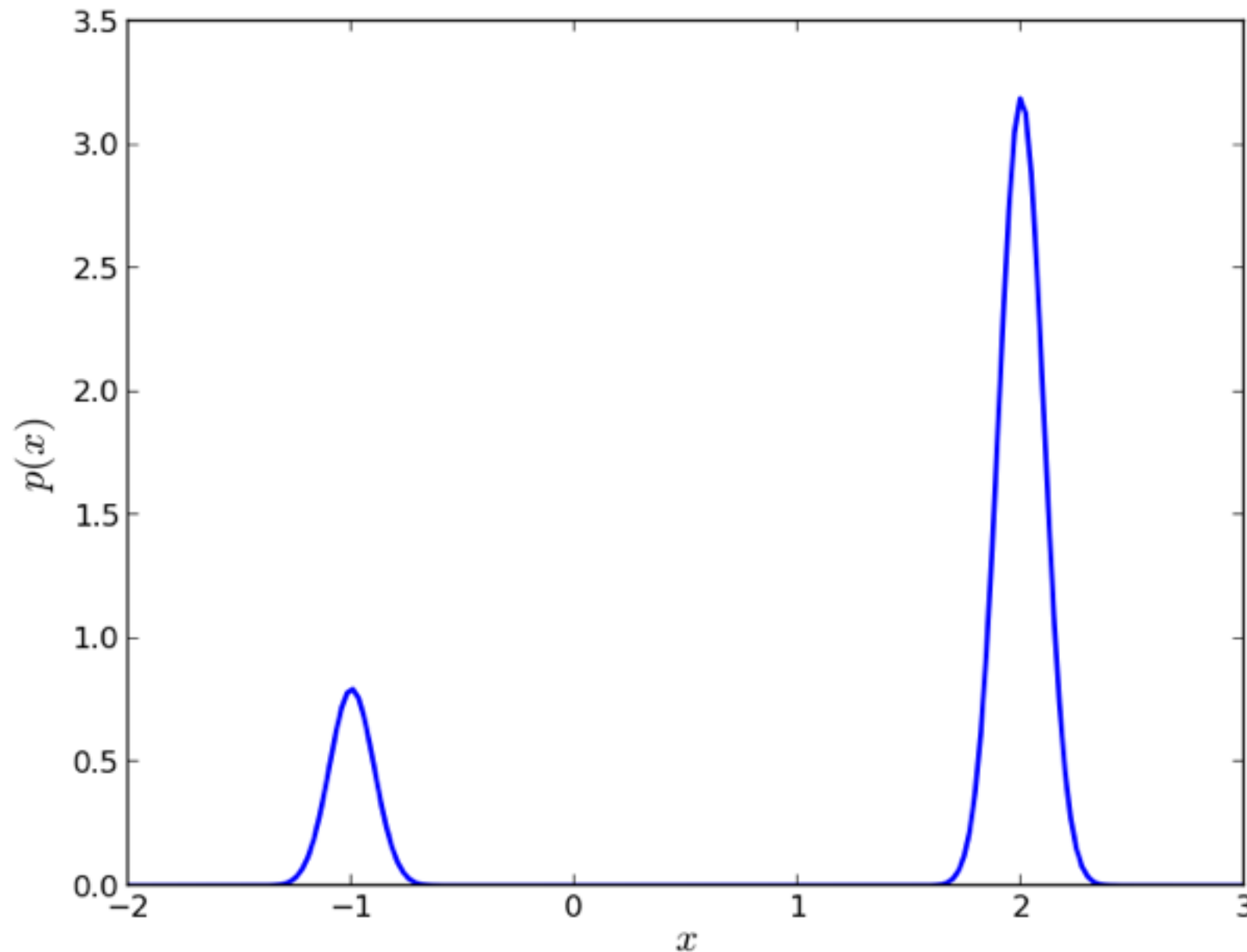
Why is it difficult?

$$p(\theta | D) = \frac{p(D | \theta)p(\theta)}{Z} \quad Z = \int p(D | \theta)p(\theta)d\theta$$

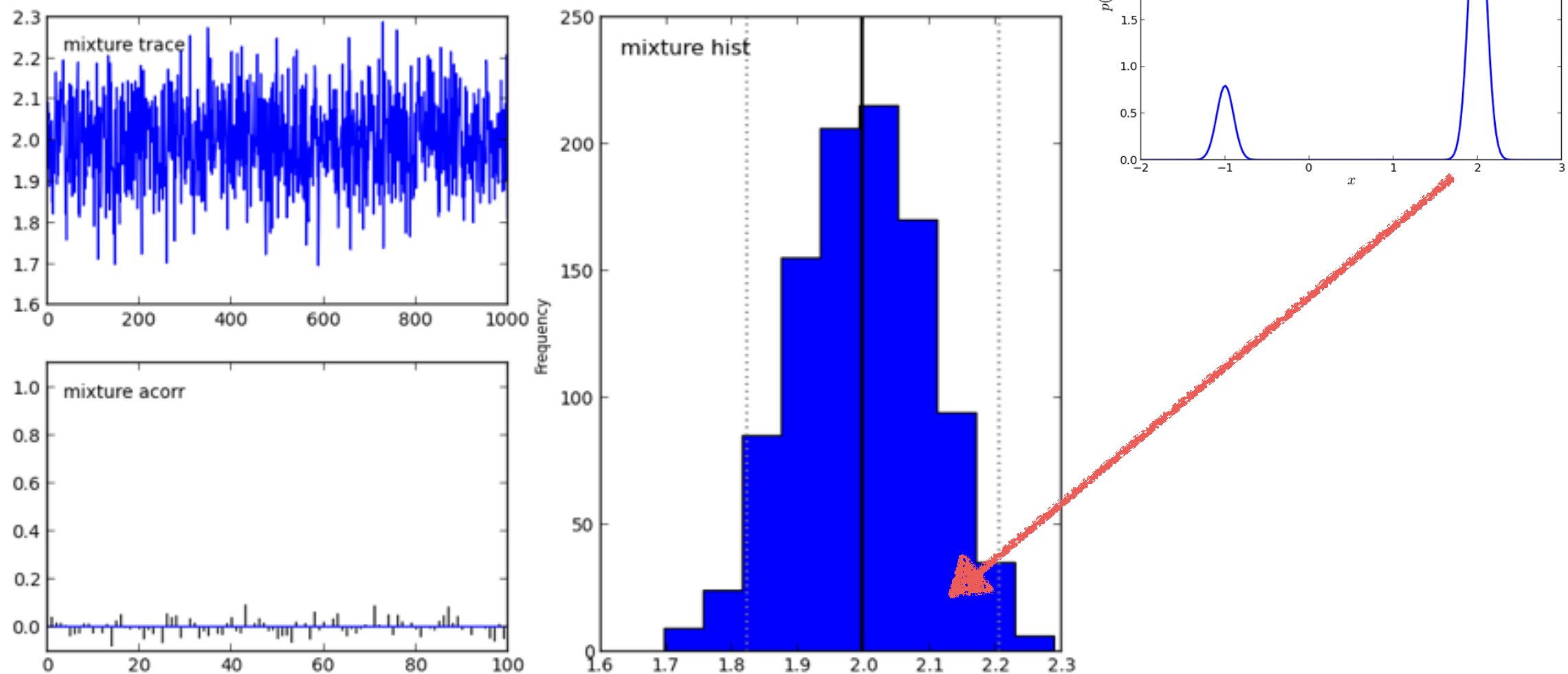
- Slow convergence (long thermalization).
- Trouble with multiple modes.
- Proposal must be hand-picked.
- The fact that the likelihood may be sharply peaked makes the evaluation of the integral is extremely difficult.



