

## Def. Metropolis-Hastings Algorithm

- start with a candidate transition matrix  $K$
- To ensure unique stationary distribution, it is sufficient to have
  - $K_{xx} > 0$ ,  $\forall x \in \mathcal{X}^n$  [aperiodic]
  - for any  $x, y \in \mathcal{X}^n$ ,  $\exists$  a path  $(x_1=x, x_2, \dots, x_m=y)$  of positive probability transitions [irreducibility]
 
$$K_{x_i x_{i+1}} > 0, \forall i \in \{1, \dots, m-1\}$$

what we want

$$(*) \quad Q_{xy} \cdot P(x) = Q_{yx} \cdot P(y)$$

what we have

$$K_{xy} P(x) \underset{\substack{\uparrow \\ \text{w.l.o.g.}}}{>} K_{yx} P(y)$$

the main trick is to remove some "probability mass" from the larger one.

$$\text{Define: } R_{xy} \triangleq \min \left\{ 1, \frac{P(y) K_{yx}}{P(x) K_{xy}} \right\} \quad \text{for each } x, y$$

$$Q_{xy} = K_{xy} \cdot R_{xy}$$

$$Q_{xx} \underset{\substack{\uparrow \\ \text{s.t.} \\ \text{probability sums to one}}}{=} 1 - \sum_{y \neq x} Q_{xy}$$

with this rejection sampling ( $\cdot R_{xy}$ )

claim:  $(Q, P(x))$  satisfy  $(*)$ .

$$\begin{aligned} \text{proof: } P(x) \cdot Q_{xy} &= P(x) \cdot R_{xy} \cdot K_{xy} \\ &= P(y) K_{yx} && \text{(if } R_{xy} \leq 1) \\ &= P(y) R_{yx} \cdot K_{yx} && \text{ \& } R_{yx} = 1 \end{aligned}$$

\* Remark: as  $P(x)$  is only used in  $R_{xy} = \min\left\{1, \frac{P(y)K_{yx}}{P(x)K_{xy}}\right\}$

we only need  $\frac{P(y)}{P(x)} = \prod_{(i,j) \in E} \frac{f_{ij}(y_i, y_j)}{f_{ij}(x_i, x_j)}$  } takes  $O(|E|)$  computation  
does not require  $\Sigma$ .

\* Remark: Do we need to store  $K$  and  $Q \in \mathbb{R}^{|\mathcal{X}|^n \times |\mathcal{X}|^n}$

- We can choose  $K$  to be simple, such as  $K = \frac{1}{|\mathcal{X}|^n} \mathbb{1}\mathbb{1}^T$

Step 1 - At time  $t$  first generate candidate  $x^{t+1/2}$  from  $K(x^t, x^{t+1/2})$

Step 2 - Accept  $x^{t+1/2}$  with probability  $R_{x^t, x^{t+1/2}}$  :  $x^{t+1} = x^{t+1/2}$

Step 3 - Otherwise reject and keep current state :  $x^{t+1} = x^t$

\* Theorem: Metropolis-Hastings Algorithm finds  $L_1$ -projection of  $K$  onto the space of reversible Markov chains with stationary distribution  $P(x)$ .

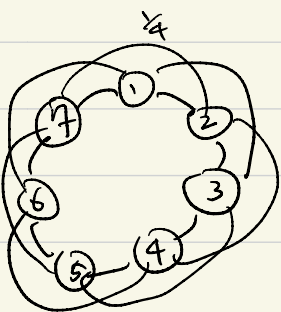
$$Q_{MH} = \arg \min_{\tilde{Q} \text{ s.t. } P^T \tilde{Q} = P^T} \sum_x \sum_{y \neq x} |P(x) \cdot K_{xy} - P(y) \tilde{Q}_{xy}|$$

\* Are we done? - the art is in choosing  $K$ .

if the spread of  $K$  is too large, then acceptance rate is low

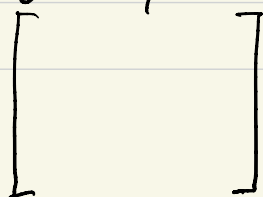
if the spread of  $K$  is too narrow, then mixing time can be large

example  $\triangleright K = \frac{1}{|\mathcal{X}|^n} \mathbb{1}\mathbb{1}^T$ ,  $R_{xy} = \min\left(1, \prod_{(i,j) \in E} \frac{f_{ij}(y_i, y_j)}{f_{ij}(x_i, x_j)}\right)$



all pairs first sampled with equal probability (as per  $K$ )

- but many candidates might be unlikely and be rejected.



