

Algorithmic Applications of Markov Chains.

Till now: Basics of Markov Chains.

- Main thm: \exists and unicity of a stat distr.
- Main thm: When distr a limiting distr?

• Notion of Ergodic chain
↑
irred, aperiodic, pos rec.

- Rate of approach of the stat distr.
mixing time, spectral gap.

within of chains that satisfy detailed balance.

All this has to be kept in Mind as a framework.
for the second part of the course.

↓
Appl & Alg.

Algorithmic Question that we will tackle is about SAMPLING from a distribution.

$$\{ \pi_i \}_{i \in S} \quad S = \text{state space.}$$

How do you sample efficiently?

Today's program.

- 1) Motivate this question. (a little). ✓
- 2) Remind some very classical methods for sampling that work well for "easy" distr. ✓
- 3) Examples of "hard" to sample distributions. ←
- 4) Markov Chain Monte Carlo method.

↓

Metropolis-Hastings algo.

In the coming weeks:

Next
lect { 1) Metropolis-Hasting to minimization of function
or a cost fct.
2) Annealing Algor.

1
weeks { 3) Coloring a graph: analyse MCMC to color a
graph.

1/2
weeks { 4) Ising Model: paradigm of "hard" to
sample distr

2
weeks { 5) last two weeks: Propp & Wilson Coupling
from the past method.
↑
special implementation
of MCMC.

etc.

Motivation for sampling from a distribution $\pi_i, i \in S$.

- For example you want to compute: $\sum_{i \in S} f(i) \pi_i$
 $= E(f(X))$.

X is r.v s.t $P(X=i) = \pi_i$.

Monte Carlo Method: you draw M samples

X_1, X_2, \dots, X_M iid from π .

Take the estimator $\frac{1}{M} \sum_{k=1}^M f(X_k)$.

By the law of large nbs: $\frac{1}{M} \sum_{k=1}^M f(X_k) \xrightarrow{M \rightarrow \infty} E(f(X))$.

Variance $\text{Var}\left\{\frac{1}{M} \sum_{k=1}^M f(X_k)\right\} = \frac{1}{M^2} \sum_{k=1}^M \text{Var}(f(X_k))$
 $= \frac{1}{M} \text{Var}(f(X))$.

$$\frac{1}{M} \sum_{k=1}^M f(X_k) = E(f(X)) + O\left(\frac{1}{\sqrt{M}}\right)$$

$\frac{1}{\sqrt{M}} \sqrt{\text{Var}(f(X))}$ ← error of estimator.

Classical Sampling Methods: (easy to sample distr).

- Most Naive one but very much used (in practice on computers):

Hyp: Efficient way to generate a $U \sim \text{Unif}[0, 1]$.

To sample from $(\pi_i)_{i \in S}$ $S = \{0, 1, 2, \dots\}$.

$$X = \begin{cases} 0 & ; & 0 \leq U \leq \pi_0 & ; & P(U=0) = \pi_0 = P(X=0) \\ 1 & ; & \pi_0 \leq U \leq \pi_0 + \pi_1 & ; & P(U=1) = \pi_1 = P(X=1) \\ \vdots & & & & \\ i & ; & \sum_{k=0}^{i-1} \pi_k \leq U \leq \sum_{k=0}^i \pi_k & ; & P(U=i) = \pi_i \\ & & & & = P(X=i). \end{cases}$$

"U Acts a die"

• Importance Sampling.

Assume we have an efficient method to sample from distr $\{\psi_i\}_{i \in S}$.

Remark:
$$\underbrace{\sum_{i \in S} f(x_i) \pi_i}_{\mathbb{E}_{X \sim \pi}(f(x))} = \underbrace{\sum_{i \in S} f(x_i) w_i \psi_i}_{\mathbb{E}_{X \sim \psi}(f(x) w(x)) \checkmark}$$

where $w_i = \frac{\pi_i}{\psi_i}$

Idea is to consider the estimator:

$$\frac{1}{M} \sum_{k=1}^M f(x_k) w(x_k) \quad \text{where } w_i = \frac{\pi_i}{\psi_i}$$

and the sum is over sample i.i.d $X_k \sim \psi$.

Variance = $\frac{1}{M} \text{Var}(f(x) w(x))$.

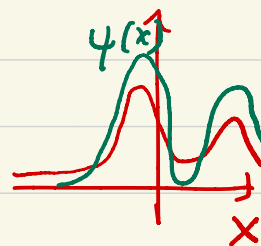
Error is order $\frac{1}{\sqrt{M}} \cdot \sqrt{\text{Var}(f(x) w(x))}$ and

you could optimize over ψ to make it smaller than for previous method.