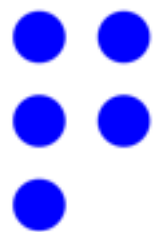
Example: Posterior with Two Modes



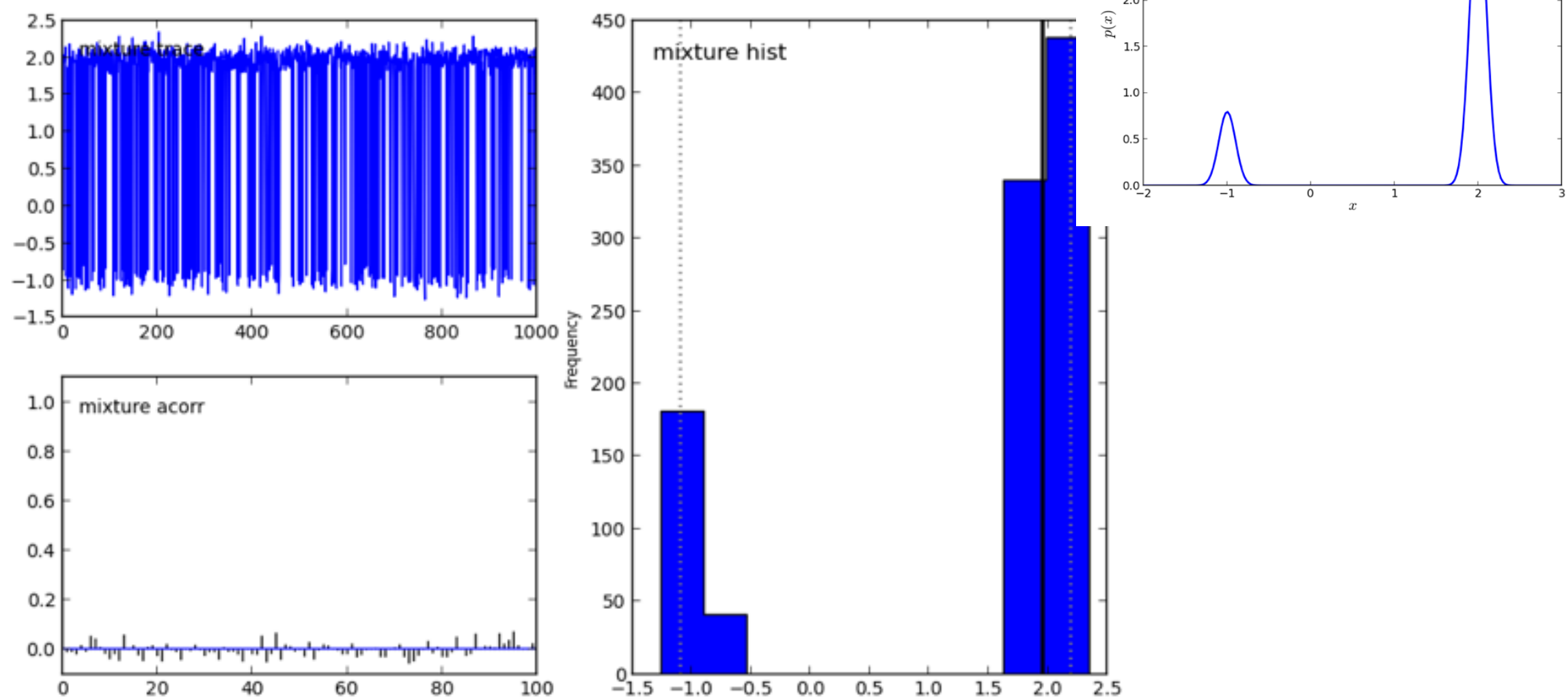
Example: Posterior with Two Modes



MCMC seems converged, but missed one mode...



Example: Posterior with Two Modes



MCMC works only if we hand-pick the right proposal...

Idea

Define a family of distributions that will take you from the prior to the posterior with increasing complexity:

$$p(\theta \mid D, \gamma) \propto \pi_\gamma(\theta) := p(D \mid \theta)^\gamma p(\theta)$$

$$0 \leq \gamma \leq 1$$

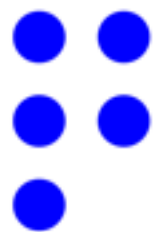
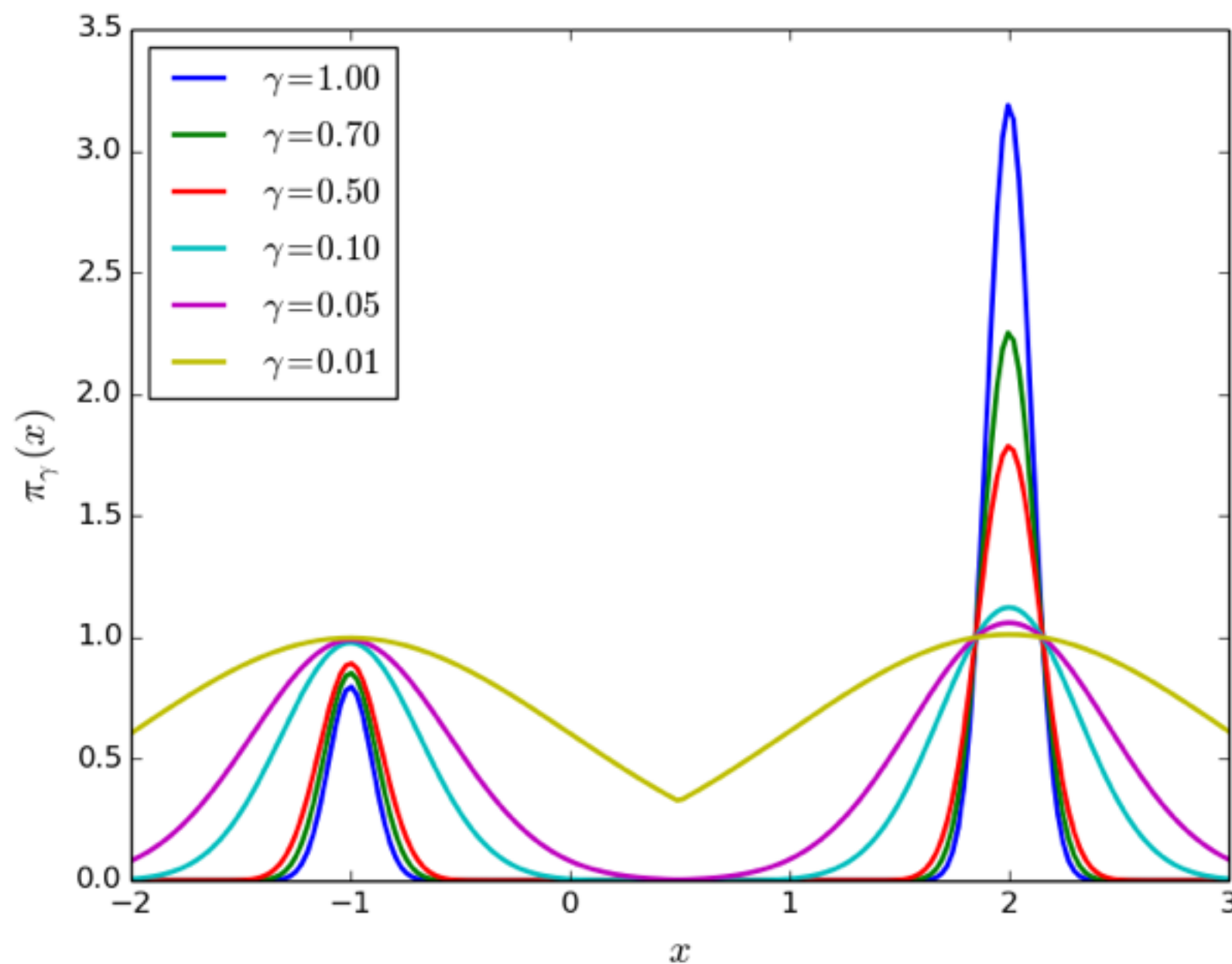
From the prior:

$$p(\theta \mid D, \gamma = 0) \propto \pi_0(\theta) = p(\theta)$$

To the posterior:

$$p(\theta \mid D, \gamma = 1) = \pi_1(\theta) = p(D \mid \theta)p(\theta)$$

Example: Posterior with Two Modes

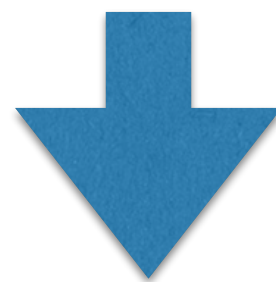


Idea

Start with a **particle approximation** to the prior:

$$p(\theta | D, \gamma = 0) = \sum_{i=1}^N w_{0,i} \delta(\theta_{0,i} - \theta)$$

$$\theta_{0,i} \sim p(\theta | D, \gamma = 0) = p(\theta) \quad \sum_{i=1}^N w_{0,i} = 1$$



Gradually increase gamma
and update particles and
weights

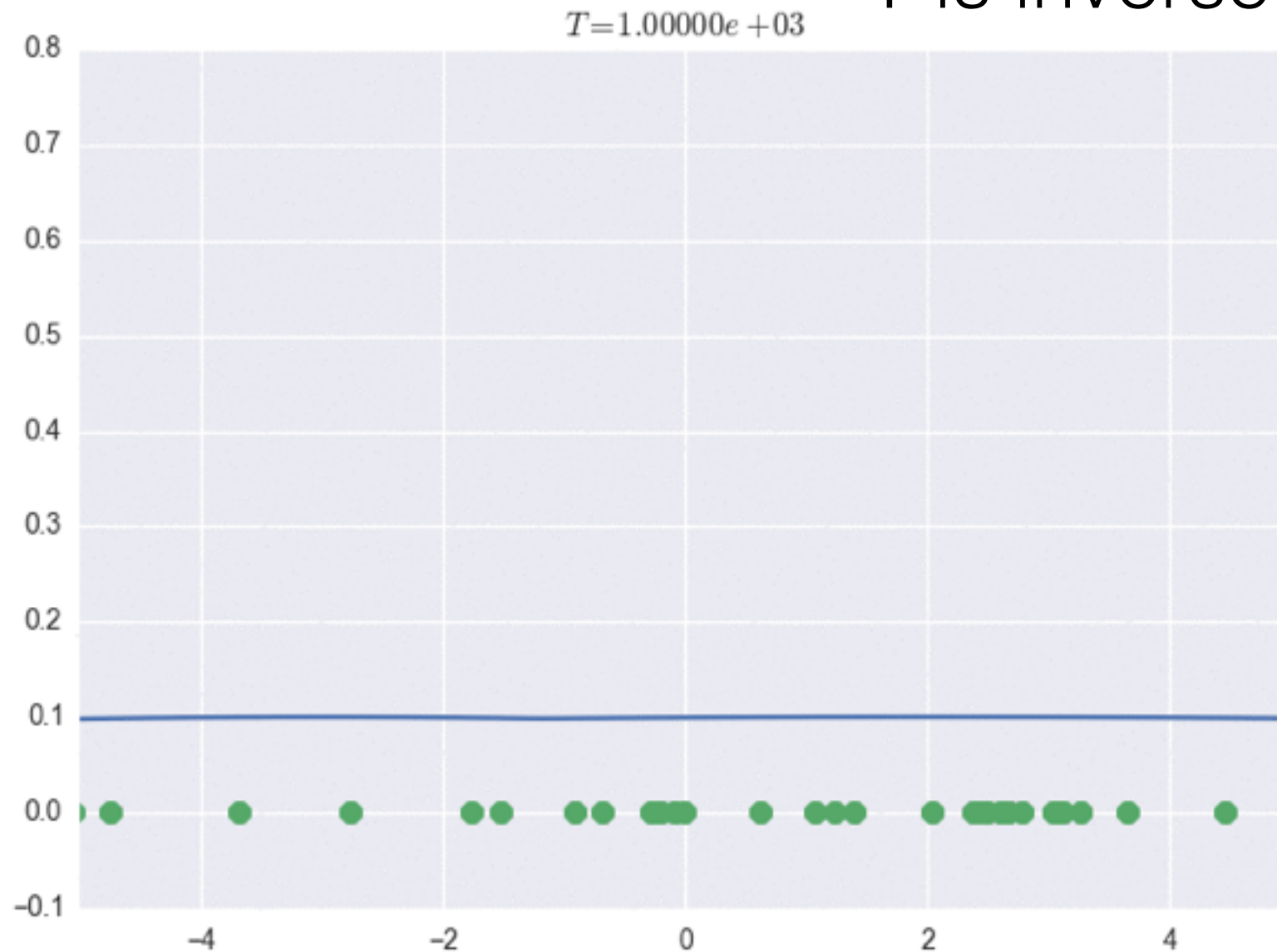
$$p(\theta | D, \gamma) = \sum_{i=1}^N w_{\gamma,i} \delta(\theta_{\gamma,i} - \theta)$$



