

Markov Chain Convergence

- Total variation distance & coupling
- MC convergence theorem
- Mixing times via $\left\{ \begin{array}{l} \text{coupling} \\ \text{strong stationary times} \end{array} \right.$
- Markov chain Monte Carlo
- Perfect sampling $\left\{ \begin{array}{l} \text{strong Doeblin cond} \\ \text{coupling from the past} \end{array} \right.$

Thm (Convergence Thm for finite MC) Given finite MC P that is irreducible, aperiodic & positive recurrent, with stationary distribution π . Then $\exists \alpha \in (0, 1]$ and $C > 0$ s.t.

$$\max_{x \in X} \|P^t(x, \cdot) - \pi\|_{TV} \leq C \alpha^t$$

We first show this for finite state chains

Pf - Since P is irreducible, aperiodic $\Rightarrow \exists r \geq 1$ s.t.

$P^r > 0$ (ie. $P^r(x, y) > 0 \forall x, y \in X$). Thus for some $\delta > 0$, we have

$$P^r(x, y) \geq \delta \pi(y) \forall x, y \in X$$

- Let $\alpha = 1 - \delta$. We can write $P^r = (1 - \alpha) \mathbb{1}^T \pi + \alpha Q$, where Q is a stochastic matrix, and we denote $\mathbb{T} = \mathbb{1}^T \pi = \begin{pmatrix} \pi \\ \pi \\ \vdots \end{pmatrix}$

- Now we claim: $\forall k \geq 1, P^{rk} = (1 - \alpha^k) \mathbb{T} + \alpha^k Q^k$

Can show this by induction - suppose true for k , then

$$P^{r(k+1)} = P^{rk} P^r = ((1 - \alpha^k) \mathbb{T} + \alpha^k Q^k) ((1 - \alpha) \mathbb{T} + \alpha Q)$$

Since $M \mathbb{T} = \mathbb{T}$
for any stochastic matrix M

$$= (1 - \alpha^k) \mathbb{T} + \alpha^k (1 - \alpha) \mathbb{T} + \alpha^{k+1} Q^{k+1}$$

$$= (1 - \alpha^{k+1}) \mathbb{T} + \alpha^{k+1} Q^{k+1}$$

- Thus $P^{rk+j} - \mathbb{T} = \alpha^k (Q^k P^j - \mathbb{T})$

Now for any $x \in X$, we have that

$$\|P^{rk+j}(x, \cdot) - \pi\|_{TV} \leq \alpha^k$$

This follows from the fact that $\|\mu - \nu\|_{TV} = \|\mu - \nu\|_1 / 2$, and since

$\|Q^k P^j(x, \cdot) - \pi\|_1 \leq 2$ for stochastic matrices Q, P □

- The total variation distance between any two probability distributions μ and ν on Ω is defined as $\|\mu - \nu\|_{TV} = \sup_{A \subset \Omega} |\mu(A) - \nu(A)|$

- Properties of $\|\mu - \nu\|_{TV}$ (or $d_{TV}(\mu, \nu)$)

i) $\|\mu - \nu\|_{TV} = \frac{1}{2} \sum_{x \in \Omega} |\mu(x) - \nu(x)| = \frac{1}{2} \|\mu - \nu\|_1$

ii) $\|\mu - \nu\|_{TV} = \sum_{x \in \Omega | \mu(x) \geq \nu(x)} (\mu(x) - \nu(x))$

iii) $\|\mu - \nu\|_{TV} = \frac{1}{2} \sup \left\{ \sum_{x \in \Omega} f(x)\mu(x) - \sum_{x \in \Omega} f(x)\nu(x) \mid \|f\|_{\infty} \leq 1 \right\}$

test function

- A coupling of 2 probability distributions F and G is a pair of r.v. (X, Y) defined on the same probability space $(\Omega, \mathcal{F}, \mathbb{P})$ s.t.

$$\forall x, y, \mathbb{P}[X \leq x] = F(x), \mathbb{P}[Y \leq y] = G(y)$$

- In other words, we 'design' X and Y so that they have the correct marginal distr (but can have any joint distr)

Eg - $F \sim \text{Ber}(p)$, $G \sim \text{Ber}(p/2)$

Coupling 1 - $X \sim \text{Ber}(p)$, $Y \sim \text{Ber}(p/2)$ independently

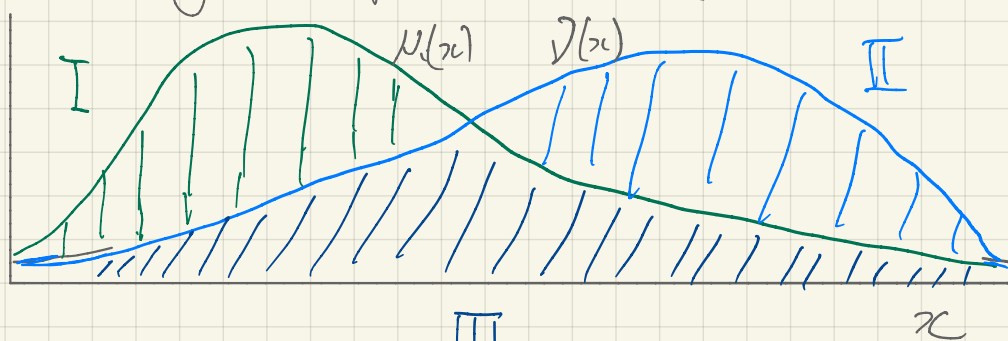
Coupling 2 - $X \sim \text{Ber}(p)$; If $X=0$ then $Y=0$, and
if $X=1$, then $Y \sim \text{Ber}(1/2)$

$$\left(\text{ie, } (X, Y) = \begin{cases} (1, 1) & \text{wp } p/2 \\ (1, 0) & \text{wp } p/2 \\ (0, 0) & \text{wp } 1-p \end{cases} \right)$$

Lemma (the 'optimal coupling') Let μ and ν be two probability distributions on Ω . Then

$$\|\mu - \nu\|_{TV} = \inf \left\{ \mathbb{P}[X \neq Y] \mid (X, Y) \text{ is a coupling of } \mu \text{ and } \nu \right\}$$

Pf - Intuitively the proof is as follows



Given $\mu(x)$ and $\nu(x)$, we can plot them as above

Now by defn, the area of regions I and II are $\|\mu - \nu\|_{TV}$, and area of III is $1 - \|\mu - \nu\|_{TV}$. Now we sample from III w.p. $1 - \|\mu - \nu\|_{TV}$ and set $X=Y$, else - sample X from I and Y from II

