

Atypical decay rates for atypical heights in RRT

Heng Ma 馬恒 (Technion)

Joint work (in progress) with

Xinxin Chen 陈昕昕

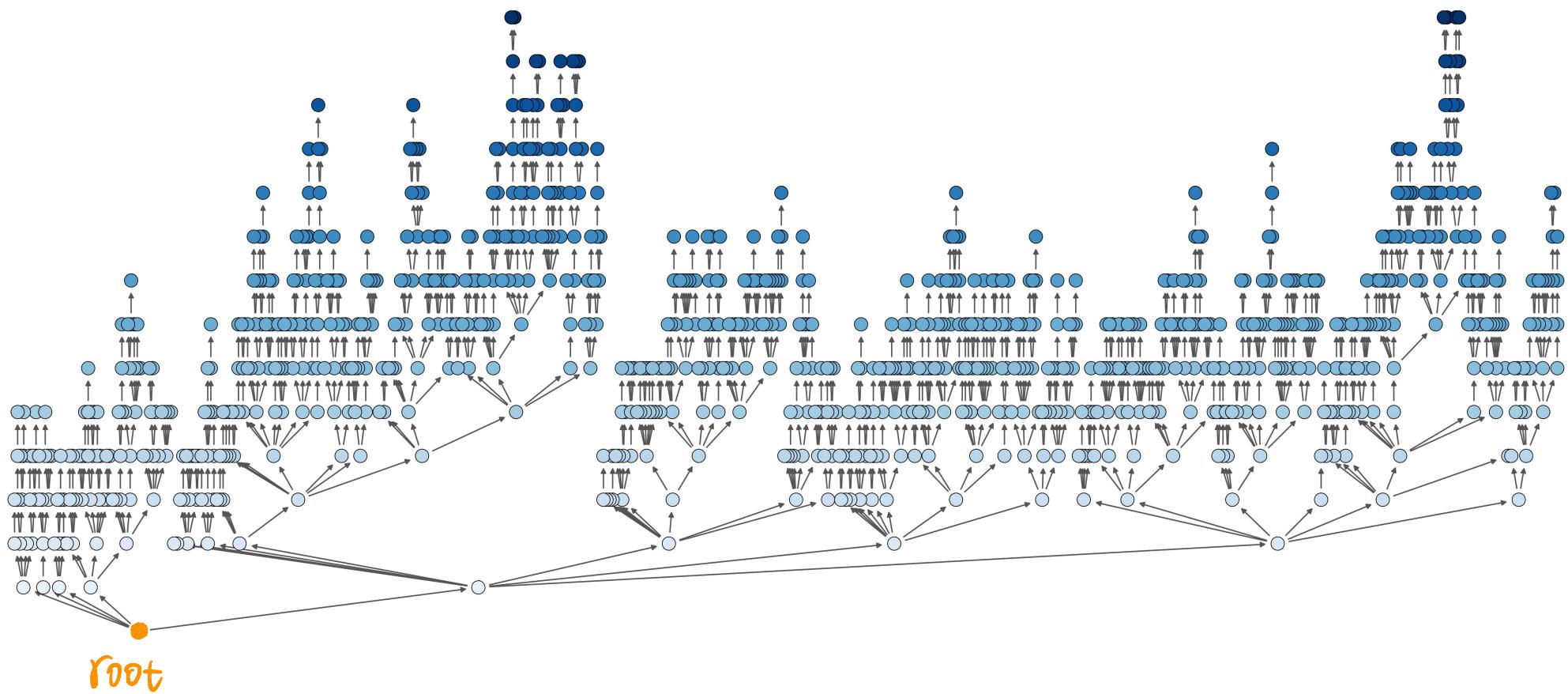
(Beijing Normal University)

Random recursive tree (RRT) T_n

Q

- Start with a root v_1
- At each step $n \geq 2$, introduce a new vertex v_n .
- We attach v_n to an existing vertex, chosen uniformly at random from $\{v_1, \dots, v_{n-1}\}$.

Also known as the "Uniform attachment tree".

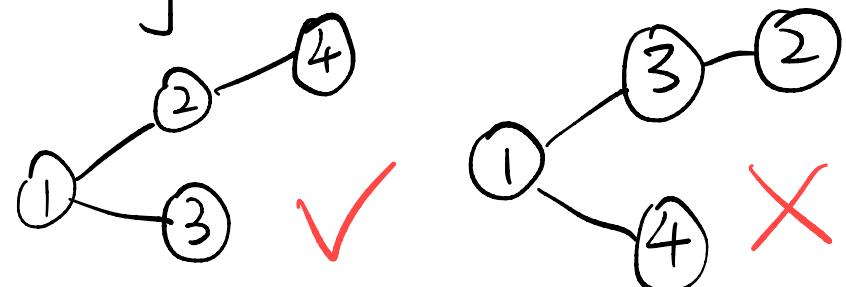


RRT on 1000 nodes

As uniform distribution on increasing trees

- A tree on $\{v_1, \dots, v_n\}$ is called an **increasing tree** if

- 1) The root has smallest label v_1
- 2) The label along any path starting from the root is increasing



- For any increasing tree T ,

$$\Pr(\mathcal{T}_n = T) = \frac{1}{(n-1)!}$$

Remark (General attachment rules)

- Start with v_1 .
- At each step $n \geq 2$, introduce a new vertex v_n .
- Attach v_n to an existing vertex $v_{I_n} \in \{v_1, \dots, v_{n-1}\}$, where

$$P(I_n = j \mid T_{n-1}) \propto \theta \cdot \underbrace{\deg_{T_{n-1}}^+(v_j)}_{\text{Out degree of } v_j} + 1 ,$$

