

Most trees are short and fat.

Louigi Addario-Berry

McGill

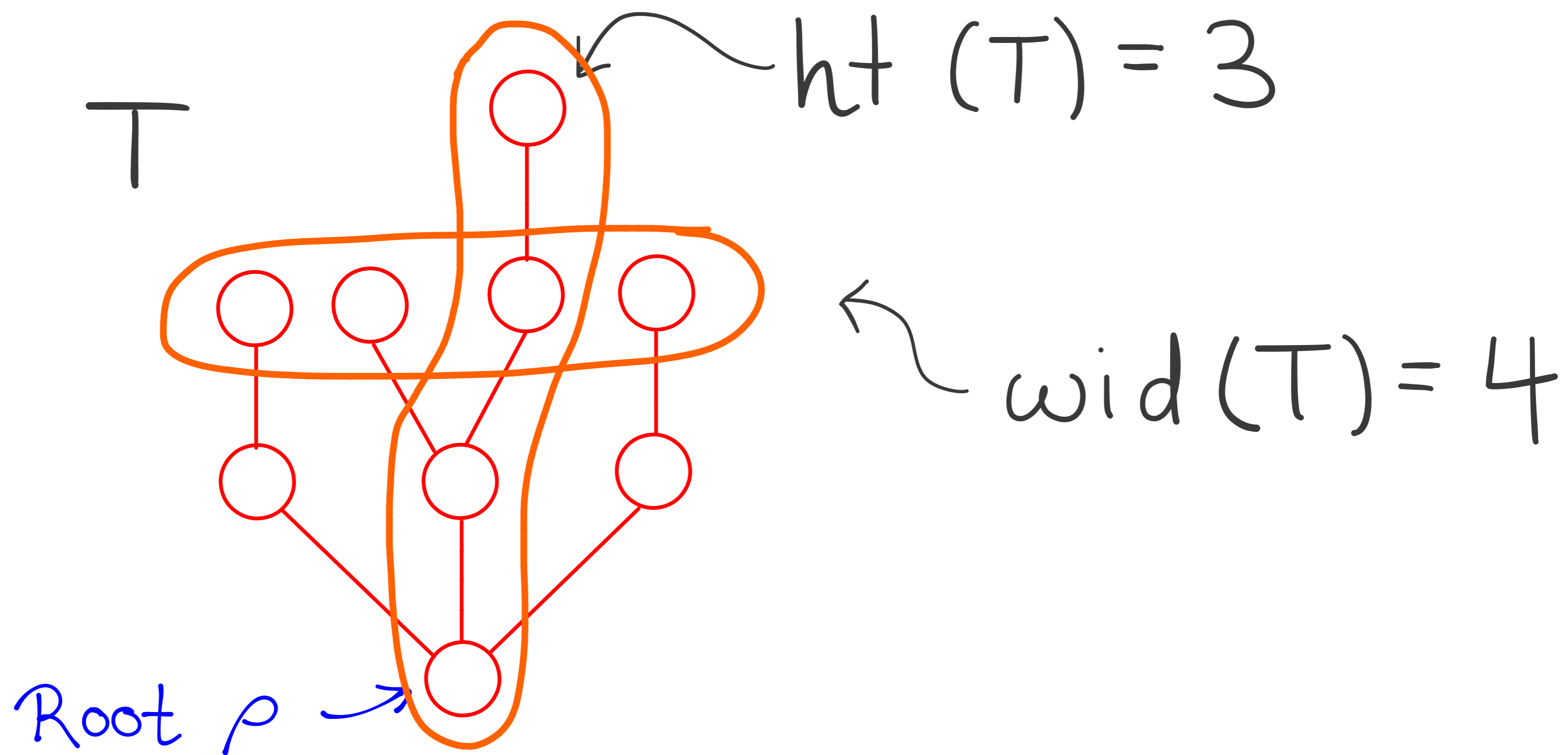
Journées de combinatoire

de BORDEAUX

January 25, 2017



Trees:

