

Quiz 6: Fri
HW 3: next Fri

COMPSCI 688: Probabilistic Graphical Models

Lecture 14: Markov Chain Monte Carlo

Dan Sheldon

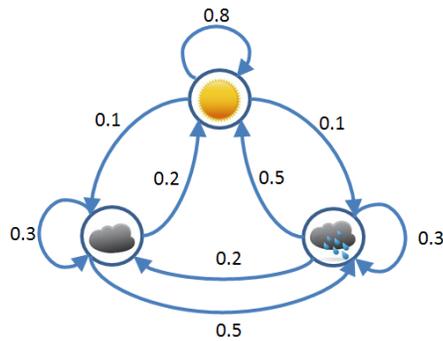
Manning College of Information and Computer Sciences
University of Massachusetts Amherst

Partially based on materials by Benjamin M. Marlin (marlin@cs.umass.edu) and Justin Domke (domke@cs.umass.edu)

Markov Chain Theory

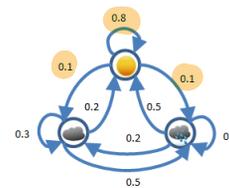
Markov Chains

A discrete Markov chain is a **set of states** with **transition probabilities** between each pair of states. **Example** (note: not a graphical model!)



Transition Matrix

- ▶ The probabilistic transitions in the state diagram can also be represented by an equivalent matrix of transition probabilities.
- ▶ The "from" states are rows and the "to" states are columns.



	To		
From	Yellow	Grey	Blue
Yellow	0.8	0.1	0.1
Grey	0.2	0.3	0.5
Blue	0.5	0.2	0.3

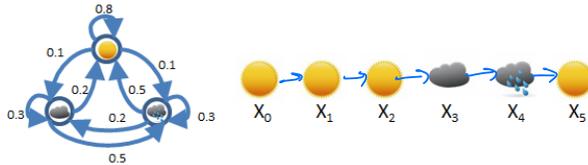
transition matrix

= T

T_{ij} = prob next state = j given current state = i

Markov Chains: Simulation and State Sequences

► To simulate a Markov chain, we draw $x_0 \sim p_0$, then repeatedly sample x_{t+1} given the current state x_t according to the transition probabilities T .



Markov Chain: Formal Definition

By repeatedly making random transitions from a starting state, we generate a *chain* of random variables $X_0, X_1, X_2, X_3, \dots$

Formally, a Markov chain is specified by:

- A set of states $\{1, 2, \dots, D\}$
- A starting distribution p_0 with $p_0(i) = P(X_0 = i)$.
- Transition probabilities $T_{ij} = P(X_{t+1} = j | X_t = i)$ for all $i, j \in \{1, 2, \dots, D\}$

A Markov chain **assumes the Markov property**:

$$P(X_t = x_t | X_0 = x_0, X_1 = x_1, \dots, X_{t-1} = x_{t-1}) = P(X_t = x_t | X_{t-1} = x_{t-1})$$

$$x_0 \rightarrow x_1 \rightarrow x_2 \rightarrow x_3 \rightarrow x_4 \rightarrow \dots$$

Markov Chain Questions

$$x_0 \rightarrow x_1 \rightarrow x_2 \rightarrow \dots \rightarrow x_t$$

p_0

Three important questions:

1. What is the joint probability of a sequence of states of length N ?
2. What is the marginal probability distribution over states after a given number of steps t ? $p(x_t)$
3. What happens to the probability distribution over states in the limit as t goes to infinity? $t \rightarrow \infty$

Markov Chain Factorization

$$x_0 \rightarrow x_1 \rightarrow x_2 \rightarrow \dots \rightarrow x_N$$

Question: What is the joint probability over the state sequence x_0, \dots, x_N ?

