



Markov
Chain Monte
Carlo
Methods

Christoph
Teichmann,
Antoine
Venant

Markov Chain Monte Carlo Methods

Christoph Teichmann Antoine Venant

Goals

Motivation

Markov
Chains

Invariant
Tricks



Markov
Chain Monte
Carlo
Methods

Christoph
Teichmann,
Antoine
Venant

Goals

Motivation

Markov
Chains

Invariant
Tricks

Goals

Markov
Chain Monte
Carlo
Methods

Christoph
Teichmann,
Antoine
Venant

Goals

Motivation

Markov
Chains

Invariant
Tricks

- Problem of Bayesian Inference
- Markov Chain Monte Carlo
- Metropolis-Hasting Technique
- Gibbs Technique



Markov
Chain Monte
Carlo
Methods

Christoph
Teichmann,
Antoine
Venant

Goals

Motivation

Markov
Chains

Invariant
Tricks

Motivation

We were discussing how to learn a language model:

- Bigram model for text
- Probabilities are hidden variables
- Dirichlet Prior for probabilities

Can we improve this model?

Assume the words are generated from hidden states

- $h \in H$ hidden tags
- $w \in L$ words
- P_w distribution over words give hidden tags (one per tag)
- P_h distribution over hidden tags given previous tag (one per tag) and initial tag
- P_w and P_h have Dirichlet Prior

$$P(w_1, w_2, \dots) = P(P_h)P(P_w)P_h(h_0) \prod_{i \in 1, \dots} P(w_i | h_i)P_h(h_i | h_{i-1})$$

- $L = \{\text{Mary, sees, something, ., John}\}$
- $H = \{1, 2\}$
- how many probability distributions/densities are we thinking about?

$$\text{Reminder: Dirichlet} = P(P_x) = \frac{\prod_{o \in \mathcal{O}} P_x(o)_{\alpha_o}^{\alpha_o}}{B(\alpha)}$$

- $L = \{\text{Mary, sees, something, ., John}\}$
- $H = \{1, 2\}$
- how many probability distributions/densities are we thinking about?
- we need 2 word given state, 2 state given state, 1 initial state + priors for each $\rightarrow 10$
- α for all parameters is 0.5 except:
- $\alpha_{\text{Mary}}^1 = 1, \alpha_{\text{sees}}^2 = 1$ (symmetry breaking)

$$\text{Reminder: Dirichlet} = P(P_x) = \frac{\prod_{o \in \mathcal{O}} P_x(o)_o^\alpha}{B(\alpha)}$$

Assume corpus: $C = \text{"Mary sees something"}$

- What is posterior probability $P(P_h^{exam} | C)$, with:
 - $P_h^{exam}(1|1) = 0.1$
 - $P_h^{exam}(2|1) = 0.9$
 - $P_h^{exam}(1|2) = 0.6$
 - $P_h^{exam}(2|1) = 0.4$
 - $P_h^{exam}(1|s) = 0.8$
 - $P_h^{exam}(2|s) = 0.2$
- What is $P(h_2 = 2)$?
- What is ...

Easy if we know the tags, e.g.:

“Mary:1 sees:2 something:1”

$$\begin{aligned}
 P(P_h^{exam} | T, C) = & \frac{\prod_{i \in \{1,2\}} P_h^{exam}(i|1)^{\alpha_i^1}}{B(\alpha^1)} \\
 & \times \frac{\prod_{i \in \{1,2\}} P_h^{exam}(i|2)^{\alpha_i^2}}{B(\alpha^2)} \\
 & \times \frac{\prod_{i \in \{1,2\}} P_h^{exam}(i|s)^{\alpha_i^s}}{B(\alpha^s)}
 \end{aligned}$$

Where $\alpha_1^s = 1.5$, $\alpha_2^s = 1.5$, $\alpha_1^1 = 1.5$ and other α s for hidden tags are still 0.5

$$\begin{aligned} P(P_h^{exam}|C) &= \sum_T P(P_h^{exam}|T, C)P(T|C) \\ &= \sum_T P(P_h^{exam}|T, C) \frac{P(C|T)P(T)}{P(C)} \end{aligned}$$

$$P(P_h^{exam} | C) = \underbrace{\sum_T}_{\text{annoying sum}} P(P_h^{exam} | T, C) \underbrace{\frac{P(C|T)P(T)}{P(C)}}_{\substack{\text{normalization factor} \\ \text{(secretly annoying sum)}}}$$

We need a generic fix for this problem!