Formally, for arbitrary $A \subset \Omega$, and $X \sim \mu, Y \sim \nu$

$$\begin{aligned} \mathbb{P}[X \neq Y] &\geq \mathbb{P}[X \in A, Y \in \bar{A}] = \mathbb{P}[X \in A] - \mathbb{P}[X \in A, Y \in A] \\ &\geq \mathbb{P}[X \in A] - \mathbb{P}[Y \in A] \\ &= \mu(A) - \nu(A) \end{aligned}$$

Thus $\|\mu - \nu\|_{TV} \leq \inf \{ \mathbb{P}[X \neq Y] \mid (X, Y) \text{ coupling of } \mu, \nu \}$

To show equality, we construct coupling (X, Y) as follows - let $U \in \{0, 1\}$ and $Z, V, W \in \Omega$ be independent r.v. defined as

$$U = \text{Ber}(1 - \|\mu - \nu\|_{TV})$$

$$Z \equiv \mathbb{P}[Z=i] = (\mu(i) \wedge \nu(i)) / (1 - \|\mu - \nu\|_{TV})$$

$$V \equiv \mathbb{P}[V=i] = (\mu(i) - \nu(i))^+ / \|\mu - \nu\|_{TV}$$

$$W \equiv \mathbb{P}[W=i] = (\nu(i) - \mu(i))^+ / \|\mu - \nu\|_{TV}$$

Now let $(X, Y) = \begin{cases} (Z, Z) & \text{if } U=1 \\ (V, W) & \text{if } U=0 \end{cases}$

Note that $V \neq W \Rightarrow \mathbb{P}[X \neq Y] = \|\mu - \nu\|_{TV}$, and

$$\mathbb{P}[X=i] = \frac{(\mu(i) \wedge \nu(i))}{(1 - \|\mu - \nu\|_{TV})} \cdot (1 - \|\mu - \nu\|_{TV}) + \|\mu - \nu\|_{TV} \frac{(\mu(i) - \nu(i))^+}{\|\mu - \nu\|_{TV}}$$

Similarly $\mathbb{P}[Y=i] = \nu(i)$

Thm (Convergence Thm for countable MC) Given P on countable X that is irreducible, aperiodic & positive recurrent, with stationary distribution π . Then $\forall x \in X$

$$\lim_{t \rightarrow \infty} \|P^t(x, \cdot) - \pi\|_{TV} = 0$$

Pf - Consider the 2-dimensional MC (X_t, Y_t) defined as $\tilde{P}[(X_{t+1}, Y_{t+1}) = (i', j') | (X_t, Y_t) = (i, j)] = P(i, i')P(j, j') \forall i, i', j, j'$

- Can view this as 2 particles starting at $X_0 = x_0$, $X_1 = x_1$ and evolving independently
- Claim - P irreducible, aperiodic \Rightarrow the joint chain (X_t, Y_t) with transition \tilde{P} is also irreducible and aperiodic (check!)
- Also $\tilde{\pi}(i, j) = \pi(i)\pi(j)$ is a stationary distribution for \tilde{P} !

\Rightarrow \tilde{P} is an ergodic MC with stat dist $\tilde{\pi}$

- Now let $\tau_c = \inf \{t \geq 0 \mid X_t = Y_t\}$ (Coupling time)

Claim - for any μ, ν , we have $\mathbb{P}_{x, y}[\tau_c < \infty] = 1$

This follows from the fact \tilde{P} is ergodic, and since τ_c is the hitting time of set $A = \{(x, x) \mid x \in X\}$

- Next, for any pair of $\text{dist}(\mu, \nu)$, we can construct coupling (X_t, Y_t) , where $X_0 \sim \mu$, $Y_0 \sim \nu$, and the chains evolve as $(\forall (i, j) \in X^2, (i', j') \in X^2)$

$$P[(X_{t+1}, Y_{t+1}) = (i', j') \mid (X_t, Y_t) = (i, j)] = \begin{cases} P(i, i')P(j, j'), & i \neq j \\ P(i, i') \mathbb{1}_{\{i=j'\}}, & i = j \end{cases}$$

- From above, we know coupling time τ_c is a.s. finite

- Now using this claim, we can prove the theorem as follows. Recall P ergodic \Rightarrow unique stationary distribution π . Now let $\nu = \pi$, and $\mu = \delta_x$ (i.e., $X_0 = x$). Then we have

$$\begin{aligned} \|P^t(x, \cdot) - \pi\|_{TV} &\leq P_{\delta_x, \pi}[X_t \neq Y_t] \\ &= P_{\delta_x, \pi}[\tau_c > t] \end{aligned}$$

$$\text{and } \lim_{t \rightarrow \infty} P_{\delta_x, \pi}[\tau_c > t] = \sum_{y \in X} \pi(y) \lim_{t \rightarrow \infty} P_{x, y}[\tau_c > t] = 0$$

- Thus $\lim_{t \rightarrow \infty} \|P^t(x, \cdot) - \pi\|_{TV} = 0 \quad \forall x \in X$



As a result, π is often referred to as the **equilibrium dist**

Mixing Times

We know use the above idea to understand how 'fast' $\mathbb{T} P^t$ converges to \mathbb{T} .

Def - For ergodic MC P with stationary distⁿ \mathbb{T} , we define $\forall t \geq 0$

$$d(t) = \max_{x \in X} \|P^t(x, \cdot) - \mathbb{T}\|_{TV}$$

$$\bar{d}(t) = \max_{x, y \in X} \|P^t(x, \cdot) - P^t(y, \cdot)\|_{TV}$$

(distance from stationarity)

• Lemma - $d(t) \leq \bar{d}(t) \leq 2d(t)$

Pf - For the upper bound, since $\|\cdot\|_{TV}$ is a norm, we have by Δ inequality -

$$\begin{aligned} \max_{x, y} \|P^t(x, \cdot) - P^t(y, \cdot)\|_{TV} &\leq \max_{x, y} (\|P^t(x, \cdot) - \mathbb{T}\|_{TV} + \|P^t(y, \cdot) - \mathbb{T}\|_{TV}) \\ &\leq 2 \max_x \|P^t(x, \cdot) - \mathbb{T}\|_{TV} = 2d(t) \end{aligned}$$

