

# CS 263: Counting and Sampling

Nima Anari



slides for

## Mixing via Transport

# Review

## Fundamental theorem

Every ergodic chain has a **unique** stationary dist  $\mu$ , and for any dist  $\nu$

$$\lim_{t \rightarrow \infty} \nu P^t = \mu.$$

# Review

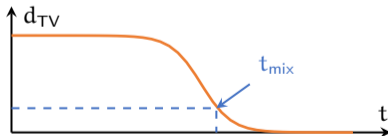
## Fundamental theorem

Every ergodic chain has a **unique** stationary dist  $\mu$ , and for any dist  $\nu$

$$\lim_{t \rightarrow \infty} \nu P^t = \mu.$$

► Mixing time  $t_{\text{mix}}(P, \epsilon, \nu)$ :

$$\min\{t \mid d_{\text{TV}}(\mu, \nu P^t) \leq \epsilon\}$$



# Review

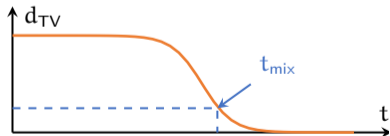
## Fundamental theorem

Every ergodic chain has a **unique** stationary dist  $\mu$ , and for any dist  $\nu$

$$\lim_{t \rightarrow \infty} \nu P^t = \mu.$$

► Mixing time  $t_{\text{mix}}(P, \epsilon, \nu)$ :

$$\min\{t \mid d_{\text{TV}}(\mu, \nu P^t) \leq \epsilon\}$$



►  $t_{\text{mix}}(P, \epsilon) = O\left(t_{\text{mix}}\left(P, \frac{1}{4}\right) \cdot \log\left(\frac{1}{\epsilon}\right)\right)$

# Review

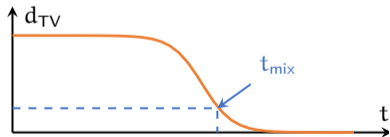
## Fundamental theorem

Every ergodic chain has a **unique** stationary dist  $\mu$ , and for any dist  $\nu$

$$\lim_{t \rightarrow \infty} \nu P^t = \mu.$$

▶ Mixing time  $t_{\text{mix}}(P, \epsilon, \nu)$ :

$$\min\{t \mid d_{\text{TV}}(\mu, \nu P^t) \leq \epsilon\}$$



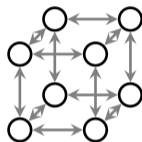
▶  $t_{\text{mix}}(P, \epsilon) = O\left(t_{\text{mix}}\left(P, \frac{1}{4}\right) \cdot \log\left(\frac{1}{\epsilon}\right)\right)$

▶ Strong stationary time:

$$\text{dist}(X_t \mid \tau = k) = \text{stationary}$$

▶  $\tau$ : all coords replaced

▶  $t_{\text{mix}}(\epsilon) \leq n \log(n/\epsilon)$



# Review

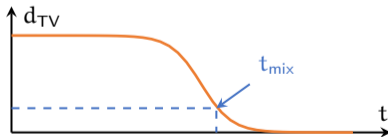
## Fundamental theorem

Every ergodic chain has a **unique** stationary dist  $\mu$ , and for any dist  $\nu$

$$\lim_{t \rightarrow \infty} \nu P^t = \mu.$$

▶ Mixing time  $t_{\text{mix}}(P, \epsilon, \nu)$ :

$$\min\{t \mid d_{\text{TV}}(\mu, \nu P^t) \leq \epsilon\}$$



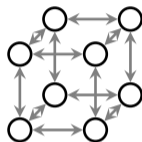
▶  $t_{\text{mix}}(P, \epsilon) = O\left(t_{\text{mix}}\left(P, \frac{1}{4}\right) \cdot \log\left(\frac{1}{\epsilon}\right)\right)$

▶ Strong stationary time:

$$\text{dist}(X_t \mid \tau = k) = \text{stationary}$$

▶  $\tau$ : all coords replaced

▶  $t_{\text{mix}}(\epsilon) \leq n \log(n/\epsilon)$



▶ **Ergodic flow**:  $Q(x, y) = \mu(x)P(x, y)$

# Review

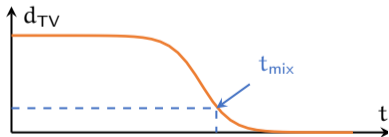
## Fundamental theorem

Every ergodic chain has a **unique** stationary dist  $\mu$ , and for any dist  $\nu$

$$\lim_{t \rightarrow \infty} \nu P^t = \mu.$$

▶ Mixing time  $t_{\text{mix}}(P, \epsilon, \nu)$ :

$$\min\{t \mid d_{\text{TV}}(\mu, \nu P^t) \leq \epsilon\}$$



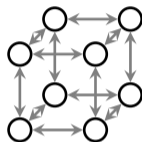
▶  $t_{\text{mix}}(P, \epsilon) = O\left(t_{\text{mix}}\left(P, \frac{1}{4}\right) \cdot \log\left(\frac{1}{\epsilon}\right)\right)$

▶ Strong stationary time:

$$\text{dist}(X_t \mid \tau = k) = \text{stationary}$$

▶  $\tau$ : all coords replaced

▶  $t_{\text{mix}}(\epsilon) \leq n \log(n/\epsilon)$



▶ **Ergodic flow**:  $Q(x, y) = \mu(x)P(x, y)$

▶ Lemma: stationary  $\leftrightarrow$  proper flow

# Review

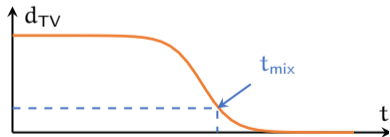
## Fundamental theorem

Every ergodic chain has a **unique** stationary dist  $\mu$ , and for any dist  $\nu$

$$\lim_{t \rightarrow \infty} \nu P^t = \mu.$$

▶ Mixing time  $t_{\text{mix}}(P, \epsilon, \nu)$ :

$$\min\{t \mid d_{\text{TV}}(\mu, \nu P^t) \leq \epsilon\}$$



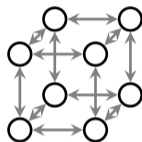
▶  $t_{\text{mix}}(P, \epsilon) = O\left(t_{\text{mix}}\left(P, \frac{1}{4}\right) \cdot \log\left(\frac{1}{\epsilon}\right)\right)$

▶ Strong stationary time:

$$\text{dist}(X_t \mid \tau = k) = \text{stationary}$$

▶  $\tau$ : all coords replaced

▶  $t_{\text{mix}}(\epsilon) \leq n \log(n/\epsilon)$



▶ **Ergodic flow**:  $Q(x, y) = \mu(x)P(x, y)$

▶ Lemma: stationary  $\leftrightarrow$  proper flow

▶ Detailed balance/time-reversible:

$$Q(x, y) = Q(y, x)$$



# Review

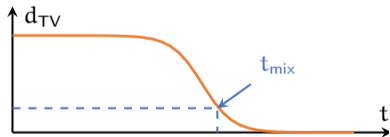
## Fundamental theorem

Every ergodic chain has a **unique** stationary dist  $\mu$ , and for any dist  $\nu$

$$\lim_{t \rightarrow \infty} \nu P^t = \mu.$$

▶ Mixing time  $t_{\text{mix}}(P, \epsilon, \nu)$ :

$$\min\{t \mid d_{\text{TV}}(\mu, \nu P^t) \leq \epsilon\}$$



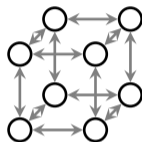
▶  $t_{\text{mix}}(P, \epsilon) = O\left(t_{\text{mix}}\left(P, \frac{1}{4}\right) \cdot \log\left(\frac{1}{\epsilon}\right)\right)$

▶ Strong stationary time:

$$\text{dist}(X_t \mid \tau = k) = \text{stationary}$$

▶  $\tau$ : all coords replaced

▶  $t_{\text{mix}}(\epsilon) \leq n \log(n/\epsilon)$



▶ **Ergodic flow**:  $Q(x, y) = \mu(x)P(x, y)$

▶ Lemma: stationary  $\leftrightarrow$  proper flow

▶ Detailed balance/time-reversible:

$$Q(x, y) = Q(y, x)$$

▶ Metropolis filter:  $P(x, y) \mapsto$

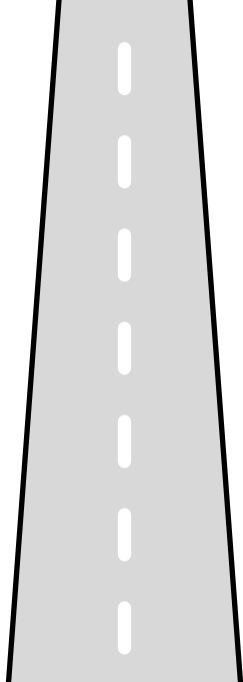
$$P(x, y) \min\left\{1, \frac{\mu(y)P(y, x)}{\mu(x)P(x, y)}\right\}$$