Markov
Chain Monte
Carlo
Methods

Christoph
Teichmann,
Antoine
Venant

Goals

Motivation

Markov
Chains

Invariant
Tricks

When you want to know what people think, you do not ask everyone, you ask a few representatives.

When you have an annoying some over T you do not consider every possible assignment of tags, you only consider a few representative ones.

Let us make our problem more general:

$$P(P_h^{exam} | C) = \sum_T P(P_h^{exam} | T, C) \frac{P(C|T)P(T)}{P(C)}$$

sum over variable $\underbrace{\sum_V}$ function of variable $\underbrace{f(V)}$ probability of variable $\underbrace{P(V)}$

Expected value problem – Many problems in Bayesian Learning/Inference can be formulated as expected value problems

The law of large numbers can be formulated as follows:

- Produce sequence V_1, V_2, \dots
- $P(V_i = v)$ given by $P(V)$ from our expected value
- Then $\lim_{n \rightarrow \infty} \sum_n \frac{1}{n} f(v_i) = \sum_V f(V)P(V)$

We will try to generate a sequence as if each v_i was drawn from

$$P(V)$$

But how do we do this?

Markov
Chain Monte
Carlo
Methods

Christoph
Teichmann,
Antoine
Venant

Goals

Motivation

Markov
Chains

Invariant
Tricks

For simple distributions (categorical) solutions exist based on pseudo-random number generators.

But we have a hard case here!



Markov
Chain Monte
Carlo
Methods

Christoph
Teichmann,
Antoine
Venant

Goals

Motivation

Markov
Chains

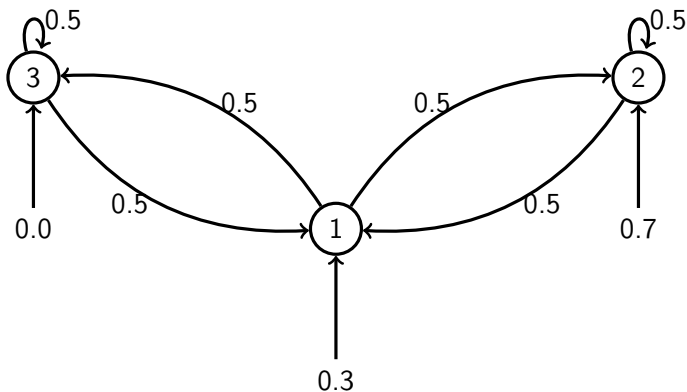
Invariant
Tricks

Markov Chains

- Produce a Markov Chain (later)
- Produce sequence V_1, V_2, \dots from Markov Chain
- Ensure certain conditions
- Then $\lim_{n \rightarrow \infty} \sum_n \frac{1}{n} f(v_i) = \sum_V f(V)P(V)$

- Set of states V – will be variables of interest, can be infinite (but we assume discrete)
- Initial Probability $S(V_0)$
- Transition Probability $T(V_i|v_{i-1})$

How to create v_1, v_2, \dots from this:



What are the Magic Requirements?

Markov
Chain Monte
Carlo
Methods

Christoph
Teichmann,
Antoine
Venant

Goals

Motivation

Markov
Chains

Invariant
Tricks

What do we need so: $\lim_{n \rightarrow \infty} \sum_n \frac{1}{n} f(v_i) = \sum_V f(V)P(V)$?