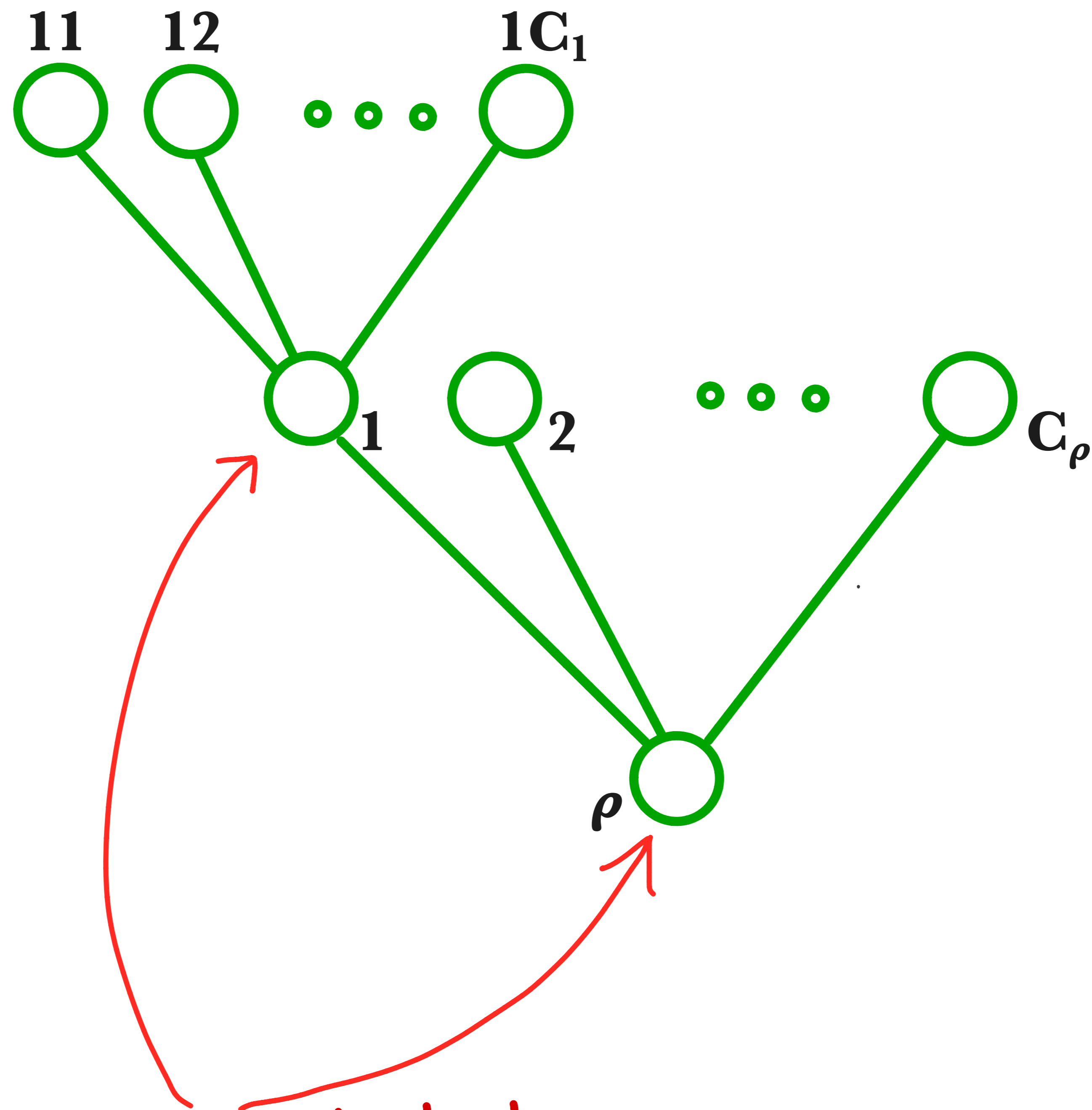


Height: Greatest distance from any node to the root } ht(T)

Width: Greatest # nodes on a single level. } wid(T)

Main Results Fix any r.v. C with $\sum_{k \geq 0} P(C=k) = 1$, write $p_k = P(C=k)$.

Let T be $GW(C)$ distributed.



children is distributed as C ,
independently at each node.

Main Results Fix any r.v. C with $\sum_{k \geq 0} \mathbb{P}(C=k) = 1$, write $p_k = \mathbb{P}(C=k)$.