- For the lower bound, note that $\forall A \subseteq X$, $\mathbb{T}(A) = \sum_{x \in X} \mathbb{T}(x) P^t(x, A)$

$$\begin{aligned} \Rightarrow \|P^t(x, \cdot) - \mathbb{T}\|_{TV} &= \sup_{A \subseteq \Omega} |P^t(x, A) - \mathbb{T}(A)| \\ &= \sup_{A \subseteq \Omega} \left| \sum_y \mathbb{T}(y) (P^t(x, A) - P^t(y, A)) \right| \\ &\leq \sup_{A \subseteq \Omega} \sum_y \mathbb{T}(y) |P^t(x, A) - P^t(y, A)| = \bar{d}(t) \quad \square \end{aligned}$$

• Lemma - $\bar{d}(s+t) \leq \bar{d}(s)\bar{d}(t) \forall s, t \geq 0$

Pr - Fix $x, y \in X$, and for any $s \geq 0$, let (X_s, Y_s) be the optimal coupling of $P^s(x, \cdot)$ and $P^t(y, \cdot)$, i.e., $X_s \sim P^s(x, \cdot)$, $Y_s \sim P^s(y, \cdot)$, $\|P^s(x, \cdot) - P^s(y, \cdot)\|_{TV} = P[X_s \neq Y_s]$

• Now $P^{s+t}(x, w) = \sum_z P[X_s = z] P^t(z, w) = E[P^t(X_s, w)]$

Similarly $P^{s+t}(y, w) = E[P^t(Y_s, w)]$

$\Rightarrow P^{s+t}(x, w) - P^{s+t}(y, w) = E[P^t(X_s, w) - P^t(Y_s, w)]$

(Note - we can do this as X_s, Y_s are on the same $(\Omega, \mathcal{F}, \mathbb{P})$)

$$\begin{aligned} \Rightarrow \|P^{s+t}(x, \cdot) - P^{s+t}(y, \cdot)\|_{TV} &= \frac{1}{2} \sum_w |P^{s+t}(x, w) - P^{s+t}(y, w)| \\ &= E \left[\frac{1}{2} \sum_w |P^t(X_s, w) - P^t(Y_s, w)| \right] \\ &= E \left[\|P^t(X_s, \cdot) - P^t(Y_s, \cdot)\|_{TV} \right] \end{aligned}$$

• However $\because P^t(X_s, \cdot) = P^t(Y_s, \cdot)$ when $X_s = Y_s$, we have

$$\begin{aligned} \|P^{s+t}(x, \cdot) - P^{s+t}(y, \cdot)\|_{TV} &\leq \bar{d}(t) E[\mathbb{1}_{\{X_s \neq Y_s\}}] \\ &= \bar{d}(t) \bar{d}(s) \quad \square \end{aligned}$$

Thus, $\bar{d}(t)$ is submultiplicative $\Rightarrow \bar{d}(ct) \leq \bar{d}(t)^c$

• Def (Mixing Time) - For any P_i and $\varepsilon > 0$

$$t_{\text{mix}}(\varepsilon) = \inf \{ t \geq 0 \mid d(t) \leq \varepsilon \}$$

$$t_{\text{mix}} = t_{\text{mix}}(1/4)$$

• Why $1/4$? Using previous lemma we have

$$d(2t_{\text{mix}}) \leq \bar{d}(2t_{\text{mix}}) \leq (\bar{d}(t_{\text{mix}}))^2 \leq (2d(t_{\text{mix}}))^2 = 2^{-1}$$

$$\Rightarrow t_{\text{mix}}(\varepsilon) \leq \lceil \log_2(1/\varepsilon) \rceil t_{\text{mix}}$$

• We now want to understand how t_{mix} behaves

Thm (Mixing from coupling) Let (X_t, Y_t) be a coupling st.

$X_0 = x, Y_0 = y$, and if $X_s = Y_s$ then $X_t = Y_t \forall t \geq s$.

Define $T_{\text{couple}} = \inf \{ t \geq 0 \mid X_t = Y_t \}$

Then we have $\bar{d}(t) \leq \mathbb{P}_{x,y} [T_{\text{couple}} > t]$

$$\mathbb{P}^t(x, z) = \mathbb{P}_{x,y} [X_t = z], \quad \mathbb{P}^t(y, z) = \mathbb{P}_{x,y} [Y_t = z]$$

\Rightarrow for coupling (X_t, Y_t) , $\|\mathbb{P}^t(x, \cdot) - \mathbb{P}^t(y, \cdot)\|_{\text{TV}} \leq \mathbb{P}_{x,y} [X_t \neq Y_t]$

Also $\mathbb{P}[X_t \neq Y_t] = \mathbb{P}[T_{\text{couple}} > t]$, and we are done!



Eg - Lazy walk on circle - $X = \{0, 1, \dots, n-1\}$

$$Z_t = \begin{cases} -1 & \text{wp } 1/4 \\ 0 & \text{wp } 1/2 \\ 1 & \text{wp } 1/4 \end{cases}, \quad X_{t+1} = (X_t + Z_{t+1}) \bmod n$$

• Check that $\Pi = \left(\frac{1}{n} \ \frac{1}{n} \ \dots \ \frac{1}{n}\right)^T$

• Coupling $(X_t, Y_t) \equiv X_0 = x, Y_0 = y$

- At time t , let $W_t = \text{Ber}(1/2)$.

• If $W_t = 1 \Rightarrow X_t = X_{t-1} \mp 1$ wp $1/2$

$W_t = 0 \Rightarrow Y_t = Y_{t-1} \mp 1$ wp $1/2$

- If $X_s = Y_s$ then $X_t = Y_t \ \forall s \geq t$

• Let $D_t =$ clockwise distance betⁿ particles at t ,

- $\tau = \min \{t \geq 0 \mid D_t \in \{0, n\}\} = \tau_{\text{couple}}$

- $\mathbb{E}_{x,y}[\tau] = \underbrace{k(n-k)}$, where $k =$ clockwise dist betⁿ x, y

from gambler's ruin - solves $\mathbb{E}[\tau|k] = 1 + \frac{1}{2}(\mathbb{E}[\tau|k+1] + \mathbb{E}[\tau|k-1])$

$$\Rightarrow \bar{d}(t) \leq \max_{x,y} \mathbb{P}_{x,y}[\tau_{\text{couple}} > t] \leq \frac{\max_{x,y} \mathbb{E}_{x,y}[\tau]}{t}$$

