

LOWER BOUNDS on SAMPLE COMPLEXITY.

LECTURE 12
28/3/2021

THEOREM.

Let H be any hypothesis class over a domain X with a finite VC-dim d .

For every $\epsilon, \delta \in (0, 1)$, for any sampling distribution D and any true labelling $f: X \rightarrow \{+1, -1\}$,

$$\mathbb{P}_{S \sim D^n} \left[\exists h \in H : \mathcal{E}_S(h, f) = 0 \wedge \mathcal{E}_D(h, f) > \epsilon \right] \leq \delta$$

whenever.

$$n \geq \frac{8d}{\epsilon} \log\left(\frac{8e}{\epsilon}\right) + \frac{4}{\epsilon} \log\left(\frac{2}{\delta}\right).$$

DEFN.

A hypothesis class H over a domain X is said to be PAC-learnable with a sample complexity $s_n: (0, 1)^2 \rightarrow \mathbb{N}$ if

for any $\epsilon, \delta \in (0, 1)$,

for any prob. dist D on X , and

any true labelling $f: X \rightarrow \{+1, -1\}$,

$$\mathbb{P}_{S \sim D^n} \left[\exists h \in H : \mathcal{E}_S(h, f) = 0 \wedge \mathcal{E}_D(h, f) > \epsilon \right] \leq \delta$$

whenever.

$$n \geq s_n(\epsilon, \delta).$$

Note: Some definitions insist that $s_n(\epsilon, \delta)$ is polynomially bounded in $1/\epsilon$ and $1/\delta$.

THEOREM (Equivalently)

A hypothesis class H with a finite VC-dimension d is PAC-learnable with sample complexity

$$S(\epsilon, d) \leq \frac{8d}{\epsilon} \log\left(\frac{8d}{\epsilon}\right) + \frac{4}{\epsilon} \log\left(\frac{2}{\epsilon}\right) \\ \in O\left(\frac{1}{\epsilon} \left(d \log \frac{1}{\epsilon} + \log \frac{1}{\epsilon}\right)\right)$$

$$\begin{array}{l} \text{VC-dim}(H) < \infty \\ \Rightarrow \\ H \text{ is PAC} \\ \text{learnable} \end{array}$$

CONVERSE

$$\text{VC-dim}(H) = \infty \stackrel{?}{\Rightarrow} \text{Not PAC-learnable}$$

TIGHTNESS

For VC-dim finite, can we reduce the sample complexity?

$VC \text{ Dim}(H) = \infty \Rightarrow \text{NOT PAC learnable}$

Given \mathcal{H} . $\forall \mathcal{D}, \forall f, \forall \epsilon, \forall \delta$
Demer (X)

Ans: Yes
 $VC\text{-dim}(H)$ is infinite
 $\Leftrightarrow H$ is PAC-learnable.

CLAIM:
 $S_H(\frac{\epsilon}{2}, \frac{\delta}{2}) \geq \frac{1}{2} VC \text{ Dim}(H)$

Given \mathcal{H} : Any hypothesis class of $VC\text{-dim } d$

Given X : Domain of \mathcal{H}
 $f(x) = +1 \forall x \in X$

$T = \{x_1, \dots, x_d\} \subseteq X$: A set shattered by \mathcal{H}
 $|T| = d$

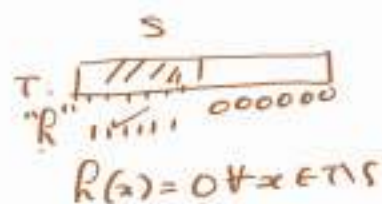
Clever Choice:

$$\mathcal{D} : P(x) = \begin{cases} 1/d & \forall x \in T \\ 0 & \forall x \notin T \end{cases}$$

$\epsilon, \delta = 1/2$

Let $n < d/2$ and $S \sim \mathcal{D}^n$

\Rightarrow At least half the points in T are not in S Training data

S
 T : 
 $R(x) = 0 \forall x \in T \cap S$

$$\Rightarrow \exists h : \mathcal{E}_S(h, f) = 0 \wedge \mathcal{E}_D(h, f) > 1/2$$

$$\Rightarrow P_{S \sim \mathcal{D}^n} [\exists h : \mathcal{E}_S(h, f) = 0 \wedge \mathcal{E}_D(h, f) > \epsilon] = 1 > \delta = 1/2$$

Hence $S_H(1/2, 1/2) \geq d/2$

(Note: S_H is decreasing in ϵ and δ
 so it is no better for smaller ϵ and δ)

Can we have a smarter algorithm?

Let it be the smart algo

$$h^* = A(S) \quad (\text{Deterministic})$$

"Best choice" among

$$\{h \in H : \mathcal{E}_S(h, f) = 0\}$$

We will beat it with same D but different f 's.

$$\text{Let } f(x) = \begin{cases} +1, & x \in X \setminus T \\ \begin{cases} +1, & \text{w.p. } 1/2 \\ -1, & \text{w.p. } 1/2 \end{cases}, & x \in T. \end{cases}$$

("Probabilistic Method")

$$f \sim F$$

$$n < d/2$$

For any set S of at most n elements from T (At least $1/2|T|$ in outside S)

$$\mathbb{E}_{f \sim F} [\mathcal{E}_D(A(S), f)] = \frac{1}{2} |T \setminus S| \geq \frac{1}{4} |T|$$

Hence $\exists f$ s.t.

$$\forall S \quad \mathcal{E}_D(A(S), f) > \frac{1}{4} \cdot \epsilon$$

$$\Rightarrow \mathbb{P}_{S \sim D^n} [\mathcal{E}_D(A(S), f) > 1/4] = 1$$

Can a randomised algorithm work?

Ans: No.

FINITE VC-dimension

- Smaller Sample Complexity?

Final Answer:

$$S_H(\epsilon, f) = \Omega\left(\frac{d}{\epsilon} \log \frac{1}{\epsilon} + \frac{1}{\epsilon} \log \frac{1}{\delta}\right)$$

so tight up to constants.