\* Def. Gibbs Sampling.

- Q {
- Step 1. sample  $i \in \{1, \dots, n\}$  uniformly at random.
  - Step 2. set  $\tilde{X}_{-i} = X_{-i}^{(t)}$
  - Step 3. sample  $\tilde{X}_i$  from  $P(\tilde{X}_i | X_{-i}^{(t)})$
  - Step 4.  $X^{(t+1)} \leftarrow \tilde{X}$

\* Remark.  $P(\tilde{X}_i | X_{-i}^{(t)}) \propto \prod_{j \in \Theta_i} f_{ij}(\tilde{X}_i, X_j^{(t)})$  is efficient to compute

\* claim.  $(Q, P(x))$  satisfy (\*).

proof for  $\tilde{x}$  that differ at only  $i$ -th coordinate from  $x$ ,

$$P(x) \cdot Q_{x\tilde{x}} = P(x) \cdot \frac{1}{n} \cdot P(\tilde{x}_i | X_{-i})$$

$$\text{Bayes} \rightarrow = P(x_i | X_{-i}) P(X_{-i}) \cdot \frac{1}{n} P(\tilde{x}_i | X_{-i})$$

$$= \underbrace{P(x_i | X_{-i})}_{Q_{\tilde{x}x}} \frac{1}{n} \underbrace{P(X_{-i}) P(\tilde{x}_i | X_{-i})}_{P(\tilde{x})}$$

otherwise  $Q_{x\tilde{x}} = 0$  if  $x$  &  $\tilde{x}$  differ more than one coordinate.

\* the resulting dynamics of the Markov chain is called **Gibbs Dynamics**.

examples > Glauber dynamics for proper k-coloring

Proper Coloring: given a graph  $G$ , and  $k$ -colors  $\{1, 2, \dots, k\}$ .

a coloring is an assignment of colors to nodes

$$X = [X_1, \dots, X_n], \quad X_i \in \{1, \dots, k\}$$

a proper coloring is a coloring where all pairs of nodes connected by an edge have different colors.

Graphical Model: 
$$P(X) = \frac{1}{\sum_k \prod_{(i,j) \in E} \mathbb{I}(X_i \neq X_j)}$$

↑  
# of proper colorings

we propose Gibbs Sampling to solve, for example,  $P(X_1 = X_2), \dots$

Gibbs Sampling: initialize with  $X_0 \in [k]^n$  that is proper

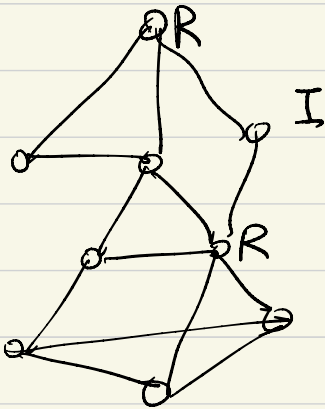
$k=3$  {R, G, B}

Repeat

Sample  $I \in [n]$  uniformly at random

Sample a color  $X_I^{(t+1)}$  uniformly s.t.

$$X_I^{(t+1)} \neq X_{\partial I}^{(t)}$$



Consider Metropolis Hastings with  $K = \frac{1}{|A|^n} \mathbb{1}\mathbb{1}^T$ ,  
each time a random configuration is proposed  
and most likely rejected. It is very slow as  
rejection rate is high.

example  $\succ$  Ising models.

$$P(x) = \frac{1}{Z} \exp \left\{ \beta \sum_{\langle ij \rangle} \theta_{ij} x_i x_j \right\}$$

$\uparrow$   
inverse temperature

$$x_i \in \{-1, 1\}$$

$\beta \uparrow$  clustered.