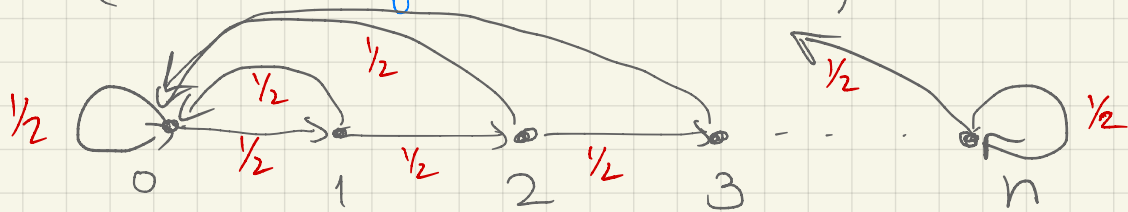
$$\leq \max_{k \in \{0, \dots, n\}} \frac{k(n-k)}{t} = \frac{n^2}{4t}$$

Thus if $t > n^2 \Rightarrow d(t) \leq \bar{d}(t) \leq 1/4$

$$\Rightarrow t_{\text{mix}} \leq n^2$$

□

Eg - (the 'winning streak' chain)



$$X_{t+1} = \begin{cases} \min(X_t + 1, n) & \text{wp } 1/2 \\ 0 & \text{wp } 1/2 \end{cases}$$

• Let $X_0 = x, Y_0 = y, Z_t \sim \text{Ber}(1/2)$

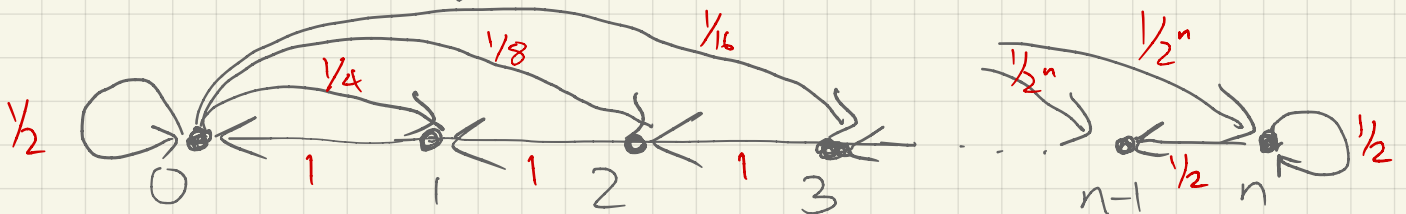
If $Z_t = 1$, then $X_t = (X_{t-1} + 1) \wedge n, Y_t = (Y_{t-1} + 1) \wedge n$

If $Z_t = 0$, then $X_t = Y_t = 0!$

$$\Rightarrow \mathbb{P}[\tau_{\text{couple}} > t] \leq 2^{-t} \Rightarrow t_{\text{mix}} \leq 2 \forall n!$$

• Note that $\pi(i) = 1/2^{i+1}$ if $i \leq n-1, 1/2^n$ if $i = n$

• Next consider the reverse chain \hat{P}



$$\pi(i) = 2^{-((i+1) \wedge n)} \Rightarrow \hat{P}(i, i-1) = 1, \hat{P}(n, n) = \hat{P}(n, n-1) = 1/2 \text{ and } \hat{P}(0, i) = \pi(i) \forall 0 \leq i \leq n$$

• If $X_t = 0 \Rightarrow X_{t+1} \sim \pi$! Now suppose $X_0 = n \Rightarrow$ MC spends $K \sim \text{Ber}(1/2)$ turns at n , then reaches 0 in $n-1$ turns, then mixes.

Also $\hat{P}^n(n, \cdot) = \pi \Rightarrow \pi_t = \pi \forall t \geq n!$

$$\text{Finally } \hat{P}^{n-1}(n, 1) = 1/2 \Rightarrow d(n-1) \geq \|\hat{P}^{n-1}(n, \cdot) - \pi\|_{TV} \geq 1/2 \Rightarrow t_{\text{mix}} = n$$

Eg - Random walk on the hypercube - $X = \{0,1\}^n$ Hamming distance of 1

- $P(x,y) = 1/n$ if x and y differ in 1 bit, else 0
- Lazy RW on hypercube - $P(z,z) = 1/2$, $P(z,y) = 1/2n$ if $d_H(z,y) = 1$
- $\pi(x) = 1/2^n$ ($\because |X| = 2^n$, by symmetry/doubly stoch)
- Alternate description for LRW - Pick index $i \in [n]$ univ, flip bit $X(i)$ w.p. $1/2$
- Coupling $(X_t, Y_t) \equiv \text{Given}(X_t, Y_t)$
 - i) Pick index $I_{t+1} \in [n]$ univ, and $Z_{t+1} = \{0,1\}$ w.p. $1/2$
 - ii) Set $X_{t+1}(I_t) = Y_{t+1}(I_t) = Z_{t+1}$
- Note - $d_H(x_0, y_0) \leq n \Rightarrow T_{\text{couple}} \leq$ 'coupon collector' on n bins
- $\Rightarrow E[T_{\text{couple}}] = n H_n \Rightarrow d(t) \leq E[T_{\text{couple}}]/t \leq n H_n/t$
- $\Rightarrow t_{\text{mix}} \leq 4n H_n$ (can be improved to $1/2 n \ln n$)

(Grand Coupling) - For any MC $X_{t+1} = f(X_t, U_{t+1})$, we can always construct a coupling $(X_{t+1}, Y_{t+1}) = (f(X_t, U_{t+1}), f(Y_t, U_{t+1}))$

- This is a generic way to construct a coupling for any chain! It gives a standard way to measure T_c

$$T_c = \min_{t \geq 0} \left\{ f(\dots f(f(X_0, U_1), U_2) \dots, U_t) = f(\dots f(f(Y_0, U_1), U_2) \dots, U_t) \right\}$$

- Note that this is not always easy to compute. However this gives us another approach to computing mixing times

Strong Stationary Times

- Suppose MC has representation $X_{t+1} = f(X_t, Z_{t+1})$ for some iid sequence Z_t

Def - (Randomized Stopping Time) - A random time \bar{T} for MC X_t is a randomized stopping time if it's a stopping time for Z_t

Def - (Stationary Time) For MC X_t with stationary dist π , a stationary time \bar{T} is a randomized stopping time (possibly dependent on x) s.t. $\mathbb{P}_x[X_{\bar{T}} = y] = \pi(y)$

- Thus, a stationary time is in a sense a signal that a chain has mixed. However, to bound mixing times, we need a slightly stronger definition.

