

# The height of Mallows trees



Based on joint work with Benoit Corsini



**McGill**

JOS-G-ADAR-O-BERRY

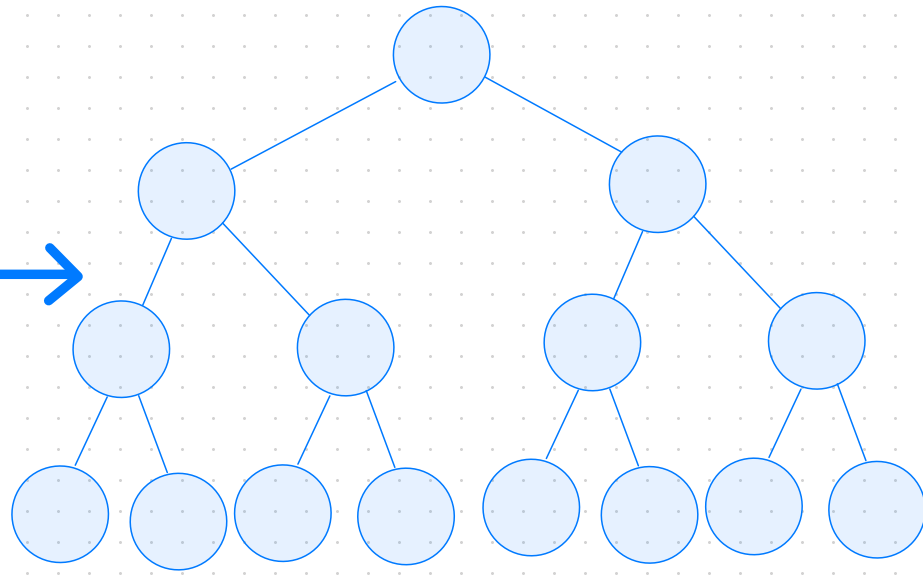


Giuseppe Penone



# Binary search trees

1 2 3 4 5 6 7 8  
(3, 5, 1, 7, 8, 4, 6, 2)



Permutation

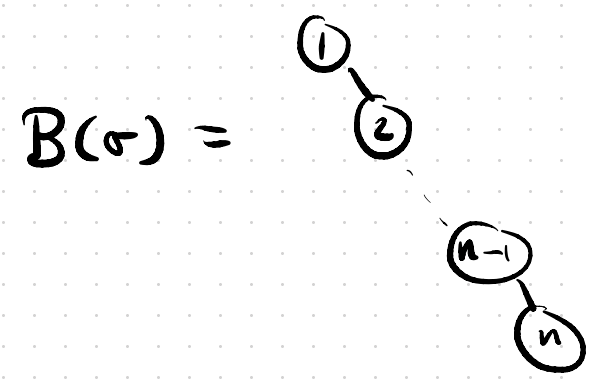


Rooted tree  
(root at top)

Insert values one-at-a-time at top, routing down through tree until empty space is reached.

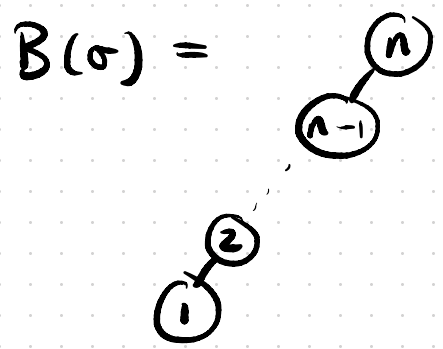
# Binary search trees

$$\sigma = \text{id} \in S_n$$



$\sigma \circ r \quad r_i = n+1-i$

$$\sigma = (n, n-1, \dots, 1) \in S_n$$



$$\sigma \in_u S_n$$

$B(\sigma) =$  Random Binary Search Tree (RBST)

## Devroye's Theorem (1986)

Let  $B_n$  be an RBST with  $n$  nodes.

