Bayesian Learning

Lecture 8 - Markov Chain Monte Carlo. Metropolis-Hastings.

Mattias Villani

Department of Statistics Stockholm University











Lecture overview

Markov Chain Monte Carlo

Metropolis-Hastings

■ MCMC - efficiency, burn-in and convergence

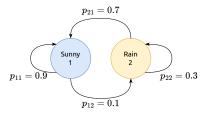
Stochastic process

- Let $S = \{1, 2, ..., K\}$ be a finite set of **states**.
 - ▶ Weather: $S = \{\text{sunny}, \text{ rain}\}.$
 - ▶ School grades: $S = \{A, B, C, D, E, F\}$
- **Stochastic process**: collection of random variables X_1, X_2, X_3, \ldots , often over time.
- A time series with categories.
- Weather: $X_1 = \text{sunny}$, $X_2 = \text{sunny}$, $X_3 = \text{rain}$, $X_4 = \text{sunny}$.
- School grades: $X_1 = C$, $X_2 = C$, $X_3 = B$, $X_4 = A$, $X_5 = B$.

Markov chains

Markov chain: probability distribution of tomorrow's state X_{t+1} depends (only) on today's state X_t :

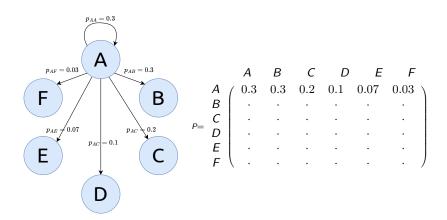
$$p_{ij} = \Pr(X_{t+1} = j | X_t = i)$$



■ Transition matrix for weather example

$$\mathbf{P} = \begin{pmatrix} p_{11} & p_{12} \\ p_{21} & p_{22} \end{pmatrix} = \begin{pmatrix} 0.9 & 0.1 \\ 0.7 & 0.3 \end{pmatrix}$$

Markov chains - Grades example



Stationary distribution

■ *h*-step transition probabilities

$$p_{ij}^{(h)} = \Pr(X_{t+h} = j | X_t = i)$$

■ h-step transition matrix by matrix power

$$\mathbf{P}^{(h)} = \mathbf{P}^h$$

- Unique equilbrium distribution $\pi = (\pi_1, ..., \pi_k)$ if chain is ergodic:
 - irreducible (possible to get to any state from any state)
 - aperiodic (does not get stuck in predictable cycles)
 - positive recurrent (expected time of returning is finite)
- Limiting long-run distribution as $h \to \infty$

$$m{P}^h
ightarrow \left(egin{array}{c} m{\pi} \ m{\pi} \ dots \ m{\pi} \end{array}
ight) = \left(egin{array}{cccc} \pi_1 & \pi_2 & \cdots & \pi_k \ \pi_1 & \pi_2 & \cdots & \pi_k \ dots & dots & dots \ \pi_1 & \pi_2 & \cdots & \pi_k \end{array}
ight)$$

Stationary distribution, cont.

Limiting long-run distribution as $h \to \infty$

$$m{P}^h
ightarrow \left(egin{array}{c} m{\pi} \ m{\pi} \ dots \ m{\pi} \end{array}
ight) = \left(egin{array}{cccc} \pi_1 & \pi_2 & \cdots & \pi_k \ \pi_1 & \pi_2 & \cdots & \pi_k \ dots & dots & dots \ \pi_1 & \pi_2 & \cdots & \pi_k \end{array}
ight)$$

Stationary distribution

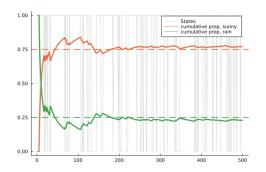
$$oldsymbol{\pi} = oldsymbol{\pi} oldsymbol{P}$$

- Draw starting point using π . All future states of the Markov Chain will have distribution π .
- Weather example:

$$\mathbf{P} = \begin{pmatrix} 0.9 & 0.1 \\ 0.3 & 0.7 \end{pmatrix}, \mathbf{P}^2 = \begin{pmatrix} 0.84 & 0.16 \\ 0.42 & 0.58 \end{pmatrix}$$
$$\mathbf{P}^5 = \begin{pmatrix} 0.77 & 0.23 \\ 0.69 & 0.31 \end{pmatrix}, \mathbf{P}^{100} = \begin{pmatrix} 0.75 & 0.25 \\ 0.75 & 0.25 \end{pmatrix}$$
$$\pi = (0.75, 0.25)$$

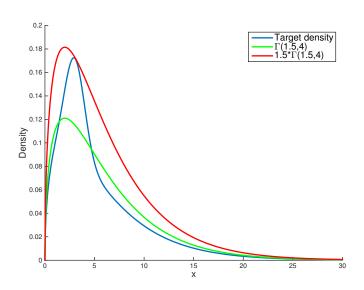
The basic MCMC idea

- Aim: simulate from a discrete distribution p(x).
- MCMC: simulate a Markov Chain with a stationary distribution that is exactly p(x).
- How to set up the transition matrix **P**? Metropolis-Hastings!



Can be extended to continuous random variables.

Rejection sampling



Random walk Metropolis algorithm

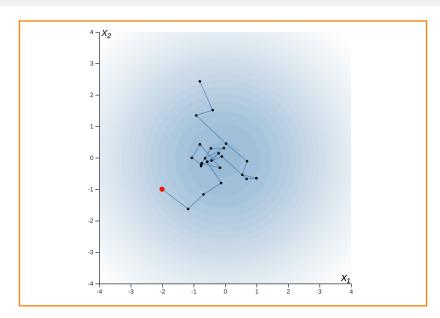
- Initialize $\theta^{(0)}$ and iterate for i = 1, 2, ...

 - **2** Compute the acceptance probability

$$\alpha = \min \left(1, \frac{p(\boldsymbol{\theta}^*|\mathbf{y})}{p(\boldsymbol{\theta}^{(i-1)}|\mathbf{y})} \right)$$

If $u \sim \mathrm{Uniform}(0,1)$ If $u \leq \alpha$ set $\boldsymbol{\theta}^{(i)} = \boldsymbol{\theta}^{\star}$ (accept and move) If $u > \alpha$ set $\boldsymbol{\theta}^{(i)} = \boldsymbol{\theta}^{(i-1)}$ (reject and stay)

Interactive - Random Walk Metropolis



Proportional form of posterior is enough!

- Assumption: we can compute $p(\theta|\mathbf{y})$ for any θ .
- Proportionality constants in posterior cancel out in

$$\alpha = \min \left(1, \frac{\rho(\boldsymbol{\theta}^{\star}|\mathbf{y})}{\rho(\boldsymbol{\theta}^{(i-1)}|\mathbf{y})} \right).$$

■ Proportional form of posterior is enough!

$$\alpha = \min \left(1, \frac{p(\mathbf{y}|\boldsymbol{\theta}^*) p(\boldsymbol{\theta}^*)}{p(\mathbf{y}|\boldsymbol{\theta}^{(i-1)}) p(\boldsymbol{\theta}^{(i-1)})} \right)$$

Random walk Metropolis - choosing a proposal

- Common choices of Σ in proposal $N(\theta^{(i-1)}, c \cdot \Sigma)$:
 - $m \Sigma = m I$ (proposes 'off the cigar')
 - $oldsymbol{\Sigma} = J_{oldsymbol{y}}^{-1}(ilde{oldsymbol{ heta}})$ (propose 'along the cigar')
 - **Adaptive**. Start with $\Sigma = I$. Update Σ from initial run.
- Set c so average acceptance probability is 25-30%.
- **■** Good proposal:
 - Easy to sample
 - **Easy to compute** α
 - \triangleright Proposals should take reasonably large steps in θ -space
 - Proposals should not be reject too often.

The Metropolis-Hastings algorithm

- Generalization when the proposal density is not symmetric.
- Initialize $\theta^{(0)}$ and iterate for i = 1, 2, ...
 - **1** Sample proposal: $oldsymbol{ heta}^{\star} \sim q\left(\cdot | oldsymbol{ heta}^{(i-1)}
 ight)$
 - 2 Compute the acceptance probability

$$\alpha = \min \left(1, \frac{p(\boldsymbol{\theta}^{\star}|\mathbf{y})}{p(\boldsymbol{\theta}^{(i-1)}|\mathbf{y})} \frac{q\left(\boldsymbol{\theta}^{(i-1)}|\boldsymbol{\theta}^{\star}\right)}{q\left(\boldsymbol{\theta}^{\star}|\boldsymbol{\theta}^{(i-1)}\right)} \right)$$

3 Draw $u \sim \mathrm{Uniform}(0,1)$ If $u \leq \alpha$ set $\boldsymbol{\theta}^{(i)} = \boldsymbol{\theta}^{\star}$ (accept and move) If $u > \alpha$ set $\boldsymbol{\theta}^{(i)} = \boldsymbol{\theta}^{(i-1)}$ (reject and stay)

The independence sampler

- Independence sampler: $q\left(\theta^{\star}|\theta^{(i-1)}\right) = q\left(\theta^{\star}\right)$.
- Proposal is independent of previous draw.
- Example: a multivariate-t distribution (we want heavy tails)

$$oldsymbol{ heta}^{\star} \sim t\left(ilde{oldsymbol{ heta}}, extstyle J_{ extstyle y}^{-1}(ilde{oldsymbol{ heta}}),
u
ight),$$

where $\tilde{\boldsymbol{\theta}}$ and $\textit{J}_{y}(\tilde{\boldsymbol{\theta}})$ are computed by numerical optimization.

- Can be very efficient, but has a tendency to get stuck.
- Make sure that $q(\theta^*)$ has heavier tails than $p(\theta|\mathbf{y})$.

The efficiency of MCMC

- How efficient is MCMC compared to iid sampling?
- If $\theta^{(1)}, \theta^{(2)}, ..., \theta^{(m)}$ are iid with variance σ^2 , then

$$\operatorname{Var}(\bar{\theta}) = \frac{\sigma^2}{m}.$$

Autocorrelated $\theta^{(1)}, \theta^{(2)}, ..., \theta^{(m)}$ generated by MCMC

$$\operatorname{Var}(\bar{\theta}) = \frac{\sigma^2}{m} \left(1 + 2 \sum_{k=1}^{\infty} \rho_k \right)$$

where $\rho_k = Corr(\theta^{(i)}, \theta^{(i+k)})$ is the autocorrelation at lag k.

■ Inefficiency factor (for large enough *K*)

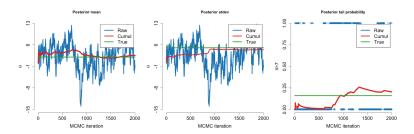
$$IF = 1 + 2\sum_{k=1}^{K} \rho_k$$

Effective sample size from MCMC

$$ESS = \frac{m}{IF}$$

Burn-in and convergence

- How long burn-in?
- How long to sample after burn-in?
- Convergence diagnostics
 - ► Raw plots of simulated sequences (MCMC trajectories)
 - Cumulative estimates plots
 - Repeated runs with different initial values
 - ▶ Potential scale reduction factor, R.



Burn-in and convergence