$\beta \downarrow$  random

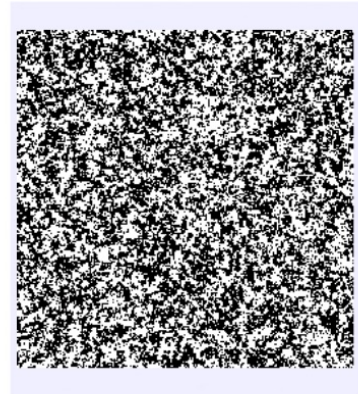
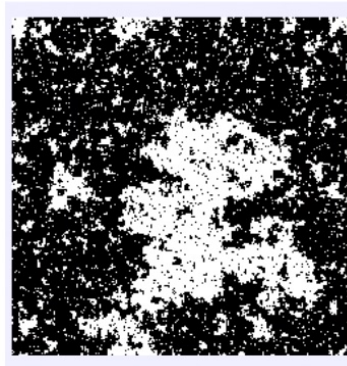
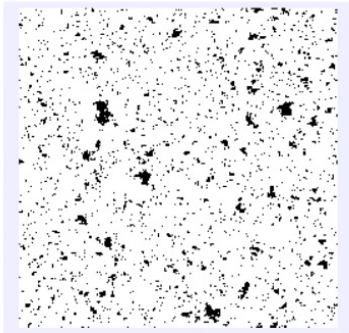


FIGURE 3.2. Glauber dynamics for the Ising model on the  $250 \times 250$  torus viewed at times  $t = 1,000$ ,  $16,500$ , and  $1,000$  at low, critical, and high temperature, respectively. Simulations and graphics courtesy of Raissa D'Souza.

## \* Mixing Time of a Markov Chain.

Def.  $\epsilon$ -mixing time of a Markov Chain  $Q$

is the smallest time  $T_{\text{mix}}(\epsilon)$  such that for all  $t > T_{\text{mix}}(\epsilon)$  and for all initial state  $z^{(0)}$  (a distribution over  $\mathcal{X}^n$ )

$$\left| (z^{(0)})^T \cdot \underbrace{Q \cdot Q \cdots Q}_t - \pi^T \right|_{\text{TV}} \leq \epsilon$$

stationary distribution

where  $|p - z|_{\text{TV}} = \frac{1}{2} \sum_{x \in \mathcal{X}^n} |p(x) - z(x)|$  is the total variation distance.

ex)  $\frac{\begin{array}{c|c} 1-p & p \\ \hline 1-z & z \\ \hline 0 & 1 \end{array} \text{ Bern}(p)}{\begin{array}{c|c} 1-p & p \\ \hline 1-z & z \\ \hline 0 & 1 \end{array} \text{ Bern}(z)} = \sum_{x \in \mathcal{X}^n} \underbrace{[p(x) - z(x)]^+}_{\max\{0, p(x) - z(x)\}}$

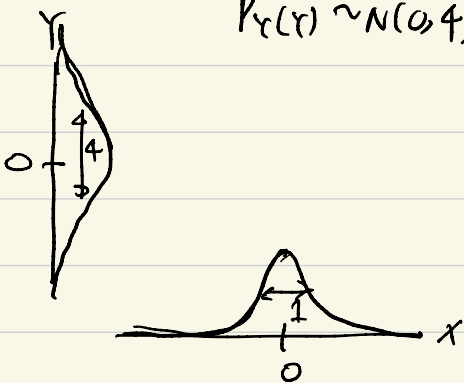
$$|p - z|_{\text{TV}} = |p - z|$$

## \* Bounding mixing time via Coupling

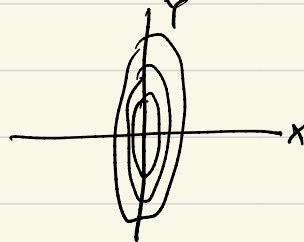
Def. a coupling of two random variables  $X$  and  $Y$  with distributions  $P_X(x)$  and  $P_Y(y)$ , is a construction of a joint probability distribution over  $(X, Y)$ , i.e.,  $P_{X,Y}(x,y)$  such that the marginals are preserved:

$$\begin{cases} \sum_y P_{X,Y}(x,y) = P_X(x) \\ \sum_x P_{X,Y}(x,y) = P_Y(y) \end{cases}$$

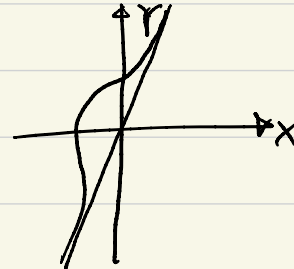
example  $\begin{cases} P_X(x) \sim N(0,1) \\ P_Y(y) \sim N(0,4) \end{cases}$