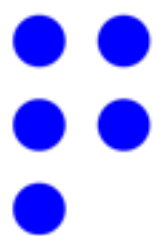
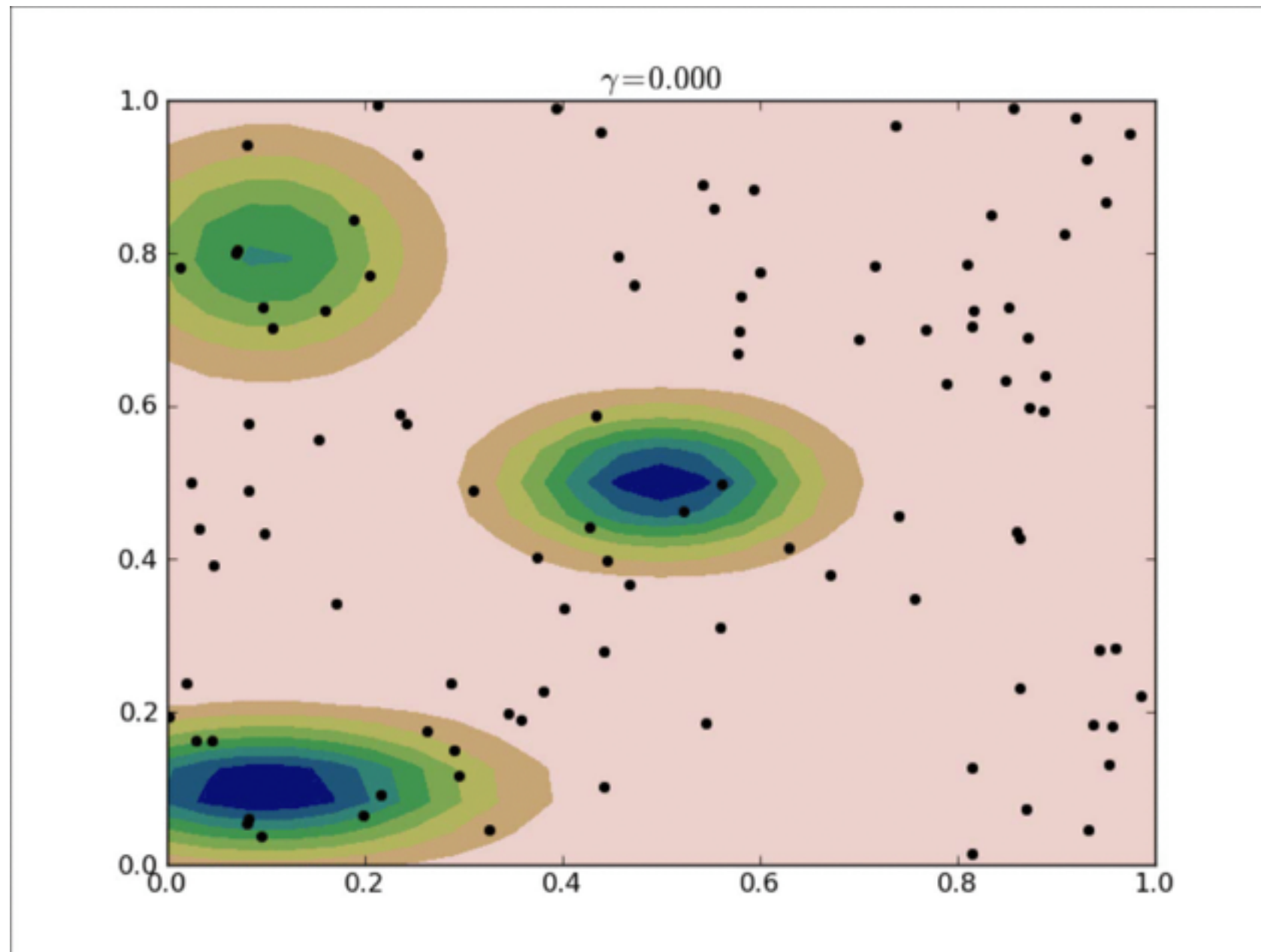
$$p(\theta | D, \gamma = 1) = \sum_{i=1}^N w_{1,i} \delta(\theta_{1,i} - \theta)$$

Example: Posterior with Two Modes

T is inverse gamma

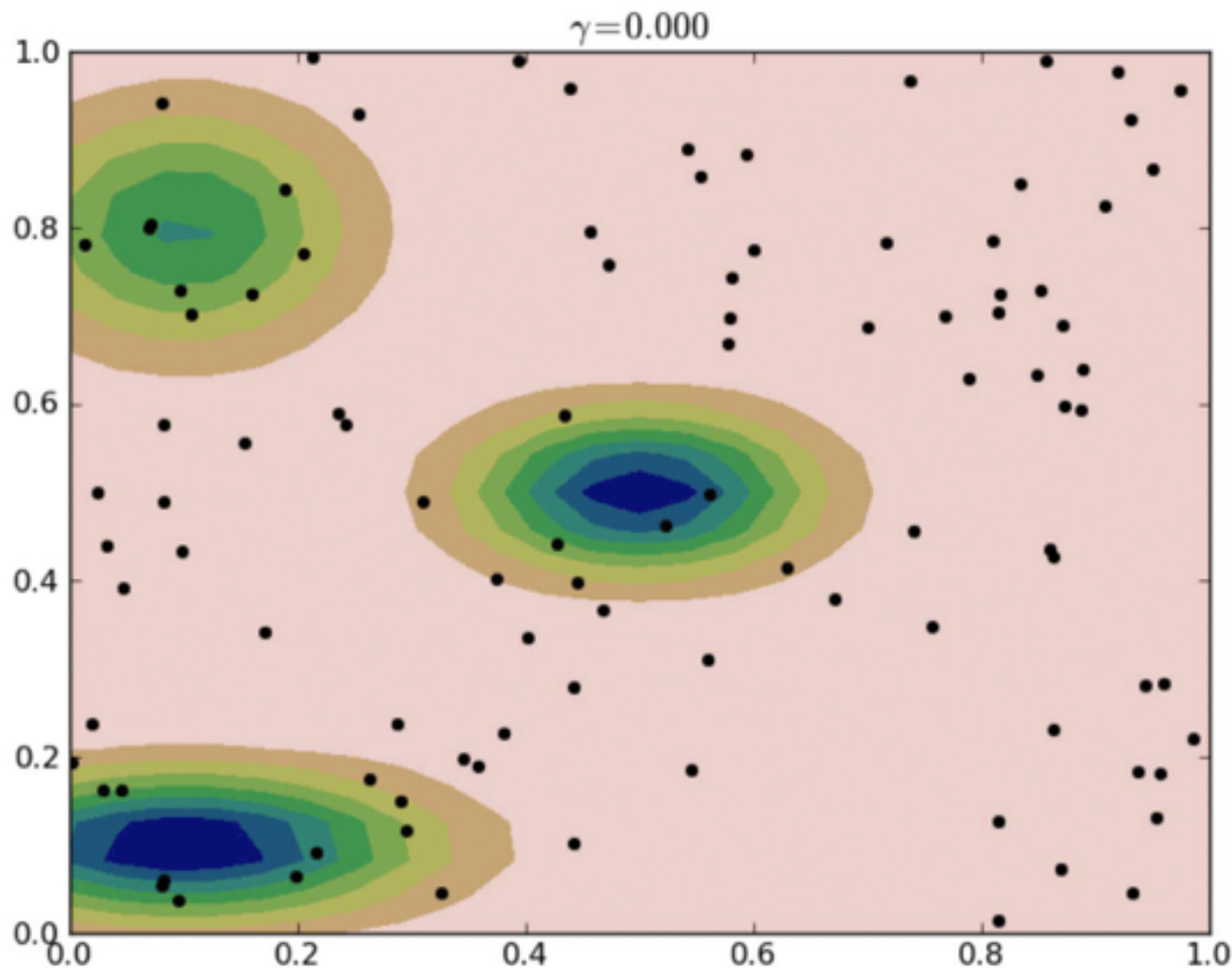


2D Three Modes



Embarrassingly Parallelizable

- Each particle can reside in a different CPU core.
- Minimal communication is required.



OK - How do we do it?

