

02477 – Bayesian Machine Learning: Lecture 9

Michael Riis Andersen

Technical University of Denmark,
DTU Compute, Department of Applied Math and Computer Science

Outline

1 Markov chain Monte Carlo

2 Gibbs sampling

3 Convergence diagnostics

4 Hierarchical models

Markov chain Monte Carlo

Bayesian inference and probabilistic modelling

■ Bayesian supervised learning in general:

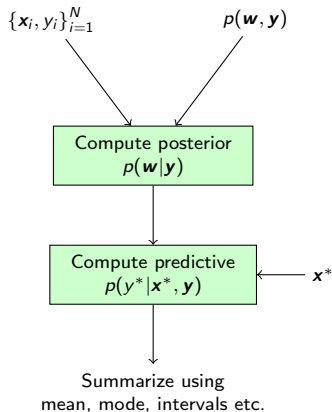
1. Joint model of data \mathbf{y} and parameters \mathbf{w} .
2. Summarize knowledge of \mathbf{w} given data \mathbf{y} .
3. Compute posterior predictive distribution.

■ Goal: separate modelling from inference:

1. Build models reflecting domain knowledge.
2. Push “inference button” and get results.

■ Inference methods:

1. Laplace approximations.
2. Markov chain Monte Carlo.
3. Variational approximations.



Monte Carlo: Posterior inference using samples

- Many posterior summaries can be phrased as expectations $\mathbb{E}_p[f(\mathbf{z})]$ for some function f

- We can compute expectations wrt. p using the *Monte Carlo estimator*

$$\bar{f} = \mathbb{E}_p[f(\mathbf{z})] = \int f(\mathbf{z})p(\mathbf{z})d\mathbf{z} \approx \frac{1}{S} \sum_{i=1}^S f(\mathbf{z}^i) \equiv \hat{f},$$

where $\mathbf{z}^i \sim p(\mathbf{z})$ for $i = 1, \dots, S$

- We showed last time that ...

1. the Monte Carlo estimator \hat{f} is *unbiased*
2. the variance of \hat{f} decreases with $1/S$ when the samples are *i.i.d.*

A zoo of sampling-based methods

■ Simple sampling methods

1. Rejection sampling
2. Ancestral sampling
3. Importance sampling
4. Transformation methods
5. Inverse transform sampling
6. ...

■ MCMC methods

1. Metropolis-Hastings
2. Gibbs Sampling
3. Slice sampling
4. Hamiltonian Monte Carlo
5. ...

MCMC using the Metropolis-Hastings algorithm

- We can use the MH to generate samples from a distribution of interest $p(\mathbf{z})$.

The Metropolis-Hastings algorithm

- Start from some initial value \mathbf{z}^1 (e.g., a sample from the prior).
- Repeat for $k = 1$ to K :
 1. Given last value \mathbf{z}^{k-1} , generate *candidate sample* using proposal distribution

$$\mathbf{z}^* \sim q(\mathbf{z}^* | \mathbf{z}^{k-1}).$$

2. Compute *acceptance probability* A_k as follows

$$A_k = \min \left(1, \frac{p(\mathbf{z}^*)q(\mathbf{z}^{k-1} | \mathbf{z}^*)}{p(\mathbf{z}^{k-1})q(\mathbf{z}^* | \mathbf{z}^{k-1})} \right).$$

3. Simulate $u_k \sim \mathcal{U}(0, 1)$ and define \mathbf{z}^k as

$$\mathbf{z}^{k+1} = \begin{cases} \mathbf{z}^* & \text{if } u_k < A_k \\ \mathbf{z}^{k-1} & \text{otherwise} \end{cases}$$

- What do we need in order to implement MH for a given model?

Markov chain Monte Carlo theory I

- Metropolis-Hastings defines a chain of samples $\mathbf{z}^0, \mathbf{z}^1, \mathbf{z}^2, \dots$ with a *Markov property*

$$p(\mathbf{z}^{k+1} | \mathbf{z}^1, \mathbf{z}^2, \dots, \mathbf{z}^k) = p(\mathbf{z}^{k+1} | \mathbf{z}^k)$$

- The *transition kernel* tells us how to iterate the chain

$$T(\mathbf{z}^{k+1} | \mathbf{z}^k) \equiv p(\mathbf{z}^{k+1} | \mathbf{z}^k)$$

- The distribution of \mathbf{z}^{k+1} is given by *sum rule*

$$p(\mathbf{z}^{k+1}) = \int T(\mathbf{z}^{k+1} | \mathbf{z}^k) p(\mathbf{z}^k) d\mathbf{z}^k$$