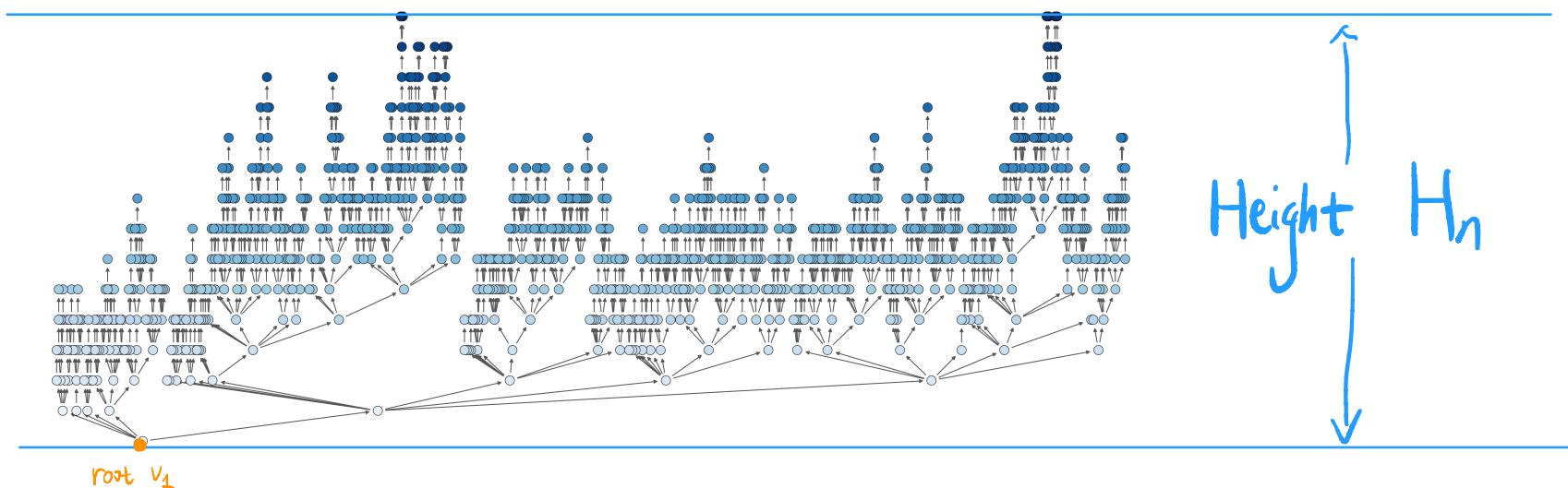
$\theta = 0$	RRT	Out degree of v_j
$\theta = 1$	preferential attachment tree	
$\theta = -\frac{1}{2}$	Binary Search tree	Plane Oriented Recursive Tree

Typical Behaviors
of RRT

The Height of RRT

- Denote by H_n the height of the tree T_n :

$$H_n := \max_{1 \leq j \leq n} d(v_1, v_j)$$



- Szymanski' 90 : with high prob. $1 \leq \frac{H_n}{\log n} \leq e$.
- Pittel '94 proved that $\lim_{n \rightarrow \infty} \frac{H_n}{\log n} = e$ almost surely,
by using a Continuous-time Embedding method
- Addario-Berry - Ford '13 :

$$\left\{ \begin{array}{l} \mathbb{E} H_n = e \log n - \frac{3}{2} \log \log n + O(1) \\ (H_n - \mathbb{E} H_n)_{n \geq 1} \text{ is tight} : \forall \lambda \in (0, \frac{1}{2e}), n \geq 1 \end{array} \right.$$

$$(H_n - \mathbb{E} H_n)_{n \geq 1} \text{ is tight} : \forall \lambda \in (0, \frac{1}{2e}), n \geq 1$$

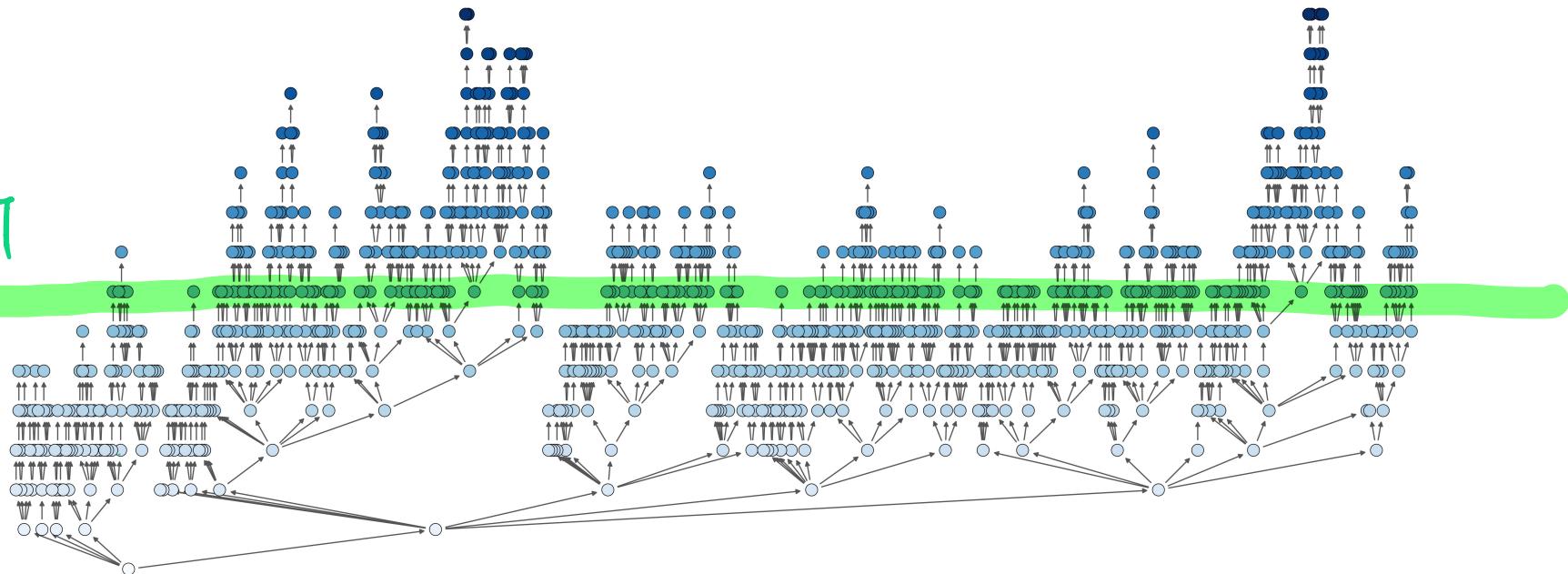
$$\mathbb{P}(|H_n - \mathbb{E} H_n| > k) \lesssim_\lambda e^{-\lambda k}$$

Level Sets of RRT

- Let $X_{n,k} :=$ number of nodes at level k in \mathcal{T}_n
 $= \sum_{j=1}^n \mathbb{1}_{\{d(v_1, v_j) = k\}}$

$(X_{n,k})_{k \geq 1}$ is also called the profile of \mathcal{T}_n in the literature.