Consider two consecutive gammas:

$$0 \leq \gamma_t < \gamma_{t+1} \leq 1$$

We have a particle approximation for the first one:

$$p(\theta \mid D, \gamma_t) = \sum_{i=1}^N w_{\gamma_t, i} \delta(\theta_{\gamma_t, i} - \theta)$$

Build a particle approximation for the second one:

$$p(\theta \mid D, \gamma_{t+1}) = \sum_{i=1}^N w_{\gamma_{t+1}, i} \delta(\theta_{\gamma_{t+1}, i} - \theta)$$

It can be done in many ways, I will just give you one...

OK - How do we do it?

$$0 \leq \gamma_t < \gamma_{t+1} \leq 1$$

$$p(\theta | D, \gamma_t) = \sum_{i=1}^N w_{\gamma_t, i} \delta(\theta_{\gamma_t, i} - \theta) \quad \rightarrow \quad p(\theta | D, \gamma_{t+1}) = \sum_{i=1}^N w_{\gamma_{t+1}, i} \delta(\theta_{\gamma_{t+1}, i} - \theta)$$

The new weights will be:

$$w_{\gamma_{t+1}, i} = \frac{w_{\gamma_t, i}}{\sum_{j=1}^N w_{\gamma_t, j}} \quad \text{where} \quad w_{\gamma_{t+1}} = w_{\gamma_t, i} \frac{\pi_{\gamma_{t+1}}(\theta_{\gamma_t, i})}{\pi_{\gamma_t}(\theta_{\gamma_t, i})}$$

Sample the thetas from:

$$\theta_{\gamma_{t+1}, i} \sim \pi_{\gamma_{t+1}}(\theta) = p(D | \theta)^{\gamma_{t+1}} p(\theta)$$

by doing a few steps of your favorite MCMC.

Where is the Normalization Constant?

$$0 \leq \gamma_t < \gamma_{t+1} \leq 1$$

$$p(\theta | D, \gamma_t) = \sum_{i=1}^N w_{\gamma_t, i} \delta(\theta_{\gamma_t, i} - \theta) \quad \rightarrow \quad p(\theta | D, \gamma_{t+1}) = \sum_{i=1}^N w_{\gamma_{t+1}, i} \delta(\theta_{\gamma_{t+1}, i} - \theta)$$
$$w_{\gamma_{t+1}, i} = \frac{W_{\gamma_{t+1}, i}}{\sum_{j=1}^N W_{\gamma_{t+1}, j}} \quad W_{\gamma_{t+1}} = w_{\gamma_t, i} \frac{\pi_{\gamma_{t+1}}(\theta_{\gamma_t, i})}{\pi_{\gamma_t}(\theta_{\gamma_t, i})}$$

We can approximate the ratio of two normalization constants by:

$$\frac{Z_{\gamma_{t+1}}}{Z_{\gamma_t}} \approx \sum_{i=1}^N w_{t, i} = \sum_{i=1}^N w_{\gamma_t, i} \frac{\pi_{\gamma_{t+1}}(\theta_{\gamma_t, i})}{\pi_{\gamma_t}(\theta_{\gamma_t, i})}$$

Where is the Normalization Constant?

We can approximate the ratio of two normalization constants by:

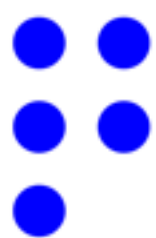
$$\frac{Z_{\gamma_{t+1}}}{Z_{\gamma_t}} \approx \sum_{i=1}^N w_{t,i} = \sum_{i=1}^N w_{\gamma_t,i} \frac{\pi_{\gamma_{t+1}}(\theta_{\gamma_t,i})}{\pi_{\gamma_t}(\theta_{\gamma_t,i})}$$

Since, the normalization constant of the prior is just:

$$Z_{\gamma_0} = Z_0 = 1$$

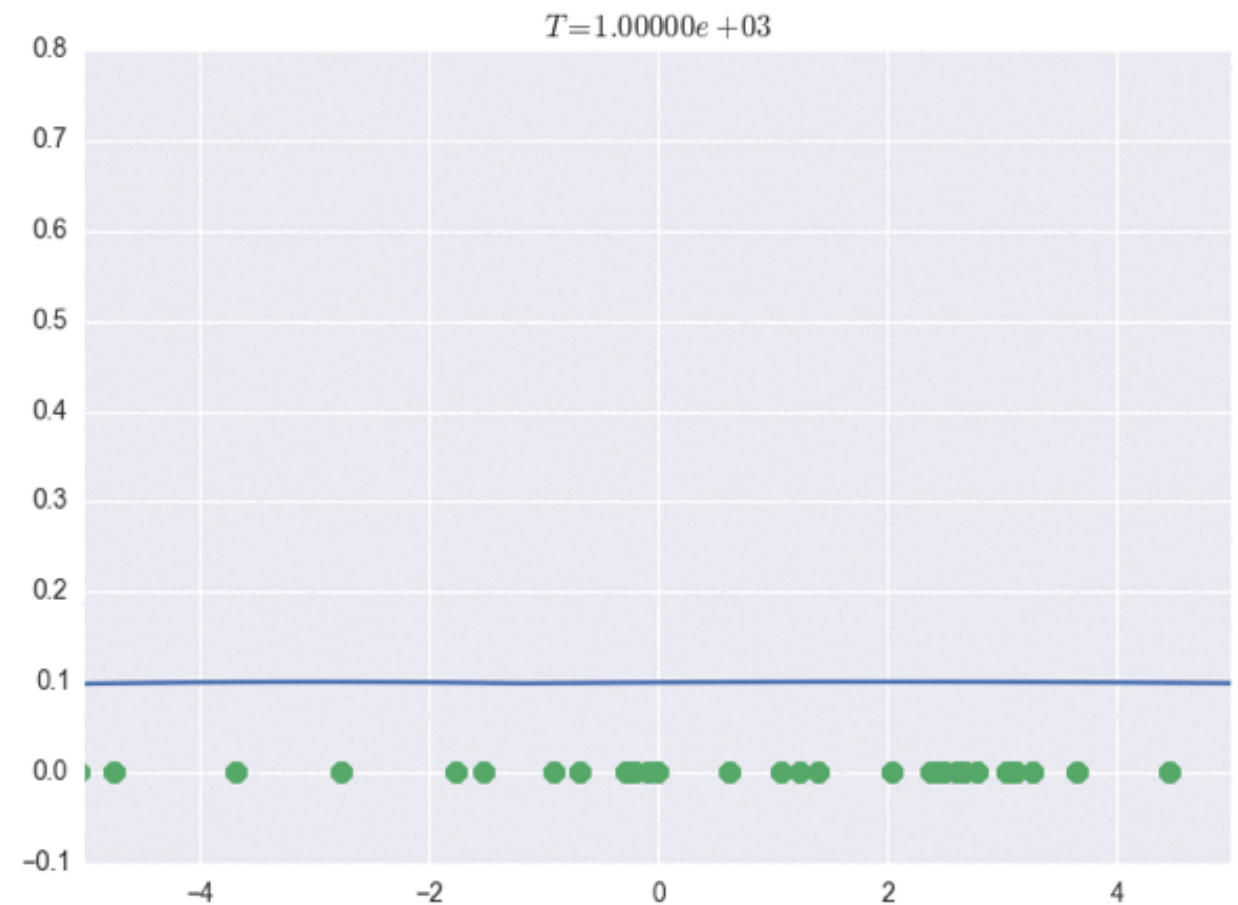
we get by induction:

$$Z = Z_1 = \prod_{t=1}^T \frac{Z_{\gamma_{t+1}}}{Z_{\gamma_t}} \approx \prod_{t=1}^T \sum_{i=1}^N w_{t,i} = \prod_{t=1}^T \sum_{i=1}^N w_{\gamma_t,i} \frac{\pi_{\gamma_{t+1}}(\theta_{\gamma_t,i})}{\pi_{\gamma_t}(\theta_{\gamma_t,i})}$$



Resampling

- Some of the particles may be trapped in regions of low probability.
- Their weight will become extremely small.
- To avoid wasting computations, we should eliminate them.