Then  $\frac{ht(B_n)}{\log n} \xrightarrow{P} c^*$

and  $E ht(B_n) = (1+o(1))c^* \log n$

where  $c^*$  is the unique solution in  $[2, \infty)$  of  $c \log\left(\frac{2}{c}\right) + c = 1$ .

## Reed's Theorem (2003)

Let  $B_n$  be an RBST with  $n$  nodes.

Then  $E ht(B_n) = c^* \log n - \beta^* \log \log n + O(1)$

$\text{Var } ht(B_n) = O(1)$

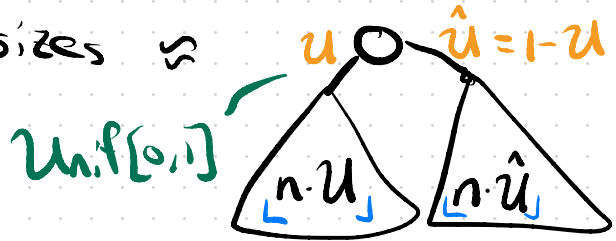
where  $\beta^* = \frac{3}{2} \frac{1}{\log(c^*/2)}$

# Proof idea (Devroye's Theorem)

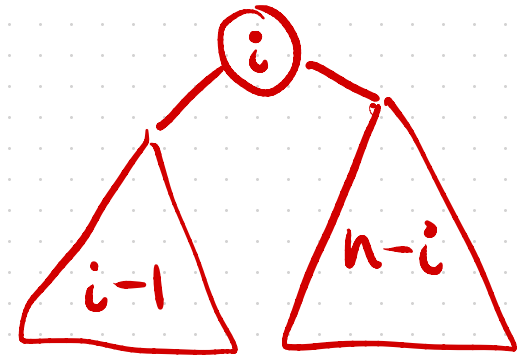
$$\mathcal{B}_n \leftarrow \sigma \in_u S_n$$

- Root value  $\sigma(i)$  is Unif. on  $\{1, 2, \dots, n\}$

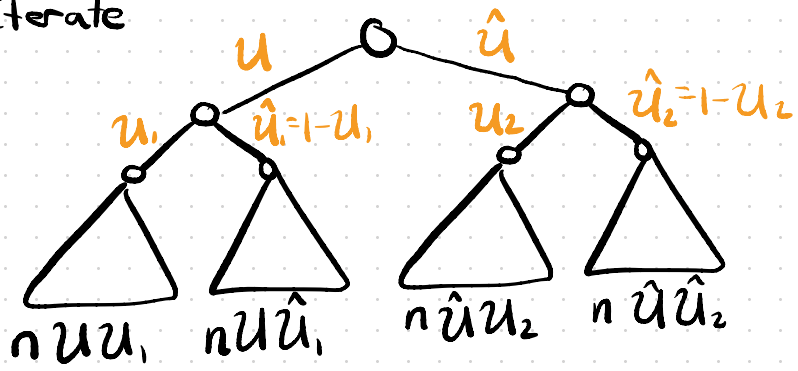
- So subtree sizes  $\approx$



$$\sigma(i) = i$$



- Iterate



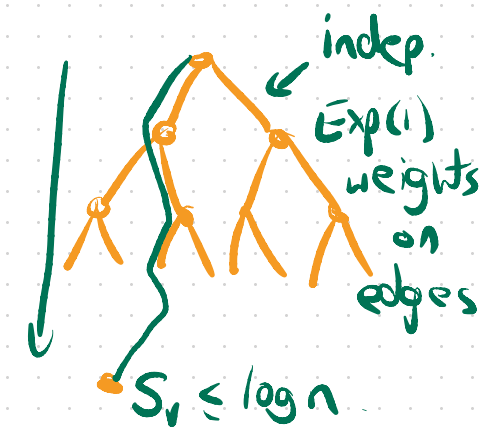
$$\text{size} \stackrel{d}{=} \lfloor n \cdot U_1 \cdot U_2 \rfloor$$

- size of subtree rooted at  $v \approx n \cdot \prod_{e \leq v} U_e$

where terms in product are indep Unif[0,1]

- Height  $\approx \max \{ht(v) : \prod_{e \leq v} U_e \geq \frac{1}{n}\}$

- Height  $\approx \max \{ ht(v) : \sum_{e \in \gamma} (-\log U_e) \leq \log n \}$   
 $=: S_v \leftarrow$  Random walk with  $\text{Exp}(1)$  steps.