# Designing Markov Chains

- ▶ Markov kernels
- ▶ Combination with time-reversal

# Mixing via Transport

- ▶ Wasserstein distance
- ▶ Path coupling

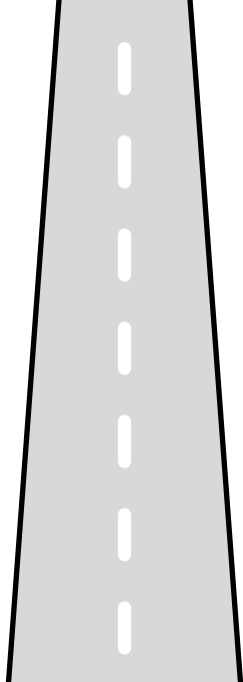


## Designing Markov Chains

- ▶ Markov kernels
- ▶ Combination with time-reversal

## Mixing via Transport

- ▶ Wasserstein distance
- ▶ Path coupling



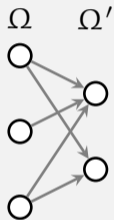
# Markov kernel

We can generalize time-reversal to

## Markov kernel

$$P \in \mathbb{R}_{\geq 0}^{\Omega \times \Omega'}$$

$$\sum_y P(x, y) = 1$$

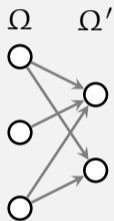


# Markov kernel

We can generalize time-reversal to

## Markov kernel

$$P \in \mathbb{R}_{\geq 0}^{\Omega \times \Omega'}$$
$$\sum_y P(x, y) = 1$$



- ▶ Markov kernels are conditional dists. Combined with dist  $\mu$  on  $\Omega$ , they give joint dist/ergodic flow:

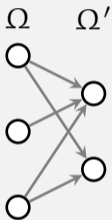
$$Q(x, y) = \mu(x)P(x, y)$$

# Markov kernel

We can generalize time-reversal to

## Markov kernel

$$P \in \mathbb{R}_{\geq 0}^{\Omega \times \Omega'}$$
$$\sum_y P(x, y) = 1$$



▶ Time-reversal:

$$Q^\circ(y, x) = Q(x, y)$$

▶ Markov kernels are conditional dists. Combined with dist  $\mu$  on  $\Omega$ , they give joint dist/ergodic flow:

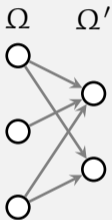
$$Q(x, y) = \mu(x)P(x, y)$$

# Markov kernel

We can generalize time-reversal to

## Markov kernel

$$P \in \mathbb{R}_{\geq 0}^{\Omega \times \Omega'}$$
$$\sum_y P(x, y) = 1$$



▶ Time-reversal:

$$Q^\circ(y, x) = Q(x, y)$$

▶ Dist on  $\Omega'$ :  $\mu^\circ = \mu P$  is marginal of  $y$  in  $Q^\circ$  or  $Q$ .

▶ Markov kernels are conditional dists. Combined with dist  $\mu$  on  $\Omega$ , they give joint dist/ergodic flow:

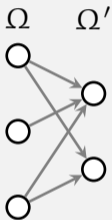
$$Q(x, y) = \mu(x)P(x, y)$$

# Markov kernel

We can generalize time-reversal to

## Markov kernel

$$P \in \mathbb{R}_{\geq 0}^{\Omega \times \Omega'}$$
$$\sum_{\mathbf{y}} P(\mathbf{x}, \mathbf{y}) = 1$$



▶ Time-reversal:

$$Q^\circ(\mathbf{y}, \mathbf{x}) = Q(\mathbf{x}, \mathbf{y})$$

▶ Dist on  $\Omega'$ :  $\mu^\circ = \mu P$  is marginal of  $\mathbf{y}$  in  $Q^\circ$  or  $Q$ .

▶ The time-reversal Markov kernel is the conditional dist of  $\mathbf{x}$  given  $\mathbf{y}$ :

$$P^\circ(\mathbf{y}, \mathbf{x}) = \frac{\mu(\mathbf{x})P(\mathbf{x}, \mathbf{y})}{\mu^\circ(\mathbf{y})}$$

▶ Markov kernels are conditional dists. Combined with dist  $\mu$  on  $\Omega$ , they give joint dist/ergodic flow:

$$Q(\mathbf{x}, \mathbf{y}) = \mu(\mathbf{x})P(\mathbf{x}, \mathbf{y})$$

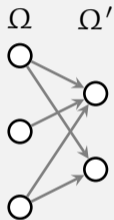


# Markov kernel

We can generalize time-reversal to

## Markov kernel

$$P \in \mathbb{R}_{\geq 0}^{\Omega \times \Omega'}$$
$$\sum_{\mathbf{y}} P(\mathbf{x}, \mathbf{y}) = 1$$



- ▶ Markov kernels are conditional dists. Combined with dist  $\mu$  on  $\Omega$ , they give joint dist/ergodic flow:

$$Q(\mathbf{x}, \mathbf{y}) = \mu(\mathbf{x})P(\mathbf{x}, \mathbf{y})$$

- ▶ Time-reversal:

$$Q^\circ(\mathbf{y}, \mathbf{x}) = Q(\mathbf{x}, \mathbf{y})$$

- ▶ Dist on  $\Omega'$ :  $\mu^\circ = \mu P$  is marginal of  $\mathbf{y}$  in  $Q^\circ$  or  $Q$ .

- ▶ The time-reversal Markov kernel is the conditional dist of  $\mathbf{x}$  given  $\mathbf{y}$ :

$$P^\circ(\mathbf{y}, \mathbf{x}) = \frac{\mu(\mathbf{x})P(\mathbf{x}, \mathbf{y})}{\mu^\circ(\mathbf{y})}$$

- ▶ Note the detailed balance equation:

$$\underbrace{\mu(\mathbf{x})P(\mathbf{x}, \mathbf{y})}_{Q(\mathbf{x}, \mathbf{y})} = \underbrace{\mu^\circ(\mathbf{y})P^\circ(\mathbf{y}, \mathbf{x})}_{Q^\circ(\mathbf{y}, \mathbf{x})}$$

# Combination with time-reversal

Design recipe:

- 1 Target dist  $\mu$  on  $\Omega$

# Combination with time-reversal

Design recipe:

- ① Target dist  $\mu$  on  $\Omega$
- ② Markov kernel  $N$  from  $\Omega$  to  $\Omega'$

# Combination with time-reversal

Design recipe:

- ① Target dist  $\mu$  on  $\Omega$
- ② Markov kernel  $N$  from  $\Omega$  to  $\Omega'$
- ③ Let  $P = NN^\circ$

# Combination with time-reversal

Design recipe:

- ① Target dist  $\mu$  on  $\Omega$
- ② Markov kernel  $N$  from  $\Omega$  to  $\Omega'$
- ③ Let  $P = NN^\circ$

## Lemma

$\mu P = \mu$  and  $P$  is time-reversible

# Combination with time-reversal

Design recipe:

- 1 Target dist  $\mu$  on  $\Omega$
- 2 Markov kernel  $N$  from  $\Omega$  to  $\Omega'$
- 3 Let  $P = NN^\circ$

## Lemma

$\mu P = \mu$  and  $P$  is time-reversible

Proof: we have  $\mu(x)P(x, z) =$

$$\sum_y \mu(x)N(x, y)N^\circ(y, z) =$$

$$\sum_y \frac{\mu(x)N(x, y)\mu(z)N(z, y)}{\mu^\circ(y)}$$

↑  
symmetric in  $x, z$

# Combination with time-reversal

Design recipe:

- 1 Target dist  $\mu$  on  $\Omega$
- 2 Markov kernel  $N$  from  $\Omega$  to  $\Omega'$
- 3 Let  $P = NN^\circ$

## Lemma

$\mu P = \mu$  and  $P$  is time-reversible

Proof: we have  $\mu(x)P(x, z) =$

$$\sum_y \mu(x)N(x, y)N^\circ(y, z) =$$

$$\sum_y \frac{\mu(x)N(x, y)\mu(z)N(z, y)}{\mu^\circ(y)}$$

↑  
symmetric in  $x, z$

## Example: Glauber dynamics



# Combination with time-reversal

Design recipe:

- 1 Target dist  $\mu$  on  $\Omega$
- 2 Markov kernel  $N$  from  $\Omega$  to  $\Omega'$
- 3 Let  $P = NN^\circ$

## Lemma

$\mu P = \mu$  and  $P$  is time-reversible

Proof: we have  $\mu(x)P(x, z) =$

$$\sum_y \mu(x)N(x, y)N^\circ(y, z) =$$

$$\sum_y \frac{\mu(x)N(x, y)\mu(z)N(z, y)}{\mu^\circ(y)}$$

↑  
symmetric in  $x, z$

## Example: Glauber dynamics



▷  $N$ : erase u.r. vertex



# Combination with time-reversal

Design recipe:

- 1 Target dist  $\mu$  on  $\Omega$
- 2 Markov kernel  $N$  from  $\Omega$  to  $\Omega'$
- 3 Let  $P = NN^\circ$