① independent:  $P_{X,Y} = P_X \cdot P_Y$



②  $Y = 2X$



example  $\succ P_X \sim \text{Bern}(p)$

$P_Y \sim \text{Bern}(q)$

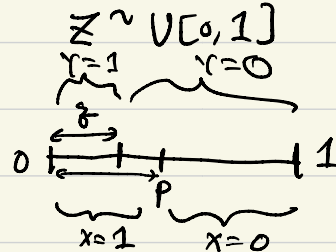
$$P_X \begin{bmatrix} 1-p \\ p \end{bmatrix}$$

$$P_X [1-p, p]$$

① indep.

$$P_{XY} \begin{bmatrix} p\bar{q} & p q \\ \bar{p}\bar{q} & \bar{p} q \end{bmatrix} \begin{matrix} 0 \quad 1 \\ 0 \quad 1 \end{matrix}$$

② construct coupling from  $U[0,1]$



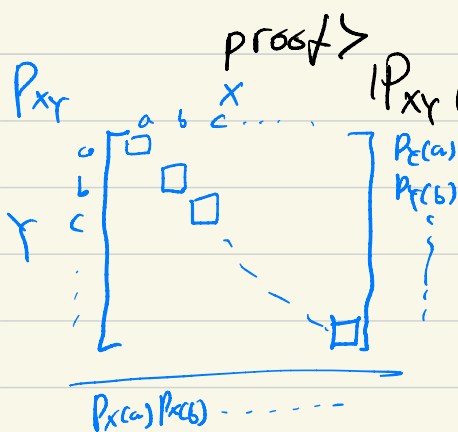
Def. optimal coupling for  $X, Y \in \mathcal{X}$

$$\min_{\text{couplings of } P_X, P_Y} P_{XY}(X \neq Y)$$

\* Coupling Lemma:

for two random variables  $X, Y$  (either continuous or discrete) in the same domain,

$$\|P_X - P_Y\|_{TV} = \min_{\text{couplings of } P_X, P_Y} P_{XY}(X \neq Y)$$



$$P_{XY}(X \neq Y) = 1 - \sum_x P_{XY}(x, x)$$

$$\geq \sum_x \{ P_X(x) - \min \{ P_X(x), P_Y(x) \} \}$$

$$= \sum_x \max \{ 0, P_X(x) - P_Y(x) \}$$

$$= \|P_X - P_Y\|_{TV}$$

\* further, exists  $P_{XY}(x,y)$  s.t. minimum is achieved.

\* Corollary:  $\|P_X - P_Y\|_{TV} \leq \|P_{X,Y}\|$  (for any coupling)

any coupling can be used to upper bound TV distance of two distributions.

ex> optimal coupling of Bernoulli distributions.

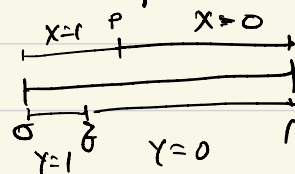
$$P_X \sim \text{Bern}(p) \quad p > q$$

$$P_Y \sim \text{Bern}(q)$$

$$TV = |p - q|$$

$$P_{X,Y} = \begin{bmatrix} 1-p & p-q \\ 0 & q \end{bmatrix}$$

you can sample  $Z \sim U(0,1]$



\* Coupling for bounding  $T_{mix}(\epsilon)$  of Gibbs Sampling.

strategy: let  $X_t, Y_t$  be the random state after  $t$  transitions as per  $Q$ , started with  $X_0, Y_0$ , respectively.

$$\|P_{X_t} - \pi\|_{TV} \leq \max_{P_{X_0}, P_{Y_0}} \|P_{X_t} - P_{Y_t}\|_{TV}$$

$$\text{coupling lemma} \rightarrow \leq \max_{P_{X_0}, P_{Y_0}} P(X_t \neq Y_t)$$

We will analyze a coupling that is

simple enough for analysis, and yet gives a tight upper bound

for any coupling



Proposed coupling of two Gibbs sampling chains. for  $x, y \in \{0, 1\}^n$   
couple two processes closely, while preserving marginal distribution

Step 1. draw  $I \in \{1, \dots, n\}$

Step 2. optimal coupling of  $P(X_I^{(t+1)} | X_{-I}^{(t)})$   $P(Y_I^{(t+1)} | Y_{-I}^{(t)})$