Write  $M_k = \min \{ S_v : ht(v) = k \}$

Then  
 Height  $\approx \min \{ k : M_k \geq \log n \}$

Law of large #s for minima in branching random walks  
 (Hammersley-Kingman-Biggins 1970's)

$$\Rightarrow \frac{M_k}{k} \xrightarrow{\text{a.s.}} \frac{1}{c^*} \Rightarrow \text{Height} \approx c^* \log n.$$

□

## Proof idea (Reed's Theorem)

Write  $M_k = \min \{ S_v : \text{ht}(v) = k \}$

Then  
Height  $\approx \min \{ k : M_k \geq \log n \}$

Reed used "truncated second moment method" to prove that

$$M_k = \frac{1}{c^*} k + \frac{1}{\beta^*} \log k + O(1) \quad \text{and} \quad \text{Var}(M_k) = O(1)$$

then used this to control the height.

Similar results independently proved by Bramson, much earlier (late 1970's) for branching Brownian motion. □

## Mallows trees

Any metric  $d$  on  $S_n$  defines a one-parameter family of distributions  $(\Psi_q, q \geq 0)$  on the permutations of  $S_n$  by

$$\Psi_q(\sigma) \propto q^{d(\sigma, id)}$$

$$\Psi_q(\sigma) = \frac{q^{d(\sigma, id)}}{\sum_{\pi \in S_n} q^{d(\pi, id)}}$$

Well-used in statistics with various distances

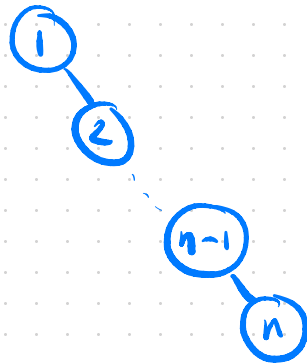
“Suppose that a set of  $k$  objects is to ranked independently by each of a panel of  $n$  judges, and that each ranking proceeds according to some criterion that presumes a true underlying ordering of the objects. Mallows (1957) introduces two one-parameter models to describe such a ranking process, the parameter reflecting the variability of the rankings about the true or modal ordering. The models are referred to as Mallows'  $\Theta$  and  $\Phi$  models, which are based on the correlation coefficients of Spearman (1904) and Kendall (1938), respectively. The  $\Phi$  model is further investigated in Feigin and Cohen (1978), and Diaconis (1982) suggests use of the Mallows' models with other distances.”

Fligner and Verducci, JRSS B (1986)

# Mallows trees

$$\Psi_q(\sigma) \propto q^{d(\sigma, \text{id})} \quad \text{dist. on } S_n$$

$$q \approx 0 \rightarrow \sigma \approx \text{id} \rightarrow B(\sigma) \approx$$



$$q \approx 1 \rightarrow \sigma \approx \text{unif. dist. on } S_n \rightarrow B(\sigma) \approx \text{RBST}$$

$q \approx \infty \rightarrow \sigma \approx \text{unif. dist. on perms. } \pi \text{ with max dist. to id.}$   
Law of  $B(\sigma)$  depends on  $d$ .