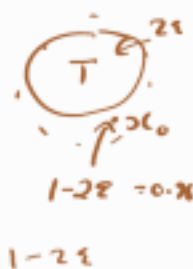
$\log 1/\epsilon$ term can be removed if we make "smart algorithms"

Idea: $S_H(\epsilon, f) \geq \Omega(d/\epsilon)$ $S_H(1/n, 1/2) \geq d/2$

$T = \{x_1, \dots, x_d\}$ shattered by H

$x_0 \in X \setminus T$

$$\mathcal{D} : P(x) = \begin{cases} 0, & x \in X \setminus (T \cup \{x_0\}) \\ 1-2\epsilon, & x = x_0 \\ \frac{2\epsilon}{d}, & x \in T. \end{cases}$$



Let $n < d/8\epsilon$

$$\begin{aligned} \mathbb{E}_{S \sim \mathcal{D}^n} [|S \cap T|] &\approx n P_{x \sim \mathcal{D}} [x \in T] \\ &= n \cdot 2\epsilon \\ &< \frac{d}{8\epsilon} \cdot 2\epsilon \\ &= d/4. \end{aligned}$$

"linearity of expectation"

$x > 0$

$$\mathbb{E}\{x > k \mathbb{E}x\} \leq \frac{1}{k}$$

$$\therefore P_{S \sim \mathcal{D}^n} [|S \cap T| \geq \frac{d}{2}] \leq 1/2 \quad (\text{Markov Inequality})$$

$$P_{S \sim \mathcal{D}^n} [|S \cap T| < \frac{d}{2}] > 1/2$$

$$|S \cap T| < \frac{d}{2} \Rightarrow \exists h : \epsilon_S(h, f) = 0 \wedge \epsilon_D(h, f) > \epsilon$$

$$\text{Hence } P_{S \sim \mathcal{D}^n} [\exists h : \epsilon_S(h, f) = 0 \wedge \epsilon_D(h, f) > \epsilon] > 1/2$$

$$\text{i.e. } S_H(\varepsilon, 1/2) \geq d/8\varepsilon \\ = \Omega(d/\varepsilon).$$

$$\frac{d}{8\varepsilon} \leq S_H(\varepsilon, 1/2) \leq \frac{8d \log(1/\varepsilon)}{\varepsilon} + \frac{4}{\varepsilon} \log\left(\frac{2}{\varepsilon}\right) \\ = \frac{8d \log(1/\varepsilon)}{\varepsilon} + \frac{8}{\varepsilon}$$

↳ "64-factor gap".

Rule of thumb: d/ε

Making full use of PAC learning

We proved that

H is PAC-learnable

if and only if

VC-dimension of H is finite

What more can one ask for?

1. Can we tolerate a non-zero in-sample error?

$$\exists h \in H : (\hat{L}_S(h, f) = 0) \wedge (\hat{L}_D(h, f) > \epsilon)$$

2. Can we use a combination of simple (small VC-dim) classifiers?

3. Anything b/n d and $d+1$?

4. Boosting

Non-Zero In-Sample Error

Goal:

$$P_{S \sim D^n} \left[\underbrace{\exists h \in H : |\mathcal{E}_D(h, f) - \mathcal{E}_S(h, f)|}_{C_h} > \epsilon \right] \leq \delta$$

H : hypothesis class
 X : domain of H
 f : true labeling
 D : sampling dist

$$= P_{S \sim D^n} \left[\bigcup_{h \in H} C_h \right] \leq \delta, \text{ where}$$

$$C_h = \left\{ \left| \mathcal{E}_D(h, f) - \mathcal{E}_S(h, f) \right| > \epsilon \right\}$$

$$\exists h \in H : \mathcal{E}_S(h, f) = 0 \wedge \mathcal{E}_D(h, f) > \epsilon$$

"this leads"

Those training data (S) which permit even one $h \in H$ which is a very different in-sample and out-of-sample error.



Step 1. Fix an arbitrary $h \in H$.

let $\mu = \mathcal{E}_D[h, f]$ (a deterministic quantity)

$X = \mathcal{E}_S[h, f]$ (random variable)

$$= \frac{1}{n} \left| \{x \in S : h(x) \neq f(x)\} \right|$$

$$= \frac{1}{n} (X_1 + X_2 + \dots + X_n),$$

where $X_i = \begin{cases} 1 & \text{if } h(x_i) \neq f(x_i) \\ 0 & \text{o/w} \end{cases}$

($S = \{x_1, \dots, x_n\}$)

So X is the average of n ^{independent} Bernoulli random variables.

Now our $\mathcal{E}X = \frac{1}{n} \sum \mathcal{E}X_i$ (linearity of expectation)

$$= \frac{1}{n} \sum_{x_i \sim D} P[h(x_i) \neq f(x_i)]$$

$$\boxed{\mathcal{E}X = \mathcal{E}_D(h, f)} = \frac{1}{n} \sum \mathcal{E}_D(h, f)$$

$$= \frac{1}{n} \sum \mu$$

$$= \mu //$$

if X is a $\{0, 1\}$ -R.V., then

$$\mathcal{E}X = P(X_i = 1) \times 1 + P(X_i = 0) \times 0$$

$$= P(X_i = 1)$$

$$= p$$

Hence

$$C_k = |X - \mu_i| > \varepsilon$$

C_k : Event that X deviates from its mean by more than ε .

Hoeffding Bounds.

Let X_1, X_2, \dots, X_n be independent $\{0,1\}$ -random variables with $P(X_i=1) = p$ ($\forall i$).

Let $X = \frac{1}{n}(X_1 + \dots + X_n)$. Then

$$EX = \frac{1}{n} \sum EX_i = p$$

$$(a) P(X > p + \varepsilon) \leq e^{-2\varepsilon^2 n}$$

$$(b) P(X < p - \varepsilon) \leq e^{-2\varepsilon^2 n}$$

$$P\{|X - p| > \varepsilon\} \leq 2e^{-2\varepsilon^2 n}$$

Hence $P_{S \sim D^n}[C_k] \leq 2e^{-2\varepsilon^2 n}$ (Exercise: Finite \mathcal{H})

Recall: $P_{S \sim D^n}[A_k] \leq (1 - \varepsilon)^n \leq e^{-\varepsilon n}$

$$A_k = (\mathcal{E}_S(h, f) = 0) \wedge (\mathcal{E}_D(h, f) > \varepsilon)$$

$$1 + x \leq e^x$$

$$A_k = (\mathcal{E}_S(h, f) = 0) \wedge (\mathcal{E}_D(h, f) > \varepsilon)$$

$$\mathcal{B} = \bigcup_{k \in \mathcal{H}} B_k$$

$$B_k = (\mathcal{E}_S(h, f) = 0) \wedge (\mathcal{E}_{S'}(h, f) > \varepsilon/2)$$

$$P(B_k/A_k) \geq \frac{1}{2}$$

$$C_k = |\mathcal{E}_D(h, f) - \mathcal{E}_S(h, f)| > \varepsilon$$

$$P(\mathcal{B}) \leq \delta/L \Rightarrow P(A) \leq \delta$$

$$D_k = |\mathcal{E}_{S'}(h, f) - \mathcal{E}_S(h, f)| > \varepsilon/2$$

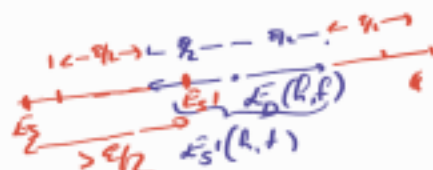
✓ CLAIM: $P(D_k/C_k) \geq 1/2$

PROOF: Let $D'_k = |\mathcal{E}_{S'}(h, f) - \mathcal{E}_D(h, f)| \leq \varepsilon/2$

