

## Information Theory

- Entropy: measure of surprise / amount of information
- Joint Entropy:  $H(X, Y)$  info conveyed by  $X, Y$
- Conditional Entropy:  $H(X|Y)$  info in  $X$  not provided by  $Y$
- Mutual Information:  $I(X; Y)$  info. shared by  $X, Y$ .

$$\text{Equations: } H(X) = \mathbb{E} \left[ \log \frac{1}{P(X)} \right] = \sum_x P(X=x) \log \frac{1}{P(X=x)}$$

$$H(X, Y) = \mathbb{E} \left[ \log \frac{1}{P_{X,Y}(x,y)} \right] = \sum_{x,y} P(X=x, Y=y) \log \frac{1}{P_{X,Y}(x,y)}$$

$$H(X|Y) = \mathbb{E} \left[ \log \frac{1}{P_{X|Y}(x|y)} \right] = \sum_y P_Y(y) H(X|Y=y)$$

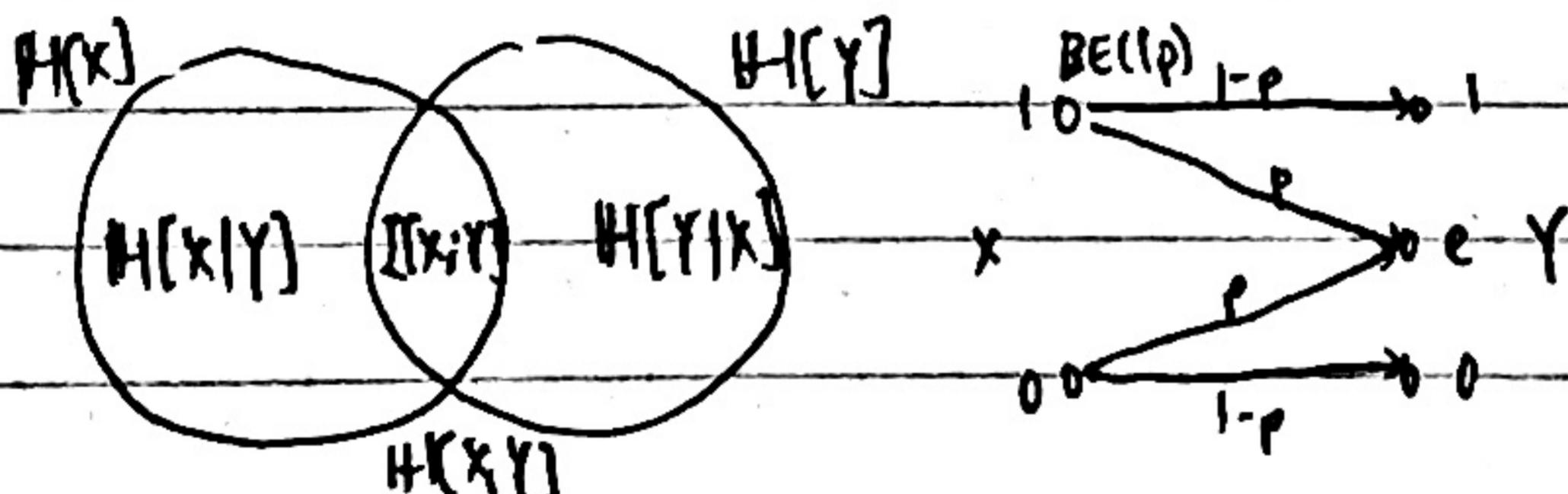
$$= \sum_y P_Y(y) \sum_x P_{X|Y}(x|y) \log \frac{1}{P_{X|Y}(x|y)}$$

$$I(X; Y) = H(X) - H(X|Y) = H(X) + H(Y) - H(X, Y) = I(Y; X)$$

$$J(X; Y) = - \sum_x \sum_y P_{X,Y}(x,y) \log \frac{P_X(x) P_Y(y)}{P_{X,Y}(x,y)}$$

## Properties:

- $H(X) \geq 0$ ;  $H$  concave (i.e.  $-H$  convex)
- $H(X)$  represents minimum number of bits per character to represent a file.  $H(X) \leq \log_2 |X|$  (maximized at dist)
- $\log X$  and  $\pi \log \frac{1}{\pi}$  are both concave (via Jensen's)
- $H(X, Y) \leq H(X) + H(Y)$  with equality iff  $X, Y$  independent
- $H(X, Y) \geq \min(H(X), H(Y))$
- [Chain Rule]  $H(X, Y) = H(Y) + H(X|Y) = H(X) + H(Y|X)$
- $H(Y|X) \leq H(Y)$  with equality iff  $X, Y$  independent.
- $I(X; Y) \geq 0$  with equality iff  $X, Y$  independent.



