

# Lecture 21

①

- Recall that the **CFN model** is a Markov Chain on a tree  $(T, P, \mu_P)$ , where  $T$  is a directed binary tree rooted at  $\rho$  and with leaf set  $[n]$ , whose edges have transition matrices

$$P_e = \begin{bmatrix} 1 - p_e & p_e \\ p_e & 1 - p_e \end{bmatrix} \text{ over the character set } \{0, 1\},$$

and  $\mu_P = (1/2 \ 1/2)$ .

[ $p_e$  is the mutation probability of the edge]

- The **tree reconstruction problem** is the following:

- given  $\underline{Z} = (Z_{[n]}^1, Z_{[n]}^2, \dots, Z_{[n]}^k)$ , that is  $k$  independent samples from the CFN model at the leaves of  $T$ , the goal is to reconstruct the unrooted, undirected tree  $T^{-\rho}$
- the strong reconstruction problem also asks for a CFN model over the leaf set  $[n]$  whose leaf character distn' is within  $\epsilon$  total variation distance <sup>from</sup> the distn' of the actual model (where the samples  $\underline{Z}$  were sampled from)

- Our goal is to reconstruct  $T^{-\rho}$  using a tree metric

- Our tree metric is inspired by the following decomposition:

$$P_e = (1 - 2p_e) \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} + 2p_e \begin{pmatrix} 1/2 & 1/2 \\ 1/2 & 1/2 \end{pmatrix}$$

interpretation: on an edge  $e = (u, v)$  :

$\begin{matrix} \theta_e \\ \vdots \\ \theta_e \end{matrix}$

$\begin{matrix} \xrightarrow{1-2p_e} \\ \xrightarrow{2p_e} \end{matrix}$

$\begin{matrix} \mathbb{F}_v = \mathbb{F}_u \text{ w/pr } 1-2p_e \\ \mathbb{F}_v \perp \mathbb{F}_u \text{ w/pr } 2p_e \end{matrix}$

consider now a path in tree  $Path(a,b) = \{e_1, \dots, e_k\}$



$$\mathbb{I}_b = \begin{cases} \mathbb{I}_a & \text{w.p. } \prod_{i=1}^k \theta_{e_i} \\ 0/1 \text{ u.r.} & \text{w.p. } 1 - \prod_{i=1}^k \theta_{e_i} \end{cases}$$

define:  
hence  $\theta(a,b) := \prod_{i=1}^k \theta_{e_i}$   
and  $p(a,b) = \frac{1 - \theta(a,b)}{2}$

Since  $\theta_{e_i}$ 's multiply, reasonable to define

$$w_e = -\log \theta_e$$

then  $\sum_{i=1}^k w_{e_i} = -\log \prod_{i=1}^k \theta_{e_i} = -\log \theta(a,b)$

- So CFN model defines tree metric  $\delta(a,b) = -\log \theta(a,b)$
- If I could estimate  $\delta(a,b)$  to within  $\frac{1}{4} w_e^{\text{min}}$ , I could solve the reconstruction problem (from last lecture)

Estimating  $\delta(a,b)$  from  $\hat{P}$

define  $\hat{P}_{ab} := \frac{1}{k} \sum_{i=1}^k \mathbb{1}_{\{z_a^i \neq z_b^i\}}$ ; then  $E[\hat{P}_{ab}] = \frac{1 - \theta(a,b)}{2} = p(a,b)$

set  $\hat{\delta}(a,b) = \begin{cases} -\log[1 - 2\hat{P}^{ab}] & \text{if } \hat{P}^{ab} < 1/2 \\ \text{equal to } +\infty & \text{if } \hat{P}^{ab} \geq 1/2 \end{cases}$

Claim: Suppose that  $\forall a, b \in [n]$

$$|p^{ab} - \hat{p}^{ab}| < \varepsilon \leq \frac{1}{2} (1 - e^{-2W_*/4}) (1 - 2p^{ab})$$

then

$$\max_{q=(a,b,c,d)} |\delta(q) - \hat{\delta}(q)| \leq 2 \max_{a,b} |\delta(a,b) - \hat{\delta}(a,b)| < \frac{1}{2} W_*$$

$$\text{where } \delta(q) = \frac{1}{2} [\delta(a,c) + \delta(b,d) - \delta(a,b) - \delta(c,d)]$$

$(a,b,c,d)$

Proof:

• observation:

$$p^{ab} + \varepsilon < p^{ab} + \frac{1}{2} - p^{ab} < \frac{1}{2}$$

• hence  $-\log(1 - 2(p^{ab} \pm \varepsilon))$  well-defined (i.e. not  $+\infty$ )

• Now:

$$|\delta(a,b) - \hat{\delta}(a,b)| = |\log(1 - 2p^{ab}) - \log(1 - 2(p^{ab} \pm \varepsilon))|$$

$$= \left| \log \frac{1 - 2p^{ab} \pm 2\varepsilon}{1 - 2p^{ab}} \right|$$

$$= \left| \log \left( 1 \pm \frac{2\varepsilon}{1 - 2p^{ab}} \right) \right| \leq \frac{1}{4} W_*$$