Level T



Fuchs-Hwang-Neininger'06:

- For $0 \leq \alpha < e$, if $k = k(n)$ satisfies $\frac{k}{\log n} \rightarrow \alpha$ then

$$\mathbb{E} X_{n,k} = n^{\alpha(1-\log \alpha) + o(1)}, \quad \frac{X_{n,k}}{\mathbb{E} X_{n,k}} \xrightarrow{d} X_\alpha$$

\tilde{X} : an i.i.d.

where X_α satisfies

$$X \stackrel{d}{=} \alpha U^\alpha X + (1-U)^\alpha \tilde{X}$$

copy of X

$U \sim \text{Unif}[0,1]$

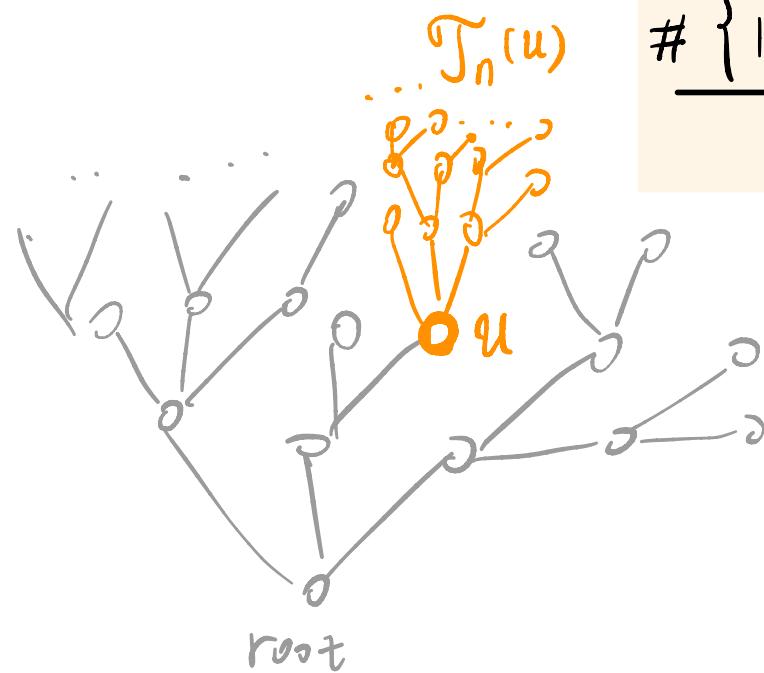
- CLT when $\alpha = 0$: If $k = k(n)$ satisfies $\frac{k}{\log n} \rightarrow 0$ then

$$\frac{(k-1)! \sqrt{2k-1}}{(\log n)^{k-\frac{1}{2}}} \left[X_{n,k} - \frac{(\log n)^k}{k!} \right] \xrightarrow[n \rightarrow +\infty]{d} \text{Normal}(0, 1)$$

Limiting Fringe Tree of RRT

Aldous' 91: Uniformly choose a node \underline{u}

in T_n . What does the subtree of u looks like?



$$\frac{\#\left\{1 \leq j \leq n : T_n(u) \subseteq \tau\right\}}{n} \xrightarrow{\text{a.s.}} P(T_M \subseteq \tau)$$

where M is a R.V. independent
of (T_n) , with $P(M \geq n) = \frac{1}{n}$.

Degree Profile of RRT

As a consequence,

$$\frac{\#\{1 \leq j \leq n : \deg_{T_n}(v_j) = k\}}{n} \xrightarrow[n \rightarrow +\infty]{\text{a.s.}} 2^{-k}$$

In particular,

proportion of leaves $\longrightarrow 50\%$

Our focus :

Rare events concerning
RRT height



RRT being atypically short

Fix $\alpha \in (0, 1)$, what is the decay rate of

$$\mathbb{P}(H_n \leq \alpha e^{\log n})$$



RRT being atypically tall

Fix $\beta > 1$, what is the decay rate of

$$\mathbb{P}(H_n \geq \beta e^{\log n})$$

Theorem 1 (Chen-M.'26+)

Define $m_\alpha(n) := n^{1-\alpha} (\log n)^{-\frac{3}{2\alpha}}$ and $w_\alpha(n)$ as follows

$$P(H_n \geq \alpha \ln n) = e^{-w_\alpha(n) m_\alpha(n)}$$

Then for $\alpha < 1$,

$$\liminf_{n \rightarrow \infty} w_\alpha(n) = +\infty$$

Moreover, for any integer $k \in \mathbb{N}$,

$$\limsup_{n \rightarrow \infty} \frac{w_\alpha(n)}{\log^{(k)}(n)} < +\infty.$$

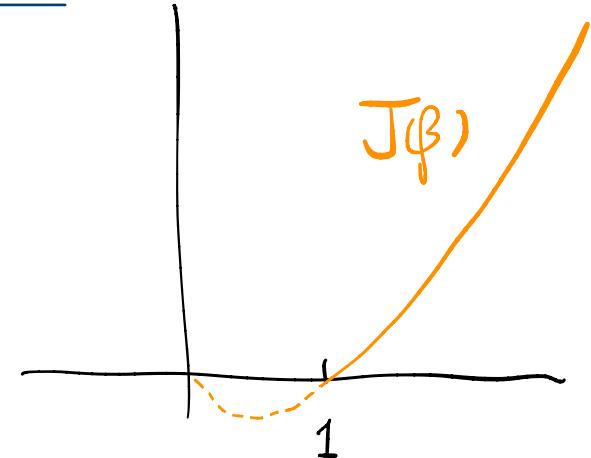
$$\begin{aligned}\log^{(k)}(x) \\ := \log^{(k-1)}(\log x)\end{aligned}$$

Theorem 2 (Chen-M.'26+)

- Define $J(\beta) := e \cdot \beta \ln \beta$

Then for $\beta > 1$, we have

$$\begin{aligned} \mathbb{P}(H_n > \beta e \log(n)) \\ = n^{-J(\beta) + o(1)} &= e^{-[J(\beta) + o(1)] \log n} \end{aligned}$$