## Fountain Code

- Encoding
  - Divide data into  $n$  chunks. Pick  $1 \leq d \leq n$  according to a distribution (degree of packet)
  - Select  $d$  random chunks of data, combine binary representations of them using XOR ( $y_i = x_1, x_2, \dots, x_d$ )
  - Transmit:  $y_i$
- Decoding
  - If degree of  $y_i \geq 1$ , assign to  $x_i$ . Else, XOR known chunks.
  - Wtch for info, peel off w/ unfinished packets... repeat.

(def  $f_k(d)$  = expected number of packets of degree  $d \geq 1$  when  $k$  chunks have been recovered.)

$$f_k(d) = \frac{n-k}{d(d-1)} \quad p(d) = \begin{cases} \frac{1}{n} & d=1 \\ \frac{1}{d(d-1)} & d=2, \dots, n \end{cases}$$

## Discrete Time Markov Chain (continued)

- For irreducible MC, some state positive recurrent  $\Leftrightarrow$  all states positive recurrent. Expected return time  $= \frac{1}{\pi_i}$ . Same for null recurrent.

## Source Coding and compression

Source Coding Theorem: Let  $(X_n)_n$  be iid and  $\varepsilon > 0$ . A source coding scheme s.t.  $\lim_{n \rightarrow \infty} \mathbb{E}[\underline{l}(X_1, \dots, X_n)] \leq H(X) + \varepsilon$  where  $\underline{l}(X_1, \dots, X_n)$  denotes length of binary encryption of  $(X_n)_n$ .  $(X_n)_n$  can be recovered from encoding with probability  $1-\varepsilon$ . Impossible to compress source  $X_1, \dots, X_n$  below its entropy  $H(X_1, \dots, X_n) = nH(X)$  bits, but possible for  $n(H(X) + \varepsilon)$ .

Typical Set:  $A_{\varepsilon}^{(n)}$  is the set of sequences in  $\mathcal{X}^n$  that fits the condition  $2^{-n(H(X) + \varepsilon)} < P(X=X_1, \dots, X_n \in A_{\varepsilon}^{(n)}) < 2^{-n(H(X) - \varepsilon)}$

Size of this set is  $\approx 2^{nH(X)} \ll |\mathcal{X}^n| = |\mathcal{X}|^n$

$\Rightarrow P((X_1, \dots, X_n) \in A_{\varepsilon}^{(n)}) > 1-\varepsilon$  for  $n$  large.

$\Rightarrow (1-\varepsilon) 2^{n(H(X) - \varepsilon)} \leq |A_{\varepsilon}^{(n)}| \leq 2^{n(H(X) + \varepsilon)}$  for  $n$  large

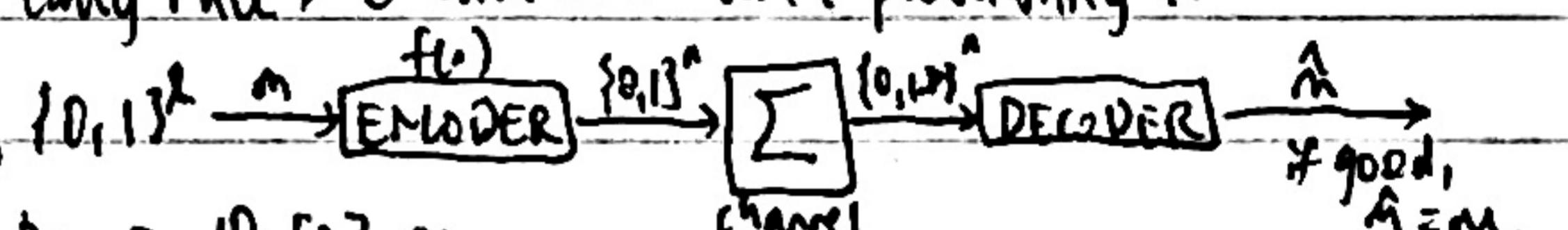
Vershov: Only need to care about encoding of things in  $A_{\varepsilon}^{(n)}$

• anything else can be handled one-off  $\stackrel{H(X)}{\text{Huffman Coding}}$ : Expected average length between  $H(X)$  and

channels and capacity

Channel Coding Theorem: Any rate below the channel capacity  $C = \max_{p \in P, x \in \mathcal{X}} I(X; Y)$  is achievable. Conversely, any sequence of codes with  $P_e(n) \rightarrow 0$  as  $n \rightarrow \infty$  has rate  $R \leq C$ .