↳ indeed:

$$\left| \log \left( 1 \pm \frac{2\varepsilon}{1 - 2p^{ab}} \right) \right| \leq \max \left\{ \log \left( 1 + \frac{2\varepsilon}{1 - 2p^{ab}} \right), -\log \left( 1 - \frac{2\varepsilon}{1 - 2p^{ab}} \right) \right\}$$

$$\leq -\log \left( 1 - \frac{2\varepsilon}{1 - 2p^{ab}} \right) \leq \frac{1}{4} W_*$$

□

Theorem 1: Let  $w_* = \min_e w(e)$  and  $W_* = \max_{a,b} \delta(a,b)$ .

Then  $k = O\left(\frac{e^{2W_*}}{(1 - e^{-w_*/4})^2} \cdot \log n\right)$  samples suffice

to get the correct tree  $T^p$  w/prob  $1 - o(1)$  as  $n \rightarrow \infty$ .

Proof: • From Chernoff bounds:

$$\mathbb{P}_r\left[|P_{ab} - \hat{P}_{ab}| \geq \varepsilon\right] \leq 2 \exp(-2\varepsilon^2 k) < \frac{1}{n^3} \tag{*}$$

choosing  $k = \Omega\left(\frac{1}{\varepsilon^2} \log n\right)$ .

• Use  $\varepsilon = \frac{1}{2} (1 - e^{-w_*/4}) \cdot e^{-W_*}$  (\*\*)

• By probability  $\geq 1 - \frac{1}{n}$ : (\*) and a union bound, it follows that w/ for all  $a, b \in [n]$ :

$$|P^{ab} - \hat{P}^{ab}| < \frac{1}{2} (1 - e^{-w_*/4}) (1 - 2p^{ab}).$$

• Hence <sup>by previous claim</sup> we'll get all quartets right w. pr.  $\geq 1 - \frac{1}{n}$  and therefore the correct tree w/prob.  $\geq 1 - \frac{1}{n}$ .

• From (\*), (\*\*) it follows that  $k = O\left(\frac{e^{2W_*}}{(1 - e^{-w_*/4})^2} \log n\right)$  suffices for this purpose.



- So assume that mutation probabilities are bounded away from 0 and  $1/2$ , i.e.  $0 < c_1 \leq p_e \leq c_2 < 1/2$ , for all edges  $e$ .

Then 
$$0 < \underbrace{-\log(1-2c_1)}_f \leq W_e \leq \underbrace{-\log(1-2c_2)}_g < +\infty \quad (f, g \text{ are constants})$$

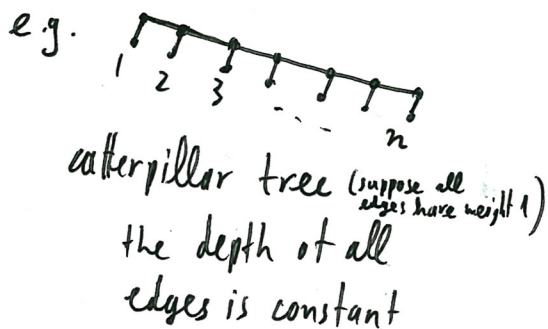
- In this case,  $\min_e W_e \geq f$  and  $W^* \leq (n+1)g$

So from previous theorem:  $k = 2^{\Omega(n)}$  suffices.

- This is not tight; it turns out that Thm 1 holds if we replace  $W$  by the weighted depth of the tree.

def (Weighted Depth): The depth of an edge is the length (under  $\delta$ ) of the shortest path between two leaves crossing  $e$ . The depth of a tree is the largest depth of an edge in the tree.

e.g.



full binary tree (suppose  $w_e = 1$ )  
the depth is at most  $2 \log_2 n$

- Ex (opt): If  $w_e = 1, \forall e$ , then the depth of an edge  $e$  in any binary tree is at most  $2 \log_2 n + 2$ .

• Hence, if  $f \leq w \leq g, \forall e$ , and we use the modified version of Thm 1, it follows that w/  $h = \text{poly}(n)$  samples from the CFN model we can reconstruct the tree w/ prob.  $1 - o(1)$  as  $n \rightarrow \infty$ .

• Can we do better than  $\text{poly}(n)$  sequence length?  
 (recall that our counting lower bound was just  $\Omega(\log n)$ )

**Steel's Conjecture:**

Let  $\theta^* = 1/2$ . Then, if  $\theta(e) \leq \theta^*, \forall e$ ,  $\text{poly}(n)$  samples are necessary for the tree reconstruction problem, while if  $\exists \theta(e) > \theta^*$ ,  $O(\log n)$  samples suffice.

More on Steel's Conjecture next time!