Let T be $\text{GW}(C)$ distributed.

Theorem ("Most trees are short & fat")

There is a universal constant $\delta > 0$ s.t.

$$\mathbb{P}(\text{ht}(T) \geq \frac{k}{1-p_1} \cdot \text{wid}(T)) \leq \exp(-\delta k).$$

Remark: If $\mathbb{E}C > 1$ then $\mathbb{P}(\sigma = \infty) > 0$, and

$$\mathbb{P}(\text{ht}(T) = \text{wid}(T) = \infty \mid \sigma = \infty) = 1.$$

Also, given that $\sigma < \infty$, the cond. dist. of T is $\text{GW}(\hat{C})$ where $\mathbb{P}(\hat{C}=1) = p_1$,

$$\mathbb{E}(\hat{C}) < 1 \text{ so can assume } \mathbb{E}C \leq 1.$$

Heuristic: GW trees satisfy $\text{wid}(T) \cdot \text{ht}(T) \cong \text{vol}(T) := \sigma$

Implies " $\text{ht} > C^2 \cdot \text{wid}$ " \cong " $\text{ht}^2 \geq C^2 \cdot \text{vol}$ " so $\mathbb{P}(\text{ht}(T) > \frac{k}{\sqrt{1-p_1}} \sqrt{\text{vol}(T)}) \leq \exp(-\delta k^2)$

Theorem: $\mathbb{P}(\text{ht}(T) \geq \frac{k}{\sqrt{1-p_1}} \text{vol}(T)^{1/2}) \leq \exp(-\delta k^2)$

Setup

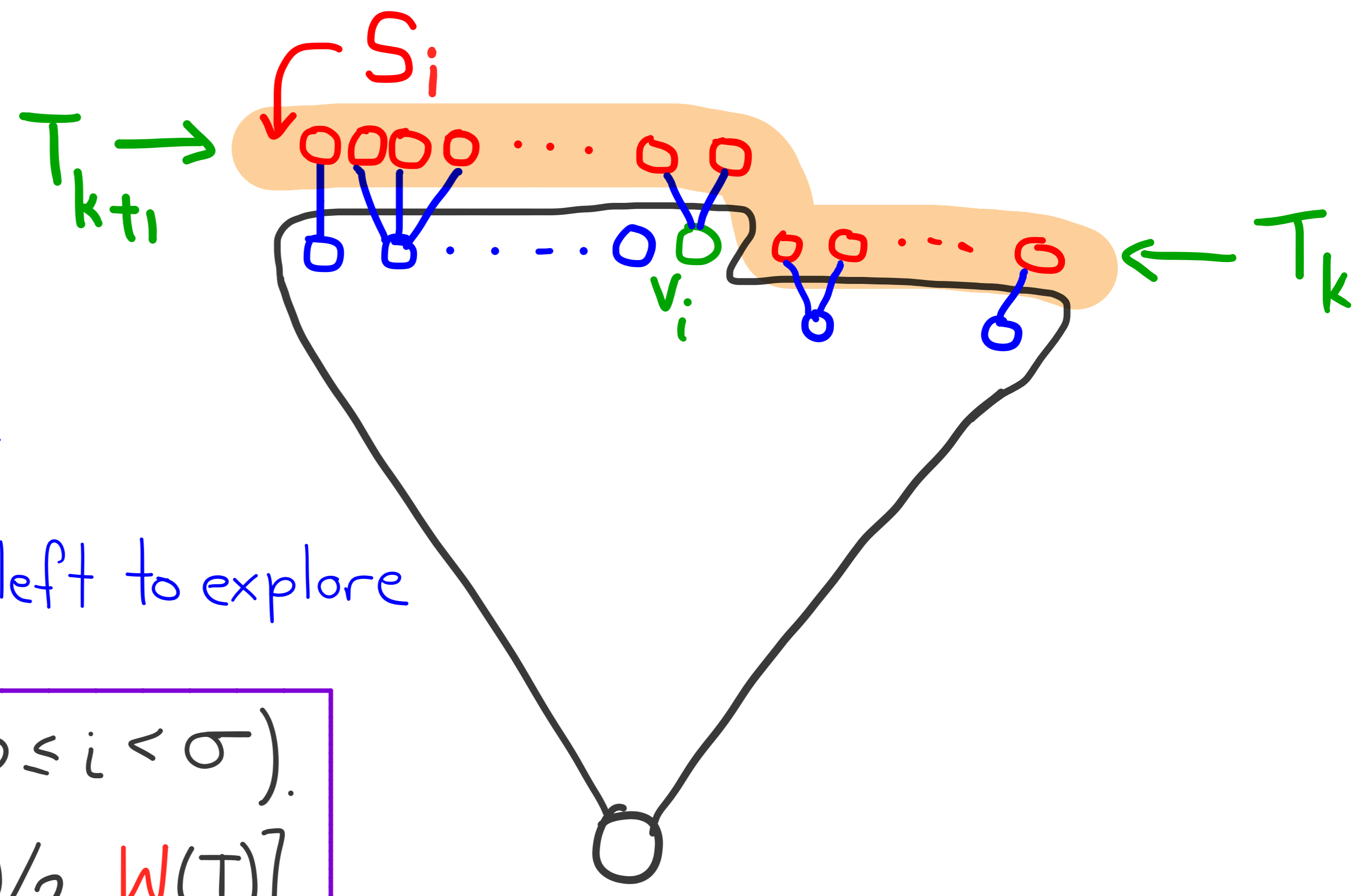
$$1 + \sum_{j=1}^i C_j = \# \text{ nodes discovered by time } i.$$