(any rate  $> C$  with error with probability 1).



Denote  $P_{e,l}$  as probability of error given  $l$  bits sent. i.e.  $P_{e,l} = \max_{m \in \{0,1\}^l} P_{e,m}$ . Rate =  $\frac{l}{n}$ . Say it is achievable if  $P_{e,l} \rightarrow 0$  as  $l \rightarrow \infty$

Define capacity = largest achievable rate.

Theorem: Capacity ( $BEC(p)$ ) =  $1-p$ . i.e.  $1-p$  bits

Key Idea:  $\frac{1}{2^n} \sum_{i=1}^n \frac{c_i}{c_i'} = \frac{1}{2^n} \sum_{i=1}^n \frac{1}{c_i'} = \frac{1}{2^n} \sum_{i=1}^n P_{e,i}$  (max. last  $n$  bits erased, reliably)

$$P_e = P(c_1' = c_1 \text{ OR } c_2' = c_2 \text{ OR } \dots) \leq \sum_{i=1}^n P(c_i' = c_i) \leq 2^n \cdot 2^{-n(1-p)} = 2^{nR} \cdot 2^{-n(1-p)} = 2^{-n(1-p)}$$

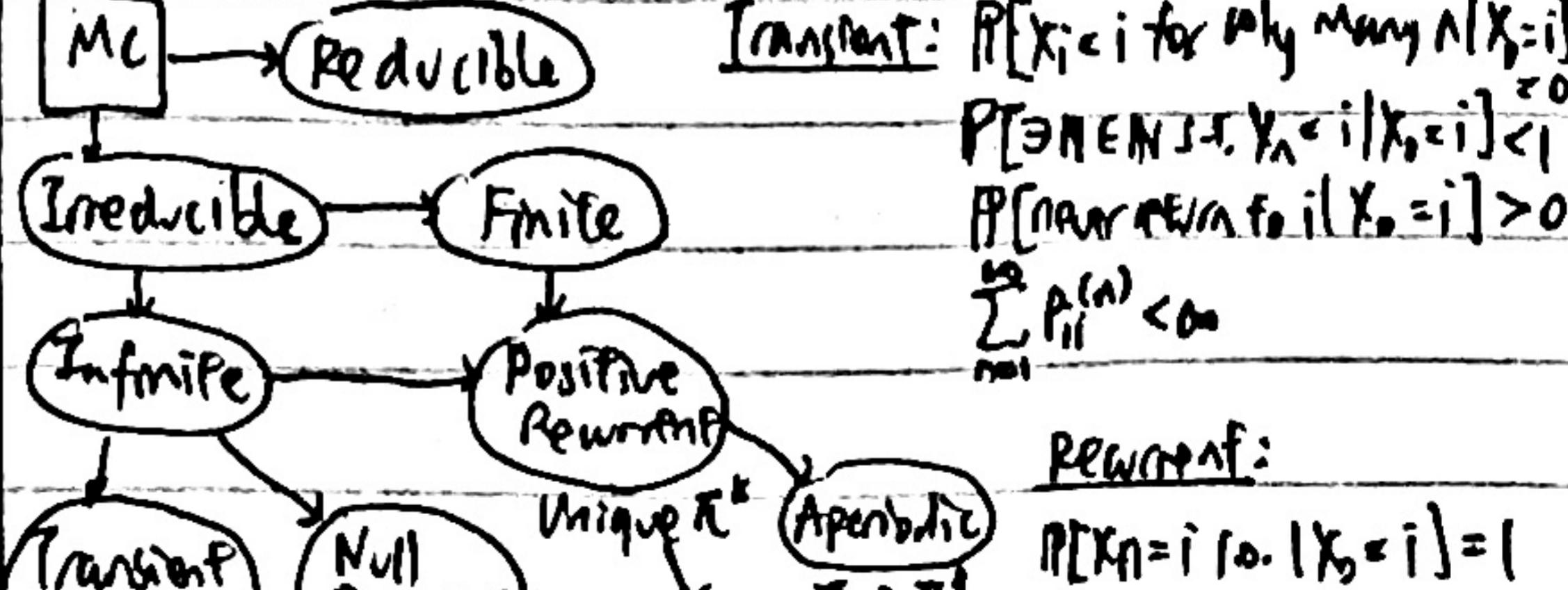
If  $R = 1-p - \varepsilon$ , then  $P_e < 2^{-n\varepsilon}$

$\mathbb{E}[l]$ : length of message in bits that you want to send over.