Def (Strong stationary time) - A strong stationary time for MC X_t is a randomized stopping time, possibly dependent on starting state x s.t. $\forall y, t: \mathbb{P}_x[\bar{T} = t, X_{\bar{T}} = y] = \mathbb{P}_x[\bar{T} = t] \pi(y)$

Why the stronger defn? Consider a rw on n -cycle, and the following \bar{T} : w.p. $1/n$, set $\bar{T} = 0$, else set $\bar{T} = \inf\{t \geq 0 \mid \text{every state } x \in X \text{ visited once}\}$. In the latter case, the terminal state is uniform over $X \setminus x_0 \Rightarrow \mathbb{P}[X_{\bar{T}} = x] = 1/n \forall x \in X \setminus x_0$. However $\bar{T} \not\perp X_{\bar{T}}$, as $\bar{T} = 0 \Rightarrow X_{\bar{T}} = x_0$!

• Lemma - For X_t ergodic MC with stationary dist π , if τ is a strong stationary time, then $\forall t \geq 0$

$$P_x[\tau \leq t, X_t = y] = P_x[\tau \leq t] \pi(y)$$

Pf - Let z_1, z_2, \dots be the sequence in the random mapping.

For any $s \leq t$

$$P_x[\tau = s, X_t = y] = \sum_{z \in X} P_x[X_t = y | \tau = s, X_s = z] P_x[\tau = s, X_s = z]$$

$\because \tau$ is a stopping time for Z_t , the event $\{\tau = s\}$ is adapted to $\sigma(z_1, z_2, \dots, z_s)$, and $X_{s+r} = \bar{f}_r(X_s, z_{s+1}, \dots, z_{s+r})$ for some fn \bar{f}_r . Also since $(z_1, \dots, z_s) \perp\!\!\!\perp (z_{s+1}, \dots, z_{s+r}) \Rightarrow$

$$P_x[X_t = y | \tau = s, X_s = z] = P_x[\bar{f}_{t-s}(z, z_{s+1}, \dots, z_t) = y | X_s = z, \sigma(z_1, \dots, z_s)] \\ = P^{t-s}(z, y)$$

\Rightarrow by definition of strong stopping times

$$P_x[X_t = y, \tau = s] = \sum_{z \in X} P^{t-s}(z, y) P_x[\tau = s] \pi(z)$$

Also since $\pi^T P = \pi^T$, we have that the RHS is $\pi(y) P_x[\tau = s]$. Now summing over all $s \leq t$, we get the result \square

- Now to use strong stationary times to bound t_{mix} , we need an additional definition.

Defn (Separation Distance) - For any x, t , given MC P , define

$$S_x(t) = \max_{y \in X} \left[1 - \frac{P^t(x, y)}{\pi(y)} \right], \quad S(t) = \max_{x \in X} S_x(t)$$

Lemma - $\|P^t(x, \cdot) - \pi\|_{TV} \leq S_x(t)$

Pf - $\|P^t(x, \cdot) - \pi\|_{TV} = \sum_{y \in X} (\pi(y) - P^t(x, y))$

$$= \sum_{\substack{y \in X \\ P^t(x, y) < \pi(y)}} \left(1 - \frac{P^t(x, y)}{\pi(y)} \right) \pi(y)$$

$\sum x_i y_i \leq \|x\|_1 \|y\|_\infty$ Hölder's Ineq $\rightarrow \leq \sum_{y \in X} \pi(y) \max_y \left(1 - \frac{P^t(x, y)}{\pi(y)} \right) \leq S_x(t)$

Lemma - If τ is a strong stationary time, then

$$S_x(t) \leq \mathbb{P}_x[\tau > t]$$

Pf - For any $x \in X$, $1 - \frac{P^t(x, y)}{\pi(y)} = 1 - \frac{\mathbb{P}_x[X_t = y]}{\pi(y)} \leq 1 - \frac{\mathbb{P}_x[X_t = y, \tau \leq t]}{\pi(y)}$

by the earlier lemma, $\mathbb{P}_x[X_t = y, \tau \leq t] = \mathbb{P}[\tau \leq t] \pi(y)$ \square

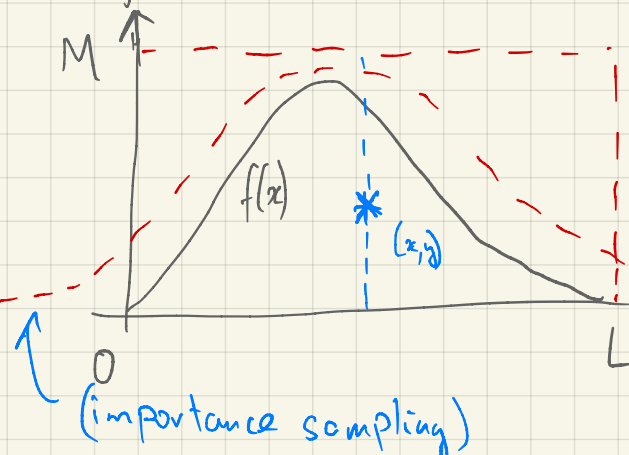
Combining we get $\forall x \in X$ and strong stationary time τ

$$\|P^t(x, \cdot) - \pi\|_{TV} \leq \mathbb{P}_x[\tau > t]$$

Markov chains & Sampling

One of the most important uses of MCs is to sample from complex distributions.

- If we want $X \sim F$ for some known F on \mathbb{R} , then given $U \sim \text{Unif}[0,1]$, can set $X = F^{-1}(U)$ (inversion method)

- If we know pdf f for X , and can bound it in a box $(E_f, \text{ in } [0,L] \times [0,M]; \text{ see figure})$, then we can generate (x,y) uniformly in the box, and set $X = x$ if $y \leq f(x)$, else sample another (x,y) ... (rejection sampling)
- 

- Suppose now we want to sample something more complex:
 - Given $G(V,E)$, sample a spanning tree uniformly at random
 - Sample $x \in X$ for some complex X , proportional to some $g(x)$ (i.e., w.p. $Zg(x)$, where $Z = \int_x g(x) dx$ may be unknown)
 - Sample a random independent set X of $G(V,E)$, proportional to $\lambda^{|X|}$

- **Markov-chain Monte Carlo (MCMC)** is a set of generic techniques that let us do this!