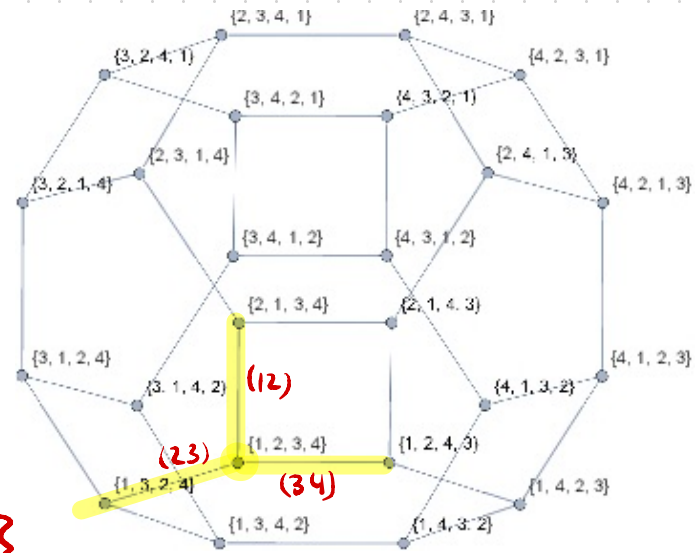
# Mallows' $\Phi$ model:

$d(\pi, \sigma) =$  distance in adjacent transposition graph

For this model

$$d(\sigma, id) = \text{inv}(\sigma)$$

$$:= \#\{i < j : \sigma(i) > \sigma(j)\}$$



The corresponding Mallows distribution on  $S_n$  is

$$\rightarrow \varphi_q(\sigma) = \frac{q^{\text{inv}(\sigma)}}{Z_{n,q}}$$

Mallows  $(n, q)$  dist.

$$Z_{n,q} = \prod_{i=1}^n \frac{1-q^i}{1-q} = i \text{ if } q=1$$

# Mallows $\Phi$ -model: previous research

Let  $\pi \sim \text{Mallows}(n, q_n)$ ,  $0 \leq q_n \leq 1$ .

Bhatnagar-Peled (2015): If  $q_n \rightarrow 1$  and  $n(1-q_n) \rightarrow \infty$  then

Length of longest increasing subseq,  $\text{LIS}_n$ , satisfies  $\frac{\text{LIS}_n}{n\sqrt{1-q_n}} \xrightarrow[\text{prob}]{\text{h.p.}} 1$

Gladkikh-Peled (2018)  $0 \leq q_n \leq 1$ .

List cycles in decreasing order of length as  $e_1^\downarrow, e_2^\downarrow, \dots$

• If  $n(1-q_n)^2 \rightarrow 0$  then  $\frac{1}{n}(|e_i^\downarrow|, i \geq 1) \xrightarrow{\text{dist}} \text{Poisson-Dirichlet}(1)$

•  $\mathbb{E}[\# \text{ cycles in } \pi] \approx (1-q_n) \cdot n + O(1)$ .

Other results by Diaconis+Ram; Mueller & Starr; Basu & Bhatnagar; Mukherjee; Gnedin & Olshanski; Evans, Grübel & Wakolbinger; Angel, Holroyd, Hutchcroft & Levy; ...

Theorem (A-B, Corsini). Fix non-negative reals  $(q_n, n \geq 1)$ .

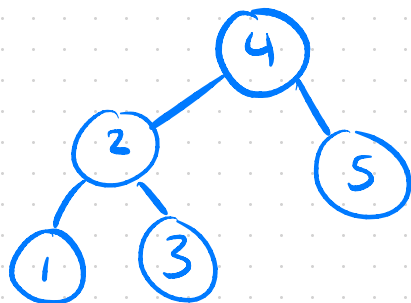
let  $\sigma_n$  be a Mallows  $(n, q_n)$  perm., let  $T_n = B(\sigma_n)$ .  
 ↪ Mallows  $(n, q_n)$ -tree.

Then height  $(T_n)$   $\rightarrow$  1 in prob. and in  $L_p$  for all  $p$ .