Given C_k , $D'_k \Rightarrow D_k/C_k$

$$A \Rightarrow B, P(A) \leq P(B)$$

$$\begin{array}{ccc} & & \geq \varepsilon \\ & \nearrow & \\ s & & s' \\ & \searrow & \\ & & \leq \varepsilon/2 \end{array}$$



$$P(D_k/C_k) \geq P(D_k'/C_k)$$

$$\geq P(|X - \mu_k| \leq \epsilon/2)$$

← Average of n Bernoulli's

$$\geq 1 - 2e^{-\frac{\epsilon^2}{2}n}$$

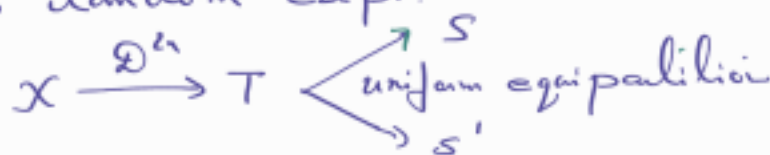
$$\geq \frac{1}{2} \quad \text{if } \underline{n > 4/\epsilon^2}$$

$$P\left\{|X - \mu_k| > \frac{\epsilon}{2}\right\} \leq 2e^{-\frac{\epsilon^2}{2}n}$$

Hence $P(D/C) \geq 1/2$ when $\quad \quad \quad (1)$

$$D = \bigcup_{L \neq H} D_k \text{ and } C = \bigcup_{L \neq H} C_k$$

New random expt.



For any $k, P_{S, S'}[D_k]$

$$= P_{T \sim D_k, (S, S') \sim \text{unif}}[D_k]$$

$$= P_{S, S'} \left[|E_S(h, f) - E_{S'}(h, f)| > \epsilon/2 \right]$$

$$\leq P_{S, S'} \left[(|E_S(h, f) - \mu| > \epsilon/4) \vee (|E_{S'}(h, f) - \mu| > \epsilon/4) \right]$$

$$\leq 2 \cdot 2e^{-\frac{2\epsilon^2}{16}n} \quad (\text{Hoeffding})$$

$$= 4e^{-\frac{\epsilon^2}{8}n}$$

—(*)

$$|X - X'| > \epsilon/2$$

$$E[X] = E[X'] = \mu = E_0$$



$$|X - \mu| > \epsilon/4$$

$$|X' - \mu| > \epsilon/4$$

Putting it all together.

$$P_{S \sim D^n} \left[\bigcup_{h \in H} C_h \right] \leq \delta$$

$\uparrow P(D/C) \geq 1/2$

$$P_{S, S' \sim D^n} \left[\bigcup_{h \in H} D_h \right] \leq \delta/2$$

\uparrow Smartest Idea.

$$P_{\substack{T \in D \\ S, S' \sim T}} \left[\bigcup_{[h]_T \in H} D_h \right] \leq \delta/2$$

\uparrow Shattering

$$g_H(2n) P_{S, S'}[D_h] \leq \delta/2$$

$\uparrow (*)$

$$g_H(2n) 4 e^{-\frac{\epsilon^2}{8} n} \leq \delta/2$$

\uparrow Saw's Lemma.

$$\binom{2n}{\leq d} e^{-\frac{\epsilon^2}{8} n} \leq \delta/8$$

$$(VCdim(H) = d)$$

$$\binom{2n}{\leq d}^d \uparrow$$

$$n \geq \frac{c}{\epsilon^2} (d \ln(1/\epsilon) + \ln(1/\delta))$$

THEOREM

$VC DIMENSION(H)$ finite

$\Leftrightarrow H$ is UNIFORM PAC LEARNABLE

(H has Uniform Convergence Property).

Sample Complexity: $O\left(\frac{1}{\epsilon^2} (d \log 1/\epsilon + \log 1/\delta)\right)$

Converse?

$$\text{VC dim}(H) = \infty$$

$\Rightarrow H$ is not PAC-learnable

$\Rightarrow H$ is not uniform PAC-learnable,

$$(\forall \epsilon(A) \leq \epsilon(C)).$$

1. Approximation vs. Generalisation Tradeoff

LECTURE 14
4/3/2021

2. Combination of Concepts.

We know

$$P_{S,D^n}[\exists h, |L_S(h,f) - L_D(h,f)| > \epsilon]$$

$$= P_{S \sim D^n}[C]$$

$$\leq 2 P_{S,S' \sim D^n}[D]$$

$$\leq 2 \cdot g_H(2n) P_{S,S' \sim D^n}[2R]$$

$$\leq 2 g_H(2n) \cdot 4 e^{-\frac{\epsilon^2}{8} n}$$

$$\leq 8(2n)^d e^{-\frac{\epsilon^2}{8} n}$$

— (1)

$$L_S = 0$$

$$|L_D - L_S| < \epsilon$$

$$L_D - L_S < \epsilon$$

$$L_D < L_S + \epsilon$$

↑
True about

Earlier we solved for n so that (1) $\leq \delta$.

We can also solve for ϵ assuming n is given.

$$8(2n)^d e^{-\frac{\epsilon^2}{8} n} \leq \delta$$

$$e^{\frac{\epsilon^2}{8} n} \geq \frac{8}{\delta} (2n)^d$$

$$\epsilon^2 \geq \frac{8}{n} \left[\ln\left(\frac{8}{\delta}\right) + d \ln 2n \right]$$

$$\epsilon \geq \sqrt{\frac{8}{n} [d \ln 2n + \ln \frac{8}{\delta}]}$$

Hence

$$= \Omega \approx \tilde{O}(\sqrt{d/n})$$

$$|L_{in} - L_D| \leq \Omega \text{ w.p. } 1 - \delta$$

(Generalisation error)

$$\mathcal{E}_D \leq \mathcal{E}_{in} + \Omega$$

\uparrow \uparrow
 Approximation Generalisation
 \mathcal{E}_{in} over.

OCCAM'S
RAZOR

Ω increases with increase in d/n ✓

\mathcal{E}_{in} is likely to decrease with d/n . ✓

Hence choice of d (i.e. choice of H)
is a **trade off**.

"Amateur ML engineers" tend to choose
a larger d/n so that \mathcal{E}_{in} (which is
staring at their face) is low.
"Overfitting"

"Professional ML engineers" are more
careful about Ω .

(Prof. Abu-Mostafa, CALTECH)

Finding the sweet spot?

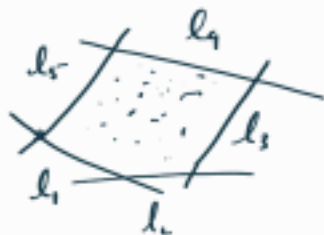
- Bias-variance trade off
- Regularisation
- Validation

(Not in this course).