$\Lambda$ : Number of channel use. Bound =  $2^{\Lambda} \cdot 2^{-\Lambda(1-p)}$ .

$$(=\max_{p \in P} I(X; Y) = \max_{p \in P} [H(X) - H(X|Y)])$$

## Discrete Time Markov Chain (DTMC)



Recurrence:

$$P[X_n=i \text{ for all } X_0=i] = 1 \quad P[\exists N \text{ s.t. } X_n=i \mid X_0=i] > 0$$

$$\sum_{n=1}^{\infty} P_{ii}^{(n)} < \infty \quad P[\text{never returns to } i \mid X_0=i] = 1$$

Transience:

$$P[X_n=i \text{ for all } X_0=i] = 0 \quad P[\exists N \text{ s.t. } X_n=i] = 1$$

Convergence: Not guaranteed

$$\sum_{n=1}^{\infty} P_{ii}^{(n)} = \infty$$

Parts:

- Irreducible, possibly as MC is positive recurrent iff  $\pi \neq e$  iffs
- State  $i$  is null recurrent if  $i$  is transient and expected return time
- State  $i$  is positive recurrent if  $i$  is recurrent and expected return time  $\pi_i < \infty$

## Continuous Time Markov Chains (CTMC)

- Properties:  $\pi Q = 0$ ;  $\sum \pi_i = 1$ ;  $Q \in \mathbb{R}^{n \times n}$ ,  $\sum q_{ij} = 0 \forall i$   
 Holding rate of state  $i = q_i = \sum_{j \neq i} q_{ij} \Rightarrow$  Transition time  $\sim \text{Exp}(q_i)$   
 $\Rightarrow$  Mean transition time  $= \frac{1}{q_i}$ . Jump chain stationary distribution  $\pi$ .
- Differential Analysis: For small time interval  $\delta$  of state  $i$ ,
- $$\mathbb{P}[0 \text{ transitions}] = 1 - q_i \delta + o(\delta), \quad \mathbb{P}[1 \text{ transition}] = q_i \delta + o(\delta)$$
- $$\mathbb{P}[\geq 2 \text{ transitions}] = o(\delta) \quad \mathbb{P}[\text{transition to } j (\neq i)] \approx q_{ij} \delta + o(\delta)$$

### Relationship to Jump chain / DTMC

- $\pi_i =$  long term fraction of time CTMC spends in state  $i$
- $p_i =$  long term fraction of number of transitions leading to state  $i$
- $\pi_i q_{ij} =$  expected number of transitions from  $i$  to  $j$  per unit time.
- $p$  may exist even if  $\pi$  does not exist.

$$\pi_i = \frac{p_i}{\sum p_i} \quad p_i = \frac{\pi_i q_i}{\sum_{j \neq i} \pi_j q_{ij}} \quad \therefore [\pi_i q_i] \text{ eigenvector of } Q \text{ with eigenvalue } 0$$

- DTMC with same stationary distribution  $\pi$ 
  - Find max holding time  $q_{\max}$
  - Divide all edges by  $q_{\max}$  in CTMC and add self-loops

### Jump Chain

$$p_{ij} = \frac{q_{ij}}{q_i}, \quad p_{ii} = 0 \quad (\text{i.e. no self loops})$$

Same CTMC can correspond to different sets of jump chains and holding times if self-loop condition is used.