- Idea - Set up MC s.t. $\Pi = \text{target dist}^n$. Run till mixing...

MCMC recipe

• $\Pi \equiv$ target distribution

$P \equiv$ sampling algorithm (typically reversible)

$Q \equiv$ Candidate-generator matrix
(Some given random walk on X)

- Given state $x \in X$, generate candidate y w.p. $Q(x, y)$

- Accept y w.p. $\alpha(x, y) = P(x, y) / Q(x, y)$, else stay at x

- Run till 'MC has mixed sufficiently'.

• We want P s.t. $P(x, y)\Pi(x) = P(y, x)\Pi(y) \forall x, y$

Q s.t. it is ergodic, easy to sample from.

Examples of P -

1) Metropolis Algorithm - $\alpha(x, y) = \min\left(1, \frac{\Pi(y)Q(y, x)}{\Pi(x)Q(x, y)}\right)$

2) Barker's Algorithm - $\alpha(x, y) = \frac{\Pi(y)Q(y, x)}{\Pi(y)Q(y, x) + \Pi(x)Q(x, y)}$

3) Gibb's Algorithm - If we want $\Pi(x^{(1)}, x^{(2)}, \dots, x^{(d)})$ on set $X = \Lambda^d$, Λ countable, we first choose I u.a.r on $[d]$, and set $x^{(I)} \rightarrow y$ w.p. $\Pi(y | x^{(1)}, \dots, x^{(I-1)}, x^{(I+1)}, \dots, x^{(d)})$

Eg - Given X , and 'energy function' (or, ^{inverse} fitness fn)

$h: X \rightarrow \mathbb{R}$, we want $\Pi(x) = e^{-h(x)} / Z$, where

$$Z = \sum_{x \in X} e^{-h(x)} \equiv \text{partition function}$$

- Suppose Q has stationary dist $\Pi(x) = 1/|X|$
(For example, choose Q to be doubly stochastic)

- Metropolis: $\alpha(x, y) = \min[1, e^{-(h(y) - h(x))}]$

$$\Rightarrow P(x, y) \Pi(x) = Q(x, y) \min[1, e^{-(h(y) - h(x))}] \cdot e^{-h(x)}, \quad \frac{Q(x, y)}{|X|} = \frac{Q(y, x)}{|X|}$$
$$P(y, x) \Pi(y) = Q(y, x) \min[1, e^{-(h(x) - h(y))}] \cdot e^{-h(y)}$$

- Barker: $\alpha(x, y) = \frac{e^{-h(y)}}{e^{-h(x)} + e^{-h(y)}}$

$$\Rightarrow P(x, y) \Pi(x) = Q(x, y) \frac{e^{-h(y)}}{e^{-h(x)} + e^{-h(y)}} \cdot e^{-h(x)}, \quad \frac{Q(x, y)}{|X|} = \frac{Q(y, x)}{|X|}$$
$$P(y, x) \Pi(y) = Q(y, x) \frac{e^{-h(x)}}{e^{-h(x)} + e^{-h(y)}} \cdot e^{-h(y)}$$

- For any reversible Q with known stationary distr $\tilde{\Pi}$ we can modify these to get desired Π

$$\leftarrow \frac{Q(x, y)}{Q(y, x)} = \frac{\tilde{\Pi}(y)}{\tilde{\Pi}(x)}$$

Eg (Glauber dynamics / Gibbs sampler)

For undirected graph $G(V, E)$, a **proper q -coloring** is a function

$\chi: V \rightarrow [q]$ (ie, element of $[q]^V$) s.t. $\chi(u) \neq \chi(w)$

$\forall (u, w) \in E$. We also denote $N(u) = \{w \mid (u, w) \in E\}$
(neighborhood of u)

- Aim: Select proper q -coloring u.a.r.

- Gibbs sampler - select $v_t \in V$ u.a.r.

- select color $\bar{j}_t \in [q]$ u.a.r. from

'allowable colors' $A_{v_t}(\chi) = \{j : \chi(w) \neq j \forall w \in N(v_t)\}$

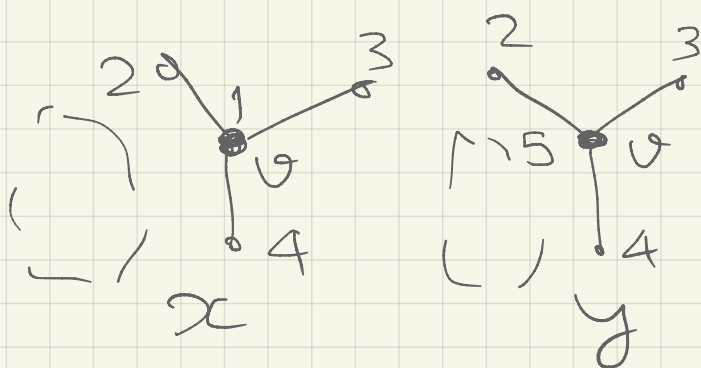
- set $X_{t+1}(v_t) = \bar{j}_t$

$$P(x, y) = \frac{1}{|V|} \cdot \frac{1}{|A_{v_t}(x)|}$$

Note though that \forall 'adjacent' configs $x, y \in [q]^V$
we have $|A_{v_t}(x)| = |A_{v_t}(y)|$.

Thus $P(x, y) = P(y, x) \forall$ adjacent $x, y \Rightarrow \pi(x)$ is uniform

$\hookrightarrow x, y$ are adjacent if $x(u) = y(u) \forall u$ except one



$$q = 7$$

$$A_{v_t}(x) = \{1, 5, 6, 7\}$$

$$A_{v_t}(y) = \{1, 5, 6, 7\}$$

Perfect Sampling

- The MCMC method thus gives us a way to sample any π via a MC. However, how do we know when to stop?