COMBINING HYPOTHESIS CLASSES

Scenarios

1



$$G = G_1 \cap G_2 \cap G_3 \cap B_4 \cap B_5$$

$$f: x_1 \wedge x_2 \wedge x_3 \wedge \bar{x}_4 \wedge \bar{x}_5$$

2



"Perceptron n/w"

Let H be a hypothesis class that is obtained as a combination of k hypothesis classes H_1, \dots, H_k .

$$\text{ie } h(x) = f(h_1(x), \dots, h_k(x))$$

where $h_i \in H_i$ and

$f: \{0,1\}^k \rightarrow \{0,1\}$ is any Boolean function.

Then $VC\text{-dim}(H) = ?$

Will depend on f (2^{2^k} f's)

An upper bound independent of f ? YES

Simplification: $\forall i, VC\text{-dim}(H_i) \leq d$.

CLAIM:
 \checkmark $VC\text{-dim}(H) \leq 2kd \ln kd$

PROOF:

Let S be a largest shattered
 by H

Let $|S| = n$.

$$\Rightarrow |H|_S = 2^n \quad (1)$$

Consider $H^* = H_1 \times H_2 \times \dots \times H_k$
 $= \{(h_1, \dots, h_k) : h_i \in H_i\}$

Obs 1. $|H|_S \leq |H^*|_S$

Obs 2. $H^*_S = H_1|_S \times H_2|_S \times \dots \times H_k|_S$

$$|H^*_S| = |H_1|_S \times \dots \times |H_k|_S$$

$$\leq \binom{n}{d_1} \times \dots \times \binom{n}{d_k}$$

$$\leq n^{d_1} \times \dots \times n^{d_k}$$

$$\leq n^{kd}$$

Hence $2^n \leq n^{kd}$

$$\uparrow$$

$$n \leq 2kd \log kd. \quad \square$$

$$(VC\text{-dim}(H) = |S|)$$

$\overset{n}{\underset{n}{}}$

$$h, h' \in H$$

$$h|_S \neq h'|_S$$

$$h(x) = f(h_1(x), \dots, h_k(x))$$

$$h'(x) = f(h'_1(x), \dots, h'_k(x))$$

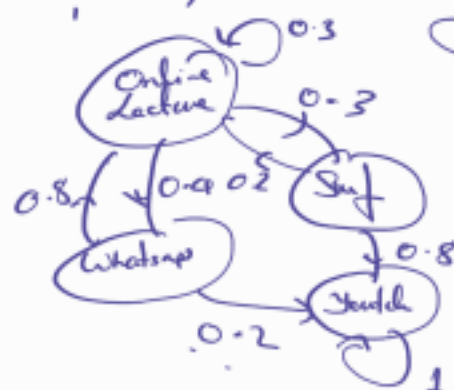
$h(x) \neq h'(x)$ needs
 at least one $h_i(x) \neq h'_i(x)$

Fact. Perceptron networks (neural networks
 with hard thresholds) have
 $VC\text{-dim} \ O(m \log m)$, where m is
 the no. of edges.



Examples:

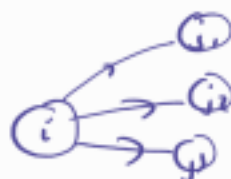
- (Web Pages, Hyperlinks)
- (Activities, Switching)



Edge weight = transition probabilities



$$P_{ij} = P[\text{Next State} = j / \text{Current State} = i]$$



Hence $\forall i, \sum_{j \in V(G)} P_{ij} = 1$ (outgoing from i).

— (PI)

$$P = \begin{bmatrix} P_{11} & \dots & P_{1n} \\ P_{21} & \dots & P_{2n} \\ \vdots & & \vdots \\ P_{n1} & \dots & P_{nn} \end{bmatrix}_{n \times n}$$

Transition Probability Matrix (TPM).

"Stochastic Matrix"

$$\begin{aligned} & m_{ij} \geq 0 \\ & \sum_{j=1}^n m_{ij} = 1 \quad (\text{row sum}) \end{aligned}$$

PI \equiv Every row of P sums to 1.

$$\begin{bmatrix} P_{11} & \dots & P_{1n} \\ \vdots & & \vdots \\ P_{n1} & \dots & P_{nn} \end{bmatrix}_{n \times n} \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}_{n \times 1} = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}_{n \times 1}$$

Hence $\mathbb{1} = [1, 1, \dots, 1]^T$ is a (right) eigen vector of P with eigen-value 1.

In fact, If P is a real non-negative matrix then

P is a stochastic matrix

$$\iff P\mathbb{1} = \mathbb{1}$$

(-P2)

Corollary: If P_1 & P_2 are $n \times n$ S.M.s,
so in $P_1 P_2$ $P = P_1 P_2$

$$\begin{aligned} P_1 P_2 \mathbb{1} &= P_1 (P_2 \mathbb{1}) \\ &= P_1 \mathbb{1} \\ &= \mathbb{1} \end{aligned}$$

$$P_{i,j} \geq 0$$

Corollary: If P is a S.M.

$\forall k \in \mathbb{N}$, P^k is a S.M.

$$P^k[i,j] = P[\text{State after } k \text{ steps} = j \mid \text{Current State} = i]$$

(Proof: Exercise)

$$\begin{aligned} P[\text{Next to Next state} = j \mid \text{Current state} = i] &= \sum_{k=1}^n P_{ik} P_{kj} \\ &= P^2[i,j] \end{aligned}$$



Prop 3.

All eigen values of P have magnitude ≤ 1 .

i.e. $Px = \lambda x \Rightarrow |\lambda| \leq 1$.

Proof: Let $x = (x_1, \dots, x_n)$ be an eigen vector corresponding to λ s.t. some $x_i = 1$ and $|x_j| \leq 1 \forall j$ (Scaling).

$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

$$x^T = [x_1 \dots x_n]$$

We have $Px = \lambda x$

$$\langle (i^{th} \text{ row of } P), x \rangle = \lambda x_i$$

$$\sum_{j=1}^n P_{ij} x_j = \lambda x_i$$

$$\lambda = \sum_{j=1}^n P_{ij} x_j$$

$$|\lambda| \leq \sum_{j=1}^n \underline{P_{ij}} |x_j|$$

$$\leq \sum_j P_{ij} 1$$

$$= 1$$

$$i: \begin{bmatrix} \leftarrow \end{bmatrix} \begin{bmatrix} x \\ \vdots \\ x \end{bmatrix} = \begin{bmatrix} \oplus \end{bmatrix}$$

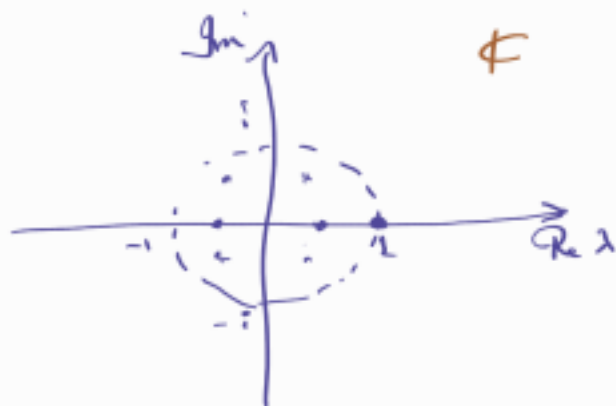
— (1) $|z_1, z_2| \leq |z_1|/|z_2|$

$\leq \dots \leq \frac{x_i}{1}$

(try it $\Rightarrow |x_j| \geq 1 \forall j$)

$\Rightarrow x_j = 1 \dots$

$\forall j: P_{ij} > 0$



Next Question:

What is the multiplicity of $\lambda = 1$?