The Effective Sample Size

The effective sample size (ESS):

$$\text{ESS}(\gamma) = \frac{1}{\sum_{i=1}^N w_{\gamma,i}^2}$$

measures the degeneracy of the particle approximation.

Notice:

$$w_{\gamma,i} = \frac{1}{N}$$

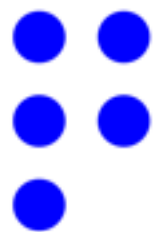


$$\text{ESS}(\gamma) = N$$

$$w_{\gamma,1} = 1, \text{ and } w_{\gamma,i} = 0, \text{ for } i = 2, \dots, N$$



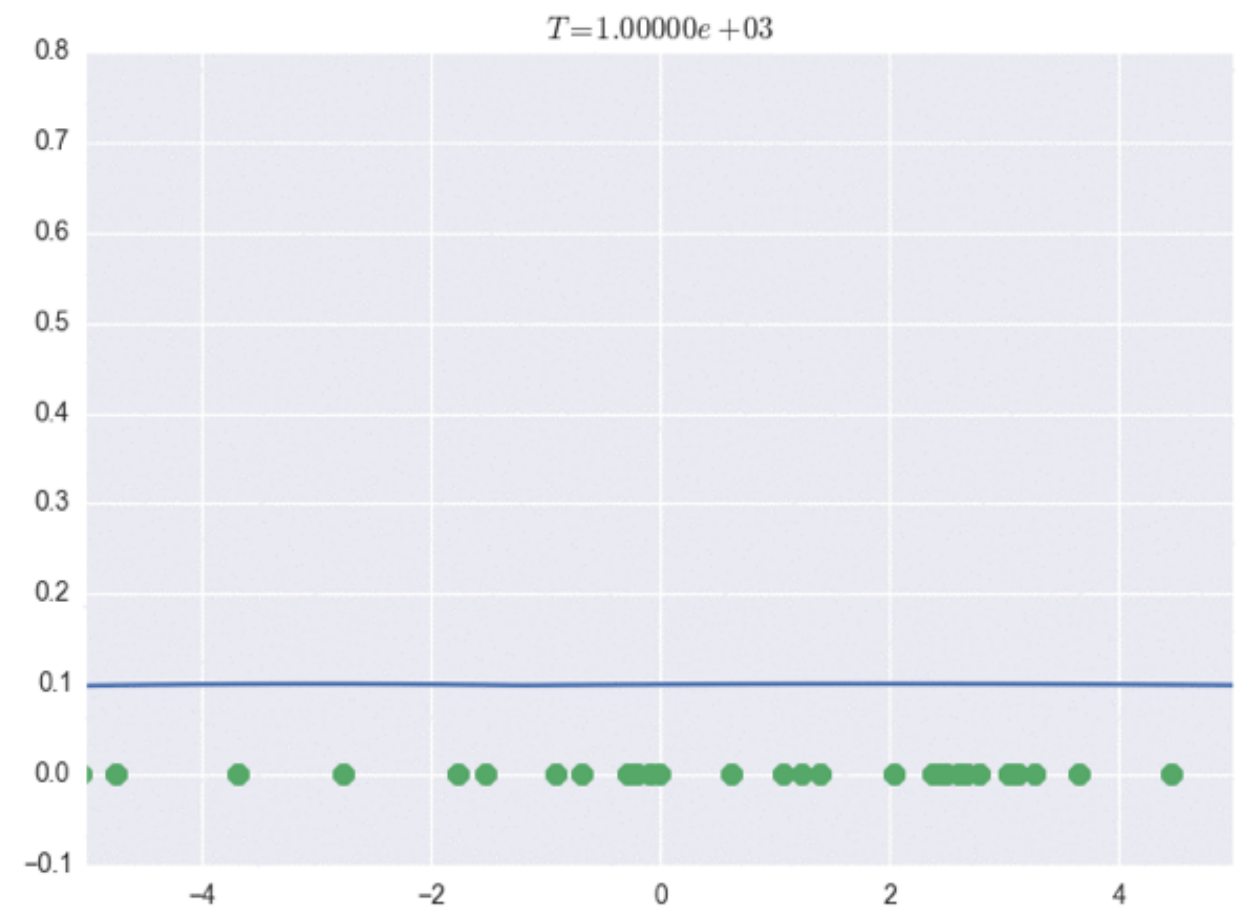
$$\text{ESS}(\gamma) = 1$$



Resampling

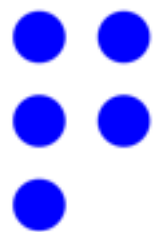
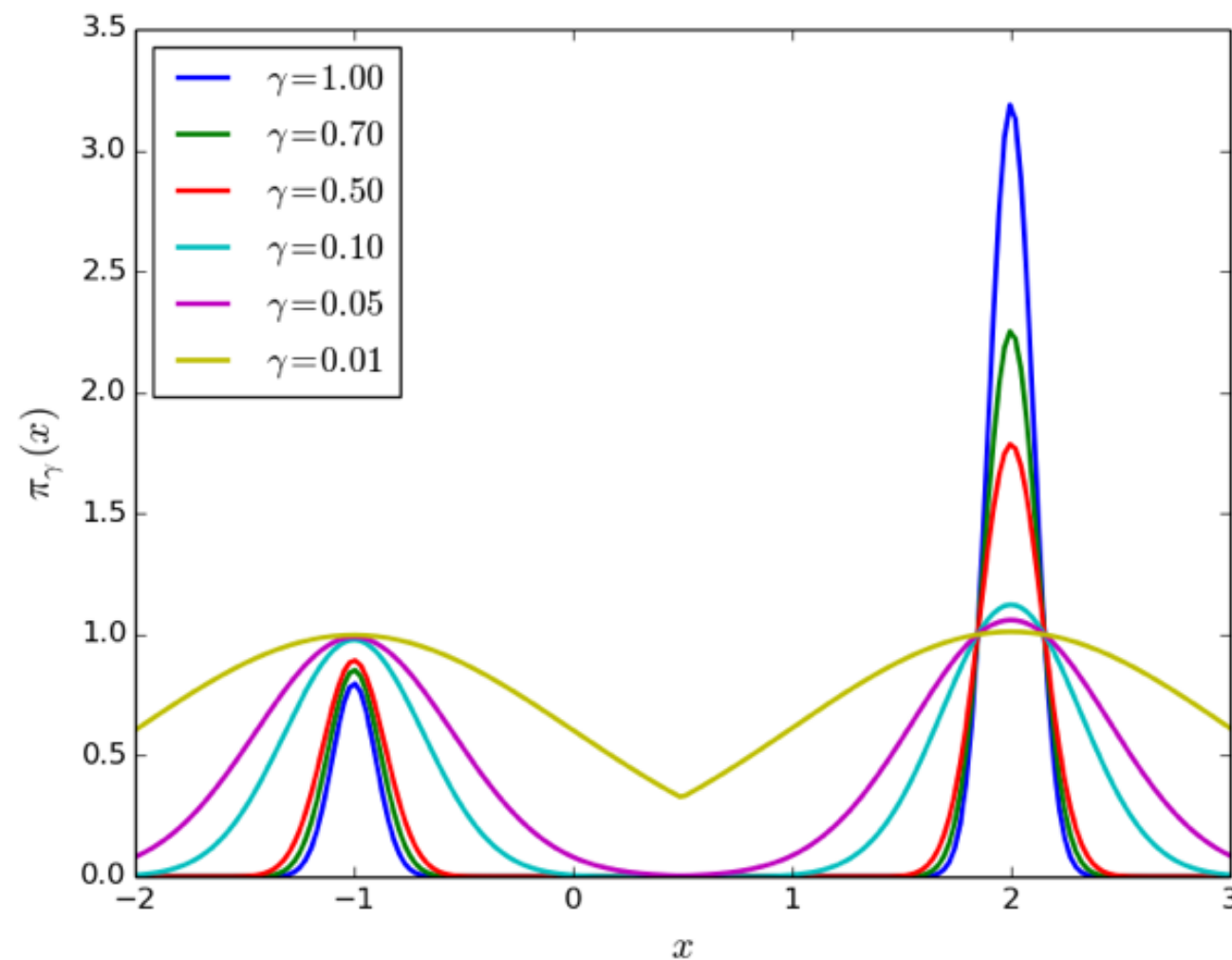
We resample every time:

$$\text{ESS}(\gamma_t) \leq 0.5N$$



Selecting the Intermediate Gammas

You can just pick them by hand:



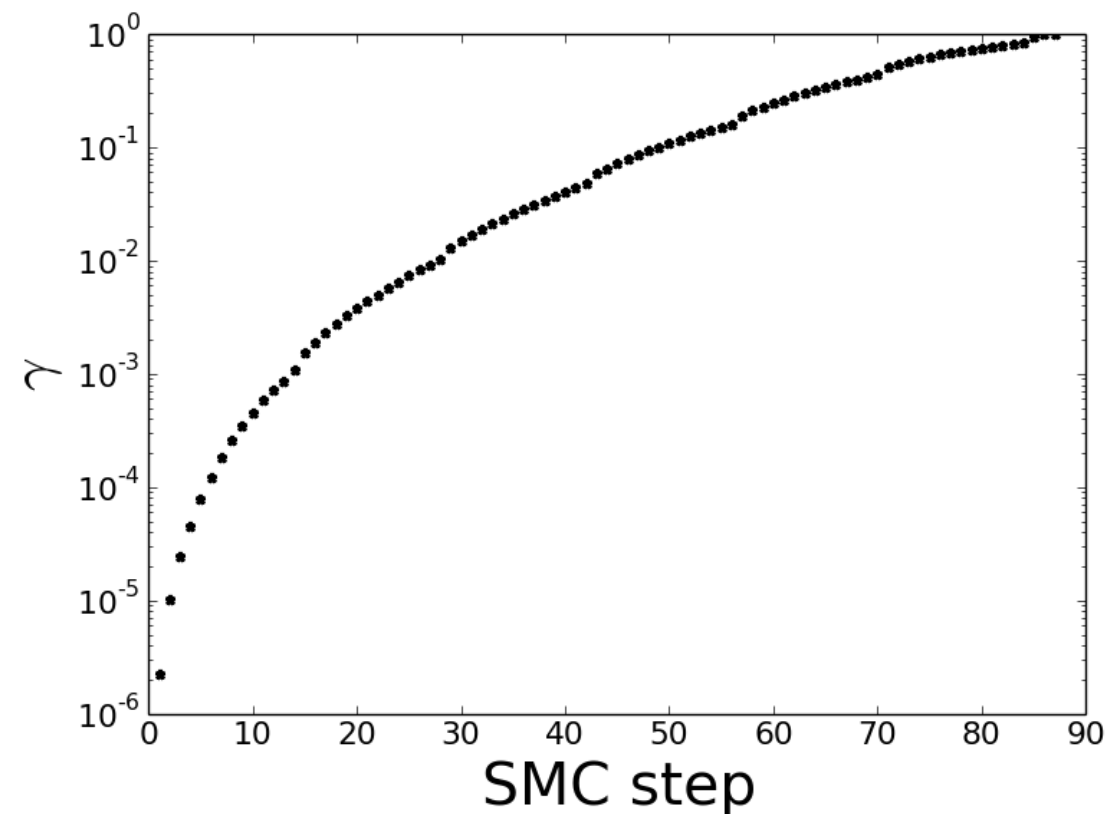
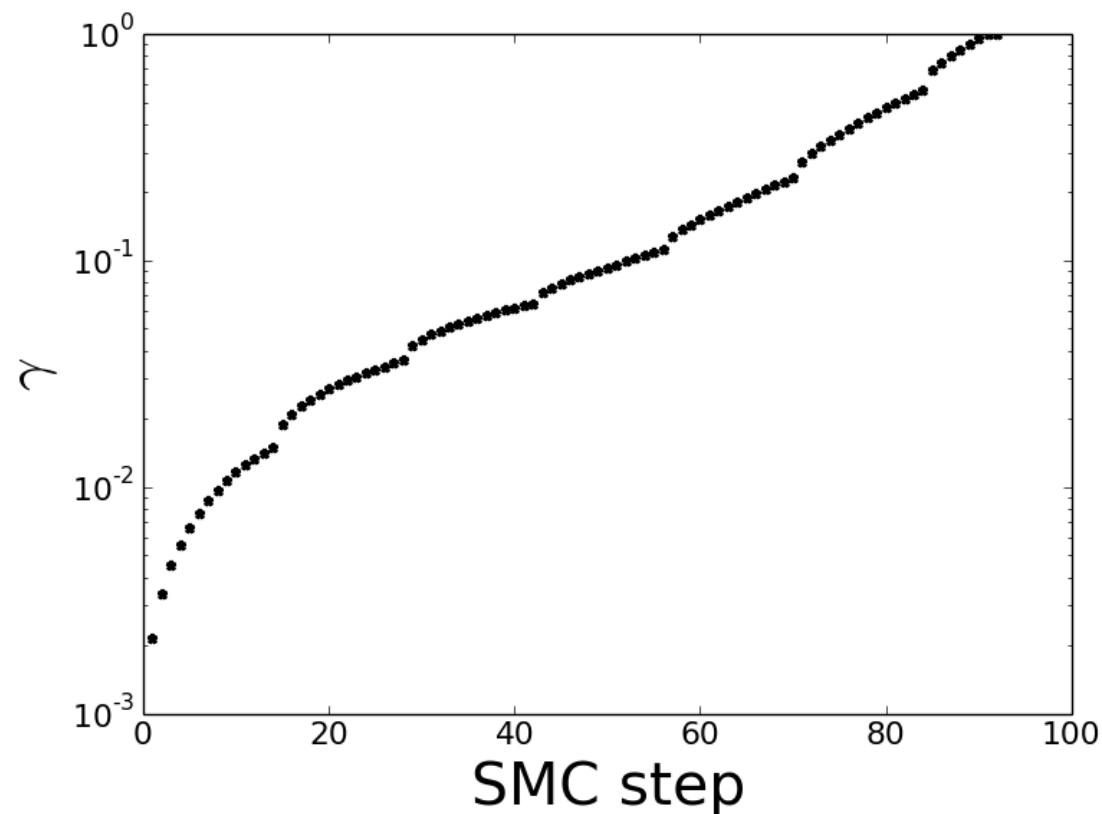
Selecting the Intermediate Gammas

Or you can pick them so that the particles do not become too degenerate from step to step:

$$\text{ESS}(\gamma_{t+1}) = 0.95\text{ESS}(\gamma_t)$$

Example (Calibration 10D)

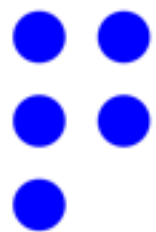
Automatic selection of inverse temperature schedule:



Bilionis, I., et al. (2015). "Crop physiology calibration in the CLM." Geoscientific Model Development **8**(4): 1071-1083.