Markov
Chain Monte
Carlo
Methods

Christoph
Teichmann,
Antoine
Venant

Goals

Motivation

Markov
Chains

Invariant
Tricks

Invariant distribution $I(V)$ of chain = $P(V)$

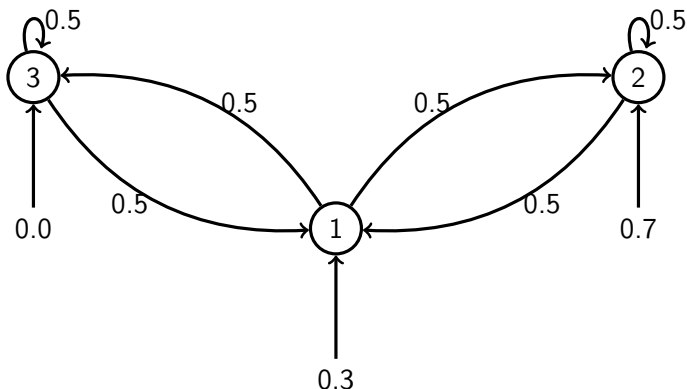
What is invariant distribution?

Invariant distribution $I(V)$ of chain = $P(V)$

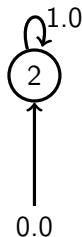
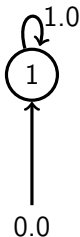
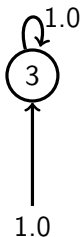
$$I(V) \text{ s.t. for all } v: I(v) = \sum_{v' \in V} T(v|v')I(v')$$

Initial probability does not matter

$$I(1) = I(2) = I(3) = \frac{1}{3} - \text{Show!}$$

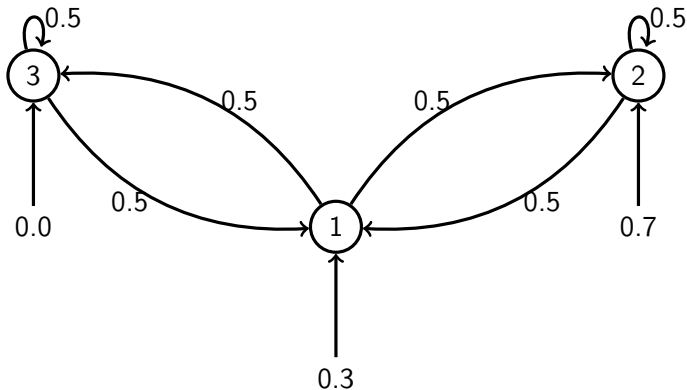


Irreducible: Same $I(v)$, no magic!



Magic Requirement 3: Recurrence

Probability of returning is 1 – always true for finite.



Markov
Chain Monte
Carlo
Methods

Christoph
Teichmann,
Antoine
Venant

Goals

Motivation

Markov
Chains

Invariant
Tricks

- Produce a Markov Chain
- Produce sequence V_1, V_2, \dots from Markov Chain
- If Markov Chain Recurrent, Irreducible, and has Invariant Distribution $I(V) = P(V)$
- Then $\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n f(v_i) = \sum_V f(V)P(V)$
See "Markov Chains" by James Norris
relevant chapters available online
<http://www.statslab.cam.ac.uk/~james/Markov/>



Markov
Chain Monte
Carlo
Methods

Christoph
Teichmann,
Antoine
Venant

Goals

Motivation

Markov
Chains

Invariant
Tricks

Invariant Tricks

Markov
Chain Monte
Carlo
Methods

Christoph
Teichmann,
Antoine
Venant

Goals

Motivation

Markov
Chains

Invariant
Tricks

- Irreducibility and Recurrence – no general trick – often easy
- Correct Invariant Distribution – super hard?

- Proposal Distribution – $p(V|v_i)$
- Must be easy to draw from!
- e.g. for example: flip one tag at random (table)
- Then accept with probability:

$$T(v_i|v_{i-1}) = \min \left(1, \frac{P(v_i)p(v_{i-1}|v_i)}{P(v_{i-1})p(v_i|v_{i-1})} \right)$$

- Otherwise $v_i = v_{i-1}$

- Proposal Distribution – $p(V|v_i)$
- Must be easy to draw from!
- e.g. for example: flip one tag at random (table)
- Then accept with probability:

$$T(v_i|v_{i-1}) = \min \left(1, \frac{P(v_i)p(v_{i-1}|v_i)}{P(v_{i-1})p(v_i|v_{i-1})} \right)$$

- Otherwise $v_i = v_{i-1}$

Show that indeed $I(V) = P(V)$!