If $\lambda = 1$, then eqn (1) above has to be an equality.

i.e. $x_j = 1$ for all j s.t. $P_{ij} \neq 0$



$$x = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

If every vertex is reachable from i , then $x_j = 1 \forall j \Rightarrow x = \mathbb{1}$

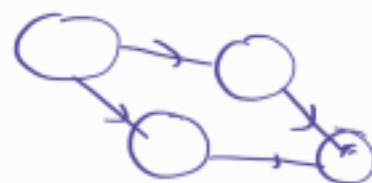
Prop 4.

If the underlying graph G is strongly connected then the eigen value 1 has multiplicity 1.

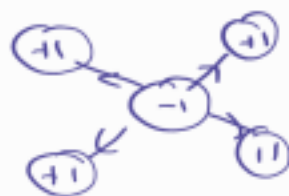
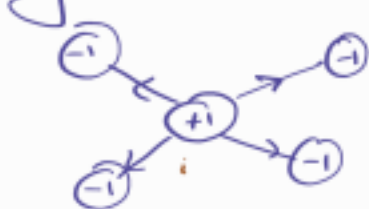
$\mathbb{1}$ is the unique eigen vector (up to scaling) for $\lambda = 1$.



Converse? Exercise:
(Hint. Sink component).



When do we get $\lambda = -1$.



iff A is "bipartite".

Strongly connected $\Rightarrow -1$ also has multiplicity 1.

Qn. Is P full rank?

Not necessarily

$$\begin{bmatrix} 1/3 & 1/3 & 1/3 \\ 1/3 & 1/3 & 1/3 \\ 1/3 & 1/3 & 1/3 \end{bmatrix}$$

Interestingly, the $n \times (n+1)$ matrix

$P - I : \mathbb{1}$ has rank n .

(Exercise).



Summary:

Let P be the TPM of a random walk on a finite directed graph G .

(1) $\sum_{j=1}^n P_{ij} = 1$. (TPM is a SM)

(2) $P \mathbf{1} = \mathbf{1}$, and hence $\lambda = 1$ is an eigenvalue.

(3) All eigenvalues of P lie in the closed unit disk around 0.

(4) $\lambda = 1$ has multiplicity 1 $\iff G$ is strongly connected.

(5) -1 is an eigenvalue of P $\iff G$ is bipartite.

Note:

All these apply to undirected graphs as well.



If A is undirected P is symmetric.

\Rightarrow Eigen values are real.

$\Rightarrow \in [-1, 1]$.

Now that we have learned so much about P , let's move on and study

P^T

P^T is more important than P .

(x_1, \dots, x_n) is called a prob. vec. if $x_i \geq 0$ & $\sum x_i = 1$.

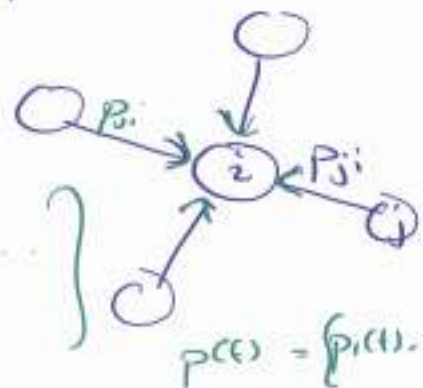
Why

Let $p = (p_1, \dots, p_n)$ be the node probabilities at current step.



Then $p' = P^T p$ gives the node probabilities at the next step.

$$p'_i = \sum_{j=1}^n P_{ji} p_j = [P^T p]_i$$



$$p(t) = (p_1(t), \dots, p_n(t))$$

$$p(t+1) = P^T p(t)$$

(if $p(t)$ is written as a col vector)

$$p(t+1) = p(t) P$$

(if $p(t)$ is written as a row vector)

P^T captures

"evolution of the random walk"

* Distribution at time t
prob. vec $p(t)$

* evolution

$$p(t+1) = P^T p(t)$$

$$p(0) = (0, 0, 1, 0, 0, \dots)$$

$$P^T p(0) = p(1) = (0, \frac{1}{2}, \frac{1}{2}, \frac{1}{2}, 0, \dots)$$

↓

$$P^T p(1) = p(2) = \dots$$

Properties of P^T

- (1) Eigen values (and multiplicities) of P and P^T are the same.

(Proof: Determinant)

$$\det(P - \lambda I) = \det((P - \lambda I)^T) \\ = \det(P - \lambda I).$$

$x = (x_1, \dots, x_n) \in \mathbb{R}^n$ is called a probability vector if

(i) $x_i \geq 0 \quad \forall i$

(ii) $\sum_{i=1}^n x_i = 1.$

- (2) x is a prob. vec $\Rightarrow P^T x$ is a prob. vec.

Proof: Let $y = P^T x$.
 $y_i \geq 0$, obvious.

$$\begin{aligned} \sum_{i=1}^n y_i &= \mathbf{1}^T y \\ &= \mathbf{1}^T P x \\ &= (\mathbf{1}^T P) x \\ &= \mathbf{1}^T x \\ &= 1. \end{aligned}$$

$$(1, \dots, 1) \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} = 1.$$

If G is strongly connected, then

- (3) P^T also has a unique eigen vector (upto scaling) for $\lambda=1$.

But is it a prob. vec?

$$P^T P = P$$

"Stationary distribution"

STATIONARY DISTRIBUTIONS of RANDOM WALKS

LECTURE 16
11/4/2021

Defn. $\pi \in \mathbb{R}^n$ is called a stationary distribution of a random walk with transition prob. matrix P if

- (i) π is a prob. vec., and
- (ii) $P^T \pi = \pi$.

$$P^T P$$

$$P^T P$$

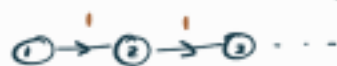
$$P^T x = 1x$$

Can we
have α to
be prob. vec?
 $P^T P = P$

Theorem 1. (Fundamental Theorem of Finite Markov Chains)

****** A random walk on every finite graph has a stationary distbn.

Obs. Finiteness is necessary.



$$V(G) = \mathbb{N}$$

$$E(G) = \{(i, i+1) : i \in \mathbb{N}\}$$

$$P^T P$$

(Many important infinite cases can be handled - but not in this course)

Standard proofs

Key idea "Fixed point theorems"

Let $f: X \rightarrow X$, then a point $x \in X$ s.t. $f(x) = x$ is called a fixed point of f .