• Strong Doeblin condition: For any MC P on finite X

$$\alpha = \sum_{y \in X} \min_{x \in X} \underbrace{[P(x,y)]}_{\alpha_y}$$

• Now we can write $P = \underbrace{\begin{pmatrix} \alpha_1 & \alpha_2 & \dots & \alpha_n \\ \alpha_1 & \alpha_2 & \dots & \alpha_n \\ \vdots & \vdots & \ddots & \vdots \\ \alpha_1 & \alpha_2 & \dots & \alpha_n \end{pmatrix}}_{\Theta} + (P - \Theta) = \alpha \underbrace{\begin{pmatrix} \alpha_1/\alpha & \alpha_2/\alpha & \dots & \alpha_n/\alpha \\ \vdots & \vdots & \ddots & \vdots \end{pmatrix}}_{I^T \Theta} + (1-\alpha) R$
↑ stochastic

$$\Rightarrow P = \alpha(I^T \Theta) + (1-\alpha)R, \quad R \equiv \text{stochastic}, \quad I^T \Theta = \begin{pmatrix} \alpha_1/\alpha & \alpha_2/\alpha & \dots \\ \vdots & \vdots & \ddots \end{pmatrix}$$

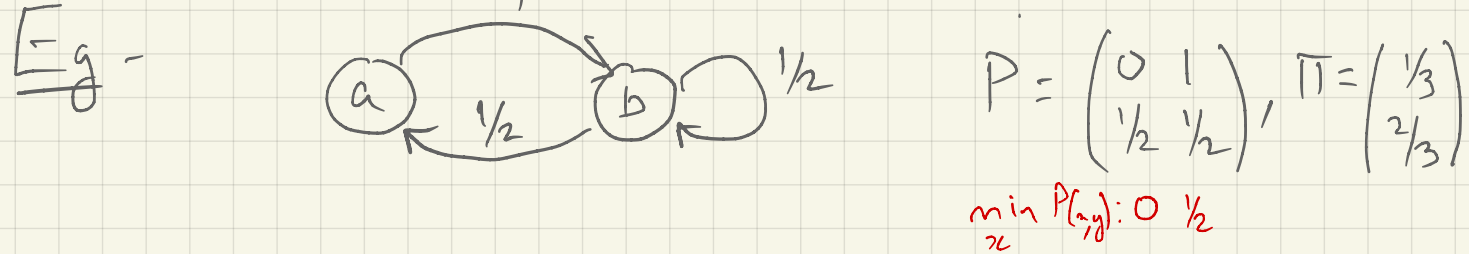
$$\text{Now } \pi^T P = \pi^T \Rightarrow \alpha(\pi^T (I^T \Theta)) + (1-\alpha)\pi^T R = \alpha \theta^T + (1-\alpha)\pi^T R$$

$$\Rightarrow \pi^T = \alpha \theta^T (I - (1-\alpha)R)^{-1} = \sum_{t=0}^{\infty} (1-\alpha)^t \alpha \theta^T R^t \quad (\text{as } |\lambda(R)| \leq 1-\alpha)$$

$$= \mathbb{E}_{N,Y} [R^{N-1}(Y, \cdot)], \quad N \sim \text{Geom}(\alpha), Y \sim \Theta, \text{ independent}$$

- Algo (Strong Doeblin Sampler) - Sample $X_0 \sim \Theta$, $N \sim \text{Geom}(\alpha)$, independent
- Output X_{N-1}

• Thm - $X_{N-1} \sim \pi$ (ie., X_{N-1} is a perfect sample from π)



- For this example - $\Theta = \begin{pmatrix} 0 \\ 1 \end{pmatrix}, \alpha = 1/2,$
- \Rightarrow Set $X_0 = \{b\}, N \sim \text{Geom}(1/2),$ output $X_{N-1} \sim \Pi!$

Eg - PageRank MC (Brin, Page, Motwani, Winograd, J. Kleinberg - HITS)

- Given graph $G(V, E),$ set $\Theta = \mathbf{1}^T / n, \alpha \in (0, 1)$ (Typically 0.2)
- $X_{t+1} = \begin{cases} \text{w.p. } \alpha, \text{ sample } v \in V \text{ u.a.r.} \\ \text{else move to random neighbor of } X_t \end{cases}$

$$P = \alpha \mathbf{1}^T \Theta + (1-\alpha) D^{-1} A$$

\leftarrow Adjacency matrix
 \leftarrow $D = \text{diag}(d_1, \dots, d_n)$
 $W \equiv$ random walk matrix

- Can set any Θ . Eg. $\Theta^T = e_i^T = \mathbb{1}_{\{v=i\}} \equiv$ Personalized PageRank for node i

- Problem - Need $\alpha > 0$, Not true for many MC.
- Idea behind CFTP - Use $X_t = f(X_{t-1}, Z_t)$ and the grand coupling

i.e. - Choose z_1, z_2, \dots s.t. $f(f(\dots f(f(x_0, z_1), z_2), z_3), \dots, z_t))$ is equal for all x_0

$\underbrace{f \circ f \circ \dots \circ f}_{z \text{ times}}(x_0)$

A more convenient way to write the 'random function' representation is as $X_{t+1} = G_t(X_t),$ where $G_t(\cdot)$ is a random function s.t. $G_t \sim \{g_i(\cdot) \text{ w.p. } p_i\}.$

Eg. for $P = \begin{pmatrix} 0 & 1 \\ 1/2 & 1/2 \end{pmatrix}$, let $g_1(x) = \begin{cases} b; & x=a \\ b; & x=b \end{cases}$, $g_2(x) = \begin{cases} b; & x=a \\ a; & x=b \end{cases}$

$X_{t+1} = \begin{cases} g_1(x_t) \text{ w.p. } 1/2 \\ g_2(x_t) \text{ w.p. } 1/2 \end{cases}$ is a random fn representation for X_t

• Now fix a labelling of time $(-\infty, \dots, -2, -1, 0, 1, 2, \dots, \infty)$, with some arbitrary 'current time' 0. Let G_i^j be the composite fn $G_{j+1} \circ G_{j+2} \circ \dots \circ G_i(\cdot)$, i.e., the evolution of the MC from $X_i \rightarrow X_j$. In particular

$$- G_0^t = G_{t-1} \circ G_{t-2} \circ \dots \circ G_1 \circ G_0 \quad (\text{forward simulation})$$

$$G_{-t}^0 = G_{-1} \circ G_{-2} \circ \dots \circ G_{-t+1} \circ G_{-t} \quad (\text{backward simulation})$$

• Similarly we can define two 'coalescence times'

$$(\text{Forward Coalescence}) \quad \tau_F = \inf \{ t \geq 0 \mid G_0^t(x) = G_0^t(y) \forall x, y \in X \}$$