$$\text{Let } S_i = 1 + \sum_{j=1}^i (C_j - 1) \\ = \# \text{ nodes in "BFS queue" at time } i$$

$$\mathbb{E} C \leq 1 \Rightarrow \mathbb{E} S_n = 1 + n(\mathbb{E} C - 1) \leq 1.$$

$\sigma = \inf\{t : S_t = 0\}$ = first time no nodes left to explore

Prop: Let $W(T) = \max(S_i, 0 \leq i < \sigma)$.
Then $\text{wid}(T) \in (W(T)/2, W(T)]$



Proof: During BFS on level k , "exploration queue" $\subset T_k \cup T_{k+1}$, and $= T_k$ at start of level k . ■

Idea: $ht(T) = \sum_{k=1}^{ht(T)} 1 = \sum_{k=1}^{ht(T)} \sum_{v \in T_k} \frac{1}{|T_k|}$.

Prop:
 $ht(T) \leq 3H(T)$.

When $v_i \in T_k$ then $S_i \approx |T_k|$ so perhaps

$$ht(T) \approx \sum_{k=1}^{ht(T)} \sum_{v_i \in T_k} \frac{1}{S_i} = \sum_{i=1}^n \frac{1}{S_i} =: H(T)?$$

[False; consider a star with n leaves. But...]

Corollary Suffices to prove
 $\mathbb{P}(H(T) \geq \frac{k}{1-p} W(T)) \leq e^{-\delta k}$,
thm. follows.

$$W(\sigma) = \max(S_i, 0 \leq i < \sigma) \quad H(\sigma) = \sum_{i=1}^{\sigma} \frac{1}{S_i} \quad \text{Aim: } \mathbb{P}(H(\sigma) \geq \frac{k}{1-p_1} W(\sigma)) \leq e^{-\delta k}$$

Key Tool: Decomposition into scales.

When $S_j \approx 2^l$ ("scale l ") for $j \in \{i, \dots, i+2^l\}$, have

$$H(i+2^l) - H(i) = \sum_{j=i+1}^{i+2^l} \frac{1}{S_j} \approx 2^l \cdot \frac{1}{2^l} = 1.$$

So bound (a) time to change scales,

(b) "# visits to scales" = $(M(l), l \geq 1)$

$$W(\sigma) = \max(S_i, 0 \leq i < \sigma) \quad H(\sigma) = \sum_{i=1}^{\sigma} \frac{1}{S_i} \quad \text{Aim: } \mathbb{P}(H(\sigma) \geq \frac{k}{1-p_1} W(\sigma)) \leq e^{-\delta k}$$

Decomposition into scales.