## Lemma

$\mu P = \mu$  and  $P$  is time-reversible

Proof: we have  $\mu(x)P(x, z) =$

$$\sum_y \mu(x)N(x, y)N^\circ(y, z) =$$

$$\sum_y \frac{\mu(x)N(x, y)\mu(z)N(z, y)}{\mu^\circ(y)}$$

↑  
symmetric in  $x, z$

## Example: Glauber dynamics



▷  $N$ : erase u.r. vertex

▷  $N^\circ$ : recolor with prob  $\alpha$

$\mu(\text{result})N(\text{result}, \text{partial})$

↑  
cancels out

# Combination with time-reversal

Design recipe:

- 1 Target dist  $\mu$  on  $\Omega$
- 2 Markov kernel  $N$  from  $\Omega$  to  $\Omega'$
- 3 Let  $P = NN^\circ$

## Lemma

$\mu P = \mu$  and  $P$  is time-reversible

Proof: we have  $\mu(x)P(x, z) =$

$$\sum_y \mu(x)N(x, y)N^\circ(y, z) =$$

$$\sum_y \frac{\mu(x)N(x, y)\mu(z)N(z, y)}{\mu^\circ(y)}$$

↑  
symmetric in  $x, z$

## Example: Glauber dynamics



▷  $N$ : erase u.r. vertex

▷  $N^\circ$ : recolor with prob  $\alpha$

$$\mu(\text{result})N(\text{result}, \text{partial})$$

↑  
cancels out

▷  $P$ : pick u.r. valid color for u.r. vert

# Combination with time-reversal

Design recipe:

- 1 Target dist  $\mu$  on  $\Omega$
- 2 Markov kernel  $N$  from  $\Omega$  to  $\Omega'$
- 3 Let  $P = NN^\circ$

## Lemma

$\mu P = \mu$  and  $P$  is time-reversible

Proof: we have  $\mu(x)P(x,z) =$

$$\sum_y \mu(x)N(x,y)N^\circ(y,z) =$$

$$\sum_y \frac{\mu(x)N(x,y)\mu(z)N(z,y)}{\mu^\circ(y)}$$

↑  
symmetric in  $x,z$

## Example: Glauber dynamics



▷  $N$ : erase u.r. vertex

▷  $N^\circ$ : recolor with prob  $\alpha$

$$\mu(\text{result})N(\text{result}, \text{partial})$$

↑  
cancels out

▷  $P$ : pick u.r. valid color for u.r. vert

▷ Note: different from Metropolis.

## Example: block dynamics



▶ N: erase  $k$  u.r. verts

▶ P: recolor  $k$  u.r. verts

u.a.r. from valid colorings

## Example: block dynamics



- ▶ N: erase  $k$  u.r. verts
- ▶ P: recolor  $k$  u.r. verts

u.a.r. from valid colorings

## Example: spanning trees (I)



- ▶ N: drop one edge u.a.r.
- ▶ P: then add edge u.a.r. from cut

### Example: block dynamics



- ▶ N: erase  $k$  u.r. verts
- ▶ P: recolor  $k$  u.r. verts

u.a.r. from valid colorings

### Example: spanning trees (II)



- ▶ N: add one edge u.a.r.
- ▶ P: drop edge u.a.r. from cycle

### Example: spanning trees (I)



- ▶ N: drop one edge u.a.r.
- ▶ P: then add edge u.a.r. from cut

### Example: block dynamics



- ▶ N: erase  $k$  u.r. verts
- ▶ P: recolor  $k$  u.r. verts

u.a.r. from valid colorings

### Example: spanning trees (I)



- ▶ N: drop one edge u.a.r.
- ▶ P: then add edge u.a.r. from cut

### Example: spanning trees (II)



- ▶ N: add one edge u.a.r.
- ▶ P: drop edge u.a.r. from cycle

- ▶ Trivial example: let  $\Omega' = \{\emptyset\}$  and N map everything to  $\emptyset$ .

### Example: block dynamics



- ▶ N: erase  $k$  u.r. verts
- ▶ P: recolor  $k$  u.r. verts

u.a.r. from valid colorings

### Example: spanning trees (I)



- ▶ N: drop one edge u.a.r.
- ▶ P: then add edge u.a.r. from cut

### Example: spanning trees (II)



- ▶ N: add one edge u.a.r.
- ▶ P: drop edge u.a.r. from cycle

- ▶ Trivial example: let  $\Omega' = \{\emptyset\}$  and N map everything to  $\emptyset$ .
- ▶ We get **ideal** Markov chain:

mixes in one step

$$P(x, y) = \mu(y)$$



### Example: block dynamics



- ▶ N: erase  $k$  u.r. verts
- ▶ P: recolor  $k$  u.r. verts

u.a.r. from valid colorings

### Example: spanning trees (I)



- ▶ N: drop one edge u.a.r.
- ▶ P: then add edge u.a.r. from cut

### Example: spanning trees (II)



- ▶ N: add one edge u.a.r.
- ▶ P: drop edge u.a.r. from cycle

▶ Trivial example: let  $\Omega' = \{\emptyset\}$  and N map everything to  $\emptyset$ .

▶ We get **ideal** Markov chain:

↑  
mixes in one step

$$P(x, y) = \mu(y)$$

▶ Algorithmically useless!

Algorithmic implementation:

```
for  $t = 0, 1, \dots$  do  
  sample  $y_t \sim N(x_t, \cdot)$   
  for  $z$  with  $N(z, y_t) > 0$  do  
     $p_z \leftarrow \mu(z)N(z, y_t)$   
  sample  $z$  with prob  $\propto p_z$   
   $x_{t+1} \leftarrow$  sample
```

Algorithmic implementation:

```
for  $t = 0, 1, \dots$  do  
  sample  $y_t \sim N(x_t, \cdot)$   
  for  $z$  with  $N(z, y_t) > 0$  do  
     $p_z \leftarrow \mu(z)N(z, y_t)$   
  sample  $z$  with prob  $\propto p_z$   
   $x_{t+1} \leftarrow$  sample
```

► Want sparse columns for  $N$ .

Algorithmic implementation:

```
for  $t = 0, 1, \dots$  do
  sample  $\mathbf{y}_t \sim \mathcal{N}(\mathbf{x}_t, \cdot)$ 
  for  $z$  with  $\mathbf{N}(z, \mathbf{y}_t) > 0$  do
     $p_z \leftarrow \mu(z)\mathbf{N}(z, \mathbf{y}_t)$ 
  sample  $z$  with prob  $\propto p_z$ 
   $\mathbf{x}_{t+1} \leftarrow$  sample
```

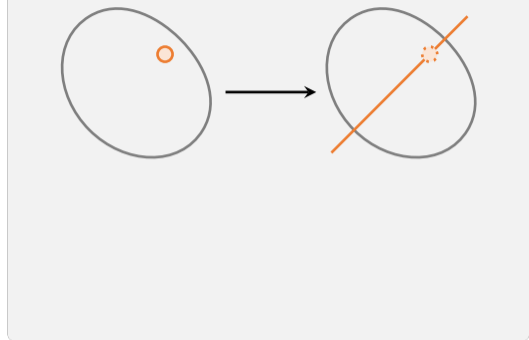
- ▶ Want sparse columns for  $\mathbf{N}$ .
- ▶ Ideally we can simulate  $\mathbf{N}$ , and its columns are not just **sparse** but **efficiently enumerable**.

Algorithmic implementation:

```

for  $t = 0, 1, \dots$  do
  sample  $\mathbf{y}_t \sim \mathcal{N}(\mathbf{x}_t, \cdot)$ 
  for  $z$  with  $\mathcal{N}(z, \mathbf{y}_t) > 0$  do
     $\mathbf{p}_z \leftarrow \mu(z)\mathcal{N}(z, \mathbf{y}_t)$ 
  sample  $z$  with prob  $\propto \mathbf{p}_z$ 
   $\mathbf{x}_{t+1} \leftarrow$  sample
  
```

- ▶ Want sparse columns for  $\mathbf{N}$ .
- ▶ Ideally we can simulate  $\mathbf{N}$ , and its columns are not just **sparse** but **efficiently enumerable**.

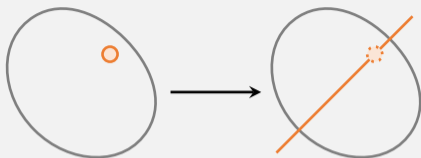


Algorithmic implementation:

```

for  $t = 0, 1, \dots$  do
  sample  $\mathbf{y}_t \sim \mathcal{N}(\mathbf{x}_t, \cdot)$ 
  for  $z$  with  $\mathcal{N}(z, \mathbf{y}_t) > 0$  do
     $\mathbf{p}_z \leftarrow \mu(z)\mathcal{N}(z, \mathbf{y}_t)$ 
  sample  $z$  with prob  $\propto \mathbf{p}_z$ 
   $\mathbf{x}_{t+1} \leftarrow$  sample
    
```

- ▶ Want sparse columns for  $\mathbf{N}$ .
- ▶ Ideally we can simulate  $\mathbf{N}$ , and its columns are not just **sparse** but **efficiently enumerable**.