Domain X
 $f: X \rightarrow X$
 $X = \{1, 2, 3\}$
 $f(x) = 4 - x$

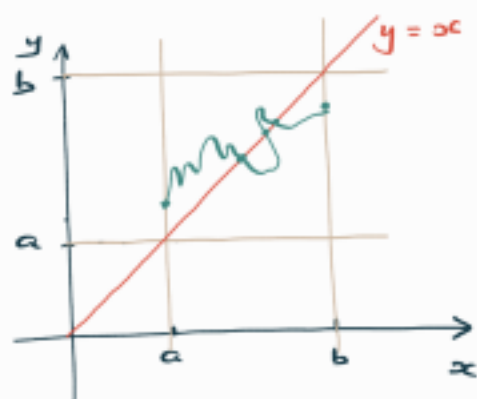
1	\mapsto	3
2	\mapsto	2
3	\mapsto	1

Example: Any continuous fn from a closed interval $[a, b]$ to itself has a fixed point.

$$X = [a, b]$$

$$f: X \rightarrow X, \text{ continuous.}$$

Let $f: [a, b] \rightarrow [a, b]$, cont



Non-example

$f: [n] \rightarrow [n]$, derangement (cyclic shift).

$$1 \mapsto 2$$

$$2 \mapsto 3$$

$$3 \mapsto 1$$

Fact 1 If X is a compact convex set and f is continuous then f has a fixed point (Brouwer's FP theorem, 1909).

Fact 2 Closed and bounded sets in \mathbb{R}^n are compact.

$X \subseteq \mathbb{R}^n$
Compact = closed & Bounded
 $[a, b]$ closed
 (a, b) not closed

What's our X ?

$X =$ set of all prob. vecs in \mathbb{R}^n
($n = |V(K)|$)

$$= \{(p_1, \dots, p_n) : p_i \geq 0, \sum p_i = 1\}$$

$$q = a + \frac{1}{n}$$

$$\rightarrow a \text{ as } n \rightarrow \infty$$

$$\bigcirc [-1]$$

$$f: X \rightarrow X$$

$$p \mapsto P^T p$$

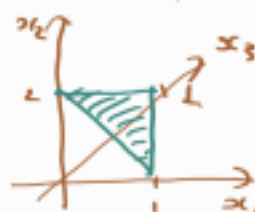
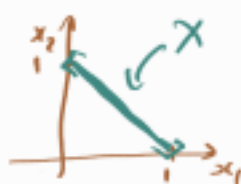
$$p \in X \Rightarrow P^T p \in X$$

Verify: X is
 ✓ (i) convex
 ✓ (ii) bounded
 ✓ (iii) closed.

$$\exists x \in X \text{ s.t. } f(x) = x, \\ \text{(i) } p, q \in X \\ \lambda p + (1-\lambda)q \in X \quad (\text{Convex})$$

Then $f: X \rightarrow X$
 $x \mapsto P^T x$

Verify: f is continuous.



$$\text{(ii) } X \text{ is bounded} \\ \|x\|_\infty = \max x_i \\ \{\|x\|_\infty : x \in X\} \leq 1 \\ \|x\|_1 = \sum_{i=1}^n |x_i| \\ \|x\|_2 = \sqrt{\sum_{i=1}^n x_i^2}$$

Done! f has a fixed point in X .

That is your stationary distribn.

$$x, x+\Delta x \\ f(x+\Delta x) - f(x) = f'(\Delta x)$$

DIRECT PROOF.

Let x be any prob. vec.

Consider the sequence of prob. vectors

$$x = [p_1, p_2, \dots, p_n] \quad x, xP, xP^2, xP^3, \dots$$

$$(x, P^T x, (P^T)^2 x, \dots)$$

(evolution of the random walk)

$x(0), x(1), \dots$, where

$$x(t) = xP^t$$

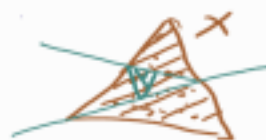
$$\text{Let } a(t) = \frac{1}{t} (x(0) + x(1) + \dots + x(t-1))$$

$$= \frac{1}{t} (x + xP + xP^2 + \dots + xP^{t-1})$$

"long-term average"

Claim 1. $\forall t$, $a(t)$ is a prob. vec.

$\therefore a(1), a(2), \dots \in X$



X bounded $\Rightarrow a(1), \dots$ contains a
convergent subseq
 $a(t_1), a(t_2), \dots$

(Proof: Pigeonhole principle).

X closed $\Rightarrow \lim_{n \rightarrow \infty} a(t_n) = a \in X$,

(i.e. a is a prob. vec.)

$\forall t$,

$$a(t)P - a(t) = \frac{1}{t} \left(\underbrace{xP}_{\in X} + \underbrace{xP^2}_{\in X} + \dots + xP^t \right) - \frac{1}{t} \left(x + xP + \dots + xP^{t-1} \right)$$

$$= \frac{1}{t} \left(\underbrace{xP^t}_{\in X} - \underbrace{x}_{\in X} \right)$$

$$\|a(t)P - a(t)\|_{\infty} \leq \frac{1}{t} \quad \text{--- (x)} \quad \left(\|x\|_{\infty} = \max_i |x_i| \right)$$

$$\|a(t_n)P - a(t_n)\|_{\infty} \leq \frac{1}{t_n} \rightarrow 0 \text{ as } n \rightarrow \infty$$

$$\Rightarrow \|aP - a\|_{\infty} = 0$$

$$aP = a$$

Hence a is a stationary distribution.

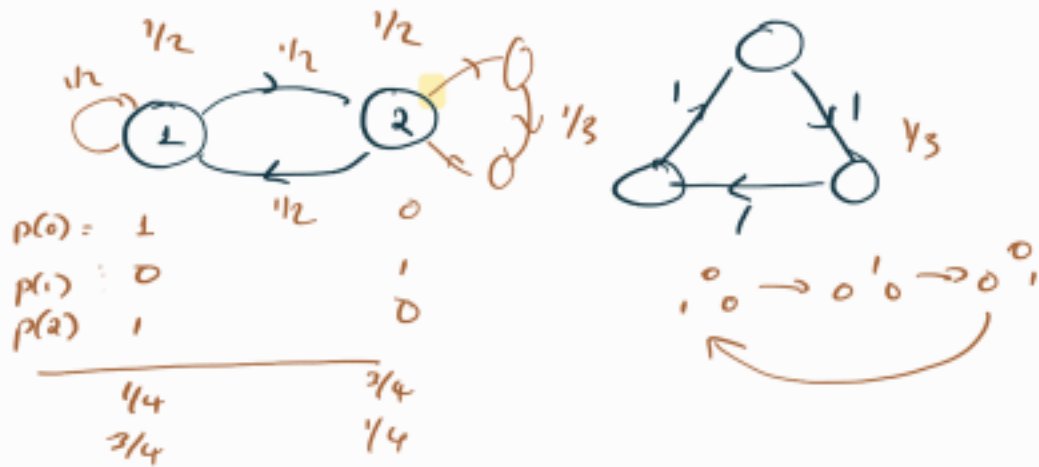
Summary

1. Every random walk on a finite graph has a stationary distribution \vec{v} ($\vec{v}^T P = \vec{v}^T$).
2. If the graph is strongly connected, the stationary distribution is unique, and $a(t) \rightarrow \vec{v}$ as $t \rightarrow \infty$.