$\pi(x)$

$\psi(x)$

x



- Rejection sampling \rightarrow proceed as in importance sampling with some acceptance/rejection probs.

$$\underbrace{\sum_{i \in S} f(x_i) \pi_i}_{\mathbb{E}_{\pi}(f(x))} = \frac{\sum_{i \in S} f(x_i) \tilde{w}_i \psi_i}{\underbrace{1/c}_{1/c}} \quad \checkmark$$

$$\tilde{w}_i = \frac{1}{c} \frac{\pi_i}{\psi_i}$$

- easy to sample from distr ψ .
- pick a sample i with probability ψ_i .
- accept the sample i with probability

here $c \geq \max_{i \in S} \left(\frac{\pi_i}{\psi_i} \right)$

$$\tilde{w}_i = \frac{1}{c} \frac{\pi_i}{\psi_i}$$

Estimator: $\frac{1}{M'} \sum_{k: X_k \text{ is accepted}} f(x_k) \quad \checkmark$

of accepted samples \rightarrow

$\parallel X_1, \dots, X_M$ are samples from ψ . You accept in total M' of them.

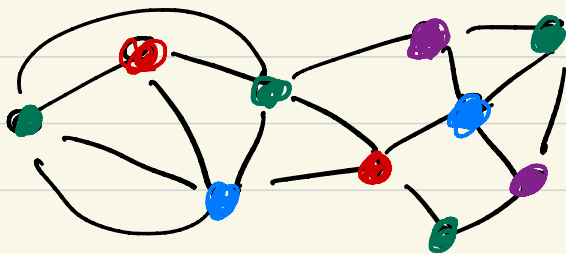
$$\frac{\sum_{i \in S} f(x_i) \psi_i \tilde{w}_i}{\sum_{i \in S} \psi_i \tilde{w}_i} = \frac{\sum_{i \in S} f(x_i) \psi_i \tilde{w}_i}{1/c} = \sum_{i \in S} f(x_i) \pi_i = \mathbb{E}_{\pi}(f(x))$$

$\psi_i \tilde{w}_i = \frac{1}{c} \pi_i \Rightarrow \sum_{i \in S} \psi_i \tilde{w}_i = \frac{1}{c}$

3) Examples of hard to sample distributions.

- Graph Theory or theoretical computer science:

Coloring an arbitrary large graph.



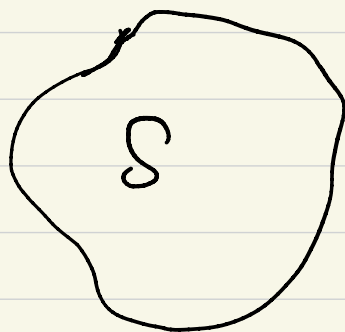
$$G = (V, E).$$

q colors at disposal $\{1, 2, \dots, q\}$.

Def: Proper Coloring of G is an assignment of colors to V s.t. if $(a, b) \in E$ then a & b don't have the same color.

Take the distribution:

Π = unif distr over the set of all possible proper colorings.



$$\Pi_{\text{proper col}} = \frac{\mathbb{1}(\text{proper col})}{\mathbb{Z}}$$

$$\mathbb{Z} = \#(\text{of proper colorings})$$

= number of proper colorings.

• You don't know the set S

• You don't know \mathbb{Z} the normalisation factor is unknown. (specially for large G).

- Ising Model. (we will come back to this in a few weeks).

$G = (V, E)$ v.v assigned to $i \in V : s_i = \pm 1$

Cost function : $H(s_1, \dots, s_{|V|}) = - \sum_{(i,j) \in E} J_{ij} s_i s_j$

\nearrow sum is over all Edges of G .
 \downarrow $J_{ij} \in \mathbb{R}^*$

distr : Ising Model distr (MRF or Gibbs distr for finite G).

$$\mathbb{P}(s_1, \dots, s_{|V|}) = \frac{e^{-\beta H(s_1, \dots, s_{|V|})}}{\mathbb{Z}}$$

state space = set of all binary assignments $(s_1, \dots, s_{|V|})$

$$\mathbb{Z} = \sum_{s_1, \dots, s_{|V|} \in \{\pm 1\}^{|V|}} e^{-\beta H(s_1, \dots, s_{|V|})}$$

sum contains $2^{|V|}$ terms. hard to compute



Markov chain Monte Carlo (MCMC) Sampling Method.

- Goal is to sample π

Idea: view π as the stationary distribution
and even in fact the limiting distribution
of a Markov chain!

Given π , we will construct:

- Markov chain s.t. π is the unique limiting distribution.
- Convergence rate (mixing time or spectral gap) of MC to π should be fast.