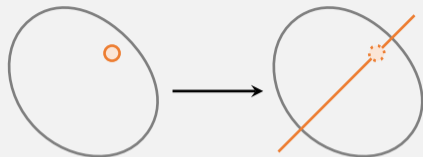
▶  $\mu$  is uniform on subset  $S$  of  $\mathbb{R}^d$

Algorithmic implementation:

```

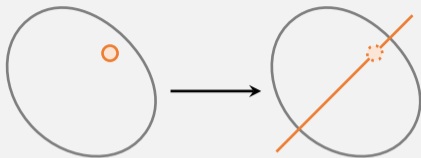
for  $t = 0, 1, \dots$  do
  sample  $y_t \sim N(x_t, \cdot)$ 
  for  $z$  with  $N(z, y_t) > 0$  do
     $p_z \leftarrow \mu(z)N(z, y_t)$ 
  sample  $z$  with prob  $\propto p_z$ 
   $x_{t+1} \leftarrow$  sample
    
```

- ▶ Want sparse columns for  $N$ .
- ▶ Ideally we can simulate  $N$ , and its columns are not just **sparse** but **efficiently enumerable**.



- ▶  $\mu$  is uniform on subset  $S$  of  $\mathbb{R}^d$
- ▶  $N: x \mapsto$  u.r. line  $\ell$  through  $x$

## Example: hit-and-run $\leftarrow$ infinite space



- ▶  $\mu$  is uniform on subset  $S$  of  $\mathbb{R}^d$
- ▶  $N: x \mapsto$  u.r. line  $\ell$  through  $x$
- ▶  $P$ : then choose u.a.r. from  $\ell \cap S$

Algorithmic implementation:

```
for  $t = 0, 1, \dots$  do
  sample  $y_t \sim N(x_t, \cdot)$ 
  for  $z$  with  $N(z, y_t) > 0$  do
     $p_z \leftarrow \mu(z)N(z, y_t)$ 
  sample  $z$  with prob  $\propto p_z$ 
   $x_{t+1} \leftarrow$  sample
```

- ▶ Want sparse columns for  $N$ .
- ▶ Ideally we can simulate  $N$ , and its columns are not just **sparse** but **efficiently enumerable**.

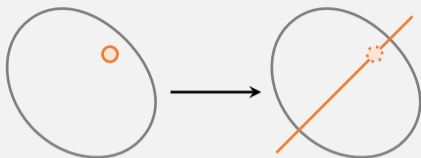


Algorithmic implementation:

```
for  $t = 0, 1, \dots$  do
  sample  $y_t \sim \mathcal{N}(x_t, \cdot)$ 
  for  $z$  with  $N(z, y_t) > 0$  do
     $p_z \leftarrow \mu(z)N(z, y_t)$ 
  sample  $z$  with prob  $\propto p_z$ 
 $x_{t+1} \leftarrow$  sample
```

- ▶ Want sparse columns for  $N$ .
- ▶ Ideally we can simulate  $N$ , and its columns are not just **sparse** but **efficiently enumerable**.

### Example: hit-and-run $\leftarrow$ infinite space



- ▶  $\mu$  is uniform on subset  $S$  of  $\mathbb{R}^d$
- ▶  $N: x \mapsto$  u.r. line  $\ell$  through  $x$
- ▶  $P$ : then choose u.a.r. from  $\ell \cap S$

### Example: restricted Gaussian

- ▶  $\mu$ : dist on  $\mathbb{R}^d$
- ▶  $N: x \mapsto y = x + g$  for  $g \sim \mathcal{N}(0, cI)$
- ▶  $P$ : then sample  $z$  w.p.  $\propto$

restricted Gaussian  $\rightarrow \mu(z)e^{-\|z-y\|^2/2c}$

# Summary

Design  $P$  time-reversible w.r.t.  $\mu$ :

$$\mu(x)P(x, y) = \mu(y)P(y, x)$$

# Summary

Design  $P$  time-reversible w.r.t.  $\mu$ :

$$\mu(x)P(x, y) = \mu(y)P(y, x)$$

## ① Metropolis filter

- ▶ Have some initial  $P$
- ▶ Modify it to

$$P(x, y) \min \left\{ 1, \frac{\mu(y)P(y, x)}{\mu(x)P(x, y)} \right\}$$

# Summary

Design  $P$  time-reversible w.r.t.  $\mu$ :

$$\mu(x)P(x, y) = \mu(y)P(y, x)$$

## 1 Metropolis filter

- ▶ Have some initial  $P$
- ▶ Modify it to

$$P(x, y) \min \left\{ 1, \frac{\mu(y)P(y, x)}{\mu(x)P(x, y)} \right\}$$

## 2 Combination with time-reversal

- ▶ Have some Markov kernel  $N$
- ▶ Form  $NN^\circ$

# Summary

Design  $P$  time-reversible w.r.t.  $\mu$ :

$$\mu(x)P(x, y) = \mu(y)P(y, x)$$

## ① Metropolis filter

- ▶ Have some initial  $P$
- ▶ Modify it to

$$P(x, y) \min \left\{ 1, \frac{\mu(y)P(y, x)}{\mu(x)P(x, y)} \right\}$$

## ② Combination with time-reversal

- ▶ Have some Markov kernel  $N$
- ▶ Form  $NN^\circ$

Question: do these guarantee irreducible/aperiodic?

## Designing Markov Chains

- ▶ Markov kernels
- ▶ Combination with time-reversal

## Mixing via Transport

- ▶ Wasserstein distance
- ▶ Path coupling

# Designing Markov Chains

- ▶ Markov kernels
- ▶ Combination with time-reversal

## Mixing via Transport

- ▶ Wasserstein distance
- ▶ Path coupling

# Transport distance

- ▶ Prevalent strategy for analyzing mixing time: **contraction**



# Transport distance

- ▶ Prevalent strategy for analyzing mixing time: **contraction**
- ▶  $d_{TV}$  is too crude; doesn't contract every step



$$d_{TV}(\nu P, \nu' P) = d_{TV}(\nu, \nu')$$

# Transport distance

- ▶ Prevalent strategy for analyzing mixing time: **contraction**
- ▶  $d_{TV}$  is too crude; doesn't contract every step



$$d_{TV}(vP, v'P) = d_{TV}(v, v')$$

- ▶ **Fix:** use a proxy for  $d_{TV}$ 
  - ▶ Transport/Wasserstein/earth-mover distance ← **today**
  - ▶ f-divergences, variance, entropy
    - ↑  
functional analysis, later

# Transport distance

- ▶ Prevalent strategy for analyzing mixing time: **contraction**
- ▶  $d_{TV}$  is too crude; doesn't contract every step



$$d_{TV}(vP, v'P) = d_{TV}(v, v')$$

- ▶ **Fix:** use a proxy for  $d_{TV}$ 
  - ▶ Transport/Wasserstein/earth-mover distance ← today
  - ▶ f-divergences, variance, entropy
    - ↑  
functional analysis, later
- ▶ Suppose  $\Omega$  is equipped with metric  $d : \Omega \times \Omega \rightarrow \mathbb{R}_{\geq 0}$ .

# Transport distance

- ▶ Prevalent strategy for analyzing mixing time: **contraction**
- ▶  $d_{TV}$  is too crude; doesn't contract every step



$$d_{TV}(\nu P, \nu' P) = d_{TV}(\nu, \nu')$$

- ▶ **Fix:** use a proxy for  $d_{TV}$ 
  - ▶ Transport/Wasserstein/earth-mover distance ← **today**
  - ▶ f-divergences, variance, entropy
    - ↑  
**functional analysis, later**
- ▶ Suppose  $\Omega$  is equipped with metric  $d : \Omega \times \Omega \rightarrow \mathbb{R}_{\geq 0}$ .

## Wasserstein distance

We define the Wasserstein distance w.r.t.  $d$  as  $\mathcal{W}(\mu, \nu) =$

$$\min \{ \mathbb{E}_{(X,Y) \sim \pi} [d(X,Y)] \mid \pi \text{ coupling} \}$$

# Transport distance

- ▶ Prevalent strategy for analyzing mixing time: **contraction**
- ▶  $d_{TV}$  is too crude; doesn't contract every step



$$d_{TV}(\nu P, \nu' P) = d_{TV}(\nu, \nu')$$

- ▶ **Fix:** use a proxy for  $d_{TV}$ 
  - ▶ Transport/Wasserstein/earth-mover distance ← **today**
  - ▶ f-divergences, variance, entropy
    - ↑  
functional analysis, later
- ▶ Suppose  $\Omega$  is equipped with metric  $d : \Omega \times \Omega \rightarrow \mathbb{R}_{\geq 0}$ .