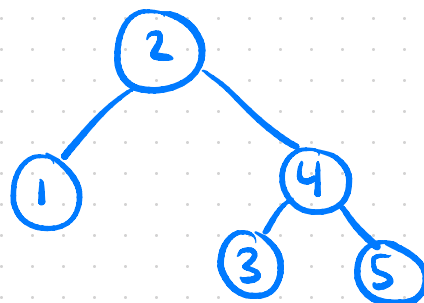
$$n \cdot \max(1 - q_n, 1 - \frac{1}{q_n}) + c \cdot \log n$$

$$n(1 - q_n) + c \cdot \log n.$$

Remark: If  $\sigma$  is Mallows  $(n, q)$  and  $r$  has  $r_i = n + 1 - i$  then  $\sigma \circ r$  is Mallows  $(n, \frac{1}{q})$ . So can assume  $q_n \leq 1 \forall n$ .



$(4, 2, 1, 5, 3)$



$(2, 4, 5, 1, 3)$

**Theorem (A-B, Corsini).** Fix reals  $(q_n, n \geq 1)$  in  $[0, 1]$ , let  $\sigma_n$  be a Mallows  $(n, q_n)$  perm., let  $T_n = B(\sigma_n)$ .

If  $nq_n \rightarrow \lambda \in (0, \infty)$  then  $n-1$ -height  $(T_n) \xrightarrow{d} \text{Poisson}(\lambda)$

If  $nq_n \rightarrow \infty$  and  $\frac{n(1-q_n)}{\log n} \rightarrow \infty$  then

$$\frac{\text{height}(T_n) - n(1-q_n) - c^* \log\left(\frac{1}{1-q_n}\right)}{\sqrt{nq_n(1-q_n)}} \xrightarrow{d} \mathcal{N}(0, 1)$$

Behaviour when  $\frac{n(1-q_n)}{\log n} = O(1)$  an open question.

When  $q_n \equiv 1$ , the RBST case, neither a Poisson nor a Normal limit obtains.

# Key proof input: insertion algorithm for Mallows trees

Fix  $q \in (0, 1)$ .

Let  $(G_i, i \geq 1)$  be i.i.d. Geometric( $1-q$ )  $P(G_i = k) = q^{k-1}(1-q)$

Let  $s_1 = G_1$ ,  $IN_1 = IN \setminus \{G_1\}$

value  $\approx \frac{1}{1-q}$

Inductively let

$s_k = k$ 'th smallest element of  $IN_{k-1}$ ,

$$\begin{aligned} IN_k &= IN_{k-1} \setminus \{s_k\} \\ &= IN \setminus \{s_1, \dots, s_k\} \end{aligned}$$

Then the rank permutation of  $(s_1, \dots, s_n)$  is Mallows( $n, q$ )-dist

$$(2, 1, 3, 1, 5, 2) \rightarrow (2, 1, 5, 3, 9, 6) \rightarrow (2, 1, 4, 3, 6, 5)$$

$$(G_1, G_2, \dots, G_6) \quad (s_1, s_2, \dots, s_6) \quad (\sigma_1, \sigma_2, \dots, \sigma_6)$$

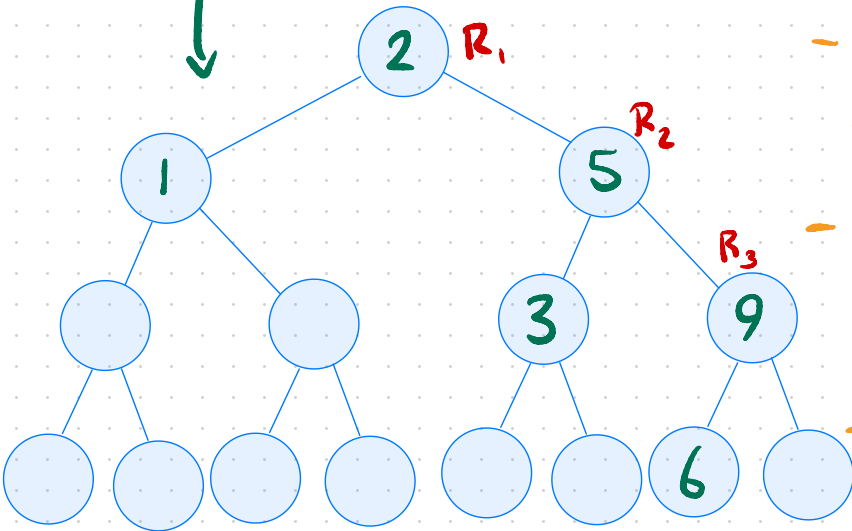


# The structure of Mallows trees

$$(2, 1, 3, 1, 5, 2) \rightarrow (2, 1, 5, 3, 9, 6) \rightarrow (2, 1, 4, 3, 6, 5)$$