Answer: by the Markov property:

$$P(X_1 = x_1, \dots, X_N = x_N | X_0 = x_0) = P(X_1 = x_1 | X_0 = x_0) \times P(X_2 = x_2 | X_1 = x_1) \times \dots \times P(X_N = x_N | X_{N-1} = x_{N-1})$$

Shorter version:

$$p(x_1, x_2, \dots, x_N | x_0) = p(x_1 | x_0) p(x_2 | x_1) \dots p(x_N | x_{N-1})$$

$$= T_{x_0 x_1} \times T_{x_1 x_2} \times \dots \times T_{x_{N-1} x_N}$$

T_{ij}

The t -Step Distribution for Fixed x_0 $x_0 \rightarrow x_1 \rightarrow \dots \rightarrow x_t$

Question: What is the marginal probability distribution after t steps given that the chain starts at x_0 ? I.e., what is $p(x_t|x_0)$?

Examples:

$$p(x_1|x_0) = T_{x_0 x_1}$$

$$p(x_2|x_0) = \sum_{x_1} p(x_1, x_2|x_0) = \sum_{x_1} p(x_1|x_0) p(x_2|x_1)$$

In general, we have the recursive expression:

$$p(x_t|x_0) = \sum_{x_{t-1}} p(x_{t-1}, x_t|x_0) = \sum_{x_{t-1}} p(x_{t-1}|x_0) p(x_t|x_{t-1})$$

$$= \sum_{x_{t-1}} p(x_{t-1}|x_0) \cdot T_{x_{t-1} x_t}$$

The t -Step Distribution for Random X_0 $X_0 \rightarrow X_1 \rightarrow \dots \rightarrow X_t$ p_0 $P(x_t)$

Question: What is the marginal probability distribution after t steps given that $X_0 \sim p_0$? I.e., what is $p(x_t)$?

By similar logic:

$$p(x_1) = \sum_{x_0} p_0(x_0) T_{x_0 x_1}$$

$$p(x_2) = \sum_{x_1} p(x_1, x_2) = \sum_{x_1} p(x_1) \cdot T_{x_1 x_2}$$

In general:

$$p(x_t) = \sum_{x_{t-1}} p(x_{t-1}) \cdot T_{x_{t-1} x_t}$$

recurrence for marginals

time t time $t-1$

t -Step Recurrence as Matrix-Vector Multiplication

The recurrences for the t -step distributions can be expressed using matrix-vector multiplication. Let p_t be the row-vector

$$p_t = [P(x_t = 1), P(x_t = 2), \dots, P(x_t = D)]$$

Then, since $T_{ij} = P(X_t = j | X_{t-1} = i)$, we can write the above recursive relationship as

$$p_t = p_{t-1} T.$$

$$[p_t(i) \dots p_t(j) \dots p_t(D)] = [p_{t-1}(i) \dots p_{t-1}(D)] \begin{bmatrix} T_{1j} \\ T_{2j} \\ \vdots \\ T_{Dj} \end{bmatrix}$$

$$p_t(j) = \sum_i p_{t-1}(i) T_{ij}$$

t-Step Distribution as Matrix Power

By unrolling the recurrence, the t -step distribution can be obtained as a matrix power

$$\begin{aligned}
 p_t &= p_{t-1}T \\
 &= (p_{t-2})T \\
 &= (p_{t-2}T)T \\
 &= (p_{t-2})TT \\
 &= (p_{t-3}T)TT \\
 &\vdots \\
 &= p_0 \underbrace{TT \dots T}_{t \text{ times}} = p_0 T^t
 \end{aligned}$$

Thus

$$p_t = p_0 T^t.$$

This also implies that T^t is the t -step transition matrix

$$(T^t)_{ij} = P(X_t = j | X_0 = i) = P(X_{s+t} = j | X_s = i)$$

One-Slide Summary So Far

- ▶ Markov chain: defined by initial distribution $p_0 \in \mathbb{R}^D$, transition matrix $T \in \mathbb{R}^{D \times D}$

$$p_0(i) = P(X_0 = i), \quad T_{ij} = P(X_t = j | X_{t-1} = i)$$

- ▶ Defines distribution of chain $X_0, X_1, X_2, \dots, X_t, \dots$ (with Markov assumption)
- ▶ Joint probability

$$p(x_1, x_2, \dots, x_N | x_0) = p(x_1 | x_0) p(x_2 | x_1) \dots p(x_N | x_{N-1})$$