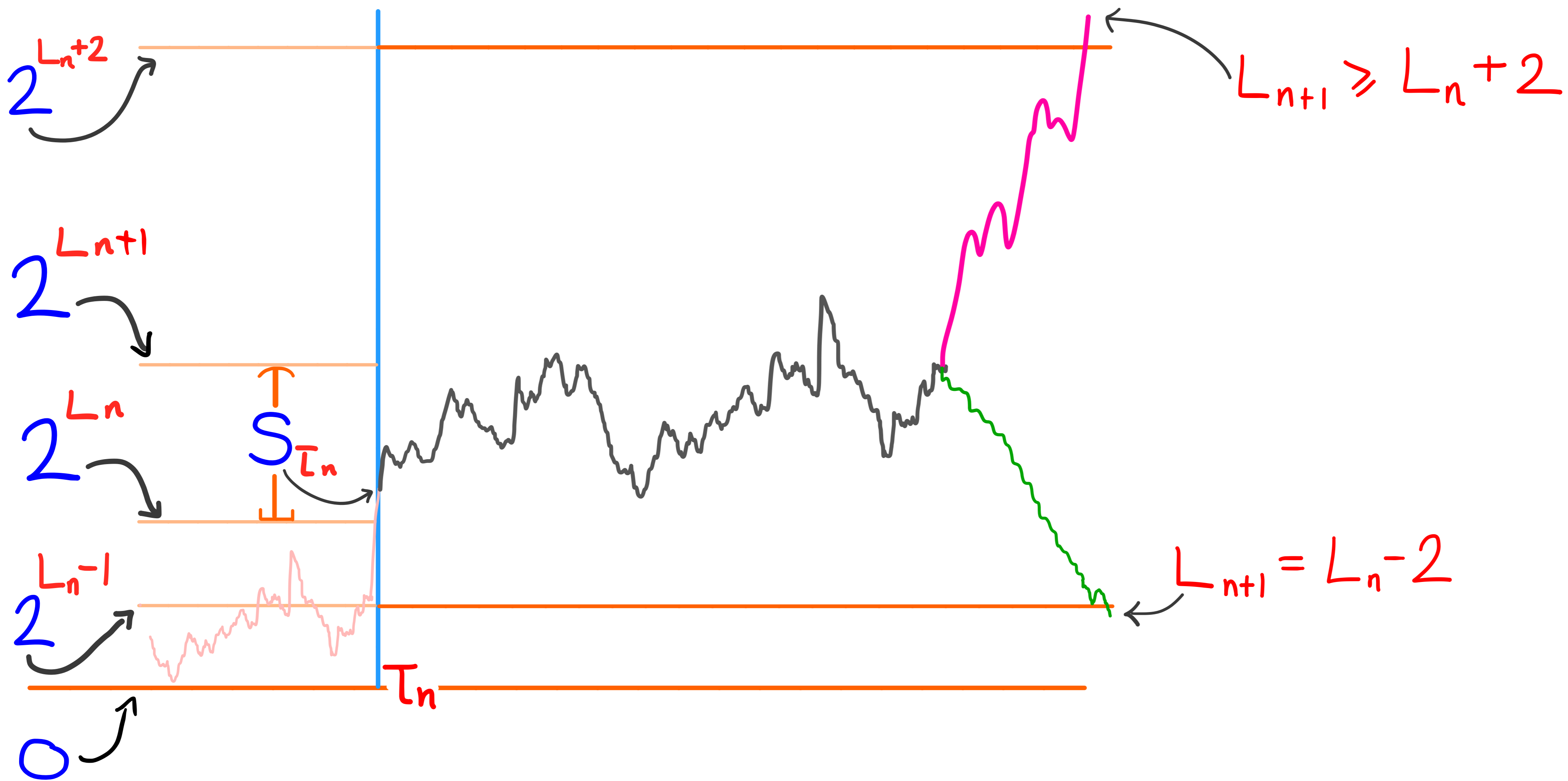
Def

$T_0 = 0 =$ initial time

$L_0 = 0 = \log_2 S_0 =$ initial scale

$$T_{n+1} = \min\{t \geq T_n : S_t \notin [2^{L_n-1}, 2^{L_n+2})\}$$

$$L_{n+1} = \sup\{k : 2^k \leq S_{T_{n+1}}\}$$



NB: $S_{T_n} < 2^{L_{n+1}}$ so $\mathbb{P}(L_{n+1} < L_n) = \mathbb{P}(\text{Hit } 2^{L_{n-1}} \text{ before } 2^{L_{n+2}}) > \frac{1}{2}$.