- A distribution $p^*(\mathbf{z})$ is said to be *invariant* or *stationary* wrt. the Markov chain if each step does not change the distribution

$$p^*(\mathbf{z}) = \int T(\mathbf{z} | \mathbf{z}') p^*(\mathbf{z}') d\mathbf{z}'$$

- We require $p^*(\mathbf{z})$ to be a limiting distribution of the chain (*independent of the initial distribution*).

$$p(\mathbf{z}^k) \rightarrow p^*(\mathbf{z}) \quad \text{for} \quad k \rightarrow \infty$$

Transition kernel for Metropolis-Hastings

- Recall the acceptance probability for Metropolis-Hastings

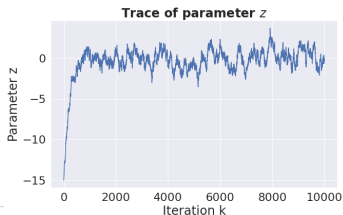
$$A(\mathbf{z}^* | \mathbf{z}^k) = \min \left[1, \frac{\rho(\mathbf{z}^*)q(\mathbf{z}^k | \mathbf{z}^*)}{\rho(\mathbf{z}^k)q(\mathbf{z}^* | \mathbf{z}^k)} \right]$$

- The transition kernel for Metropolis-Hastings

$$T(\mathbf{z}' | \mathbf{z}) = \begin{cases} q(\mathbf{z}' | \mathbf{z})A(\mathbf{z}' | \mathbf{z}) & \text{if } \mathbf{z}' \neq \mathbf{z} \\ q(\mathbf{z} | \mathbf{z})A(\mathbf{z} | \mathbf{z}) + \int q(\mathbf{z}'' | \mathbf{z}) [1 - A(\mathbf{z}'' | \mathbf{z})] d\mathbf{z}'' & \text{if } \mathbf{z}' = \mathbf{z} \end{cases}$$

- The big picture

1. We initialize \mathbf{z}^1 .
2. We iterate $\mathbf{z}^{k+1} | \mathbf{z}^k \sim T(\mathbf{z}^{k+1} | \mathbf{z}^k)$ (*warm-up phase*).
3. Eventually the distribution of \mathbf{z}^k will converge to the target distribution p^* (*sampling phase*).



Markov chain Monte Carlo theory II

$$p^*(z) = \int T(z|z')p^*(z')dz'$$

- We require $p^*(z)$ to be a limiting distribution of the chain (*independent of the initial distribution*).

$$p(z^k) \rightarrow p^*(z) \quad \text{for} \quad k \rightarrow \infty$$

- Conditions required for a Markov chain to have a limiting distribution equal to the unique stationary distribution:

(A1) *Irreducible*: all states z can be reached.

Counter example: $z^{k+1} = z^k + |e_k|$, where $e_k \sim \mathcal{N}(0, 1)$

(A2) *Aperiodic*: no deterministic cycles.

Counter example: $z^1 = 1, z^2 = 2, z^3 = 3, z^4 = 1, z^5 = 2, \dots$

(A3) *Positive recurrent*: the chain has positive probability of returning to any given state (+ finite expected return-time).

$$P(z^k \in A | z^0) > 0 \quad \text{for all sets } A \quad \text{where} \quad P(A) > 0$$

- A chain that satisfies (A1)–(A3) is said to be *ergodic* and ensures $p(z^k) \rightarrow p^*(z)$.

Markov chain Monte Carlo theory III

- These conditions are generally hard to check in practice.
- Simpler condition: if a chain satisfies the *detailed balance condition*, then p^* is its stationary distribution

$$T(z'|z)p^*(z) = T(z|z')p^*(z').$$

- MH with reasonable proposal distributions satisfies detailed balance (see slide 33).
- It's often not a question of convergence or not, but rather how fast we converge for a given proposal distribution.
- For example, all (non-degenerate) Gaussian proposals lead to detailed balance, but convergence time may vary drastically depending on the proposal variance.
- More theory and details: Monte Carlo Statistical Methods by Robert and Casella.



Pros and cons of Metropolis-Hastings

Pros

- Strong mathematical guarantees: If we sample long enough, the iterates \mathbf{z}^k will converge to the exact target distribution

$$p(\mathbf{z}^k) \rightarrow p^*(\mathbf{z}) \quad \text{for} \quad k \rightarrow \infty$$

- Easy to implement.
- Easy to prototype and evaluate different models.

Cons

- May have to sample “infinitely” long for difficult distributions.
- Acceptance ratio can be low.
- Slow for large datasets.
- Proposal distribution may require tuning.

Questions: True or false?