For these reasons almost always one considers MC that is ergodic and satisfies detailed balance

But also the construction should be such that we don't need to know too much about computing whatever is hard in π_i 's.

[In the 40's group in Los Alamos found a way to achieve this construction.]

General MCMC algorithm.

($i \in S$ state space)

It constructs of MC that has π as limiting distr:

1. Select an easy-to-simulate Markov Chain with transition probabilities ψ_{ij} for $i, j \in S$. We require that the chain ψ is aperiodic & irreducible. We also require

that $\psi_{ij} > 0 \iff \psi_{ji} > 0$.



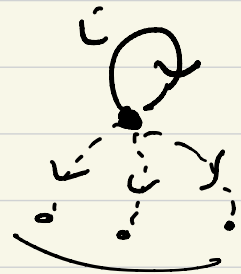
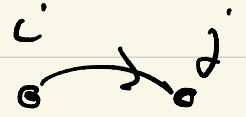
ψ is called the base chain or the proposal chain

2. Design acceptance probabilities: $a_{ij} =$ the prob that the transition $i \rightarrow j$ of the base chain or proposal chain is accepted.

3. Construct the new chain (MCMC) with transition probability matrix P :

$$\begin{cases} P_{ij} = \psi_{ij} a_{ij} & i \neq j \\ P_{ii} = 1 - \sum_{k \neq i} \psi_{ik} a_{ik} \quad \checkmark \end{cases}$$

accepted move.

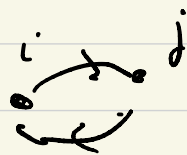


rejected move.

Remark:

$$P_{ii} = \psi_{ii} + \sum_{k \neq i} \psi_{ik} (1 - a_{ik}) = 1 - \sum_{k \neq i} \psi_{ik} a_{ik}.$$

We have not specified a_{ij} yet.



Theorem: [Metropolis - Hastings]

$$\text{Set } a_{ij} = \min\left(1, \frac{\pi_j \psi_{ji}}{\pi_i \psi_{ij}}\right) \quad \textcircled{*}$$

Then the chain with matrix P (above) is ergodic with stat distr π .

Proof: ψ irred and $\psi_{ij} > 0 \Leftrightarrow \psi_{ji} > 0$. Thus

$$a_{ij} > 0 \text{ and also } p_{ij} = \psi_{ij} a_{ij} > 0.$$

\Rightarrow chain P is irred.

Moreover $p_{ii} > 0$ for some $i \Rightarrow P$ aperiodic.

One can check that detailed balance is satisfied (with Metropolis-Hastings)

π is the stat distr. (so it exists).

$\Rightarrow P$ is per recurrent.

Therefore P irred, a per, per rec [ERGODIC]; limiting stat distr unique; and we know it's π \blacksquare

Check that detailed balance is satisfied:

$$\pi_i P_{ij} \stackrel{?}{=} \pi_j P_{ji}$$



$$\pi_i P_{ij} = \pi_i \psi_{ij} \min\left(1, \frac{\pi_j \psi_{ji}}{\pi_i \psi_{ij}}\right)$$

Metro-Hast rule

$$= \min\left(\pi_i \psi_{ij}, \pi_j \psi_{ji}\right)$$

symmetric under $i \leftrightarrow j$

$$= \pi_j P_{ji}$$

□

① Remark: Why is this a nice rule Metropolis-Hastings and why does it help to sample from "hard" disk.

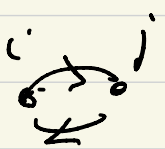
$$\frac{\mathbb{1}(\text{proper coloring})}{Z} ; \frac{e^{-\beta \sum_{i,j \in E} J_{ij} s_i s_j}}{Z}$$

$$\text{in } a_{ij} = \min \left(1, \frac{\pi_j \psi_{ji}}{\pi_i \psi_{ij}} \right)$$

The Z will always simplify in $\left(\frac{\pi_j}{\pi_i} \right)$

we never have here to compute Z .

Also we will see that $\frac{\pi_j}{\pi_i}$ often can be computed very easily. (even by hand) sometimes.



② Remark: Interpretation of the rule. Suppose $\psi_{ij} = \psi_{ji}$

$$\Rightarrow a_{ij} = \min \left(1, \frac{\pi_j}{\pi_i} \right) = \begin{cases} 1 & \text{if } \underline{\pi_j} > \pi_i \\ \frac{\pi_j}{\pi_i} & \text{if } \underline{\pi_j} < \pi_i \end{cases}$$

Metropolis original rule.

state j is more probable than i so if we start at i we should certainly go to j .

But sometimes you should move to lower probability state otherwise you could be stuck.