## Wasserstein distance

We define the Wasserstein distance w.r.t.  $d$  as  $\mathcal{W}(\mu, \nu) =$

$$\min \{ \mathbb{E}_{(X, Y) \sim \pi} [d(X, Y)] \mid \pi \text{ coupling} \}$$

## Example: total variation

If we use  $d(x, y) = \mathbb{1}[x \neq y]$ :  $\mathcal{W} = d_{TV}$

# Transport distance

- ▶ Prevalent strategy for analyzing mixing time: **contraction**
- ▶  $d_{TV}$  is too crude; doesn't contract every step



$$d_{TV}(\nu P, \nu' P) = d_{TV}(\nu, \nu')$$

- ▶ **Fix:** use a proxy for  $d_{TV}$ 
  - ▶ Transport/Wasserstein/earth-mover distance ← **today**
  - ▶ f-divergences, variance, entropy
    - ↑  
functional analysis, later
- ▶ Suppose  $\Omega$  is equipped with metric  $d : \Omega \times \Omega \rightarrow \mathbb{R}_{\geq 0}$ .

## Wasserstein distance

We define the Wasserstein distance w.r.t.  $d$  as  $\mathcal{W}(\mu, \nu) =$

$$\min \{ \mathbb{E}_{(X, Y) \sim \pi} [d(X, Y)] \mid \pi \text{ coupling} \}$$

## Example: total variation

If we use  $d(x, y) = \mathbb{1}[x \neq y]$ :  $\mathcal{W} = d_{TV}$

## Example: Hamming

$$\Omega = [q]^n \quad d(x, y) = |\{i \mid x_i \neq y_i\}|$$

# Transport distance

- ▶ Prevalent strategy for analyzing mixing time: **contraction**
- ▶  $d_{TV}$  is too crude; doesn't contract every step



$$d_{TV}(\nu P, \nu' P) = d_{TV}(\nu, \nu')$$

- ▶ **Fix:** use a proxy for  $d_{TV}$ 
  - ▶ Transport/Wasserstein/earth-mover distance ← **today**
  - ▶ f-divergences, variance, entropy
    - ↑  
**functional analysis, later**
- ▶ Suppose  $\Omega$  is equipped with metric  $d : \Omega \times \Omega \rightarrow \mathbb{R}_{\geq 0}$ .

## Wasserstein distance

We define the Wasserstein distance w.r.t.  $d$  as  $\mathcal{W}(\mu, \nu) =$

$$\min \{ \mathbb{E}_{(X,Y) \sim \pi} [d(X,Y)] \mid \pi \text{ coupling} \}$$

## Example: total variation

If we use  $d(x,y) = \mathbb{1}[x \neq y]$ :  $\mathcal{W} = d_{TV}$

## Example: Hamming

$$\Omega = [q]^n \quad d(x,y) = |\{i \mid x_i \neq y_i\}|$$

$$\mu = \text{unif on } \{(\bullet, \bullet, \bullet), (\bullet, \bullet, \bullet)\}$$

$$\nu = \text{unif on } \{(\bullet, \bullet, \bullet), (\bullet, \bullet, \bullet), (\bullet, \bullet, \bullet)\}$$

# Transport distance

- ▶ Prevalent strategy for analyzing mixing time: **contraction**
- ▶  $d_{TV}$  is too crude; doesn't contract every step



$$d_{TV}(\nu P, \nu' P) = d_{TV}(\nu, \nu')$$

- ▶ **Fix:** use a proxy for  $d_{TV}$ 
  - ▶ Transport/Wasserstein/earth-mover distance ← **today**
  - ▶ f-divergences, variance, entropy
    - ↑  
**functional analysis, later**
- ▶ Suppose  $\Omega$  is equipped with metric  $d : \Omega \times \Omega \rightarrow \mathbb{R}_{\geq 0}$ .

## Wasserstein distance

We define the Wasserstein distance w.r.t.  $d$  as  $\mathcal{W}(\mu, \nu) =$

$$\min \{ \mathbb{E}_{(X,Y) \sim \pi} [d(X,Y)] \mid \pi \text{ coupling} \}$$

## Example: total variation

If we use  $d(x,y) = \mathbb{1}[x \neq y]$ :  $\mathcal{W} = d_{TV}$

## Example: Hamming

$$\Omega = [q]^n \quad d(x,y) = |\{i \mid x_i \neq y_i\}|$$

$$\mu = \text{unif on } \{(\bullet, \bullet, \bullet), (\bullet, \bullet, \bullet)\}$$

$$\nu = \text{unif on } \{(\bullet, \bullet, \bullet), (\bullet, \bullet, \bullet), (\bullet, \bullet, \bullet)\}$$

$$\mathcal{W}(\mu, \nu) = \frac{1}{3} \cdot 0 + \frac{1}{6} \cdot 3 + \frac{1}{2} \cdot 2 = 1.5$$



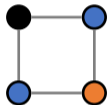
# Coloring

▶ Input: graph  $G$  and  $q \in \mathbb{N}$

# Coloring

- ▶ Input: graph  $G$  and  $q \in \mathbb{N}$
- ▶ Goal: sample proper colorings

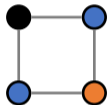
adjacent verts colored differently



# Coloring

- ▶ Input: graph  $G$  and  $q \in \mathbb{N}$
- ▶ Goal: sample proper colorings

↑  
adjacent verts colored differently

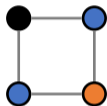


- ▶ NP-hard to even find one!

# Coloring

- ▶ Input: graph  $G$  and  $q \in \mathbb{N}$
- ▶ Goal: sample proper colorings

adjacent verts colored differently



- ▶ NP-hard to even find one!
- ▶ Easy when  $q \geq \Delta + 1$

maximum degree

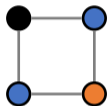
**for** *each vertex*  $v$  **do**

┌ pick a color from  
└  $[q] - \{\text{neighbors' colors}\}$

# Coloring

- ▶ Input: graph  $G$  and  $q \in \mathbb{N}$
- ▶ Goal: sample proper colorings

↑  
adjacent verts colored differently



- ▶ NP-hard to even find one!
- ▶ Easy when  $q \geq \Delta + 1$

↑  
maximum degree

**for each vertex  $v$  do**

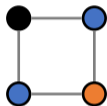
┌ pick a color from  
└  $[q] - \{\text{neighbors' colors}\}$

- ▶ **Open:** approx sample/count when  $q \geq \Delta + 1$

# Coloring

- ▶ Input: graph  $G$  and  $q \in \mathbb{N}$
- ▶ Goal: sample proper colorings

↑  
adjacent vertices colored differently



- ▶ NP-hard to even find one!
- ▶ Easy when  $q \geq \Delta + 1$

↑  
maximum degree

**for each vertex  $v$  do**

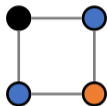
┌ pick a color from  
└  $[q] - \{\text{neighbors' colors}\}$

- ▶ **Open:** approx sample/count when  
 $q \geq \Delta + 1$
- ▶ **Open:** Metropolis/Glauber when:  
 $q \geq \Delta + 2$

# Coloring

- ▶ Input: graph  $G$  and  $q \in \mathbb{N}$
- ▶ Goal: sample proper colorings

adjacent verts colored differently



- ▶ NP-hard to even find one!
- ▶ Easy when  $q \geq \Delta + 1$

maximum degree

**for each vertex  $v$  do**

- pick a color from  
 $[q] - \{\text{neighbors' colors}\}$

- ▶ **Open:** approx sample/count when  
 $q \geq \Delta + 1$
- ▶ **Open:** Metropolis/Glauber when:  
 $q \geq \Delta + 2$
- ▶ Best-known [Chen-Delcourt-Moitra-Perarnau-Postle'18]:

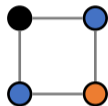
$$q \geq \left(\frac{11}{6} - \epsilon\right)\Delta$$

some tiny constant

# Coloring

- ▶ Input: graph  $G$  and  $q \in \mathbb{N}$
- ▶ Goal: sample proper colorings

adjacent verts colored differently



- ▶ NP-hard to even find one!
- ▶ Easy when  $q \geq \Delta + 1$

maximum degree

**for** each vertex  $v$  **do**

┌ pick a color from  
└  $[q] - \{\text{neighbors' colors}\}$

- ▶ **Open:** approx sample/count when  $q \geq \Delta + 1$
- ▶ **Open:** Metropolis/Glauber when:  $q \geq \Delta + 2$
- ▶ Best-known [Chen-Delcourt-Moitra-Perarnau-Postle'18]:

$$q \geq \left(\frac{11}{6} - \epsilon\right)\Delta$$

some tiny constant

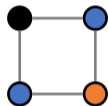
- ▶ We will show  $q \geq 4\Delta + 1$  works.



# Coloring

- ▶ Input: graph  $G$  and  $q \in \mathbb{N}$
- ▶ Goal: sample proper colorings

↑  
adjacent vertices colored differently



- ▶ NP-hard to even find one!
- ▶ Easy when  $q \geq \Delta + 1$

↑  
maximum degree

**for each vertex  $v$  do**

┌ pick a color from  
└  $[q] - \{\text{neighbors' colors}\}$

- ▶ **Open:** approx sample/count when  $q \geq \Delta + 1$
- ▶ **Open:** Metropolis/Glauber when:  $q \geq \Delta + 2$
- ▶ Best-known [Chen-Delcourt-Moitra-Perarnau-Postle'18]:

$$q \geq \left(\frac{11}{6} - \epsilon\right)\Delta$$

↑  
some tiny constant

- ▶ We will show  $q \geq 4\Delta + 1$  works.
- ▶ Then we will improve to  $q \geq 2\Delta + 1$ .