Quiz via DTU Learn:

Lecture 9: Metropolis-Hastings (12 questions)

Check you knowledge

Gibbs sampling

Gibbs sampling

- When using Metropolis-Hastings,
 1. we have to choose a proposal distribution (and sometimes tune it), and
 2. it may suffer from low acceptance rates.
- *Gibbs sampling* works by iteratively updating each coordinate of \mathbf{z} by sampling from the posterior conditionals $p(z_i | \mathbf{z}_{-i})$ (\mathbf{z}_{-i} means the entire vector except index i).

The Gibbs Sampler

- Initialize all parameter values $\{z_i^0\}_{i=1}^D$
- Repeat for $k = 1$ to K :
 - Sample $z_1^k \sim p(z_1 | z_2^{k-1}, z_3^{k-1}, \dots, z_D^{k-1})$.
 - Sample $z_2^k \sim p(z_2 | z_1^k, z_3^{k-1}, \dots, z_D^{k-1})$.
 - Sample $z_3^k \sim p(z_3 | z_1^k, z_2^k, z_4^{k-1}, \dots, z_D^{k-1})$
 - Sample ...
 - Sample $z_D^k \sim p(z_D | z_1^k, z_2^k, z_3^k, \dots, z_{D-1}^k)$.

Example: Gaussian linear model I

- Suppose we want to derive a Gibbs sampler for the following target distribution

$$y|w \sim \mathcal{N}(y|w_1x_1 + w_2x_2, \sigma^2)$$

$$w_1 \sim \mathcal{N}(w_1|0, \kappa^2)$$

$$w_2 \sim \mathcal{N}(w_2|0, \kappa^2)$$

- The posterior distribution is proportional to the joint density $p(y, w_1, w_2)$

$$\begin{aligned} p(w_1, w_2|y) &= \frac{p(y|w_1, w_2)p(w_1)p(w_2)}{p(y)} \\ &\propto p(y|w_1, w_2)p(w_1)p(w_2) \\ &= \mathcal{N}(y|w_1x_1 + w_2x_2, \sigma^2)\mathcal{N}(w_1|0, \kappa^2)\mathcal{N}(w_2|0, \kappa^2) \end{aligned}$$

- *Gibbs sampling* requires us to derive the *posterior conditionals* $p(w_1|y, w_2)$ and $p(w_2|y, w_1)$

- *General technique* for identifying $p(w_i|y, \mathbf{w}_{-i})$:

1. Write up the log joint density.
2. Identify all quantities that depends on w_i and ignore the rest.
3. Identify the distribution $p(w_i|y, \mathbf{w}_{-i})$ from its *functional form*.

Example: Gaussian linear model II

- Write out the logarithm of the joint density

$$p(w_1, w_2|y) \propto \mathcal{N}(t|w_1x_1 + w_2x_2, \sigma^2)\mathcal{N}(w_1|0, \kappa^2)\mathcal{N}(w_2|0, \kappa^2)$$

- Recall the expression for a Gaussian density

$$\mathcal{N}(x|m, v) = \frac{1}{\sqrt{2\pi v}} \exp\left(-\frac{(x-m)^2}{2v}\right)$$

- Let's write it out the log density and identify all terms that depend on w_1 or w_2

$$\begin{aligned}\log p(w_1, w_2|y) &= -\frac{1}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} (y - w_1x_1 - w_2x_2)^2 + \\ &\quad -\frac{1}{2} \log(2\pi\kappa^2) - \frac{1}{2\kappa^2} w_1^2 - \frac{1}{2} \log(2\pi\kappa^2) - \frac{1}{2\kappa^2} w_2^2 + K \\ &= -\frac{1}{2\sigma^2} (y - w_1x_1 - w_2x_2)^2 - \frac{1}{2\kappa^2} w_1^2 - \frac{1}{2\kappa^2} w_2^2 + K' \\ &= -\frac{1}{2\sigma^2} (y^2 + (w_1x_1 + w_2x_2)^2 - 2y(w_1x_1 - w_2x_2)) - \frac{1}{2\kappa^2} w_1^2 - \frac{1}{2\kappa^2} w_2^2 + K' \\ &= -\frac{1}{2\sigma^2} (w_1x_1 + w_2x_2)^2 + \frac{1}{\sigma^2} y(w_1x_1 - w_2x_2) - \frac{1}{2\kappa^2} w_1^2 - \frac{1}{2\kappa^2} w_2^2 + K'' \\ &= -\frac{1}{2\sigma^2} (w_1^2x_1^2 + w_2^2x_2^2 + 2w_1x_1w_2x_2) + \frac{1}{\sigma^2} y(w_1x_1 - w_2x_2) - \frac{1}{2\kappa^2} w_1^2 - \frac{1}{2\kappa^2} w_2^2 + K''\end{aligned}$$