### Equations

$$\pi_j \sum_{k \neq j} q_{ik} = \sum_{k \neq j} \pi_k q_{kj} \quad \sum_{k \neq j} \pi_k = 1 \quad \pi_j = 0 \text{ for transient states}$$

Flow out      Flow in

$$\pi_j \geq 0 \text{ for positive reward}$$

$$\pi_i q_{ij} = \pi_j q_{ji} \quad (\text{CTMC reversible} \Leftrightarrow \text{jump chain reversible})$$

$$\mathbb{E}[S_i] = \frac{1}{q_i} + \sum_j \frac{q_{ij}}{q_i} \mathbb{E}[S_j] \quad (\text{first step equations})$$

$$\mathbb{P}[X(t) = j | X(s) = i] = \mathbb{P}[X(t-s) = j | X(0) = i] \quad \forall s < t, i, j \in \mathcal{X}. \quad (\text{Time homogeneous})$$

### Poisson process PP( $\lambda$ )

Notation:  $N(t) = N([0, t])$ ,  $N([a, b]) = \# \text{ arrivals in } [a, b]$ ,  $N(0) = 0$   
 Arrival times  $T_n \sim \text{Erlang}(n; \lambda)$ ; interarrival times  $S_n \sim \text{Exp}(\lambda)$  i.i.d  
 $N(a, b) \sim \text{Poisson}(\lambda(b-a))$  time homogeneous, memoryless  
 independence of disjoint intervals.

### Differential Analysis

For small time interval  $\delta$ ,  $\mathbb{P}[0 \text{ arrival}] = 1 - \lambda \delta + o(\delta)$

$$\mathbb{P}[1 \text{ arrival}] = \lambda \delta + o(\delta)$$

### Facts & Techniques

- Poisson process runs backwards & still Poisson
- Given  $N(t)$ , the  $N(t)$  arrivals are uniformly distributed in  $[0, t]$  same order regardless as  $n$  i.i.d. uniform  $(0, t)$  R.V.s
- Poisson splitting (splitting processes independent)
- Poisson merging (arrivals from each source is  $\sum_{i=1}^n$ )
- Reduce to first arrival time; interarrival time, n arrivals in fixed interval  $\Rightarrow$  use binomial/geometric distribution.
- Random incidence paradox

$$f_{T_1, \dots, T_n}(t_1, \dots, t_n) = \lambda^n e^{-\lambda t_n}$$

$$N(\lambda) = \sum_{i=1}^n N(\lambda_i) - N(\lambda) = \sum_{i=1}^n U_i \text{ where } U_i \sim \text{Poisson}(\lambda_i)$$

$$\therefore \frac{N(\lambda)}{n} \xrightarrow{n \rightarrow \infty} \lambda$$

$S_N = X_1 + \dots + X_N$  where  $X_i \sim \text{Exp}(\lambda_i)$ ,  $X_i \sim \text{Exp}(\lambda)$  i.i.d  $\Rightarrow S_N \sim \text{Exp}(\lambda \sum_{i=1}^n \lambda_i)$ .

$$X \sim \text{Poisson}(\lambda), Y \sim \text{Poisson}(\mu), X+Y \sim \text{Binomial}(n, \frac{\lambda}{\lambda+\mu})$$

$$\mathbb{P}[T \leq t + \varepsilon | T > t] = \lambda \varepsilon + o(\varepsilon)$$

## Metropolis Hastings Algorithm

Target distribution  $p(x)$ . Want to create an aperiodic Markov chain with stationary distribution  $\pi(x) = p(x)$ .  $\Rightarrow$  allows generation of random sample even if  $p$  is not

### Assumptions:

- Can compute  $f(x)$ : a direct proportional estimate of  $p(x)$   
 i.e.  $p(x) = \frac{f(x)}{\sum f(y)}$  ( $x$  is current state)
- Can compute  $g(x, \cdot)$ : a proposal distribution for next state

Algorithm: For each time step  $t$ ,

- Propose next candidate state  $y$  according to  $g(x, \cdot)$
- Accept  $y$  with probability  $A(x, y) = \min(1, \frac{f(y)g(x, y)}{f(x)g(y, x)})$
- If accept, move to  $y$ , else stay in  $x$ .

### Related Extensions:

- Stationary distribution of MC equals  $p(x)$
- Taking every  $k$ th step space reduces dependence between samples
- Letting chain propagate a bit lets distribution converge to stationary
- Any MH MC can be made aperiodic by giving each state a self-loop with same probability.
- Burn-in time: how long initial distribution converges to stationary distribution.
- Disadvantage for bimodal distributions

### Convergence

General procedure: • Markov/Chobyshev; reduce to WLLN/SLLN  
 • Convergence in distribution (using MGF)

Theorems/Traits: •  $y \xrightarrow{P} y, z \xrightarrow{P} z \Rightarrow y+z \xrightarrow{P} y+z$