► Strategy: show  $\mathcal{W}$  contracts.

► Strategy: show  $\mathcal{W}$  contracts.

### Lemma

When  $q \geq 4\Delta + 1$ , for Metropolis P:

$$\mathcal{W}(v_P, v'_P) \leq \left(1 - \frac{1}{\text{poly}(n)}\right) \mathcal{W}(v, v')$$

► Strategy: show  $\mathcal{W}$  contracts.

### Lemma

When  $q \geq 4\Delta + 1$ , for Metropolis P:

$$\mathcal{W}(\nu^P, \nu'^P) \leq \left(1 - \frac{1}{\text{poly}(n)}\right) \mathcal{W}(\nu, \nu')$$

► This is enough because

$$\mathcal{W}(\nu^{P^t}, \mu) \leq \left(1 - \frac{1}{\text{poly}}\right)^t \mathcal{W}(\nu, \mu)$$

↑  
upper bounds  $d_{TV}$                       ↑  
at most  $n$

► Strategy: show  $\mathcal{W}$  contracts.

### Lemma

When  $q \geq 4\Delta + 1$ , for Metropolis P:

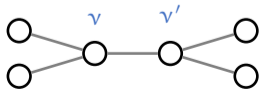
$$\mathcal{W}(vP, v'P) \leq \left(1 - \frac{1}{\text{poly}(n)}\right) \mathcal{W}(v, v')$$

► This is enough because

$$\mathcal{W}(vP^t, \mu) \leq \left(1 - \frac{1}{\text{poly}}\right)^t \mathcal{W}(v, \mu)$$

↑ upper bounds  $d_{TV}$                       ↑ at most  $n$

► Note: unlike  $d_{TV}$ , weak contraction is **NOT** guaranteed.



▷ Strategy: show  $\mathcal{W}$  contracts.

### Lemma

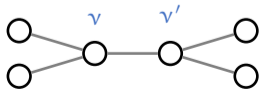
When  $q \geq 4\Delta + 1$ , for Metropolis P:

$$\mathcal{W}(\nu P, \nu' P) \leq \left(1 - \frac{1}{\text{poly}(n)}\right) \mathcal{W}(\nu, \nu')$$

▷ This is enough because

$$\underbrace{\mathcal{W}(\nu P^t, \mu)}_{\text{upper bounds } d_{TV}} \leq \left(1 - \frac{1}{\text{poly}}\right)^t \underbrace{\mathcal{W}(\nu, \mu)}_{\text{at most } n}$$

▷ Note: unlike  $d_{TV}$ , weak contraction is **NOT** guaranteed.



▷ Use **coupling**:

- ▷ Sample  $X_0, X'_0$  from optimal coupling of  $\nu, \nu'$ .
- ▷ Evolve to get  $X_1, X'_1$  while minimizing  $\mathbb{E}[d(X_1, X'_1)]$

- ▶ Strategy: show  $\mathcal{W}$  contracts.

### Lemma

When  $q \geq 4\Delta + 1$ , for Metropolis P:

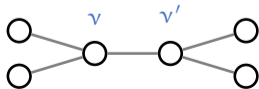
$$\mathcal{W}(\nu P, \nu' P) \leq \left(1 - \frac{1}{\text{poly}(n)}\right) \mathcal{W}(\nu, \nu')$$

- ▶ This is enough because

$$\mathcal{W}(\nu P^t, \mu) \leq \left(1 - \frac{1}{\text{poly}}\right)^t \mathcal{W}(\nu, \mu)$$

↑
↑  
 upper bounds  $d_{TV}$                       at most  $n$

- ▶ Note: unlike  $d_{TV}$ , weak contraction is **NOT** guaranteed.

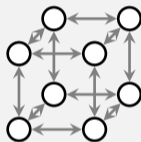


- ▶ Use **coupling**:

- ▶ Sample  $X_0, X'_0$  from optimal coupling of  $\nu, \nu'$ .
- ▶ Evolve to get  $X_1, X'_1$  while minimizing  $\mathbb{E}[d(X_1, X'_1)]$

### Warmup example: hypercube

- ▶  $\Omega = \{0, 1\}^n$
- ▶ Pick u.r.  $i \in [n]$
- ▶ Replace coord  $i$  with  $\text{Ber}(\frac{1}{2})$



- ▷ Strategy: show  $\mathcal{W}$  contracts.

### Lemma

When  $q \geq 4\Delta + 1$ , for Metropolis P:

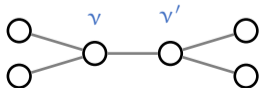
$$\mathcal{W}(\nu P, \nu' P) \leq \left(1 - \frac{1}{\text{poly}(n)}\right) \mathcal{W}(\nu, \nu')$$

- ▷ This is enough because

$$\mathcal{W}(\nu P^t, \mu) \leq \left(1 - \frac{1}{\text{poly}}\right)^t \mathcal{W}(\nu, \mu)$$

↑
↑  
 upper bounds  $d_{TV}$ 
at most  $n$

- ▷ Note: unlike  $d_{TV}$ , weak contraction is **NOT** guaranteed.

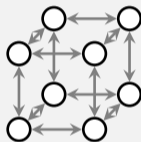


- ▷ Use **coupling**:

- ▷ Sample  $X_0, X'_0$  from optimal coupling of  $\nu, \nu'$ .
- ▷ Evolve to get  $X_1, X'_1$  while minimizing  $\mathbb{E}[d(X_1, X'_1)]$

### Warmup example: hypercube

- ▷  $\Omega = \{0, 1\}^n$
- ▷ Pick u.r.  $i \in [n]$
- ▷ Replace coord  $i$  with  $\text{Ber}(\frac{1}{2})$
- ▷ Pick **same**  $i$  and **same**  $\text{Ber}(\frac{1}{2})$





- ▷ Strategy: show  $\mathcal{W}$  contracts.

### Lemma

When  $q \geq 4\Delta + 1$ , for Metropolis P:

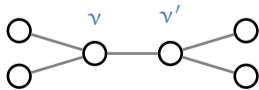
$$\mathcal{W}(\nu P, \nu' P) \leq \left(1 - \frac{1}{\text{poly}(n)}\right) \mathcal{W}(\nu, \nu')$$

- ▷ This is enough because

$$\mathcal{W}(\nu P^t, \mu) \leq \left(1 - \frac{1}{\text{poly}}\right)^t \mathcal{W}(\nu, \mu)$$

↑
↑  
 upper bounds  $d_{TV}$ 
at most  $n$

- ▷ Note: unlike  $d_{TV}$ , weak contraction is **NOT** guaranteed.



- ▷ Use **coupling**:

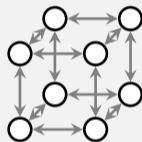
- ▷ Sample  $X_0, X'_0$  from optimal coupling of  $\nu, \nu'$ .
- ▷ Evolve to get  $X_1, X'_1$  while minimizing  $\mathbb{E}[d(X_1, X'_1)]$

### Warmup example: hypercube

- ▷  $\Omega = \{0, 1\}^n$

- ▷ Pick u.r.  $i \in [n]$

- ▷ Replace coord  $i$  with  $\text{Ber}(\frac{1}{2})$



- ▷ Pick **same**  $i$  and **same**  $\text{Ber}(\frac{1}{2})$

- ▷ If  $d(X_0, X'_0) = k$ , then

$$\mathbb{E}[d(X_1, X'_1) \mid X_0, X'_0] = k - \frac{k}{n}$$

- ▷ Strategy: show  $\mathcal{W}$  contracts.

### Lemma

When  $q \geq 4\Delta + 1$ , for Metropolis P:

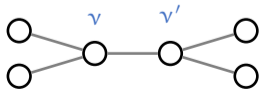
$$\mathcal{W}(\nu P, \nu' P) \leq \left(1 - \frac{1}{\text{poly}(n)}\right) \mathcal{W}(\nu, \nu')$$

- ▷ This is enough because

$$\mathcal{W}(\nu P^t, \mu) \leq \left(1 - \frac{1}{\text{poly}}\right)^t \mathcal{W}(\nu, \mu)$$

↑ upper bounds  $d_{TV}$ 
↑ at most  $n$

- ▷ Note: unlike  $d_{TV}$ , weak contraction is **NOT** guaranteed.

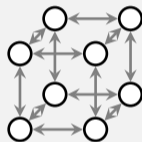


- ▷ Use **coupling**:

- ▷ Sample  $X_0, X'_0$  from optimal coupling of  $\nu, \nu'$ .
- ▷ Evolve to get  $X_1, X'_1$  while minimizing  $\mathbb{E}[d(X_1, X'_1)]$

### Warmup example: hypercube

- ▷  $\Omega = \{0, 1\}^n$
- ▷ Pick u.r.  $i \in [n]$
- ▷ Replace coord  $i$  with  $\text{Ber}(\frac{1}{2})$



- ▷ Pick **same**  $i$  and **same**  $\text{Ber}(\frac{1}{2})$
- ▷ If  $d(X_0, X'_0) = k$ , then

$$\mathbb{E}[d(X_1, X'_1) \mid X_0, X'_0] = k - \frac{k}{n}$$

- ▷  $\mathcal{W}(\nu P, \nu' P) \leq (1 - 1/n) \mathcal{W}(\nu, \nu')$

- ▷ Strategy: show  $\mathcal{W}$  contracts.