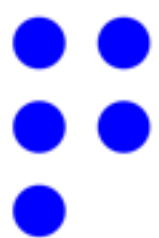
The algorithm

1. Build initial particle approximation (sample the prior).
2. Find the next gamma by fixing the change in the ESS.
3. Resample if ESS falls below threshold.
4. Draw samples at new gamma starting at the old particles.



Why bother?

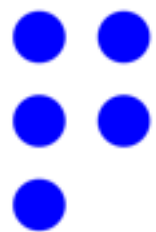
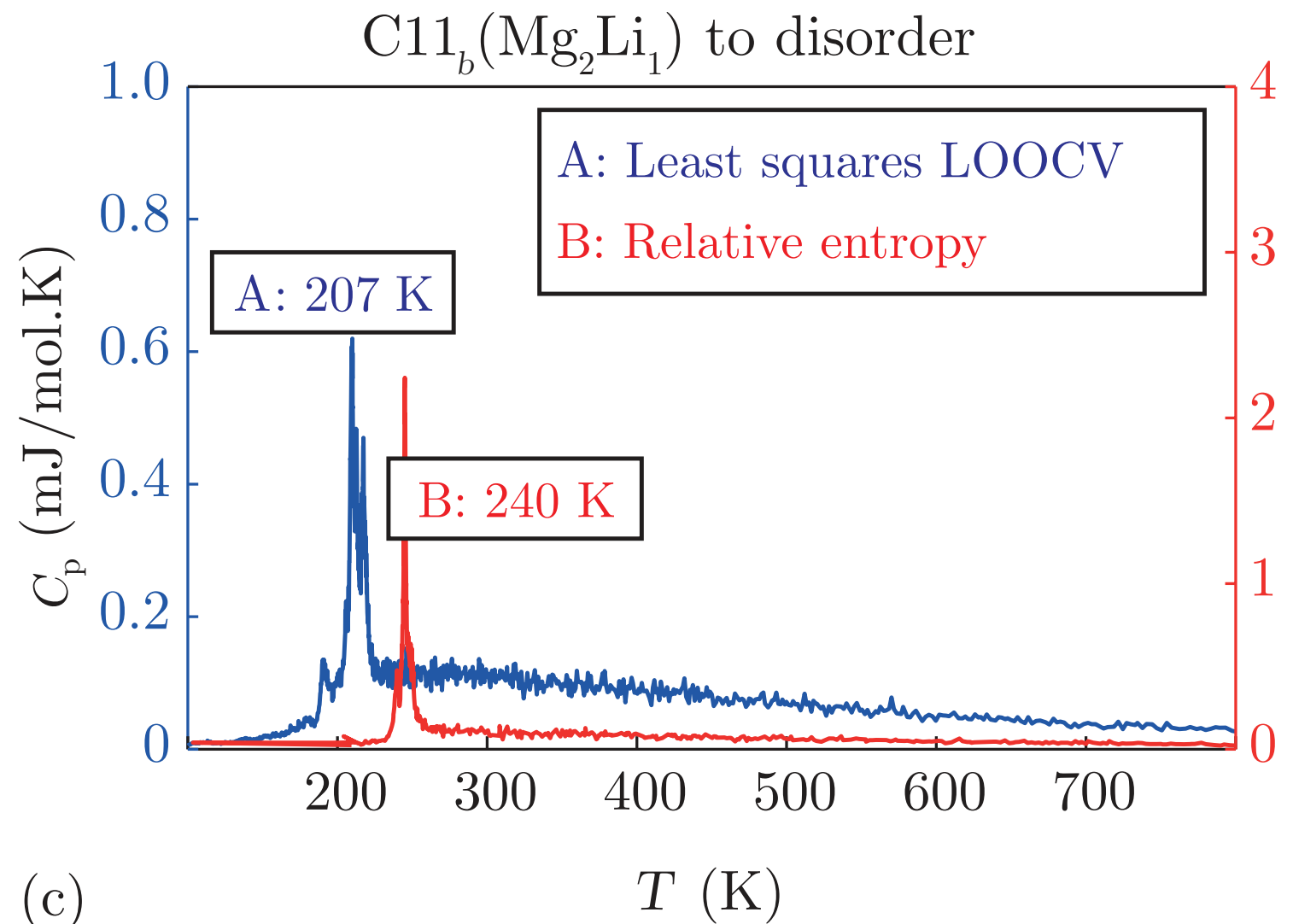
- Very easy to implement
- Embarrassingly parallelizable
- Makes use of existing sampling code
- Captures multiple modes
- Don't have to fix gamma schedule.
- Can adjust parameters of stochastic samplers.



Example (Calibration 16KD)

Scanning the temperature to find where transition occurs in binary alloys (x 10K faster than thermodynamic integration)

Kristensen, J., et al. (2013).
"Relative entropy as model
selection tool in cluster
expansions." Physical Review B
87(17).



PYSMC CODE

This repository Search Pull requests Issues Gist

PredictiveScienceLab / pysmc Unwatch 5 Unstar 22 Fork 2

Sequential Monte Carlo working on top of pymc — Edit

61 commits 5 branches 1 release 1 contributor

Branch: master pysmc / +

Fixed minor typo.

ebillionis authored on Dec 3, 2014 latest commit 46f9a5cb2c

doc	Trying to fix step methods.	2 years ago
examples	Progress bar issue fixed.	a year ago
pysmc	Progress bar issue fixed.	a year ago
.gitignore	Finished tutorial SMC simple model. Fixed a bug on pysmc.SMC (printin...	2 years ago
LICENSE.txt	Preparing for pypi	10 months ago
README.md	Update README.md	a year ago
setup.cfg	Preparing for pypi	10 months ago
setup.py	Fixed minor typo.	10 months ago

README.md

pysmc

Sequential Monte Carlo working on top of pymc.

The complete documentation can be found here: <http://ebillionis.github.io/pysmc/>

Code Issues 1 Pull requests 0 Wiki Pulse Graphs Settings

HTTPS clone URL [https://github.com](https://github.com/PredictiveScienceLab/pysmc)

You can clone with HTTPS, SSH, or Subversion.

Clone in Desktop Download ZIP

© 2015 GitHub, Inc. Terms Privacy Security Contact Help Status API Training Shop Blog About Pricing