Example: Gaussian linear model III

- We just arrived at

$$\log p(w_1, w_2 | y) = -\frac{1}{2\sigma^2}(w_1^2 x_1^2 + w_2^2 x_2^2 + 2w_1 x_1 w_2 x_2) + \frac{1}{\sigma^2} y(w_1 x_1 - w_2 x_2) - \frac{1}{2\kappa^2} w_1^2 - \frac{1}{2\kappa^2} w_2^2 + K''$$

- Let's compare that to a generic Gaussian distribution:

$$\begin{aligned}\log \mathcal{N}(w_1 | m, v) &= -\frac{1}{2} \log(2\pi v) - \frac{1}{2v} (w_1 - m)^2 \\ &= -\frac{1}{2} \log(2\pi v) - \frac{1}{2v} (w_1^2 + m^2 - 2w_1 m) = -\frac{1}{2v} w_1^2 + \frac{1}{v} m w_1 + C\end{aligned}$$

- Recall: *The functional form* of the logarithm of a Gaussian density is a quadratic function.

- Identify the distribution $p(w_1 | y, w_2)$ based on the functional dependence on w_1 :

$$\begin{aligned}\log p(w_1 | y, w_2) &= -\frac{1}{2\sigma^2}(w_1^2 x_1^2 + w_2^2 x_2^2 + 2w_1 x_1 w_2 x_2) + \frac{1}{\sigma^2} y(w_1 x_1 - w_2 x_2) - \frac{1}{2\kappa^2} w_1^2 - \frac{1}{2\kappa^2} w_2^2 + K'' \\ &= -\frac{1}{2\sigma^2}(w_1^2 x_1^2 + 2w_1 x_1 w_2 x_2) + \frac{1}{\sigma^2} y w_1 x_1 - \frac{1}{2\kappa^2} w_1^2 + K''' \\ &= -\frac{1}{2} \left(\frac{1}{\sigma^2} x_1^2 + \frac{1}{\kappa^2} \right) w_1^2 + \left(\frac{1}{\sigma^2} y x_1 - \frac{1}{\sigma^2} x_1 w_2 x_2 \right) w_1 + K'''\end{aligned}$$

- We conclude that the distribution $\log p(w_1 | y, w_2)$ must be a Gaussian, because *its functional form is quadratic wrt. w_1* .

Example: Gaussian linear model IV

- A generic Gaussian distribution

$$\ln \mathcal{N}(w_1 | m, v) = -\frac{1}{2v} w_1^2 + \frac{1}{v} m w_1 + C$$

- We know $p(w_1 | y, w_2) = \mathcal{N}(w_1 | m_1, v_1)$ is Gaussian, so all we need is a mean and variance

$$\log p(w_1 | y, w_2) = -\frac{1}{2} \left(\frac{1}{\sigma^2} x_1^2 + \frac{1}{\kappa^2} \right) w_1^2 + \left(\frac{1}{\sigma^2} y x_1 - \frac{1}{\sigma^2} x_1 w_2 x_2 \right) w_1 + K'''$$

- *Comparing the coefficients* for the *second order term* w_1^2 , we get the variance

$$v_1 = \left(\frac{1}{\sigma^2} x_1^2 + \frac{1}{\kappa^2} \right)^{-1}$$

- and by comparing coefficients for the *first order* term w_1 , we get the mean

$$\frac{m_1}{v_1} = \left(\frac{1}{\sigma^2} y x_1 - \frac{1}{\sigma^2} x_1 w_2 x_2 \right) \iff m_1 = \frac{v_1}{\sigma^2} (y x_1 - x_1 w_2 x_2)$$

- *By symmetry*, we get $p(w_2 | y, w_1) = \mathcal{N}(w_2 | m_2, v_2)$

$$v_2 = \left(\frac{1}{\sigma^2} x_1^2 + \frac{1}{\kappa^2} \right)^{-1} \iff m_2 = \frac{v_2}{\sigma^2} (y x_2 - x_2 w_1 x_1)$$

Example: Gaussian linear model V

- Initialize w_1 and w_2 .

1. Sample w_1 conditioned w_2 ,

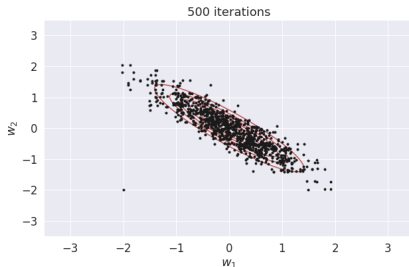
$$w_1 \sim p(w_1|y, w_2) = \mathcal{N}(w_1|m_1, v_1).$$

2. Sample w_2 conditioned w_1 ,

$$w_2 \sim p(w_2|y, w_1) = \mathcal{N}(w_2|m_2, v_2).$$

3. Repeat.

- Example shows typical *staircase behavior* of Gibbs samplers due sampling from the *posterior conditionals*



Why does Gibbs sampling work?