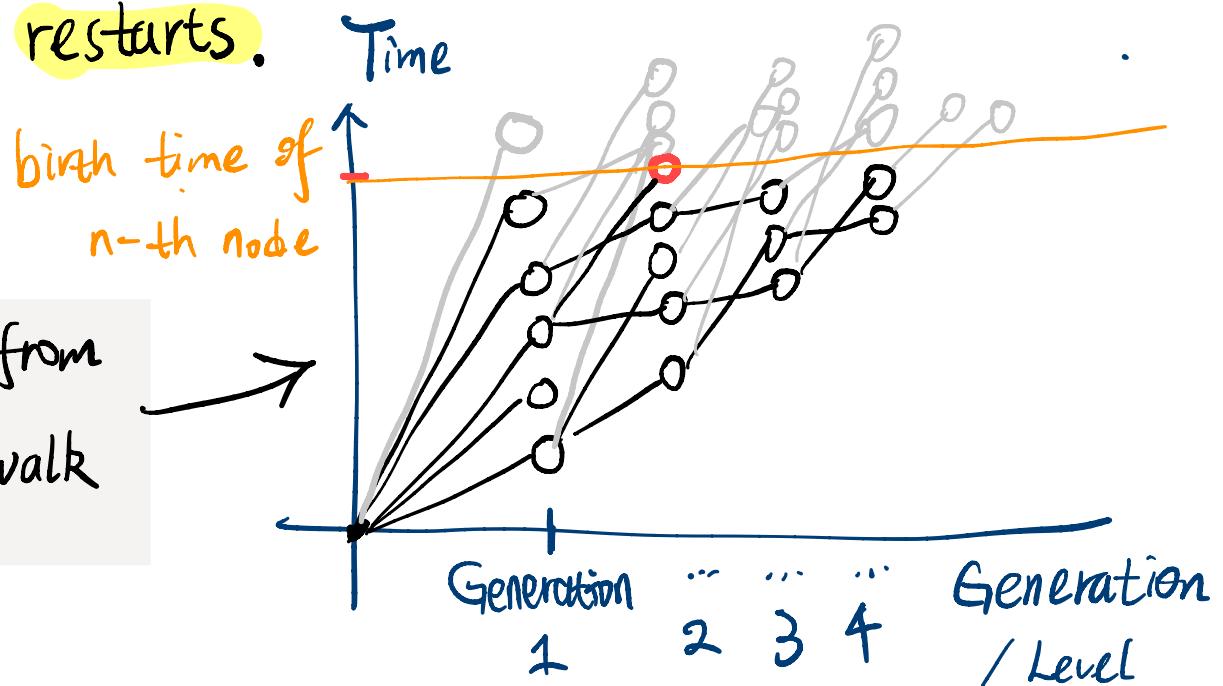


A powerful approach

Pittel's Continuous-time embedding

- Each node has an **independent** alarm clock ,
exponentially distributed with rate 1
- When it rings , this node give birth to a new child
and its clock **restarts**.

Extracting a RRT from
a branching random walk

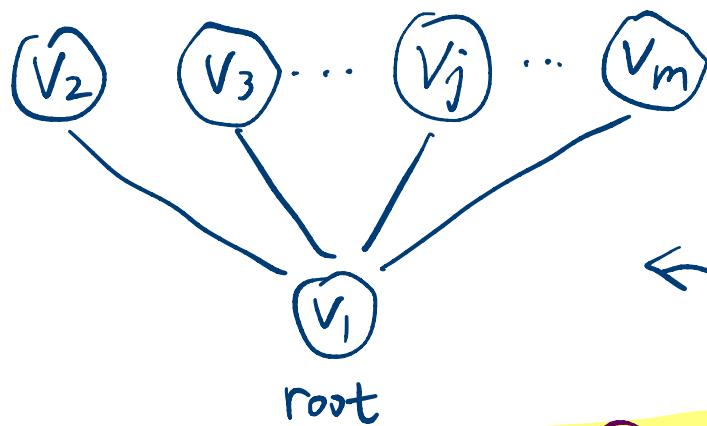


Proof ideas of Thm 1

Strategy to make
the tree atypically short

A crude strategy:

Make the tree as short as possible at an early stage , then let it grow "freely".

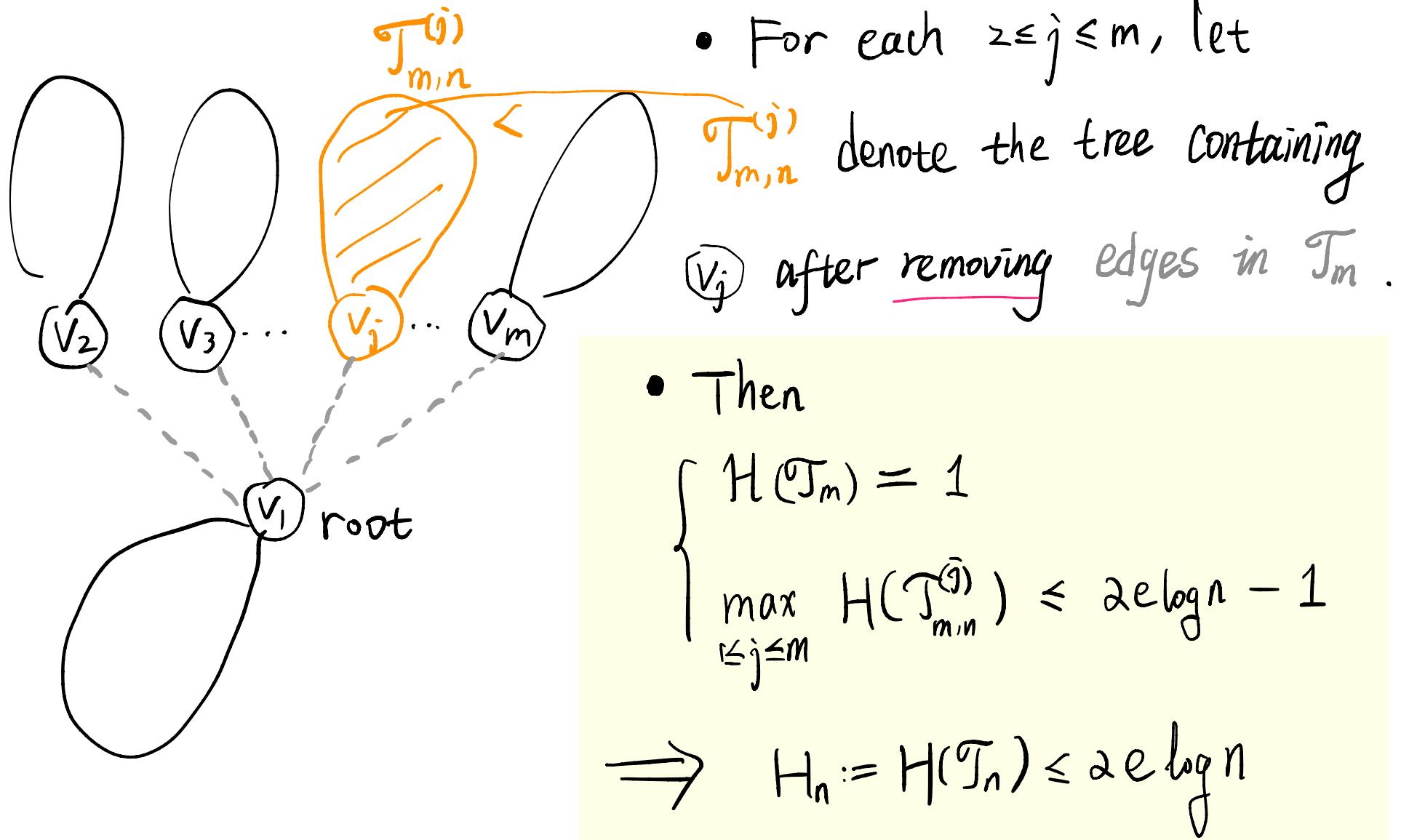


That is, fix $m \geq 2$,
we ask $\underline{H(T_m) = 1}$

Cost: $P(H(T_m) = 1)$

$$= \frac{1}{(m-1)!} = e^{-m \ln m + O(m)}$$

Make the tree as short as possible at an early stage, then let it grow "freely".