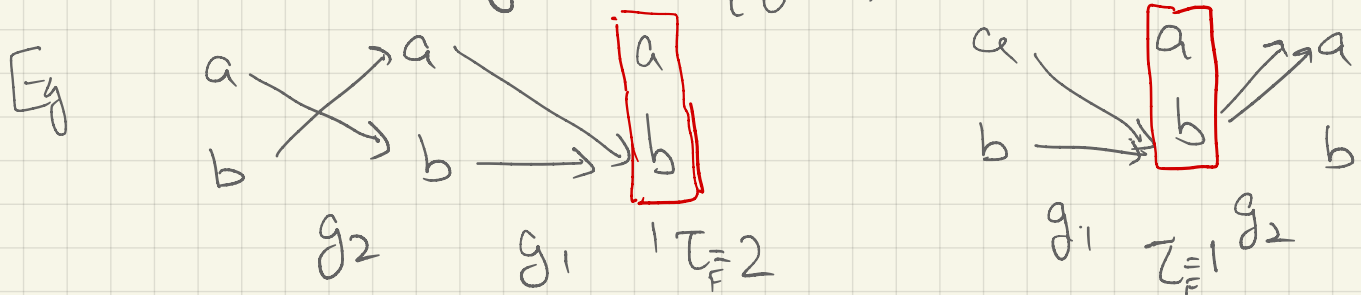
$$(\text{Reverse Coalescence}) \quad \tau_R = \inf \{ t \geq 0 \mid G_{-t}^0(x) = G_{-t}^0(y) \forall x, y \in X \}$$

- In other words, τ_F is the stopping time when all initial states coalesce in the forward simulation; τ_R is the stopping time when all states at time $-\tau_R$ coalesce at 0 in the reverse simulation. We see an example below.

• Note $\tau_F \sim \tau_R$ (i.e., they have the same distrⁿ)

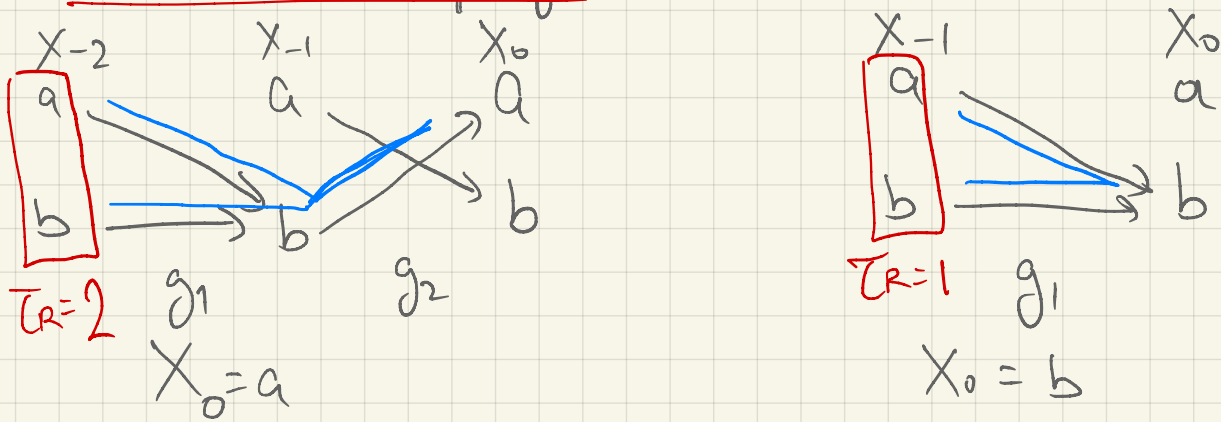
Eg - Consider $P = \begin{pmatrix} 0 & 1 \\ 1/2 & 1/2 \end{pmatrix}$, $g_1(x) = \begin{cases} b & ; x=a \\ b & ; x=b \end{cases}$, $g_2(x) = \begin{cases} b & ; x=a \\ a & ; x=b \end{cases}$

• Forward coupling - $G_t = \begin{cases} g_1 & \text{wp } 1/2 \\ g_2 & \text{wp } 1/2 \end{cases}$ -



• $\tau_F \sim \text{Geom}(1/2)$, $X_{\tau_F} = b$

• Backward coupling



Again $\tau_R \sim \text{Geom}(1/2)$. However, $X_0 = a$; τ_R even
 b ; τ_R odd

$$\Rightarrow \mathbb{P}[X_0 = b] = \frac{1}{2} + \frac{1}{2^3} + \dots = \frac{1/2}{1 - 1/4} = \frac{2}{3} = \pi(b)!$$

Thus we observe -

i) $\tau_F \sim \tau_R$

ii) $G_0^{\tau_F}(x) = b$ (not $\sim \pi$), $G_{-\tau_R}^0(x) \sim \pi$

This procedure to compute $G_{-\tau_R}^0$ is called coupling from the past (CFTP)

Thm (CFTP - Propp & Wilson) - Assuming τ_R is finite w.p.1, then the constant value $Z_{-\infty}^0 = G_{-\tau_R}^0(x)$ has distribution

$$Z_{-\infty}^0 \sim \pi$$

Pf - Let $\tau_R = \inf_{t \geq 0} \{ G_{-t}^0(x) = G_{-t}^0(x') \forall x, x' \in X \}$

$$\tau_R^1 = \inf_{t \geq 0} \{ G_{-t}^1(x) = G_{-t}^1(x') \forall x, x' \in X \}$$

Since τ_R is finite w.p.1 $\Rightarrow \tau_R^1$ is also finite w.p.1, and $Z_{-\infty}^0$ and $Z_{-\infty}^1$ are well defined.

- Now we couple G_{-t}^0 and \hat{G}_{-t}^1 to use the same $G_R \forall k$, i.e.

$$\hat{G}_{-t}^1 = G_0 \circ G_{-1} \circ G_{-2}^0 \dots \circ G_{-t+1} \circ G_{-t}$$

$$G_{-t}^0 = G_{-1} \circ G_{-2}^0 \dots \circ G_{-t+1} \circ G_{-t}$$

- Let $\hat{Z}_{-\infty}^1 = \hat{G}_{-\tau_R^1}^1(x)$, $Z_{-\infty}^0 = G_{-\tau_R}^0(x)$

Then $\hat{Z}_{-\infty}^1 \sim Z_{-\infty}^0$ (by coupling), and $\hat{Z}_{-\infty}^1 = G_0(Z_{-\infty}^0)$

$\therefore G$ is a random fn representation for π , this means $\hat{Z}_{-\infty}^1 \sim Z_{-\infty}^0 \sim \pi$ ($\because \pi$ is the unique distr st if $X \sim \pi$, then $G(x) \sim \pi$). □