- ▶ Recurrence for t -step distribution: $p(x_t) = \sum_{x_{t-1}} p(x_{t-1}) T_{x_{t-1} x_t}$
- ▶ Recurrence as matrix-vector multiplication. Let $p_t \in \mathbb{R}^D$ with $p_t(i) = P(X_t = i)$. Then

$$p_t = p_{t-1} T$$

- ▶ **Next:** what happens as $t \rightarrow \infty$?

Limiting Distribution

Given: π

What happens as t becomes large? Does p_t converge to a some *limiting distribution* π ? That is, is there some π such that the following is true?

$$\lim_{t \rightarrow \infty} p_t = \pi \quad (\text{limiting distribution})$$

The algorithmic idea of Markov chain Monte Carlo is:

- ▶ Suppose π is hard to sample from directly
- ▶ If we can **design a Markov chain** such that $\lim_{t \rightarrow \infty} p_t = \pi$, then we can draw samples by simulating the Markov chain for many time steps
- ▶ It's remarkable that this could be possible, but it can be done for very general target distributions!
- ▶ We need to reason about limiting distributions their properties

Stationary Distribution

$$\lim_{t \rightarrow \infty} p_t = \pi$$

Suppose a chain converges exactly, so that $p_t = p_{t+1} = \pi$. Since $p_{t+1} = p_t T$, this implies

$$\pi = \pi T \quad (\text{stationary distribution})$$

- ▶ we call any such π a *stationary distribution* of the Markov chain
- ▶ If you start from π and run the chain for any number of steps, the distribution is unchanged.
- ▶ If π is a limiting distribution, it is a stationary distribution
- ▶ (Linear algebra connection: π is an *eigenvector* of T with *eigenvalue* 1. Useful for computing stationary distributions.)

Stationary and Limiting Distributions

We reason about *limiting distributions* via *stationary distributions*:

- ▶ If a Markov chain: (1) converges, and (2) has a *unique* stationary distribution π , then it converges to π .
- ▶ When can we guarantee (1) and (2)? What could go wrong?

What Could Go Wrong: Periodicity

A Markov chain can fail to converge by being periodic:



Spse $X_0=1 \quad X_1=2 \quad X_2=1 \quad X_3=2$

$$p_0 = [1 \ 0] \quad p_1 = [0 \ 1] \quad p_2 = [1 \ 0] \quad p_3 = [0 \ 1]$$

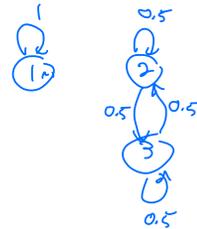
$$p_0 = [d, 1-d] \quad p_1 = [1-d, d] \quad [d, 1-d] \quad \dots$$



On the other hand, if $p_0 = [\frac{1}{2}, \frac{1}{2}]$ it does converge.

What Could Go Wrong: Reducibility

A Markov chain can fail to have a unique stationary distribution by being reducible:



$$\pi_d = \left[d, \frac{1-d}{2}, \frac{1-d}{2} \right] \text{ is stationary}$$

for every $d \in [0, 1]$

Regularity

A Markov chain is **regular** if there exists a t such that, for all i, j pairs,

$$Pr(X_{s+t}=j | X_s=i) = (T^t)_{ij} > 0,$$

- ▶ Recall that T^t is the t -step transition probability matrix. This means it is possible to get *from* any state i to any state j in exactly t steps.
- ▶ A regular Markov chain cannot be periodic or reducible (why?), and guarantees the desired computational property

Theorem: A regular Markov chain has a unique stationary distribution π and $\lim_{t \rightarrow \infty} p_t = \pi$ for all starting distributions p_0 .

(We can sample from the unique stationary distribution by simulating the chain.)