Observation 1: Conditionally on $\{\mathcal{T}_m = \mathcal{T}, |\mathcal{T}_{m,n}| = n_j, 1 \leq j \leq m\}$

$\mathcal{T}_{m,n}^{(j)}$ are independent RRT on n_j nodes.

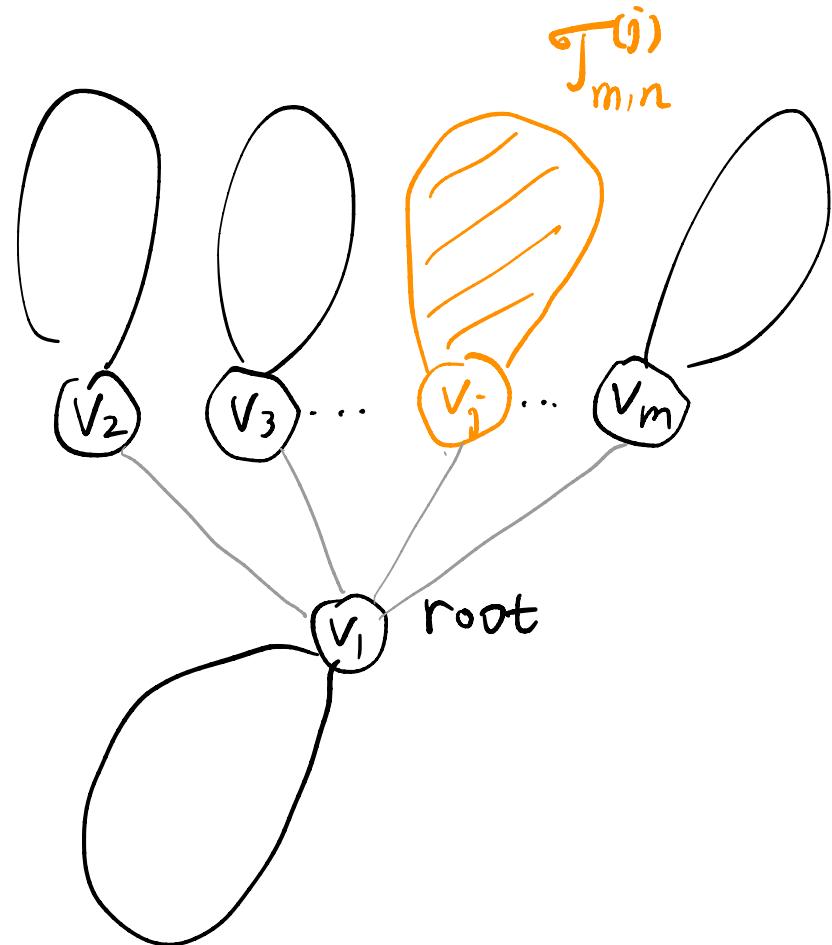
Observation 2:

$$(|\mathcal{T}_{m,n}^{(j)}| : 1 \leq j \leq m)$$

has the Uniform distribution on

$$\left\{ (n_j)_{j=1}^m : n_j \in \mathbb{N}; \sum_{j=1}^m n_j = n \right\}.$$

In particular, $\mathbb{E} |\mathcal{T}_{m,n}^{(j)}| = \frac{n}{m}$.



- Need a lower bound for

$$P\left(\max_{1 \leq j \leq m} H(T_{m,n}^{(j)}) \leq 2e \log n - 1\right)$$

- Suppose $|T_{m,n}^{(j)}| \leq 100 \frac{n}{m}$, then

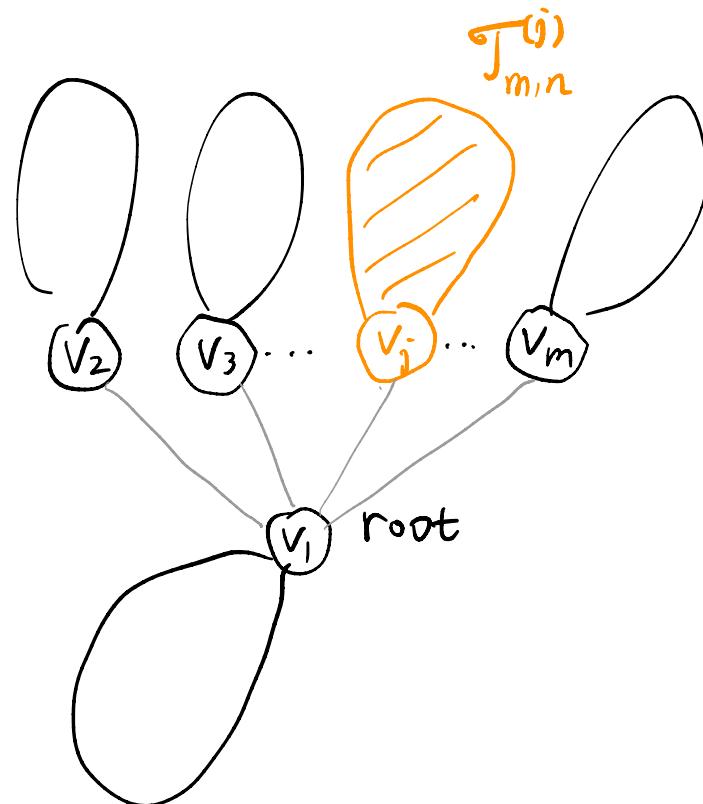
$$H(T_{m,n}^{(j)}) \stackrel{st}{\leq} H_{100 \frac{n}{m}} = \mathbb{E}[H_{\frac{n}{m}}] + O_p(1)$$

- Choose m carefully so that

$$\mathbb{E}[H_{\frac{n}{m}}] = 2e \log n - \text{large constant}$$

and hence

$$P(H_{100 \frac{n}{m}} \leq 2e \log n - 1) \geq \frac{1}{2}.$$



$$m = C_{st} \times m_2(n)$$

All together,

$$\mathbb{P}(H_n \leq 2e \ln n) \geq \underbrace{\mathbb{P}(H(T_m) = 1)}_{\frac{1}{(m-1)!}} = e^{-m \ln m + O(m)}$$

$$\times \mathbb{P}(|T_{m,n}^{(j)}| \leq 100 \frac{n}{m} \quad \forall 1 \leq j \leq m)$$