Qn. If $P(0)$ is an arb. starting distn
does $P(n) = P(0)P^n \rightarrow \sqrt{1}$ as $n \rightarrow \infty$?

Ans: Not always, but
Yes in most cases.

Some no cases.



Yes iff gcd of lengths of all directed cycles in α is 1!

Summary (Once again)

1. Every random walk on a finite graph } has a stationary distribution $\vec{\pi}$
2. If the graph is strongly connected } the stationary distribution is unique, and $a(t) \rightarrow \vec{\pi}$ as $t \rightarrow \infty$
3. If the graph G is strongly connected and $\gcd(\text{cycle lengths}) = 1$ } For any prob. vec $P(0)$ $P(t) \rightarrow \vec{\pi}$ as $t \rightarrow \infty$

Applications

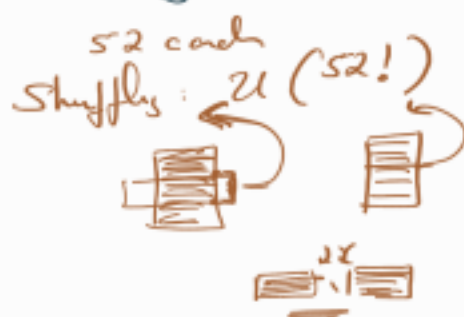
Analysing natural phenomena

- Bioconic motion
- Financial markets
- Text / Speech
- Genetics
- Web

Generating samples according to some target distribution ($\sim D$)

- Card shuffling
- Metropolis-Hasting algo
- Gibbs sampling

(p_1, p_2, \dots, p_n)
 $n = \# \text{web page}$
 $\uparrow \quad 1/n$



MARKOV CHAIN MONTÉ CARLO METHODS

LECTURE 17

17/4/2021

GOAL: Sample a point x
according to a distribution D on X .
 $x \sim D$.

E.g.: σ is a permutation of $[n]$
chosen uniformly at random.

$x \in \mathbb{R}^n$ is sampled from an
 n -dimensional gaussian

Key Idea: (π : distn over X)

Let the domain X be finite ($|X|=n$)

Design a directed graph G and
transition probabilities P s.t

π is a stationary distribution of P

π will be
a n -length
prob. dist

$$\pi = (\pi_1, \dots, \pi_n) \\ \text{s.t. } \pi_i \geq 0, \\ \sum_{i=1}^n \pi_i = 1$$

Def 6: $P \rightarrow \pi$

Def 7: $\pi \rightarrow P$

Desirable properties:

- (1) G is strongly connected. : unique st. dist
- (2) $\gcd(\text{cycle lengths}) = 1$. : $P^t \rightarrow \pi$
- (3) Low degrees
- (4) Symmetries (regular for example).
- (5) Rapid mixing (convergence to stationary distribution).

Easy Case: π is uniform on X . $\pi_i = 1/n$.

G : undirected ^{connected non-bipartite} graph with $V(G) = X$



: k -regular ($k \geq 2$)

: $P_{ij} = 1/k \quad \forall i, j$

: Expander.

$$P^T = P$$

$\mathbf{1}$ is eigenvec of P
& P^T (since $P^T = P$)

$\therefore (\frac{1}{n}, \frac{1}{n}, \dots, \frac{1}{n}) = \pi$
is the st. dist of P

Non-uniform π

A simple sufficient condition.

1) π is a prob. vec. & P is a st. mat s.t.,

$$\overline{P^T \pi = \pi}$$

$$\forall i, j. \quad \pi_i P_{ij} = \pi_j P_{ji} \quad (*)$$



then π is a stationary distn of P



Proof:

$$\text{let } \sigma = P^T \pi$$

$$\text{then } \sigma_i = \sum_{j=1}^n P_{ji} \pi_j$$

$$= \sum_{j=1}^n P_{ij} \pi_i \quad (by *)$$

$$= \pi_i$$

□

$$P^T \begin{bmatrix} \longleftrightarrow \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}$$

METROPOLIS - HASTING

Input: \bar{v} , a prob. dist on $[n]$
(an n -length prob. vec.)

Design of A

Pick a "good" connected undirected graph H and replace each undirected edge with 2 opposite arcs and add a self loop at each node to get A .



Obs. A is strongly connected,
 $\gcd(\text{cycle-lengths}) = 1$.

Automatically satisfied
for missing edges of H
and self-loops

Design of P

(Aim: $\forall i, j$ $\bar{v}_i P_{ij} = \bar{v}_j P_{ji}$)

Attempt 1.

$$P_{ij} = \bar{v}_j \quad \forall j \in N^+(i)$$

$$\begin{aligned} \text{Then, } \bar{v}_i P_{ij} &= \bar{v}_i \bar{v}_j \\ \bar{v}_j P_{ji} &= \bar{v}_j \bar{v}_i \end{aligned} \quad \checkmark$$

But is $\sum_{j=1}^n P_{ij} = 1$?

$$\sum_{j=1}^n P_{ij} = \sum_{j \in N^+(i)} \bar{v}_j \ll 1. \quad \because H \text{ is sparse.}$$

We can fix this by changing P_{ii}

$$P_{ij} = \begin{cases} \bar{v}_j & , j \in N^+(i) \setminus \{i\} \\ 1 - \sum_{k \in N^+(i) \setminus \{i\}} \bar{v}_k & , j = i \end{cases}$$

Issue: too slow a walk.

Why not scale?



- Subject to
1. Equal scaling for P_{ij} and P_{ji}
 2. $\sum_{j=1}^n P_{ij} \leq 1$.
(Defect $\rightarrow P_{ii}$).

Let α = maximum out-degree of G .
not counting the self-loop.