- The acceptance probability in Metropolis-Hastings algorithm:

$$A_k = \min \left[1, \frac{p(\mathbf{z}^*)q(\mathbf{z}^k|\mathbf{z}^*)}{p(\mathbf{z}^k)q(\mathbf{z}^*|\mathbf{z}^k)} \right]$$

- Writing $\mathbf{z}^k = \{\mathbf{z}_i^k, \mathbf{z}_{-i}^k\}$, the proposal distribution for the Gibbs sampler is

$$q(\mathbf{z}_i^*|\mathbf{z}^k) = p(\mathbf{z}_i^*|\mathbf{z}_{-i}^k). \quad (1)$$

- We need the following fact

$$p(\mathbf{z}) = p(z_i, \mathbf{z}_{-i}) = p(z_i|\mathbf{z}_{-i})p(\mathbf{z}_{-i}). \quad (2)$$

- Plugging in the proposal

$$A_k = \min \left[1, \frac{p(\mathbf{z}^*)p(\mathbf{z}_i^k|\mathbf{z}_{-i}^*)}{p(\mathbf{z}^k)p(\mathbf{z}_i^*|\mathbf{z}_{-i}^k)} \right] \quad (\text{Plugging in the proposal in eq. (1)})$$

$$= \min \left[1, \frac{p(\mathbf{z}_i^*|\mathbf{z}_{-i}^*)p(\mathbf{z}_{-i}^k)p(\mathbf{z}_i^k|\mathbf{z}_{-i}^*)}{p(\mathbf{z}_i^k|\mathbf{z}_{-i}^k)p(\mathbf{z}_{-i}^k)p(\mathbf{z}_i^*|\mathbf{z}_{-i}^k)} \right] \quad (\text{Using eq. (2)})$$

$$= \min \left[1, \frac{p(\mathbf{z}_i^*|\mathbf{z}_{-i}^k)p(\mathbf{z}_{-i}^k)p(\mathbf{z}_i^k|\mathbf{z}_{-i}^k)}{p(\mathbf{z}_i^k|\mathbf{z}_{-i}^k)p(\mathbf{z}_{-i}^k)p(\mathbf{z}_i^*|\mathbf{z}_{-i}^k)} \right] \quad (\text{Using } \mathbf{z}_{-i}^* = \mathbf{z}_{-i}^k)$$

$$= 1$$

- A Gibbs sampler is a special case of MH, where the proposed candidate is always accepted.

Questions: True or false?

Quiz via DTU Learn:

Lecture 9: Gibbs (5 questions)

Check you knowledge

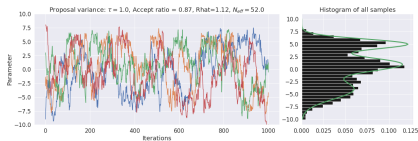
Convergence diagnostics

Are we there yet?

- MCMC theory states that samplers converge to the true target distribution as $S \rightarrow \infty$.
- *Intuitive heuristic for assessing stationarity*: Run *multiple* chains from *different initial conditions*. After K iterations, we compare the distributions for each chain. If they are different, the chains have not yet reached the stationary distribution.
- Let B denote the *between-chain variance* and let W denote the *within-chain variance*, then \hat{R} -statistic (or the *potential scale reduction factor*) for chains of length N is defined by

$$\hat{R}^2 = \frac{S - 1}{S} + \frac{1}{S} \frac{B}{W}$$

- If $B = W$, then $\hat{R} = 1$. If $B > W$, then $\hat{R} > 1$. In practice, we say that the *chains have mixed* if $\hat{R} < 1.1$



For more details and motivation, see p. 284 in Bayesian Data Analysis (<http://www.stat.columbia.edu/~gelman/book/BDA3.pdf>)

How accurate is MCMC? Quantifying the error

- Recall for i.i.d. samples, the error of the MC estimator decreases with rate $1/\sqrt{S}$:

$$\mathbb{V}[\hat{f}] = \frac{1}{S} \mathbb{V}[f(\theta)]$$

- We can *estimate the variance* based on the samples and use this to *quantify the Monte Carlo error*. Let $\widehat{\text{sd}}(f(\theta))$ be the standard deviation of the MCMC samples,

$$\text{MCSE} = \frac{1}{\sqrt{S}} \widehat{\text{sd}}(f(\theta))$$

- MCMC methods produces *highly correlated* samples. The *autocorrelation function* ρ_t measures the correlation between two samples θ^i and θ^{i+t}

$$\rho_t = \frac{1}{\sigma^2} \int (\theta^i - \mu)(\theta^{i+t} - \mu) p(\theta) d\theta$$

- The *effective sample size (ESS)* S_{eff} takes this correlation into account

$$S_{\text{eff}} = \frac{S}{\sum_{t=-\infty}^{\infty} \rho_t} = \frac{S}{1 + 2 \sum_{t=1}^{\infty} \rho_t}$$

- and then

$$\text{MCSE} = \frac{1}{\sqrt{S_{\text{eff}}}} \widehat{\text{sd}}(f(\theta))$$

Example: MCMC diagnostics

- The \hat{R} – statistic

$$\hat{R}^2 = \frac{N-1}{N} + \frac{1}{N} \frac{B}{W}$$

- The effective sample size

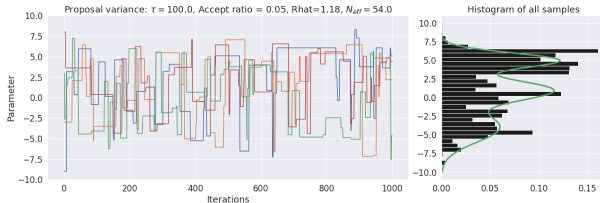
$$S_{\text{eff}} = \frac{S}{1 + 2 \sum_{t=1}^{\infty} \rho_t}$$

- The Monte Carlo
Standard Error

$$\text{MCSE} = \frac{1}{\sqrt{S_{\text{eff}}}} \widehat{\text{sd}}(f(\theta))$$

- Next week

1. A few words on Hamiltonian Monte Carlo
2. We will do a short discussion on pros and cons of MCMC in practice
3. Start discussing variational inference.



Hierarchical models

Revisiting the Poisson regression

- "*Being Bayesian*" usually refers to treating quantities of interest as *random variable* and reason using the rules of *probability theory*
- When we talked about "*fully Bayesian*" inference, we refer the setting, where we have *prior distributions on all parameters, including hyperparameters*
- Before the holidays, we worked with a fully Bayesian Poisson regression model via MCMC

$$y_n | \mu_n \sim \text{Poisson}(\mu_n),$$

$$\mu_n = \exp(f_n)$$

$$f_n = \mathbf{w}^T \mathbf{x}_n$$

$$\mathbf{w} | \kappa \sim \mathcal{N}(0, \kappa^2 \mathbf{I})$$

$$\kappa \sim \mathcal{N}_+(0, 1)$$

with the following joint distribution

$$p(\mathbf{y}, \mathbf{w}, \kappa) = \prod_{n=1}^N p(y_n | \mathbf{w}) p(\mathbf{w} | \kappa) p(\kappa),$$

- This is an example of a *hierarchical model*

Hierarchical modelling

- *Hierarchical* or *multi-level* models are one of the key strength of the Bayesian framework
- Suppose we have a model $p(\mathcal{D}|\theta)$ with data \mathcal{D} , parameters θ and hyperparameters ξ



with the following joint distribution

$$p(\mathcal{D}, \theta, \xi) \propto p(\mathcal{D}|\theta)p(\theta|\xi)p(\xi)$$

- Useful when you want to
 1. make inference more robust
 2. reason probabilistically about data, parameters and hyperparameters
 3. model hierarchical structure in data (random effects models)
 4. squeeze out every bit of predictive performance of a model/dataset

- Degrees of "Bayesianity" according to Murphy1

Method	Definition
Maximum likelihood	$\hat{\theta} = \arg \max_{\theta} p(\mathcal{D} \theta)$
MAP	$\hat{\theta} = \arg \max_{\theta} p(\mathcal{D} \theta)p(\theta \xi)$
ML-II	$\hat{\xi} = \arg \max_{\xi} \int p(\mathcal{D} \theta)p(\theta \xi)d\theta$
MAP-II	$\hat{\xi} = \arg \max_{\xi} \int p(\mathcal{D} \theta)p(\theta \xi)p(\xi)d\theta$
Full Bayes	$p(\theta, \xi \mathcal{D}) \propto p(\mathcal{D} \theta)p(\theta)p(\xi)$

Example

- Recall the model for Bayesian linear regression

$$p(\mathbf{w}|\alpha) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I}) \quad (\text{prior})$$

$$p(\mathbf{y}|\mathbf{w}, \beta) = \mathcal{N}(\mathbf{y}|\Phi\mathbf{w}, \beta^{-1}\mathbf{I}) \quad (\text{likelihood})$$

$$p(\mathbf{w}|\mathbf{y}, \alpha, \beta) = \mathcal{N}(\mathbf{w}|\mathbf{m}, \mathbf{S}) \quad (\text{posterior})$$

$$p(\mathbf{y}|\alpha, \beta) = \mathcal{N}(\mathbf{y}|\mathbf{0}, \beta^{-1}\mathbf{I} + \alpha^{-1}\Phi\Phi^T) \quad (\text{marginal likelihood})$$