One can show it's greater than e^{-cm}

$$\times \underbrace{\mathbb{P}(H_{100 \frac{n}{m}} \leq 2e \log n - 1)}_m$$

by our choice of m , it is greater than 2^{-m}

$$\geq e^{-c m \ln m} \geq e^{-c m_2(n) \log n}$$

□

All together,

$$\mathbb{P}(H_n \leq 2e \ln n) \geq \boxed{\mathbb{P}(H(J_m) = 1)} \frac{1}{(m-1)!} = e^{-m \ln m + O(m)}$$

$$\times \mathbb{P}\left(|J_{m,n}^{(j)}| \leq 100 \frac{n}{m} \quad \forall 1 \leq j \leq m\right)$$

One can show it's greater than $\frac{m}{m}$

$$\times \mathbb{P}(H_{100 \frac{n}{m}} \leq 2e \log n - 1)$$

by our choice of m , it is greater than

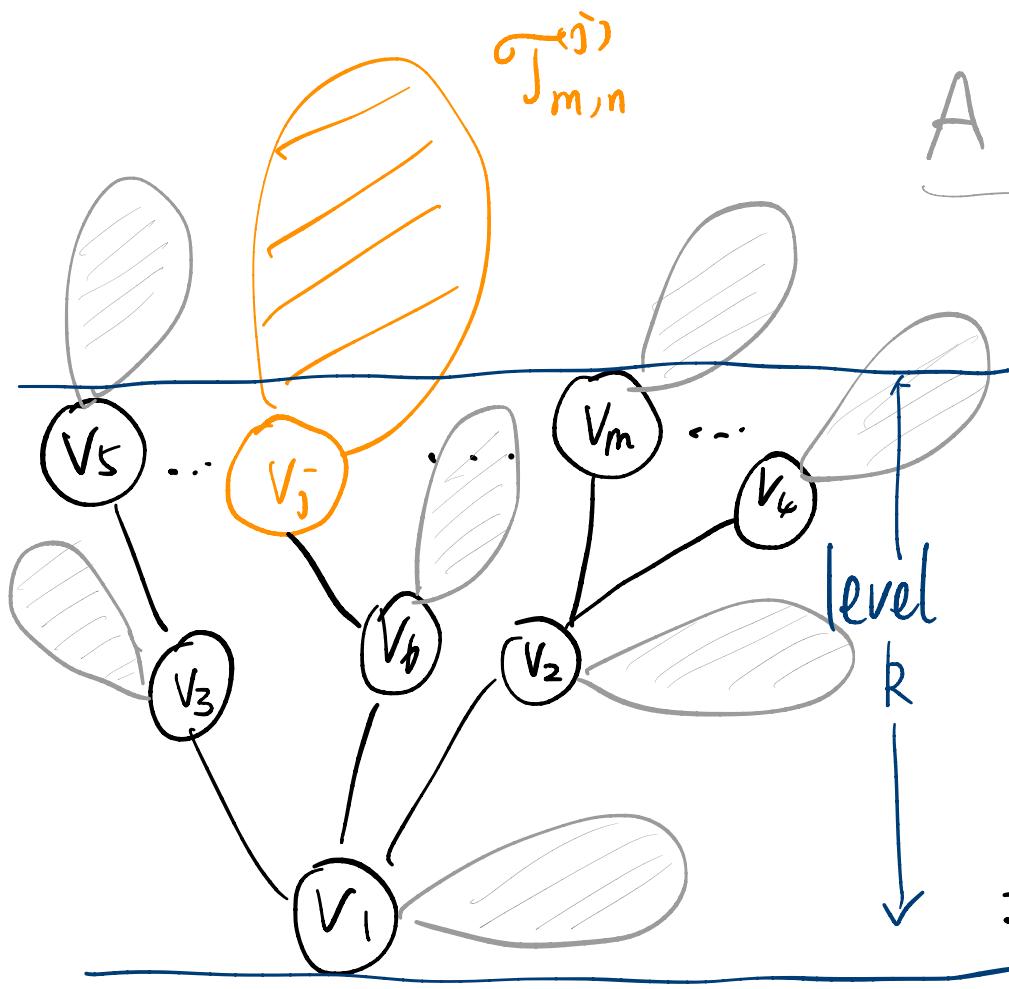
$$H(J_m) = 1$$

$$\geq e^{-c m \ln m}$$

$$\geq e^{-c m_2(n) \log n}$$

□

Just need
to improve
the crude
requirement



A slight improvement :

Fix any $k \in \mathbb{N}$. Then

$$H(\mathcal{T}_m) = k$$

$$\max_{1 \leq j \leq m} H(\mathcal{T}_{m,n}^{(j)}) \leq 2e \log n - k$$

$$H_n := H(\mathcal{T}_n) \leq 2e \log n$$

- $P(H(\mathcal{T}_m) = k) = \frac{A_k(m)}{(m-1)!}$ # Increasing trees with height k . on m nodes
- $\geq e^{-m \log^{(k)} m} = e^{-m \log^{(k)} n}$

□

Remarks

- We have a intuitive enumeration showing that

$$A_k(n) \geq \exp(n \log n - n \log^{(k)} n + O(n)) \quad (*)$$

- It is well known that, the generating function

$$F_k(z) = \sum_{n=1}^{\infty} \frac{A_k(n)}{n!} z^n = \sum_{n=1}^{\infty} \frac{1}{n} P(H_n \leq k) z^n$$

satisfies

$$\begin{cases} F_k'(z) = \exp(F_{k-1}(z)) \\ F_1(z) = e^z - 1; \quad F_k(0) = 0 \end{cases}$$

In principle, $(*)$ should be proved by analyzing $F_k(z)$

But it seems hard (for us) to do this ...