$$(G_1, G_2, \dots, G_6) \quad (S_1, S_2, \dots, S_6) \quad (\sigma_1, \sigma_2, \dots, \sigma_6)$$

Insert the values  $(s_1, \dots, s_n)$  into a BST



- List the record values in  $(S_n, n \geq 1)$  as  $(R_i, i \geq 1)$ , set  $R_0 = 0$

- The record increments

$(R_i - R_{i-1}, i \geq 0)$  are IID  $\text{Geom}(1-q)$

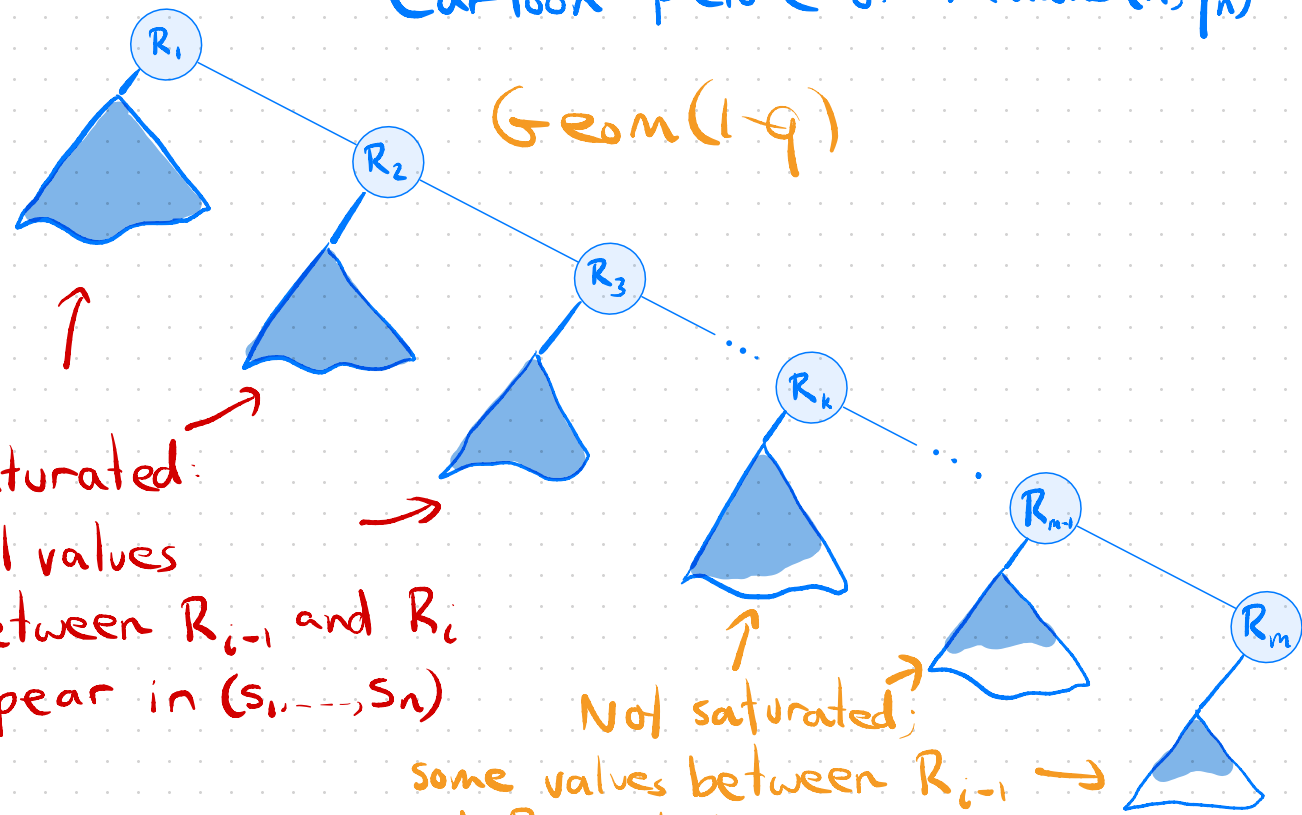
- If  $R_i - R_{i-1} = s$  then for  $n$  large, the left subtree hanging