Summary: Markov Chain Theory

- ▶ **t -step distribution:** Distribution of X_t , obtained by repeated multiplication with transition matrix: $p_t = p_0 T^t$
- ▶ **Limiting distribution:** the distribution of $\lim_{t \rightarrow \infty} p_t$, if it exists
- ▶ **Stationary distribution:** a distribution π such that $\pi T = \pi$. If you start from π and run the chain for any number of steps, the distribution is unchanged. Every limiting distribution is a stationary distribution.
- ▶ **Regularity:** if there is a t such that $(T^t)_{ij} > 0$ for all i, j , a Markov chain is regular. It is possible to get from any state i to any state j in exactly t steps.
- ▶ **Convergence to stationary distribution:** if T is regular, the chain converges to a unique stationary distribution π for any starting distribution.

Understanding MCMC

High-Level Idea

Suppose we want to sample from p , but can't do so directly. Instead, we can

- ▶ **Design a Markov chain** that has p as a stationary distribution
- ▶ Run it for a long time to get a sequence of states x_1, x_2, \dots, x_S
- ▶ Approximate an expectation as

$$\mathbb{E}_{p(X)}[f(X)] \approx \frac{1}{S} \sum_{t=1}^S f(x_t).$$

x_t

If we run the chain long enough, the approximation will be good! We can often make the following guarantees:

- ▶ Asymptotically correct: $\lim_{S \rightarrow \infty} \frac{1}{S} \sum_{t=1}^S f(x_t) = \mathbb{E}_{p(X)}[f(X)]$
- ▶ Variance decreases like $1/S$
- ▶ The chain converges exponentially quickly to the stationary distribution, so bias decreases quickly. (But in practice, we almost never know the rate!)

$bias = C \cdot 0.9999999^t$

Some concerns:

- ▶ X_1, X_2, \dots are not true samples from p , especially early in the chain
 - ▶ X_1, X_2, \dots, X_S are not independent
 - ▶ How to create a Markov chain with p as a stationary distribution?
 - ▶ How to make sure that p is the only stationary distribution? (check regular)
 - ▶ How long to run the chain?
 - ▶ How to initialize the chain?
 - ▶ What is the best Markov chain?
- } practical

MCMC for Multivariate Distributions

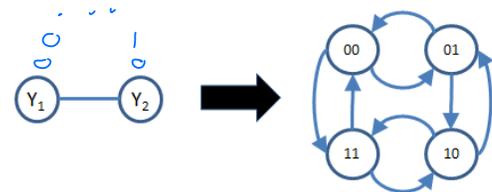
(x_1, \dots, x_D)

- ▶ To sample from a multivariate distribution $p(\mathbf{x})$ for $\mathbf{x} \in \mathbb{R}^D$, an MCMC algorithm generates a sequence of states

○○○ 101 110 ... $2^3 = 8$ states
 $x_1, x_2, x_3, \dots, x_S$

- ▶ Each $\mathbf{x}_t = (x_{t1}, \dots, x_{tD})$ is a full vector — with a setting for each variable
- ▶ The state space of the Markov chain is the full domain $\mathbf{x} \in \text{Val}(\mathbf{X})$. E.g., with D binary variables, the Markov chain has 2^D states. $T = 2^D \times 2^D$
- ▶ Because state spaces are huge, MCMC algorithms specify rules for random transitions between states without materializing the full transition matrix.

Example: Binary MRF



MRF: Two Binary-Valued Random Variables

Markov Chain: One Random Variable with Four states

Detailed Balance \implies Stationary

Theorem: If T satisfies detailed balance with respect to π then π is a stationary distribution of T .

Proof: Let $\pi' = \pi T$ be the result of running the Markov chain for 1 iteration. Then

$$\begin{aligned}\pi'(x') &= \sum_x \pi(x) T(x'|x) \quad \text{detailed balance} \\ &= \sum_x \pi(x') T(x|x') \\ &= \pi(x') \underbrace{\sum_x T(x|x')}_{=1} \\ &= \pi(x')\end{aligned}$$