Proof ideas of Thm 1

Constraints on
many sub-structures

(1) Reselect m so that

$$\mathbb{E}[H_{\frac{1}{4}\frac{n}{m}}] = 2e \log n + \text{large constant}$$

and hence

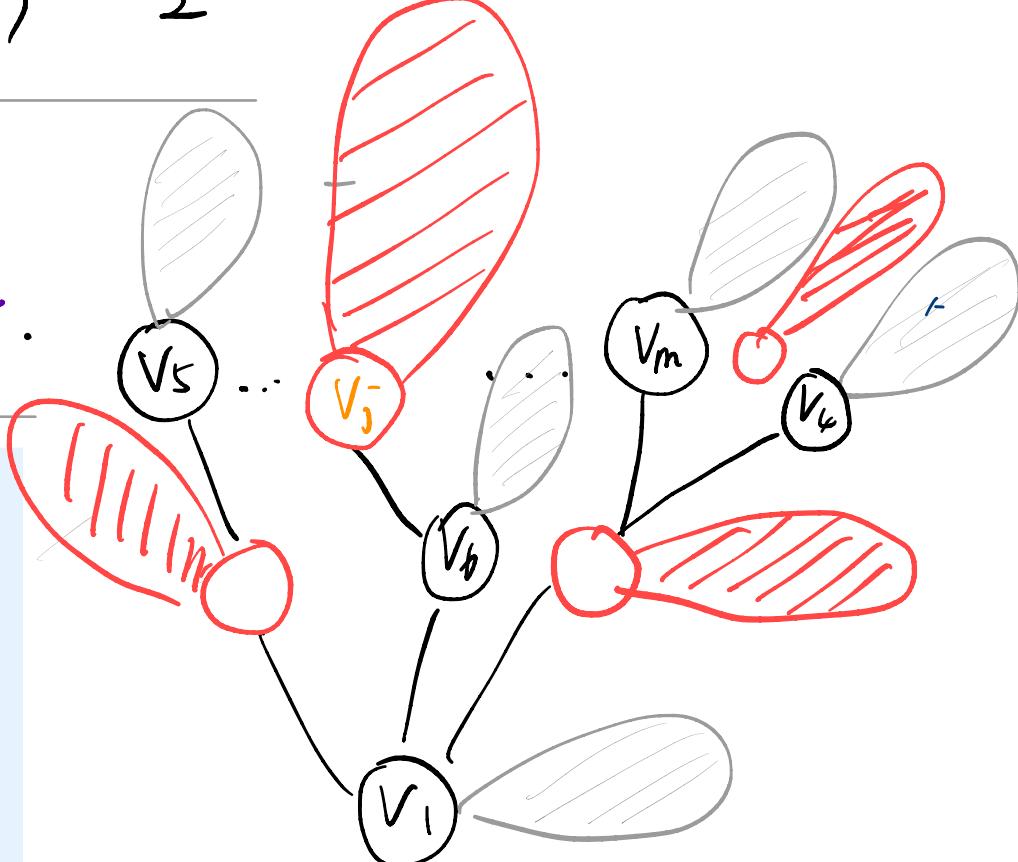
$$\mathbb{P}(H_{\frac{1}{4}\frac{n}{m}} \leq 2e \log n) < \frac{1}{2}$$

(2) Let BIG_{min} be the collection of

all $T_{m,n}^{(j)}$'s with $|T_{m,n}^{(j)}| \geq \frac{1}{4} \frac{n}{m}$.

Combining (2) with (1) :

$$\begin{aligned} \mathbb{P}(H_n \leq 2e \log n) \\ \leq \mathbb{E}[2^{-|\text{BIG}_{min}|}] \end{aligned}$$



$$\cdot \mathbb{P}(H_n \leq 2e\log n) \leq \mathbb{E}[2^{-|\text{BIG}_{m,n}|}]$$

CLAIM: $\mathbb{P}(|\text{BIG}|_{m,n} < \gamma m) \leq e^{-C_\gamma^* m}$

$$\leq 2^{-cm} + e^{-c'm} \leq e^{-c''m}$$

□

Recall that $(|\mathcal{T}_{m,n}^{(j)}|)_{j=1}^m \stackrel{\text{Unif}}{\sim} \left\{ (n_j)_1^m : n_j \in \mathbb{N}, \sum_j n_j = n \right\}$

$\xi_j := \mathbf{1}\{\mathcal{T}_{m,n}^{(j)} \text{ is small}\}$ are negative correlated

$$\mathbb{P}\left(\sum_j \xi_j \geq (1-\gamma)m\right) \leq \mathbb{P}\left(\sum_j \xi_j^{\text{IID}} \geq (1-\gamma)m\right)$$

$$\leq \exp\left\{-\frac{D(1-\gamma||p_0)m}{\gamma}\right\}$$

[Rmk: $\lim_{\gamma \downarrow 0} C_\gamma^* = +\infty$] ←

Where can be improved ?