from the  $i$ th node along the right path will have size  $s-1$ .

# The structure of Mallows trees

## Cartoon picture of Mallows $(n, q_n)$ -tree

Geom  $(1-q)$



Saturated:  
All values  
between  $R_{i-1}$  and  $R_i$   
appear in  $(s_1, \dots, s_n)$

Not saturated;  
some values between  $R_{i-1}$   
and  $R_i$  yet to arrive

$m = \#$  records  
by time  $n$ .

# Height of the Mallows tree

Two contributions

- Rightmost path
- Left subtree

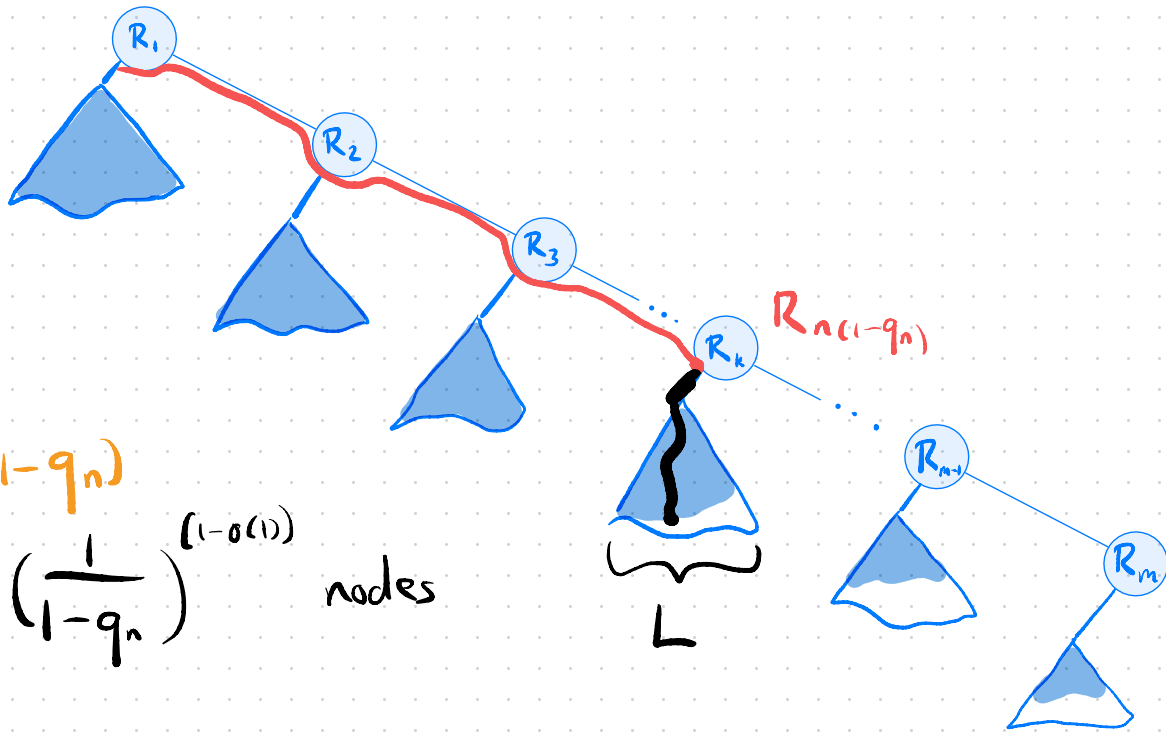
Rightmost path:

Stop at dist.  $n(1-q_n)$

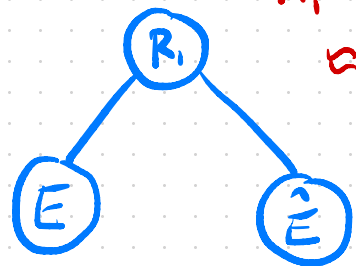
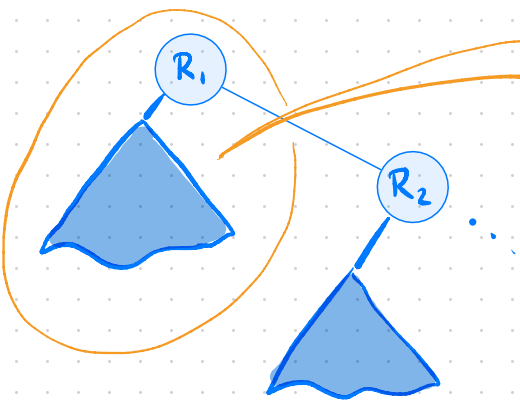
- Left subtree has  $\left(\frac{1}{1-q_n}\right)^{(1-o(1))}$  nodes

with high prob.

Prove left subtree has height  $(1+o_p(1))c^* \log |L|$  like in Devroye's thm  
(For CLT need finer control over fluctuations of left subtree heights)



The left subtrees (when  $q = 1 - O(1)$ )



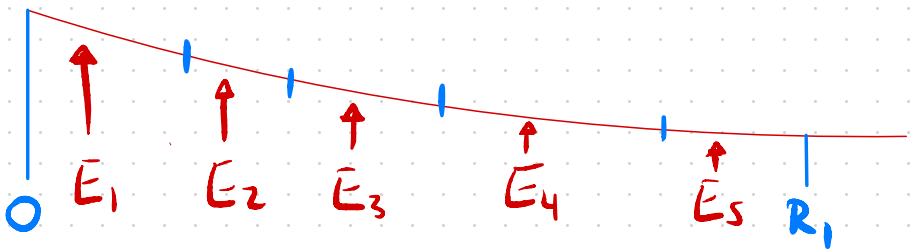
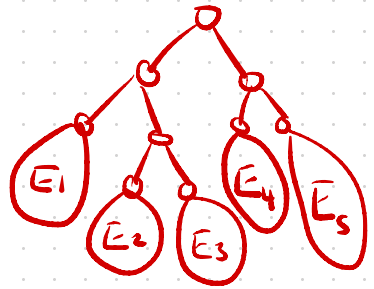
$$R_1 = \text{Geom}(1-q) \approx \frac{1}{1-q} \cdot \text{Exp}(1)$$

$$(1-q)E \approx \text{Exp}(1)$$

$$\hat{E} = R_1 - 1 - E$$

conditioned to be  $\leq R_1$

After  $k-1$  insertions:



## Open questions & research directions

- Let  $L(q)$  be a left subtree in a large Mallows( $n, q$ ) tree.

Prove Reed's theorem for  $\text{height}(L(q))$

$$\mathbb{E} \text{height}(L(q)) = c^* \log \frac{1}{1-q} + \beta^* \log \log \frac{1}{1-q} + O(1)$$

$$\text{Var}(\text{height}(L(q))) = O(1) \text{ uniformly in } 0 \leq q < 1.$$

- Extend the range of the CLT, figure out when the CLT fails
- Carry out this research for other Mallows/Diaconis models.

# Reference :

LAB+Benoît Corsini, The height of Mallows trees. Annals of Probability, to appear.  
<https://arxiv.org/pdf/2007.13728.pdf>