- Fully Bayesian inference on the full joint distribution

$$p(\alpha, \beta, \mathbf{w}|\mathbf{y}) \propto p(\mathbf{y}|\mathbf{w}, \beta)p(\mathbf{w}|\alpha)p(\alpha, \beta)$$

- Fully Bayesian inference on the marginalized joint distribution

$$p(\alpha, \beta|\mathbf{y}) \propto p(\mathbf{y}|\alpha, \beta)p(\alpha, \beta)$$

- Making predictions via MCMC

$$p(y^*|\mathbf{y}) = \mathbb{E}_{p(\alpha, \beta|\mathbf{y})} [p(y^*|\mathbf{y}, \alpha, \beta)] \approx \frac{1}{S} \sum_{i=1}^S p(y^*|\mathbf{y}, \alpha^{(i)}, \beta^{(i)})$$

for $\alpha^{(i)}, \beta^{(i)} \sim p(\alpha, \beta|\mathbf{y})$

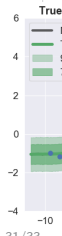
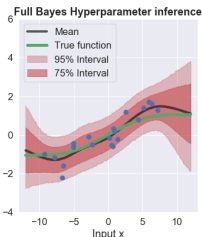
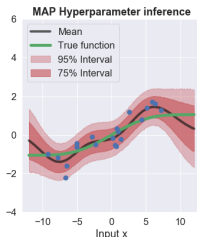
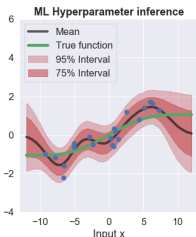
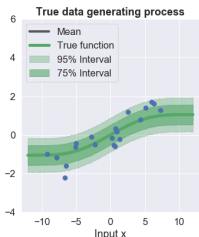
Fully Bayesian Gaussian process regression: Example

- Example with $N = 20$ data points and additive Gaussian noise
- Gaussian process regression with squared exponential kernel
- We impose a *weakly informative* prior on the lengthscale to
 1. rule out really short length scales (smaller than the grid size)
 2. rule out really long length scales (larger than span data of data)
- Joint distribution

$$p(\mathbf{y}, \mathbf{f}, \ell, \sigma, \kappa) = p(\mathbf{y}|\mathbf{f}, \sigma)p(\mathbf{f}|\ell, \kappa)p(\kappa)p(\ell)$$

- Marginalized joint distribution

$$p(\mathbf{y}, \ell, \sigma, \kappa) = p(\mathbf{y}|\sigma, \ell, \kappa)p(\kappa)p(\ell)$$



Medical example

- Suppose you work for a medical company and are asked to analyze data from a drug evaluation on rats prior to human trials

- Suppose the drug was administered to N rats, where y rats ended up developing tumors.

- We could use a Beta-Binomial model to estimate the probability of developing tumors θ

$$p(\theta|y) \propto \text{Bin}(y|N, \theta)\text{Beta}(\theta|\alpha_0, \beta_0)$$

- What if the company tested the drugs on J different types of rats, where (y_j, N_j) denote the data for the j 'th group. How to analyze the data?

1. We could fit J *individual* models

$$p(\theta_j|y_j) \propto \text{Bin}(y_j|N_j, \theta_j)\text{Beta}(\theta_j|\alpha_0, \beta_0)$$

2. We could *pool* all the data such that $N_p = \sum_j N_j$ and $y_p = \sum_j y_j$

$$p(\theta_p|y_p) \propto \text{Bin}(y_p|N_p, \theta)\text{Beta}(\theta_p|\alpha_0, \beta_0)$$

- *Individual* models may perform poorly if N_j is small, but the *pooled* model might perform poorly if the different groups exhibit different behaviors
- *Bayesian hierarchical models* allows us to borrow statistical strength from groups with lots of data to help groups with less data

Medical example cont.