$$W(\sigma) = \max(S_i, 0 \leq i < \sigma) \quad H(\sigma) = \sum_{i=1}^{\sigma} \frac{1}{S_i} \quad \text{Aim: } \mathbb{P}(H(\sigma) \geq \frac{k}{1-p} W(\sigma)) \leq e^{-\delta k}$$

Key Tool: Decomposition into scales.

When $S_j \approx 2^l$ ("scale l ") for $j \in \{i, \dots, i+2^l\}$, have

$$H(i+2^l) - H(i) = \sum_{j=i+1}^{i+2^l} \frac{1}{S_j} \approx 2^l \cdot \frac{1}{2^l} = 1.$$

So bound (a) time to change scales,

(b) "# visits to scales" = $(M(l), l \geq 1)$

(a) Thm (Lévy; Doeblin; Kolmogorov; Rogozin; Le Cam; Esséen; Kesten):

With $p = \max p_i$, have

$$\max_k \mathbb{P}(S_n = k) \leq \frac{Cp}{\sqrt{n(1-p)}} \quad C > 0 \text{ universal.}$$

"Any random walk spreads out over $\geq \sqrt{n}$ values by time n ". Here $\sqrt{n} \approx 2^l$.

Thus if $L_k = l$ then we expect $\tau_{k+1} - \tau_k \lesssim 4^l$ so

$$H(\tau_{k+1}) - H(\tau_k) = \sum_{j=\tau_k+1}^{\tau_{k+1}} \frac{1}{S_j} \lesssim 4^l \cdot \frac{1}{2^l} = 2^l \approx 2^{L_k}.$$

More precisely, obtain $\mathbb{P}(H_{\tau_{k+1}} - H_{\tau_k} \geq \frac{c}{1-p} l \mid L_k = l) \leq e^{-\delta c}$

$$W(\sigma) = \max(S_i, 0 \leq i < \sigma) \quad H(\sigma) = \sum_{i=1}^{\sigma} \frac{1}{S_i} \quad \text{Aim: } \mathbb{P}(H(\sigma) \geq \frac{k}{1-p}, W(\sigma)) \leq e^{-\delta k}$$

Key Tool: Decomposition into scales.

When $S_j \approx 2^l$ ("scale l ") for $j \in \{i, \dots, i+2^l\}$, have

$$H(i+2^l) - H(i) = \sum_{j=i+1}^{i+2^l} \frac{1}{S_j} \approx 2^l \cdot \frac{1}{2^l} = 1.$$

So bound (a) time to change scales,

(b) "# visits to scales" = $(M(l), l \geq 1)$

(b) $M(l) = \# \text{ visits to scale } l = \#\{i : L_i = l\}$

$N(l, i) = \text{Duration of } i\text{th visit to scale } l$

$N(l) = N(l, 1) + \dots + N(l, M(l)) = \text{Total duration at scale } l.$

Fact: Given that $M(l) \neq 0$, $M(l)$ dominated by sum of 2 $\text{Geom}(\frac{1}{2})$ r.v.s; $\Rightarrow \mathbb{P}(M(l) > k \mid M(l) > 0) \leq 2^{-k/2}.$

Proof: visits to scale l entail upcrossings of $[2^{l-1}, 2^{-l})$ or of $[2^{l+1}, 2^{l+2})$.

Both are hard since walk has non-positive drift. \blacksquare

$$W(\sigma) = \max(S_i, 0 \leq i < \sigma) \quad H(\sigma) = \sum_{i=1}^{\sigma} \frac{1}{S_i} \quad \text{Aim: } \mathbb{P}(H(\sigma) \geq \frac{k}{1-p} W(\sigma)) \leq e^{-\delta k}$$

Key Tool: Decomposition into scales.