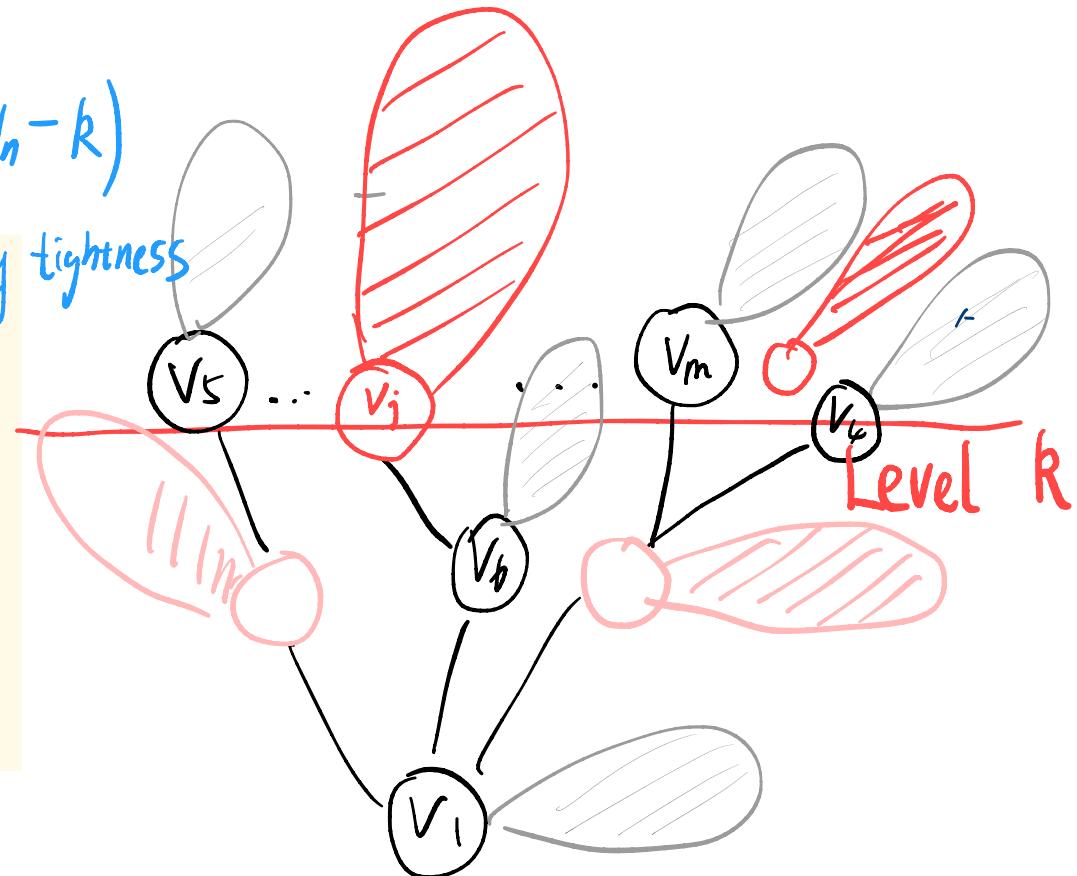
- Fix $k \in \mathbb{N}_+$

- Let $BIG_{m,n}^{(k)} = \left\{ 1 \leq j \leq m : \boxed{d(v_1, v_j) \geq k} \right. \\ \left. T_{m,n}^{(j)} \text{ is big; i.e., } |T_{m,n}^{(j)}| \geq \frac{1}{4} \frac{n}{m} \right\}$

Then

$$\delta(k) = \sup_{n \geq 1} P(H_n \leq \mathbb{E} H_n - k) \xrightarrow[k \rightarrow \infty]{\text{by tightness}} 0$$

$$P(H_n \leq \alpha \mathbb{E} H_n) \leq E[\delta(k) | BIG_{m,n}^{(k)}|]$$



Then for any $\varepsilon \in (0, 1]$. $\exists C_\varepsilon > 0$ s.t. for $n \geq 1$

$$\mathbb{P}\left(\sum_{l=0}^k X_{n,l} > \varepsilon n\right) \leq e^{-C_\varepsilon n \log^{(k)}(n)}$$

\nwarrow size of level l in T_n

- Consequently, $\mathbb{P}(|\text{BIG}_{m,n}^{(k)}| < \frac{\gamma}{2} m)$

$$\leq \mathbb{P}(|\text{BIG}|_{m,n} < \gamma m) + \mathbb{P}\left(\sum_{l=0}^k X_{m,l} > \frac{\gamma}{2}\right)$$

$$\lesssim e^{-C_\gamma^* m} + e^{-\Theta(m \log^{(k)} m)}$$

$$\text{BIG}_{m,n}^{(k)} = \left\{ 1 \leq j \leq m : d(v_1, v_j) \geq k \right. \\ \left. \text{J}_{m,n}^{(j)} \text{ is big, } |\text{J}_{m,n}^{(j)}| \geq \frac{1}{4} \frac{n}{m} \right\}$$

Then

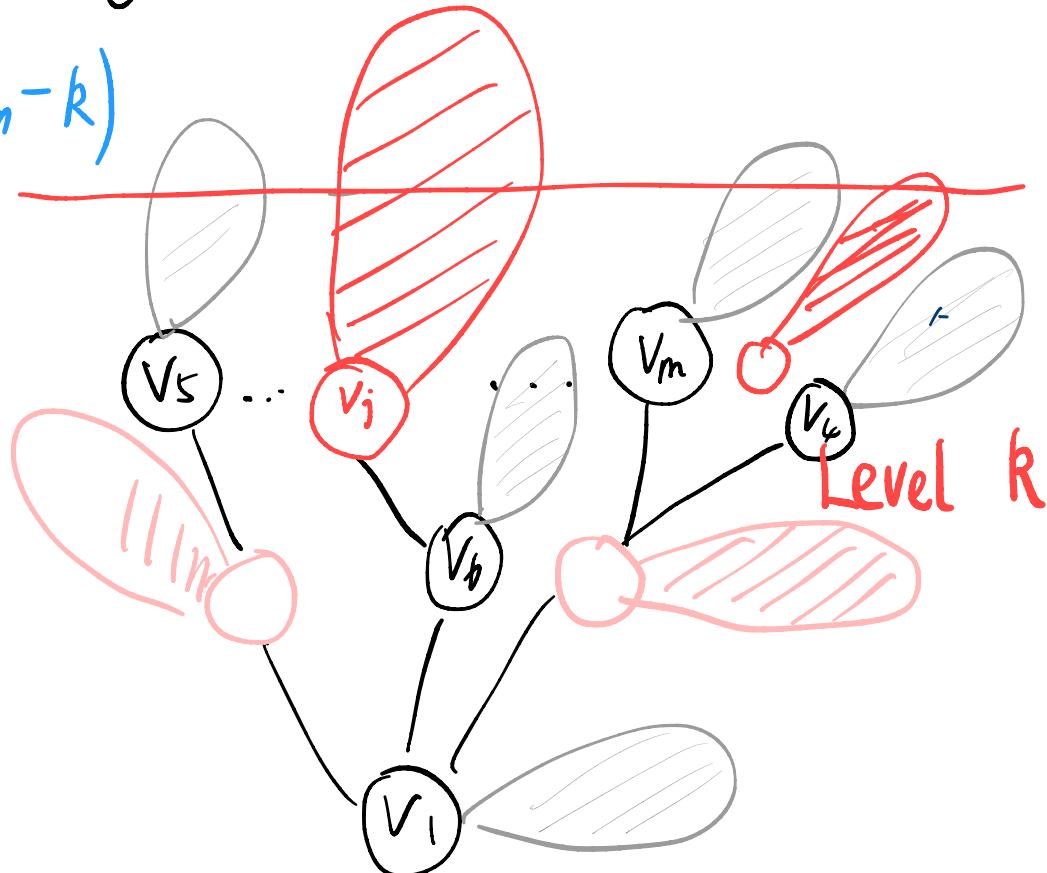
$$\delta(k) = \sup_{n \geq 1} P(H_n \leq \{E H_n - k\})$$

$$\Pr(H_n \leq \alpha \ln n) \leq \mathbb{E} \left[\delta(k) \Big| \text{BIG}_{m,n}^{(k)} \right]$$

$$\leq \delta(k)^{\gamma m} + e^{-C_r^* m}$$

$\leq e^{-Km}$ for any $K > 0$, provided γ small then

k large enough. 



Future directions

- Identifying the decay rate of $\omega_2(n)$
- LDP on the height of Preferential attachment tree
binary search tree and other random tree models .
-

Thanks for
your attention //

o

As a shortest-path tree

- K_n = complete graph on $\{v_1 \dots v_n\}$
- assign **IID exponential weights** (mean 1) to each edge

The tree T_n , formed by the shortest-path from v_1 to all other nodes $\{v_2 \dots v_n\}$, is a RRT.