### Lemma

When  $q \geq 4\Delta + 1$ , for Metropolis P:

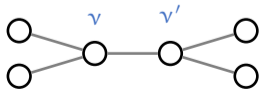
$$\mathcal{W}(\nu P, \nu' P) \leq \left(1 - \frac{1}{\text{poly}(n)}\right) \mathcal{W}(\nu, \nu')$$

- ▷ This is enough because

$$\mathcal{W}(\nu P^t, \mu) \leq \left(1 - \frac{1}{\text{poly}(n)}\right)^t \mathcal{W}(\nu, \mu)$$

↑
↑  
 upper bounds  $d_{TV}$ 
at most  $n$

- ▷ Note: unlike  $d_{TV}$ , weak contraction is **NOT** guaranteed.

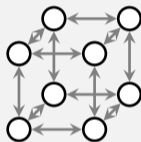


- ▷ Use **coupling**:

- ▷ Sample  $X_0, X'_0$  from optimal coupling of  $\nu, \nu'$ .
- ▷ Evolve to get  $X_1, X'_1$  while minimizing  $\mathbb{E}[d(X_1, X'_1)]$

### Warmup example: hypercube

- ▷  $\Omega = \{0, 1\}^n$
- ▷ Pick u.r.  $i \in [n]$
- ▷ Replace coord  $i$  with  $\text{Ber}(\frac{1}{2})$



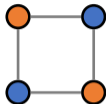
- ▷ Pick **same**  $i$  and **same**  $\text{Ber}(\frac{1}{2})$
- ▷ If  $d(X_0, X'_0) = k$ , then

$$\mathbb{E}[d(X_1, X'_1) \mid X_0, X'_0] = k - \frac{k}{n}$$

- ▷  $\mathcal{W}(\nu P, \nu' P) \leq (1 - 1/n) \mathcal{W}(\nu, \nu')$
- ▷  $t_{\text{mix}}(\epsilon) \leq n \log n + n \log(1/\epsilon)$

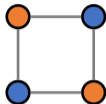
Take the **Metropolis chain** for colorings:

- ▶ Pick u.r. vertex  $v$
- ▶ Pick u.r. color  $c$
- ▶ Color  $v$  with  $c$  if valid



Take the **Metropolis chain** for colorings:

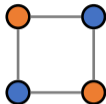
- ▶ Pick u.r. vertex  $v$
- ▶ Pick u.r. color  $c$
- ▶ Color  $v$  with  $c$  if valid



Coupling:

Take the **Metropolis chain** for colorings:

- ▶ Pick u.r. vertex  $v$
- ▶ Pick u.r. color  $c$
- ▶ Color  $v$  with  $c$  if valid

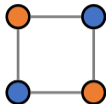


Coupling:

- ▶ Pick **same**  $v$  and **same**  $c$

Take the **Metropolis chain** for colorings:

- ▶ Pick u.r. vertex  $v$
- ▶ Pick u.r. color  $c$
- ▶ Color  $v$  with  $c$  if valid

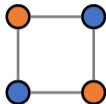


Coupling:

- ▶ Pick **same**  $v$  and **same**  $c$
- ▶ If  $d(X_0, X'_0) = k$ , then  $d(X_1, X'_1)$  is:
  - ▶  $k - 1$  (lucky)
  - ▶  $k + 1$  (unlucky)
  - ▶  $k$  (neutral)

Take the **Metropolis chain** for colorings:

- ▶ Pick u.r. vertex  $v$
- ▶ Pick u.r. color  $c$
- ▶ Color  $v$  with  $c$  if valid



Coupling:

- ▶ Pick **same**  $v$  and **same**  $c$
- ▶ If  $d(X_0, X'_0) = k$ , then  $d(X_1, X'_1)$  is:
  - ▶  $k - 1$  (lucky)
  - ▶  $k + 1$  (unlucky)
  - ▶  $k$  (neutral)

$$\mathbb{P}[\text{lucky}] \geq \underbrace{(k/n)}_{\text{pick differing } v} \cdot \underbrace{(q - 2\Delta)}_{\text{c available to both}} / q$$

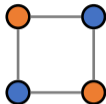
pick differing  $v$

$c$  available to both



Take the **Metropolis chain** for colorings:

- ▶ Pick u.r. vertex  $v$
- ▶ Pick u.r. color  $c$
- ▶ Color  $v$  with  $c$  if valid



▶  $\mathbb{P}[\text{unlucky}] \leq 2k\Delta/qn$

$c$  color of differing neighbor in  $X_0$  or  $X'_0$

Coupling:

- ▶ Pick **same**  $v$  and **same**  $c$
- ▶ If  $d(X_0, X'_0) = k$ , then  $d(X_1, X'_1)$  is:
  - ▶  $k - 1$  (lucky)
  - ▶  $k + 1$  (unlucky)
  - ▶  $k$  (neutral)

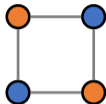
▶  $\mathbb{P}[\text{lucky}] \geq (k/n) \cdot (q - 2\Delta)/q$

pick differing  $v$

$c$  available to both

Take the **Metropolis chain** for colorings:

- ▶ Pick u.r. vertex  $v$
- ▶ Pick u.r. color  $c$
- ▶ Color  $v$  with  $c$  if valid



▶  $\mathbb{P}[\text{unlucky}] \leq 2k\Delta/qn$

$c$  color of differing neighbor in  $X_0$  or  $X'_0$

▶ We get  $\mathbb{E}[d(X_1, X'_1) \mid X_0, X'_0] \leq$

$$k - \frac{k(q-2\Delta)}{qn} + \frac{2k\Delta}{qn} = k \cdot \left(1 - \frac{q-4\Delta}{qn}\right)$$

Coupling:

- ▶ Pick **same**  $v$  and **same**  $c$
- ▶ If  $d(X_0, X'_0) = k$ , then  $d(X_1, X'_1)$  is:
  - ▶  $k - 1$  (lucky)
  - ▶  $k + 1$  (unlucky)
  - ▶  $k$  (neutral)

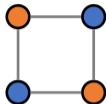
▶  $\mathbb{P}[\text{lucky}] \geq (k/n) \cdot (q - 2\Delta)/q$

pick differing  $v$

$c$  available to both

Take the **Metropolis chain** for colorings:

- ▶ Pick u.r. vertex  $v$
- ▶ Pick u.r. color  $c$
- ▶ Color  $v$  with  $c$  if valid



▶  $\mathbb{P}[\text{unlucky}] \leq 2k\Delta/qn$

$c$  color of differing neighbor in  $X_0$  or  $X'_0$

▶ We get  $\mathbb{E}[d(X_1, X'_1) \mid X_0, X'_0] \leq$

$$k - \frac{k(q-2\Delta)}{qn} + \frac{2k\Delta}{qn} = k \cdot \left(1 - \frac{q-4\Delta}{qn}\right)$$

▶ As long as  $q \geq 4\Delta + 1$ , we have **contraction**. 😊

Coupling:

- ▶ Pick **same**  $v$  and **same**  $c$
- ▶ If  $d(X_0, X'_0) = k$ , then  $d(X_1, X'_1)$  is:
  - ▶  $k - 1$  (lucky)
  - ▶  $k + 1$  (unlucky)
  - ▶  $k$  (neutral)

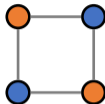
▶  $\mathbb{P}[\text{lucky}] \geq (k/n) \cdot (q - 2\Delta)/q$

pick differing  $v$

$c$  available to both

Take the **Metropolis chain** for colorings:

- ▶ Pick u.r. vertex  $v$
- ▶ Pick u.r. color  $c$
- ▶ Color  $v$  with  $c$  if valid



▶  $\mathbb{P}[\text{unlucky}] \leq 2k\Delta/qn$

$c$  color of differing neighbor in  $X_0$  or  $X'_0$

▶ We get  $\mathbb{E}[d(X_1, X'_1) \mid X_0, X'_0] \leq$

$$k - \frac{k(q-2\Delta)}{qn} + \frac{2k\Delta}{qn} = k \cdot \left(1 - \frac{q-4\Delta}{qn}\right)$$

▶ As long as  $q \geq 4\Delta + 1$ , we have **contraction**. 😊

▶ We get

$$t_{\text{mix}}(\epsilon) = O\left(\frac{q}{q-4\Delta} \cdot n \log(n/\epsilon)\right)$$

Coupling:

- ▶ Pick **same**  $v$  and **same**  $c$
- ▶ If  $d(X_0, X'_0) = k$ , then  $d(X_1, X'_1)$  is:
  - ▶  $k - 1$  (lucky)
  - ▶  $k + 1$  (unlucky)
  - ▶  $k$  (neutral)

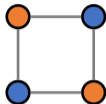
▶  $\mathbb{P}[\text{lucky}] \geq (k/n) \cdot (q - 2\Delta)/q$

pick differing  $v$

$c$  available to both

Take the **Metropolis chain** for colorings:

- ▶ Pick u.r. vertex  $v$
- ▶ Pick u.r. color  $c$
- ▶ Color  $v$  with  $c$  if valid



▶  $\mathbb{P}[\text{unlucky}] \leq 2k\Delta/qn$

$c$  color of differing neighbor in  $X_0$  or  $X'_0$

▶ We get  $\mathbb{E}[d(X_1, X'_1) \mid X_0, X'_0] \leq$

$$k - \frac{k(q-2\Delta)}{qn} + \frac{2k\Delta}{qn} = k \cdot \left(1 - \frac{q-4\Delta}{qn}\right)$$

▶ As long as  $q \geq 4\Delta + 1$ , we have **contraction**. 😊

▶ We get

$$t_{\text{mix}}(\epsilon) = O\left(\frac{q}{q-4\Delta} \cdot n \log(n/\epsilon)\right)$$

▶ Exercise: analyze Glauber this way.