- $(X_n)_n$  be nonnegative and integer-valued,  $X \neq 0$   
 $\text{Then } (X_n)_n \xrightarrow{P} X \text{ iff } \lim_n \mathbb{P}[X_n = k] = \mathbb{P}[X = k] \text{ for } k = 0, 1, 2, \dots$
- $|X| \leq |X - \mathbb{E}(X)| + |\mathbb{E}(X)|$  ( $\epsilon$ -inequality)
- If  $(X_n)_n \xrightarrow{d} c$  where  $c$  constant, then  $(X_n)_n \xrightarrow{P} c$ .
- If  $\forall \varepsilon > 0$ ,  $\sum_m \mathbb{P}[|X_n - X| > \varepsilon] < \infty$ , then  $(X_n)_n \xrightarrow{a.s.} X$
- For  $\varepsilon > 0$ , define  $A_m = \{|X_n - X| < \varepsilon \forall n \geq m\}$ . Then  $(X_n)_n \xrightarrow{a.s.} X$  if and only if  $\forall \varepsilon > 0$   $\lim_m \mathbb{P}[A_m] = 1$ .
- If  $(X_n)_n, (Y_n)_n$  defined on same sample space and  $(X_n)_n \xrightarrow{a.s.} X$  and  $(Y_n)_n \xrightarrow{a.s.} Y$  then  $(X_n + Y_n)_n \xrightarrow{a.s.} X + Y$ .
- If  $(X_n)_n$  and  $(Y_n)_n$  defined on same sample space and  $(X_n)_n \xrightarrow{P} X$  and  $(Y_n)_n \xrightarrow{P} Y$  then  $(X_n + Y_n)_n \xrightarrow{P} X + Y$ .

### Continuous Mapping Theorem:

- Let  $f$  be continuous function.
- If  $(X_n)_n \xrightarrow{a.s.} X$  then  $(f(X_n))_n \xrightarrow{a.s.} f(X)$
- If  $(X_n)_n \xrightarrow{P} X$  then  $(f(X_n))_n \xrightarrow{P} f(X)$
- If  $(X_n)_n \xrightarrow{d} X$  then  $(f(X_n))_n \xrightarrow{d} f(X)$

### Borel-Cantelli Lemma:

- Let  $A_1, A_2, \dots$  be sequence of events in probability space
- If  $\sum_{i=1}^{\infty} \mathbb{P}[A_i] < \infty$  then  $\mathbb{P}[A_i \text{ i.o.}] = 0$ .
- If  $\sum_{i=1}^{\infty} \mathbb{P}[A_i] = \infty$  and  $(A_n)_n$  independent,  $\mathbb{P}[A_i \text{ i.o.}] = 1$ .

Convergence Theorem: Suppose  $(X_n)_n \xrightarrow{a.s.} X$

- If  $0 \leq x_1 \leq \dots \leq x_n \leq \dots$  then  $(\mathbb{E}[X_n])_n \rightarrow \mathbb{E}[X]$  (Monotone)
- If  $\exists$  R.V.  $Y \geq 0$  with  $\mathbb{E}[Y] < \infty$  and  $|X_n|, |Y| \leq Y \forall n$   
 $\text{then } (\mathbb{E}[X_n])_n \rightarrow \mathbb{E}[X]$ . (Dominated)

### Classics

Gambler's ruin:  $P = \frac{1}{2} \Rightarrow \mathbb{P}[W_n] = \frac{n}{n+1}, \mathbb{E}[\text{stopping}] = n_1 \cdot n_2$

$$\xrightarrow{n_1 \quad n_2} P \neq \frac{1}{2} \Rightarrow \mathbb{P}[W_n] = \frac{1 - (\frac{n_2}{n_1})^n}{1 - (\frac{n_2}{n_1})^{n_1+n_2}}$$

$$\mathbb{E}[\text{stopping}] = \frac{n_1 \cdot n_2}{1 - (\frac{n_2}{n_1})^{n_1+n_2}} = \frac{n_1 \cdot n_2}{1 - (\frac{n_2}{n_1})^{n_1+n_2}} = \frac{n_1 \cdot n_2}{1 - (\frac{n_2}{n_1})^{n_1+n_2}} = \frac{n_1 \cdot n_2}{1 - (\frac{n_2}{n_1})^{n_1+n_2}}$$

$$H(p) = -p \log p + (1-p) \log \frac{1}{1-p}$$

$$H(p) = -p \log p + (1-p) \log \frac{1}{1-p}$$