- Let $\mathbf{y} = \{y_j\}_{j=1}^J$ and $\boldsymbol{\theta} = \{\theta_j\}_{j=1}^J$, then

$$p(\mathbf{y}, \boldsymbol{\theta}, \alpha, \beta) = \prod_{j=1}^J \text{Bin}(y_j | N_j, \theta_j) \prod_{j=1}^J \text{Beta}(\theta_j | \alpha, \beta) p(\alpha, \beta)$$

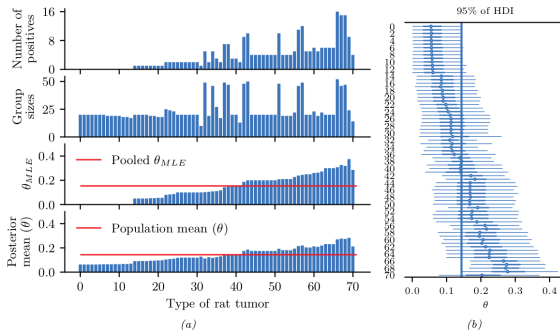


Figure 3.15: Data and inferences for the hierarchical binomial model fit using HMC. Generated by `hierarchical_binom_rats.ipynb`.

The detailed balance condition I

- If a chain satisfies the *detailed balance condition*, then p^* is its stationary distribution

$$T(z'|z)p^*(z) = T(z|z')p^*(z')$$

- Integrating both sides wrt. z' yields the stationary distribution p^*

$$\int T(z'|z)p^*(z)dz' = p^*(z) = \int T(z|z')p^*(z')dz'$$

- Does the Metropolis-Hastings satisfy this condition? Let's check
- Recall the acceptance probability for Metropolis-Hastings

$$A(z^*|z^k) = \min \left[1, \frac{p(z^*)q(z^k|z^*)}{p(z^k)q(z^*|z^k)} \right]$$

- The transition kernel for Metropolis-Hasting

$$T(z'|z) = \begin{cases} q(z'|z)A(z'|z) & \text{if } z' \neq z \\ q(z|z)A(z|z) + \int q(z''|z) [1 - A(z''|z)] dz'' & \text{if } z' = z \end{cases}$$

The detailed balance condition II

- The acceptance probability

$$A(\mathbf{z}^*, \mathbf{z}^k) = \min \left[1, \frac{p(\mathbf{z}^*)q(\mathbf{z}^k|\mathbf{z}^*)}{p(\mathbf{z}^k)q(\mathbf{z}^*|\mathbf{z}^k)} \right]$$

- The transition kernel

$$T(\mathbf{z}'|\mathbf{z}) = \begin{cases} q(\mathbf{z}'|\mathbf{z})A(\mathbf{z}'|\mathbf{z}) & \text{if } \mathbf{z}' \neq \mathbf{z} \\ q(\mathbf{z}|\mathbf{z})A(\mathbf{z}|\mathbf{z}) + \int q(\mathbf{z}''|\mathbf{z}) [1 - A(\mathbf{z}''|\mathbf{z})] d\mathbf{z}'' & \text{if } \mathbf{z}' = \mathbf{z} \end{cases}$$

- If $p(\mathbf{z}^k)q(\mathbf{z}^*|\mathbf{z}^k) > p(\mathbf{z}^*)q(\mathbf{z}^k|\mathbf{z}^*)$, then $A(\mathbf{z}^*|\mathbf{z}^k) < 1$ and $A(\mathbf{z}^k|\mathbf{z}^*) = 1$

- To jump from \mathbf{z}^k to \mathbf{z}^* , we first need to propose it and then accept it

$$T(\mathbf{z}^*|\mathbf{z}^k) = q(\mathbf{z}^*|\mathbf{z}^k)A(\mathbf{z}^*|\mathbf{z}^k) = q(\mathbf{z}^*|\mathbf{z}^k) \frac{p(\mathbf{z}^*)q(\mathbf{z}^k|\mathbf{z}^*)}{p(\mathbf{z}^k)q(\mathbf{z}^*|\mathbf{z}^k)} = \frac{p(\mathbf{z}^*)q(\mathbf{z}^k|\mathbf{z}^*)}{p(\mathbf{z}^k)}$$

- It follows

$$p(\mathbf{z}^k)T(\mathbf{z}^*|\mathbf{z}^k) = p(\mathbf{z}^*)q(\mathbf{z}^k|\mathbf{z}^*)$$

The detailed balance condition III

- We just arrived at

$$p(\mathbf{z}^k)T(\mathbf{z}^*|\mathbf{z}^k) = p(\mathbf{z}^*)q(\mathbf{z}^k|\mathbf{z}^*)$$

- What about the opposite direction?

$$T(\mathbf{z}^k|\mathbf{z}^*) = q(\mathbf{z}^k|\mathbf{z}^*)A(\mathbf{z}^k|\mathbf{z}^*) = q(\mathbf{z}^k|\mathbf{z}^*)$$

- Combining the two and we are done

$$p(\mathbf{z}^k)T(\mathbf{z}^*|\mathbf{z}^k) = p(\mathbf{z}^*)T(\mathbf{z}^k|\mathbf{z}^*)$$