Coupling:

- ▶ Pick **same**  $v$  and **same**  $c$
- ▶ If  $d(X_0, X'_0) = k$ , then  $d(X_1, X'_1)$  is:
  - ▶  $k - 1$  (lucky)
  - ▶  $k + 1$  (unlucky)
  - ▶  $k$  (neutral)

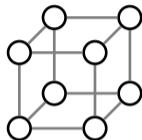
▶  $\mathbb{P}[\text{lucky}] \geq (k/n) \cdot (q - 2\Delta)/q$

pick differing  $v$

$c$  available to both

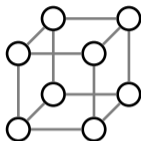
▶ Hamming distance is special.

- ▶ Hamming distance is special.
- ▶ There is a **sparse** graph s.t.  $d(x, y)$  is shortest path from  $x$  to  $y$ .



$x \sim y$  when  $x_i \neq y_i$  for one  $i$

- ▶ Hamming distance is special.
- ▶ There is a **sparse** graph s.t.  $d(x, y)$  is shortest path from  $x$  to  $y$ .

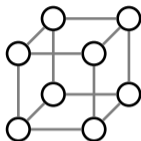


$x \sim y$  when  $x_i \neq y_i$  for one  $i$

- ▶ In general, if  $d$  is shortest path metric derived from a (possibly weighted) graph, we can use **path coupling** [Bubley-Dyer].



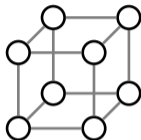
- ▶ Hamming distance is special.
- ▶ There is a **sparse** graph s.t.  $d(x, y)$  is shortest path from  $x$  to  $y$ .



$x \sim y$  when  $x_i \neq y_i$  for one  $i$

- ▶ In general, if  $d$  is shortest path metric derived from a (possibly weighted) graph, we can use **path coupling** [Bubley-Dyer].
- ▶ **Idea:** only couple starting states  $X_0, X'_0$  that are **adjacent**.

- ▶ Hamming distance is special.
- ▶ There is a **sparse** graph s.t.  $d(x, y)$  is shortest path from  $x$  to  $y$ .



$x \sim y$  when  $x_i \neq y_i$  for one  $i$

- ▶ In general, if  $d$  is shortest path metric derived from a (possibly weighted) graph, we can use **path coupling** [Bubley-Dyer].
- ▶ **Idea:** only couple starting states  $X_0, X'_0$  that are **adjacent**.

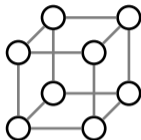
## Path coupling lemma

Suppose for all adjacent  $X_0 \sim X'_0$  we can couple  $X_1, X'_1$  s.t.

$$\mathbb{E}[d(X_1, X'_1)] \leq (1 - c)d(X_0, X'_0).$$

Then  $\mathcal{W}(\nu P, \nu' P) \leq (1 - c) \mathcal{W}(\nu, \nu')$ .

- ▶ Hamming distance is special.
- ▶ There is a **sparse** graph s.t.  $d(x, y)$  is shortest path from  $x$  to  $y$ .



$x \sim y$  when  $x_i \neq y_i$  for one  $i$

- ▶ In general, if  $d$  is shortest path metric derived from a (possibly weighted) graph, we can use **path coupling** [Bubley-Dyer].
- ▶ **Idea:** only couple starting states  $X_0, X'_0$  that are **adjacent**.

## Path coupling lemma

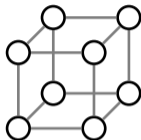
Suppose for all adjacent  $X_0 \sim X'_0$  we can couple  $X_1, X'_1$  s.t.

$$\mathbb{E}[d(X_1, X'_1)] \leq (1 - c)d(X_0, X'_0).$$

Then  $\mathcal{W}(\nu P, \nu' P) \leq (1 - c) \mathcal{W}(\nu, \nu')$ .

Proof:

- ▶ Hamming distance is special.
- ▶ There is a **sparse** graph s.t.  $d(x, y)$  is shortest path from  $x$  to  $y$ .



$x \sim y$  when  $x_i \neq y_i$  for one  $i$

- ▶ In general, if  $d$  is shortest path metric derived from a (possibly weighted) graph, we can use **path coupling** [Bubley-Dyer].
- ▶ **Idea:** only couple starting states  $X_0, X'_0$  that are **adjacent**.

## Path coupling lemma

Suppose for all adjacent  $X_0 \sim X'_0$  we can couple  $X_1, X'_1$  s.t.

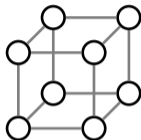
$$\mathbb{E}[d(X_1, X'_1)] \leq (1 - c)d(X_0, X'_0).$$

Then  $\mathcal{W}(\nu P, \nu' P) \leq (1 - c) \mathcal{W}(\nu, \nu')$ .

Proof:

- ▶ Take arbitrary  $X_0, X'_0$ .

- ▶ Hamming distance is special.
- ▶ There is a **sparse** graph s.t.  $d(x, y)$  is shortest path from  $x$  to  $y$ .



$x \sim y$  when  $x_i \neq y_i$  for one  $i$

- ▶ In general, if  $d$  is shortest path metric derived from a (possibly weighted) graph, we can use **path coupling** [Bubley-Dyer].
- ▶ **Idea:** only couple starting states  $X_0, X'_0$  that are **adjacent**.

## Path coupling lemma

Suppose for all adjacent  $X_0 \sim X'_0$  we can couple  $X_1, X'_1$  s.t.

$$\mathbb{E}[d(X_1, X'_1)] \leq (1 - c)d(X_0, X'_0).$$

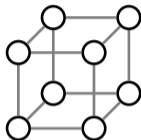
Then  $\mathcal{W}(\nu P, \nu' P) \leq (1 - c) \mathcal{W}(\nu, \nu')$ .

Proof:

- ▶ Take arbitrary  $X_0, X'_0$ .
- ▶ Let shortest path be

$$X_0 = v_0 \rightarrow v_1 \rightarrow \cdots \rightarrow v_k = X'_0$$

- ▶ Hamming distance is special.
- ▶ There is a **sparse** graph s.t.  $d(x, y)$  is shortest path from  $x$  to  $y$ .



$x \sim y$  when  $x_i \neq y_i$  for one  $i$

- ▶ In general, if  $d$  is shortest path metric derived from a (possibly weighted) graph, we can use **path coupling** [Bubley-Dyer].
- ▶ **Idea:** only couple starting states  $X_0, X'_0$  that are **adjacent**.

## Path coupling lemma

Suppose for all adjacent  $X_0 \sim X'_0$  we can couple  $X_1, X'_1$  s.t.

$$\mathbb{E}[d(X_1, X'_1)] \leq (1 - c)d(X_0, X'_0).$$

Then  $\mathcal{W}(\nu P, \nu' P) \leq (1 - c) \mathcal{W}(\nu, \nu')$ .

Proof:

- ▶ Take arbitrary  $X_0, X'_0$ .
- ▶ Let shortest path be

$$X_0 = v_0 \rightarrow v_1 \rightarrow \dots \rightarrow v_k = X'_0$$

- ▶ By triangle ineq  $\mathcal{W}(\mathbb{1}_{X_0} P, \mathbb{1}_{X'_0} P) \leq$

$$\sum_i \mathcal{W}(\mathbb{1}_{v_i} P, \mathbb{1}_{v_{i+1}} P) \leq (1 - c) \sum_i d(v_i, v_{i+1}) = (1 - c)d(X_0, X'_0)$$

Triangle inequality holds because couplings can be **stitched** together!

$$\nu_0 \xrightarrow{\pi_{0,1}} \nu_1 \xrightarrow{\pi_{1,2}} \dots \xrightarrow{\pi_{k-1,k}} \nu_k$$

Exercise: there is joint dist with marginals  $\pi_{i,i+1}$ !