When $S_j \approx 2^l$ ("scale l ") for $j \in \{i, \dots, i+2^l\}$, have

$$H(i+2^l) - H(i) = \sum_{j=i+1}^{i+2^l} \frac{1}{S_j} \approx 2^l \cdot \frac{1}{2^l} = 1.$$

So bound (a) time to change scales,

(b) "# visits to scales" = $(M(l), l \geq 1)$

(a) $\tau_{k+1} - \tau_k \sim (2^{L_k})^2$ so $\mathbb{P}(H_{\tau_{k+1}} - H_{\tau_k} \geq \frac{c}{1-p} \cdot 2^l \mid L_k = l) \leq e^{-\delta c}$

(b) $M(l) = \# \text{ visits to scale } l$ $\mathbb{P}(M(l) > k \mid M(l) > 0) \leq 2^{-k/2}$

(a)+(b) \Rightarrow Total contribution of level l to height is

$\begin{cases} 0 & \text{if } M(l) = 0. \\ O(2^l) & \text{with exp. tails if } M(l) > 0. \end{cases}$

So $H(\sigma) \leq 2^{\max(l: M(l) \neq 0)} \leq \underline{\underline{4 \cdot W(\sigma)}}$

Remarks

- Stronger results if add info. about tails of degrees.

Ex. • If $\mathbb{P}(C \geq k) = \Theta(t^{-\alpha})$, $\alpha \in (1, 2)$,

then $\mathbb{P}(\text{ht}(T) > A \cdot m \cdot \text{wid}(T)^{\alpha-1}) \leq 2^{-\delta m}$

• If $\text{Var}(C) = \infty$ then $\mathbb{P}(\text{ht}(T) > A \cdot m \cdot \text{wid}(T)) \leq e^{-m \cdot f(m)}$; $f(m) \xrightarrow{m \rightarrow \infty} \infty$

FALSE

e.g. if $p_0 > 0$, $p_2 > 0$ then

$$\mathbb{P}(\text{ht}(T) > m \cdot \text{wid}(T)) \geq \mathbb{P}\left(\begin{array}{c} \circ \\ | \\ \circ_1 - \circ_2 - \dots - \circ_m \\ | \\ \circ \end{array}\right)$$
$$= p_2^m p_0^{m+1} > e^{-C \cdot m}.$$

But this is the only problem.

Thm: If $\text{Var}(C) = \infty$ then for any $\varepsilon(m) \rightarrow 0$ there is $f(m) \rightarrow \infty$ such that

$$\mathbb{P}(\text{ht}(T) > m \cdot \text{wid}(T), \text{ht}(T) < \varepsilon(m) \cdot \text{vol}(T)) \leq \exp(-m \cdot f(m))$$

Remarks

- Stronger results if add info. about tails of degrees.

Ex. • If $P(C \geq k) = \Theta(t^{-\alpha})$, $\alpha \in (1, 2)$,
then $P(\text{ht}(T) > A \cdot m \cdot \text{wid}(T)^{\alpha-1}) \leq 2^{-\delta m}$

- Conjecture: All this works even conditional on size of tree: $P(\text{ht}(T) > A \cdot m \cdot \text{wid}(T) \mid \sigma = n) \leq 2^{-\delta m}$.
- Conjecture: Binary trees are the tallest.

Consider random trees $T_{\vec{n}}$ with a fixed degree seq $\vec{n} = (n_i, i \geq 1)$

Here $n_i = \# \text{ nodes of deg } i$. With $\sum n_i = n$, then $\sum i n_i = 2(n-1)$.

To stochastically maximize $\text{ht}(T_{\vec{n}})$ among sequences with $n_0 = k$, $n_1 = 0$,
choose the seq. $(k, 0, k-1, 0, \dots)$

I can prove binary trees stochastically maximize the depth of a random node, but not (yet) for the deepest node.

IMPACASTORINA - PRINCEP D

IMPA-NWL - NET

Thank you!



- Claimed theorem with dependence only on p_i proved it with dependence on $p = \max p_i$.

Fix: requires more careful "dispersion" bound for our setting.

Idea: If subcritical then $p_{\max} = p_0$ or p_i ; if p_0 close to 1 then either very subcritical or make large jumps.)