- Proposal Distribution – $p(V|v_i)$
- Must be easy to draw from!
- e.g. for example: flip one tag at random (table)
- Then accept with probability:

$$T(v_i|v_{i-1}) = \min \left(1, \frac{P(v_i)p(v_{i-1}|v_i)}{P(v_{i-1})p(v_i|v_{i-1})} \right)$$

- Otherwise $v_i = v_{i-1}$

No problem with

$$P(v) = \frac{f(v)}{\text{super complicated normalizer}}$$



Main Problem: Where to get $p(v_i|v_{i-1})$

Markov
Chain Monte
Carlo
Methods

Christoph
Teichmann,
Antoine
Venant

Goals

Motivation

Markov
Chains

Invariant
Tricks

- Bad $p \rightarrow$ lot of proposals that will be rejected
- Lot of rejection \rightarrow takes forever to converge
- Often need just the right p for a given problem

$$T(v_i|v_{i-1}) = \min \left(1, \frac{P(v_i)p(v_{i-1}|v_i)}{P(v_{i-1})p(v_i|v_{i-1})} \right)$$

- Assume that each $V = \langle t_1, \dots, t_n \rangle$
- E.g., hidden tags in our language model
- Pick a position i at random (or systematically \rightarrow harder to prove)
- We want to change only t_i
- Pick it according to $P(t_i | t_1, \dots, t_{i-1}, t_{i+1}, \dots, t_n)$

Show that this is case of Metropolis-Hastings

- Assume that each $V = \langle t_1, \dots, t_n \rangle$
- E.g., hidden tags in our language model
- Pick a position i at random (or systematically \rightarrow harder to prove)
- We want to change only t_i
- Pick it according to $P(t'_i | t_1, \dots, t_{i-1}, t_{i+1}, \dots, t_n)$

$$P(t'_i | t_1, \dots, t_{i-1}, t_{i+1}, \dots, t_n) = \frac{P(t_1, \dots, t'_i, \dots, t_n)}{\sum_{\bar{t}_i} P(t_1, \dots, \bar{t}_i, \dots, t_n)}$$

- Assume that each $V = \langle t_1, \dots, t_n \rangle$
- E.g., hidden tags in our language model
- Pick a position i at random (or systematically \rightarrow harder to prove)
- We want to change only t_i
- Pick it according to $P(t'_i | t_1, \dots, t_{i-1}, t_{i+1}, \dots, t_n)$

$$P(t'_i | \text{rest}) = \frac{f(t_1, \dots, t'_i, \dots, t_n)}{\text{super complicated normalizer} \sum_{\bar{t}_i} P(t_1, \dots, \bar{t}_i, \dots, t_n)}$$

- Assume that each $V = \langle t_1, \dots, t_n \rangle$
- E.g., hidden tags in our language model
- Pick a position i at random (or systematically \rightarrow harder to prove)
- We want to change only t_i
- Pick it according to $P(t'_i | t_1, \dots, t_{i-1}, t_{i+1}, \dots, t_n)$

$$P(t'_i | \text{rest}) = \frac{f(t_1, \dots, t'_i, \dots, t_n)}{\sum_{\bar{t}_i} f(t_1, \dots, \bar{t}_i, \dots, t_n)}$$

super 2
complicated
normalizer

- For many problems values of t_i very dependent
- Need to really make sure that Irreducible
- Likely to get stuck in certain locations for long time
- Two box example (table)



Markov
Chain Monte
Carlo
Methods

Christoph
Teichmann,
Antoine
Venant

Goals

Motivation

Markov
Chains

Invariant
Tricks

Let us work through Gibbs sampling for our example.