(2) is ensured if $\forall i, j \quad P_{ij} \leq 1/\alpha$. (i+j) - 1/2

$$\text{Scaling for } P_{ij} = \frac{1/\alpha}{\max\{P_{ij}, P_{ji}\}}$$

$\alpha = \frac{1}{\beta}$
Scaling $\frac{1/\alpha}{\beta}$

Hence

$$\begin{aligned} \forall j \in \mathcal{N}(i) \setminus \{i\}, P_{ij} &= \frac{1/\alpha}{\max\{P_{ij}, P_{ji}\}} \times P_{ji} \\ &= \frac{1}{\alpha} \times \min\left\{\frac{1}{P_{ji}}, \frac{1}{P_{ii}}\right\} \times P_{ji} \\ &= \frac{1}{\alpha} \min\{1, P_{ji}/P_{ii}\} \end{aligned}$$

Metropolis-Hasting ($G \rightarrow H$)

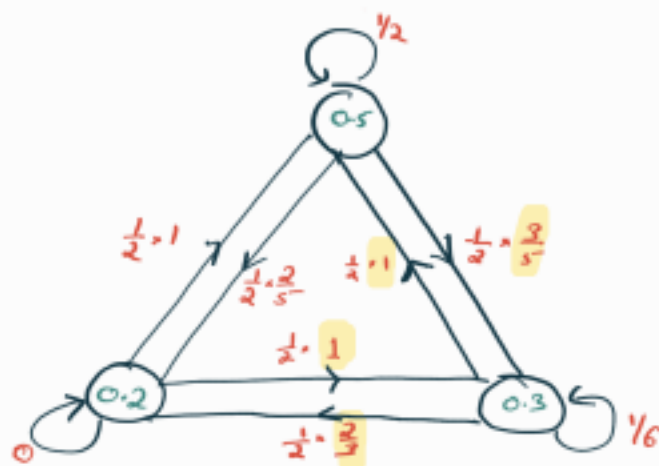
$$\forall i \quad P_{ij} = \begin{cases} \frac{1}{\alpha} \min\{1, P_{ji}/P_{ii}\} & , j \in \mathcal{N}(i) \setminus \{i\} \\ 1 - \sum_{j \in \mathcal{N}(i) \setminus \{i\}} P_{ij} & , j = i \end{cases}$$

Sanity check.

$$\begin{aligned}
 \bar{v}_i, p_{ij} &= \frac{\bar{v}_i}{n} \min \{1, \bar{v}_j / \bar{v}_i\} \\
 &= \frac{1}{n} \min \{\bar{v}_i, \bar{v}_j\} \\
 &= \frac{\bar{v}_j}{n} \min \{1, \frac{\bar{v}_i}{\bar{v}_j}\} \\
 &= \bar{v}_j, p_{ji} \quad \checkmark
 \end{aligned}$$

$$\forall i, \sum_{j \in \omega^+(i)} p_{ij} = 1 \quad \checkmark$$

Example: $\bar{v} = (0.5, 0.3, 0.2)$



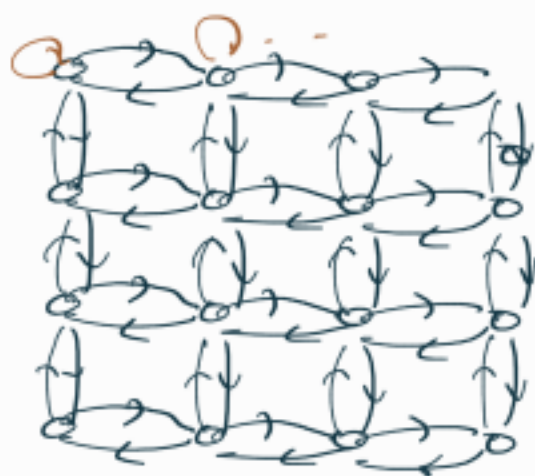
Walker's rule

When at node i

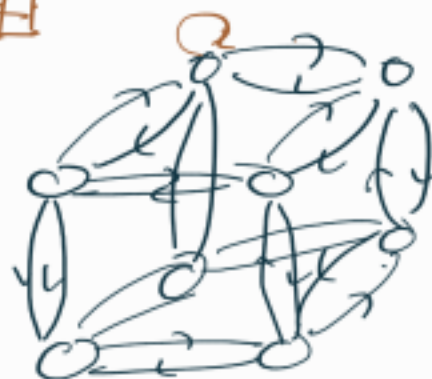
- Select each out edge (i,j) w.p. $1/n$ $(14j)$
- If $(\bar{v}_j \geq \bar{v}_i)$ move to node j (w.p. 1)
- Else move to node j w.p. \bar{v}_j / \bar{v}_i

Usually,

Graph: d -dimensional lattice $[m]^d$



$$d=2, m=4$$



$$d=3, m=2$$

- $n = m^d \geq |X|$

- Almost $2d$ -regular
(except boundary vertices),
 $\alpha = 2d$.

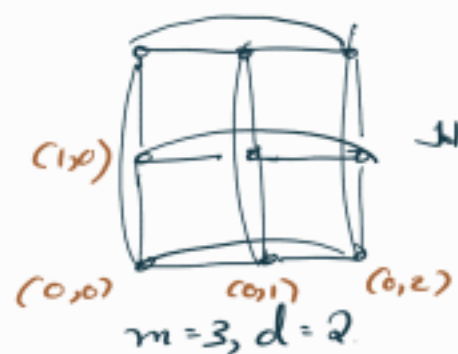
GIBBS SAMPLING

Domain: $X = [m]^d$
 $= \{(x_1, \dots, x_d) : x_i \in [m]\}$

Target: p , a distn on X .

$$H = \underbrace{K_m \square K_m \square \dots \square K_m}_{d \text{ times}}$$

(Hamming graph)



$V(H) = X$
 $\{x, y\}$ is an edge of H iff
 x & y differ in exactly
one co-ordinate

Obs. Much denser than lattice
 $(M-H)$.

Let $x = (x_1, x_2, \dots, x_d)$
 $y = (y_1, x_2, \dots, x_d)$, $y_1 \neq x_1$
 $(x \text{ and } y \text{ are adjacent in } H)$

$$P_{xy} = \frac{1}{d} \pi(y_1/x_2, \dots, x_d)$$

$$\left[\pi(y_1/x_2, \dots, x_d) = \frac{\pi(y_1, x_2, \dots, x_d)}{\sum_{z=1}^m \pi(z, x_2, \dots, x_d)} \right]$$

$$\begin{aligned} \text{Hence } \pi(x) P_{xy} &= \pi(y) P_{yx} \quad \forall x, y. \\ &= \frac{\pi(x_1, x_2, \dots, x_d) \pi(y_1, x_2, \dots, x_d)}{\sum_{z=1}^m \pi(z, x_2, \dots, x_d)}. \end{aligned}$$

$$\sum_{y \in N^+(x)} P_{xy} = \underbrace{\frac{1}{d} + \frac{1}{d} + \dots + \frac{1}{d}}_{d \text{ times}} = 1.$$

Walker's Rule

- When at node (x_1, \dots, x_d)
- Pick a dim $k \in [d]$ w.p. $1/d$ each.
- (Follow attempt 1 on dimension k)
Move to $(x_1, \dots, x_{k-1}, y_k, x_{k+1}, \dots, x_d)$
w.p. $\pi(y_k / x_{k-1}, \dots, x_{k+1}, \dots, x_d)$

Example:

