# CONCENTRATION PROPERTIES OF EXTREMAL PARAMETERS IN RANDOM DISCRETE STRUCTURES\*

#### **Michael Drmota**

Inst. of Discrete Mathematics and Geometry

Vienna University of Technology, A 1040 Wien, Austria

michael.drmota@tuwien.ac.at

www.dmg.tuwien.ac.at/drmota/

\* supported by the Austrian Science Foundation FWF, grant S9600.

Fourth Colloquium on Mathematics and Computer Science, Nancy, 18-22. 9. 2006

# Outline of the Talk

- Types of Concentration
- Chromatic Number of Random Graphs
- Travelling Salesman Problem
- Longest Increasing Subsequence in Random Permutations
- Diameter and Maximum Degree in Random Graphs
- Height and Maximum Degree in Random Trees
- Height of Scale-Free Trees

# Outline of the Talk

- Martingales  $\rightarrow$  concentration inequality
- Talagrand's *convex distance*  $\rightarrow$  concentration inequality
- Poisson transform, analytic methods
- Martingales, moment methods, ...
- Generating functions, analytic methods
- Generating functions, analytic methods

## Conclusions

- Concentration (almost) always appears.
- Smaller extremal parameters are more concentrated than larger ones.
- Concentration is *easy* to prove (compared to the precise position of the mean).

### **Recent Books on Concentration**

S. Boucheron, G. Lugosi, and O. Bousquet, *Concentration inequalities*, Lecture Notes in Computer Science 3176, Springer, Berlin, 2004.

D. P. Dubhashi and A. Panconesi, *Concentration of measure for the analysis of randomized algorithms*, draft available at http://www.dsi.uniroma1.it/~ale.

G. Lugosi, *Concentration-of-measure inequalities*, draft available at http://www.econ.upf.es/~lugosi.

C. McDiarmid, *Concentration*, Probabilistic methods for algorithmic discrete mathematics, 195–248, Algorithms Combin., 16, Springer, Berlin, 1998.

 $X_n \ldots$  non-negative random variable with  $\mathbf{E} X_n \to \infty$ 

**Concentration:** 

$$\lim_{n \to \infty} \mathbb{P}\left\{ \left| \frac{X_n}{a(n)} - 1 \right| \ge \epsilon \right\} = 0$$

for all  $\epsilon > 0$  and some sequence a(n) with  $a(n) \rightarrow \infty$ 

Equivalently 
$$X_n/a(n) \xrightarrow{d} \delta_1$$
, usually  $a(n) = \mathbb{E} X_n$ .

 $X_n \ldots$  non-negative random variable with  $| \mathbf{E} X_n \to \infty |$ 

#### **Type 1:** No Concentration:

$$\frac{X_n}{\mathbf{E}\,X_n} \not \longrightarrow \delta_1$$

Typically:

$$\frac{X_n}{\mathbf{E} X_n} \xrightarrow{\mathsf{d}} Y \dots \text{ not concentrated at 1}$$

and  $\mathbf{E} X_n^2 \sim c \cdot (\mathbf{E} X_n)^2$  for some c > 1.



**Type 2:** Weak Concentration:

For all  $\epsilon > 0$  there exists K > 0 such that

$$\lim_{n \to \infty} \mathbb{P}\left\{ \left| \frac{X_n - a(n)}{b(n)} \right| \ge K \right\} \le \epsilon.$$

with 
$$a(n) \to \infty$$
,  $b(n) \to \infty$ , and  $b(n) = o(a(n))$ 

Usually one takes  $a(n) = \mathbb{E} X_n$  and  $b(n) = (\mathbb{V}X_n)^{1/2}$ 

If  $\mathbb{E} X_n^2 \sim (\mathbb{E} X_n)^2$  the Chebyshev's inequality implies weak concentration.

Typically

$$\frac{X_n - \mathbf{E} X_n}{\sqrt{\operatorname{Var} X_n}} \xrightarrow{\mathsf{d}} Y.$$

**Type 2:** Weak Concentration:

E.g. Central Limit Theorem

$$\frac{X_n - \mathbf{E} X_n}{\sqrt{\operatorname{Var} X_n}} \to N(0, 1).$$



**Type 3: Strong Concentration:** 

For all  $\epsilon > 0$  there exists K > 0 with

 $\limsup_{n \to \infty} \mathbb{P}\{|X_n - a(n)| \ge K\} \le \epsilon$ 

for some sequence a(n) with  $a(n) \to \infty$ .

Usually  $a(n) = \mathbb{E} X_n$  or a(n) = median of  $X_n$  and one has bounded centralized moments:

$$\mathbb{E} |X_n - \mathbb{E} X_n|^d = O(1) \qquad (d \ge 1).$$

**Type 3:** Strong Concentration:

Typically: travelling wave F(x)

$$\mathbb{P}\{X_n \le k\} = F(k - m(n)) + o(1)$$

(m(n) is close to the median of  $X_n$ )



**Type 4:** Very Strong Concentration:

Concentration on two (or finitely many values):

$$\mathbb{P}\{m(n) \le X_n \le m(n) + L\} = 1 + o(1)$$

with  $m(n) \rightarrow \infty$  and some fixed L

#### Definition

Let n be a positive integer and p a real number with  $0 \le p \le 1$ .

The random graph G(n, p) is a probability space over the set of graphs on the vertex set  $\{1, 2, ..., n\}$  determined by

 $\mathbb{P}\{\{i,j\}\in G\}=p$ 

for all possible (undirected) edges  $\{i, j\}$  with  $1 \le i < j \le n$  with these events mutually independent.

٠

٠

**Random Graph** 

**Random Graph** 

٠



**Random Graph** 



**Random Graph** 



**Random Graph** 



**Random Graph** 



**Random Graph** 



**Random Graph** 



**Random Graph** 



٠

#### Definition

A regular k-coloring of the vertices of a graph G is a coloring of the vertices with k colors such that adjacent vertices have different colors.

The chromatic number  $\chi(G)$  of a graph G is the smallest number k such that there exists a regular k-coloring of the vertices of G

**Notation:** We use the notion *almost always* as an abbreviation for the property that the probability that a certain condition holds converges to 1 as the *size* of the problem goes to the infinity.

Chromatic Number = 3



Theorem (Bollobás, Frieze, Grimmet, McDiarmid)

(i) If  $C_0/n \le p = p(n) \le (\log n)^{-7}$  (for a proper constant  $C_0 > 0$ ) then almost always

 $\frac{np}{2\log(np)-2\log\log(np)+1} \le \chi(G(n,p)) \le \frac{np}{2\log(np)-40\log\log(np)}.$ 

(ii) If  $(\log n)^{-2} \le p = p(n) \le c$  (for some arbitrary c < 1) then almost always

 $\frac{n}{2\log_b n - \log_b\log_b n} \le \chi(G(n,p)) \le \frac{n}{\log_b n - 6\log_b\log_b n},$ where b = 1/(1-p).

(iii) If  $p = p(n) > n^{-\delta}$  for every  $\delta > 0$  (and sufficiently large n) but  $p = p(n) \le c$  (for some arbitrary c < 1) then almost always  $\chi(G(n,p)) = \frac{n}{2\log_b n - 2\log_b \log_b n + O(1/p)}$ .

Theorem (Łuczak, Alon and Krivelevich 1997)

Fix some  $\varepsilon > 0$ . For every sequence p = p(n) there exists a function h(n) such that almost always

(i) if  $p \ge n^{-\frac{1}{2}-\varepsilon}$  then  $\chi(G(n,p)) \sim h(n)$ , and

(ii) if  $p \le n^{-\frac{1}{2}-\varepsilon}$  then  $h(n) \le \chi(G(n,p)) \le h(n)+1$ .

**Theorem** (Shamir and Spencer 1987)

$$\mathbb{P}\{|\chi(G(n,p)) - \mathbb{E}\left(\chi(G(n,p))\right)| > \lambda\sqrt{n-1}\} < 2e^{-\lambda^2/2}.$$

**Remark.** This theorem is weaker than the previous one (for  $p \le n^{-\frac{1}{2}-\epsilon}$ ) but the basis for further considerations.

**Definition.** A martingale is a sequence of random variables  $Y_0, Y_1 \dots, Y_n$ on a probability space  $(\Omega, \mathcal{F}, \mathbb{P})$  with

$$\mathbb{E}\left(Y_{k+1}|\mathcal{F}_k\right) = Y_k,$$

where  $\mathcal{F}_0 = \{\emptyset, \Omega\} \subseteq \mathcal{F}_1 \subseteq \cdots \subseteq \mathcal{F}_n = \mathcal{F}$  is an increasing sequence of  $\sigma$ -fields.

**Theorem** (Azuma's Inequality) Suppose that  $Y_0, Y_1 \dots, Y_n$  is a martingale with constant  $Y_0$  and that

$$|Y_{k+1} - Y_k| \le c_k$$

for some some constants  $c_k$  ( $0 \le k < n$ ). Then, for every t > 0,

$$\mathbb{P}\{|Y_n - \mathbb{E} Y_n| \ge t\} \le 2 \exp\left(-\frac{t^2}{2\sum_{k=1}^n c_k^2}\right).$$

**Theorem** (McDiarmid's Independent Bounded Difference inequality)

Let  $X_1, \ldots, X_n$  be independent random variables, with  $X_k$  taking values in a set  $\Omega_k$ . Suppose that a function  $f : \Omega_1 \times \cdots \times \Omega_n \to \mathbb{R}$  satisfies the property that

$$|f(x_1,\ldots,x_n)-f(y_1,\ldots,y_n)|\leq c_k$$

if  $(x_1, \ldots, x_n)$  and  $(y_1, \ldots, y_n)$  differ only at the k-th coordinate, that is  $x_j = y_j$  for  $j \neq k$ .

Then, the random variable  $Y = f(X_1, \ldots, X_n)$  satisfies, for any  $t \ge 0$ ,

$$\mathbb{P}\{|Y - \mathbb{E} Y| \ge t\} \le 2 \exp\left(-\frac{t^2}{2\sum_{k=1}^n c_k^2}\right).$$

#### Proof

 $\mathcal{F}_k \ldots \sigma$ -field generated by  $X_1, \ldots, X_k$ 

 $Y_k = \mathbb{E}(f(X_1, \dots, X_n) | \mathcal{F}_k), k = 0, 1, \dots, n \text{ is a martingale}$ 

 $|f(x_1,...,x_n) - f(y_1,...,y_n)| \le c_k$  implies that  $|Y_{k+1} - Y_k| \le c_k$ .

Hence, Azuma's inequality applies.

#### Vertex Exposure Martingale

 $A_k = \{\{j,k\} : 1 \le j < k\} \dots$  edges that connect k with j < k.

 $X_k = (\mathbb{I}_{[e \in G(n,p)]} : e \in A_k) \dots$  rand. vector of indicators of edges in  $A_k$ .

 $f \dots$  graph theoretical function (for example, the chromatic number).

 $\mathcal{F}_k \ldots \sigma$ -field generated by  $X_1, \ldots, X_k$ 

 $Y_k = \mathbb{E}(f(G(n,p))|\mathcal{F}_k)$  ... vertex exposure martingale

# (It can be interpreted as the conditional expectation of f with partial information on the first k vertices and their internal edges.)

**Remark.**  $|f(x_1, \ldots, x_n) - f(y_1, \ldots, y_n)| \le c_k$  with  $x_k, y_k \in A_k$  says that  $|f(H_1) - f(H_2)| \le c_k$  if  $H_1, H_2$  are subgraphs of the complete graph on the vertices  $\{1, 2, \ldots, n\}$  such that the symmetric difference of the edge sets of  $H_1$  and  $H_2$  is contained in  $A_k$ .

If one adds a vertex to a graph then the chromatic number changes at most by 1. (Here we use the vertex k.)

 $\implies$  This condition is satisfied for the chromatic number with  $c_k = 1$ .

 $\mathbf{X} = (X_1, X_2, \dots, X_n) \dots$  *n*-tuple of random point selected uniformly and independently in the unit square  $[0, 1]^2$ 

Length of the minimum (travelling salesman) tour:

$$TSP(\mathbf{X}) = \min_{\pi \in S_n} \sum_{j=1}^n |X_{\pi(j)} - X_{\pi(j+1)}|$$

Theorem (Beardwood, Halton and Hammersley 1959)

$$\frac{\mathsf{TSP}(\mathbf{X})}{\sqrt{n}} \to \beta_2 \quad \text{in prob.}$$

for some  $\beta_2 > 0$ .

**Remark:** Up to now there is no known analytic expression for  $\beta_2$ .





**Notation:** M(Y) ... median of r.v. Y

**Theorem** (Rhee and Talagrand)

$$\mathbb{P}\left\{|\mathsf{TSP}(\mathbf{X}) - \mathsf{M}(\mathsf{TSP}(\mathbf{X}))| \ge t\right\} < 4e^{-t^2/c}.$$

for some constant c > 0.

**Corollary.** All central moments of TSP(X) are bounded.

(However, the exact location of the mean is unknown.)
#### **Talagrand's Inequality**

 $\Omega_1, \Omega_2, \ldots, \Omega_n \ldots$  probability spaces,  $\Omega = \Omega_1 \times \cdots \times \Omega_n$ 

 $\mathbf{X} = (X_1, X_2, \dots, X_n)$  ... independent random variables,  $X_k$  taking values in  $\Omega_k$ .

Weighted Hamming distance related to  $\alpha = (\alpha_1, \ldots, \alpha_n)$  with  $\alpha_k \ge 0$ :

$$d_{\alpha}(\mathbf{x}, \mathbf{y}) = \sum_{x_i \neq y_i} \alpha_i$$

**Talagrand's convex distance** 

$$d_T(\mathbf{x}, A) = \sup_{\alpha \ge 0, \|\alpha\| = 1} \inf_{\mathbf{y} \in A} d_\alpha(\mathbf{x}, \mathbf{y})$$

between  $\mathbf{x} \in \Omega$  and  $A \subset \Omega$ .

**Talagrand's Inequality** 

$$\left| \mathbb{P}\{\mathbf{X} \in A\} \cdot \mathbb{P}\{d_T(\mathbf{X}, A) \ge t\} < e^{-t^2/4}. \right|$$

#### Theorem

 $f \dots$  real valued function on  $\Omega = \Omega_1 \times \cdots \times \Omega_n$ 

For every  $\mathbf{x} \in \Omega$  there exists a non-negative unit *n*-vector  $\alpha$  and a constant c > 0 such that for all  $\mathbf{y} \in \Omega$ 

$$f(\mathbf{x}) \leq f(\mathbf{y}) + c d_{\alpha}(\mathbf{x}, \mathbf{y}).$$

Then, for every random *n*-tuple  $\mathbf{X} = (X_1, \dots, X_n)$  of independent random variables  $X_k$  taking values in  $\Omega_k$  we have

$$\left|\mathbb{P}\left\{|f(\mathbf{X}) - \mathsf{M}(f(\mathbf{X}))| \ge t\right\} \le 4 e^{-t^2/(4c^2)}.$$

#### Proof

 $A_a := \{ \mathbf{y} \in \Omega : f(\mathbf{y}) \le a \}$ 

By assumption for every  $\mathbf{x} \in \Omega$  there exists a non-negative unit *n*-vector  $\alpha$  such that for all  $\mathbf{y} \in A_a$ :

$$f(\mathbf{x}) \leq f(\mathbf{y}) + c d_{\alpha}(\mathbf{x}, \mathbf{y}) \leq a + c d_{\alpha}(\mathbf{x}, \mathbf{y}).$$

By taking the miminum over all  $\mathbf{y} \in A_a$  we, thus, get

$$f(\mathbf{x}) \leq a + c d_{\alpha}(\mathbf{x}, A_a) \leq a + c d_T(\mathbf{x}, A_a).$$

Hence

$$f(\mathbf{x}) \ge a + t \implies d_T(\mathbf{x}, A_a) \ge t/c.$$

 $\implies \mathbb{P}\{f(\mathbf{X}) \le a\} \cdot \mathbb{P}\{f(\mathbf{X}) \ge a+t\} \le \mathbb{P}\{\mathbf{X} \in A_a\} \cdot \mathbb{P}\{d_T(\mathbf{x}, A_a) \ge t/c\} \\ \le e^{-t^2/(4c^2)}.$ 

$$a = \mathsf{M}(f(\mathbf{X})), \ \mathbb{P}\{f(\mathbf{X}) \le a\} = \frac{1}{2}:$$
$$\mathbb{P}\{f(\mathbf{X}) \ge \mathsf{M}(f(\mathbf{X})) + t\} \le 2e^{-t^2/(4c^2)}.$$

 $a = \mathsf{M}(f(\mathbf{X})) - t:$ 

$$\mathbb{P}\{f(\mathbf{X}) \le \mathsf{M}(f(\mathbf{X})) - t\} \le 2e^{-t^2/(4c^2)}.$$

#### Lemma

For every  $\mathbf{x} \in ([0,1]^2)^n$  there exists non-negative unit vector  $\alpha$  and a constant c > 0 such that for all  $\mathbf{y} \in ([0,1]^2)^n$ 

### $\mathsf{TSP}(\mathbf{x}) \leq \mathsf{TSP}(\mathbf{y}) + c d_{\alpha}(\mathbf{x}, \mathbf{y}).$

(Elementary proof that uses an approximate minumum tour to construct  $\alpha$ .)

#### Remark

This method can be applied to several other problem, for example to the minimal Steiner tree problem etc.

 $S_n$  ... the set of permutations of the numbers  $\{1, 2, ..., n\}$ (We assume that every permutation in  $S_n$  is equally likely.)

For  $\pi \in S_n$  we say that  $\pi(i_1), \pi(i_2), \ldots, \pi(i_k)$  is an increasing subsequence in  $\pi$  if  $i_1 < i_2 < \cdots < i_k$  and  $\pi(i_1) < \pi(i_2) < \cdots < \pi(i_k)$ .

 $L_n = L_n(\pi)$  ... length of the longest increasing subsequence.

Ulam's Problem:  $\mathbb{E}L_n \sim ?$ 

Ulam's conjecture:  $\mathbb{E} L_n \sim c\sqrt{N}$  for some constant c > 0.

$$\begin{pmatrix} 1 & 2 & 3 & 4 & 5 & 6 \\ 3 & 1 & 5 & 6 & 2 & 4 \end{pmatrix}$$

$$\begin{pmatrix} 1 & 2 & 3 & 4 & 5 & 6 \\ 3 & 1 & 5 & 6 & 2 & 4 \end{pmatrix}$$

Erdős Szekeres 1935:  $c \geq \frac{1}{2}$ 

Logan and Shepp 1977:  $c \ge 2$ 

Vershik and Kerov 1977: c = 2.

(Alternate proofs are due to Aldous and Diaconis, Seppäläinen, and Johansson).

Frieze 1991:  $L_n$  is concentrated

Bollobás and Brightwell 1992, Talagrand 1995:  $\mathbb{V}L_n = O(\sqrt{N})$ .

Odlyzko and Rains 2000: order of  $\mathbb{V}L_n$  should be  $N^{1/3}$ .

Baik, Deift, and Johansson 1999: complete solution.

Theorem (Baik, Deift, and Johansson 1999)

Let  $S_n$  be the group of permutations of n numbers with uniform distribution and  $L_n$  the longest increasing subsequence. Then there exists a random variable Y such that

$$\frac{L_n - 2\sqrt{n}}{n^{1/6}} \stackrel{\mathsf{d}}{\longrightarrow} Y.$$

Furthermore, we have convergence of all moments.

#### Remark.

The limiting distribution Y is exactly the same at the limiting distribution of the largest eigenvalue in random Hermitian matrices. However, it seems that there is no direct connection between these two problems.

**Tracy-Widom distribution:**  $F(t) = \mathbb{P}\{Y \le t\}$ 

u(x) ... solution of the Painlevé II equation

$$u'' = 2u^3 + xu, \quad u(x) \sim Ai(x) \quad (as \ x \to \infty);$$

Ai(x) denotes the Airy function.

$$F(t) = \exp\left(\int_t^\infty (x-t)^2 u(x)^2 \, dx\right),\,$$

#### **Proof Method**

Basically one determines the asymptotic behaviour of the Poisson transform

$$\phi_k(\lambda) = \sum_{n=0}^{\infty} \frac{e^{-\lambda} \lambda^n}{n!} \mathbb{P}\{L_n \le k\}$$

that can be represented as

$$\phi_k(\lambda) = \frac{e^{-\lambda}}{(2\pi)^k k!} \int_{[-\pi,\pi]^k} \exp\left(2\sqrt{\lambda} \sum_{j=1}^k \theta_j\right) \prod_{1 \le j < \ell \le k} \left|e^{i\theta_j} - e^{i\theta_\ell}\right| \, d\theta_1 \cdots d\theta_k.$$

One has to use the theory of orthogonal polynomials on the unit circle, sophisticated Riemann-Hilbert problem techniques and certain properties on eigenvalues of random matrices.

#### Definition

The diameter diam(G) of a graph connected G is the largest distance between two nodes in G. If G is not connected then diam(G) =  $\infty$ .

The maximum degree of an (undirected) graph G will be denoted by  $\Delta(G)$ .

**Diameter** = 3



Maximum Degree = 3



G(n, p)-Random Graphs

Theorem (Burtin, Bollobás)

(i) If  $(pn)/\log n \to \infty$  and  $\log n/\log(pn) \to \infty$  then almost always

diam(
$$G(n,p)$$
) ~  $\frac{\log n}{\log(pn)}$ .

(ii) Let c be a positive constant and p = p(n) and d = d(n), an integer  $\geq 2$ , be related by  $p^d n^{d-1} = \log(n^2/c)$ . Further suppose that  $(pn)/(\log n)^3 \to \infty$ . Then

$$\lim_{n \to \infty} \mathbb{P}\{\operatorname{diam}(G(n,p)) = d\} = e^{-c/2}$$

and

$$\lim_{n \to \infty} \mathbb{P}\{ \text{diam}(G(n, p)) = d + 1\} = 1 - e^{-c/2}.$$

Notation

$$\lambda_k = \lambda_k(n) := n \binom{n-1}{k} p(n)^k (1-p(n))^{n-k-1}$$

٠

 $k = k(n) \ge np$  be such that the quantity max $\{\lambda_k, 1/\lambda_k\}$  is minimal.

**Theorem** (Bollobas) Suppose that  $p = p(n) = o(\log n/n)$ .

(i) If 
$$0<\liminf\lambda_k\leq\limsup\lambda_k<\infty$$
 then, as  $n\to\infty,$  
$$\mathbb{P}\{\Delta=k(n)\}=1-e^{-\lambda_k}+o(1)$$
 and

$$\mathbb{P}\{\Delta = k(n)\} = e^{-\lambda_k} + o(1).$$

(ii) If  $\lim \lambda_k = \infty$  then

$$\mathbb{P}\{\Delta = k(n)\} = 1 + o(1).$$

(iii) If  $\lim \lambda_k = 0$  then

$$\mathbb{P}\{\Delta = k(n) - 1\} = 1 + o(1).$$

(iv) If there is a function D(n) with  $\mathbb{P}\{\Delta = D(n)\} = 1 + o(1)$  as  $n \to \infty$  then  $p = p(n) = o(\log n/n)$ .

**Theorem** (Bollobas, Riordan and Selby)

Suppose that 0 is fixed and <math>q = 1 - p.

(i) For every real number y we have

$$\mathbb{P}\left\{\Delta \leq pn + \sqrt{2pqn\log n} \left(1 - \frac{\log\log n}{4\log n} + \frac{y - 2\sqrt{\pi}}{2\log n}\right)\right\}$$
$$= \exp\left(-e^{-y}\right) + o(1).$$

(ii) Almost always we have

$$\left|\Delta - pn - \sqrt{2pqn\log n} + \log\log n\sqrt{\frac{pqn}{8\log n}}\right| \le \log\log n\sqrt{\frac{n}{\log n}}.$$

(iii) For every real number b there exists c(b) such that

$$\mathbb{P}\{\Delta < pn + b\sqrt{npq}\} = (c(b) + o(1))^n$$

Barabási-Albert model (for real-world graphs, internet etc.):

- Randomly growing graph
- A new node is joint to an existing one with probability proportional to the degree.

This definition is not unambigous!!!!

#### Scale-Free Random Graphs

A *power law* is a distribution Z with tail of the form  $\mathbb{P}\{Z = d\} \sim c \cdot d^{-k}$  (for some k > 1).

If a (random graph) that has an power law as (asymptotic) degree distribution is called **scale-free**.

Bollobás and Riordan:  $G_m^n$  multi-graph

 $m = 1 (G_1^n)$ :

- Initial node 1 with a loop.
- $\bullet$  at step k we add one node that is connected to  $j \leq k$  with propability

$$rac{\deg_{G_1^{k-1}}(j)}{2k-1} \qquad ext{if } j < k, \ rac{1}{2k-1} \qquad ext{if } j = k.$$

 $m \geq 1$ 

•  $G_m^n$  is constructed from  $G_1^{mn}$  by identifying the nodes  $\{(\ell - 1)m + 1, (\ell - 1)m + 2, ..., \ell m\}$   $(1 \le \ell \le n)$  of  $G_1^{mn}$  to a new node  $\ell$  (and all edged within the nodes  $\{(\ell - 1)m + 1, (\ell - 1)m + 2, ..., \ell m\}$  are now loops of the new node  $\ell$ )

٠

٠



٠



٠



٠



٠

 $G_2^n$ 



٠

 $G_2^n$ 



Scale-Free Random Graphs:  $G_m^n$  is scale free, the tail of the degree distribution is of the form  $c \cdot d^{-3}$ .

**Theorem** (Bollobás and Riordan 2004)

Suppose that  $m \ge 2$ . Then for every  $\epsilon > 0$ 

$$\lim_{n \to \infty} \mathbb{P}\left\{ (1-\epsilon) \frac{\log n}{\log \log n} \le \operatorname{diam}(G_m^n) \le (1+\epsilon) \frac{\log n}{\log \log n} \right\} = 1.$$

# Height and Maximum Degree of Random Trees

- Galton-Watson Trees
- Pólya Trees
- *m*-Ary Search Trees
- Recursive Trees
- Scale Free Trees
- Tries
- Digital Search Trees

 $\xi$  ... non-negative integer valued random variable,  $\mathbb{E}\ \xi=1,\ 0<\mathbb{V}\ \xi=\sigma^2<\infty.$ 

 $(Z_k)_{k>0}$  Galton-Watson branching process:  $Z_0 = 1$ ,

$$Z_{k} = \sum_{j=1}^{Z_{k-1}} \xi_{j}^{(k)},$$

where the  $(\xi_j^{(k)})_{k,j}$  are iid random variables distributed as  $\xi$ .

A Galton-Watson branching processes can be represented by ordered (finite or infinite) rooted trees T.

$$y_n = \mathbb{P}\{|T| = n\}, \ y(x) = \sum_{n \ge 1} y_n x^n, \ \varphi(t) = \mathbb{E} \ t^{\xi}:$$
  
 $\implies y(x) = x\varphi(y(x)).$ 

 $\mathcal{T}_n$  ... set of rooted trees T of size |T| = n,  $\nu(T)$  ... probability that T occurs in Galton-Watson branching process:

$$\nu_n(T) := \frac{\nu(T)}{y_n}$$

is a probability distribution on  $\mathcal{T}_n$ .

 $\sigma^2 < \infty$ :

$$y_n \sim \frac{d}{\sqrt{2\pi\sigma}} n^{-3/2} \qquad (n \equiv 1 \mod d),$$

where  $d = \gcd\{i > 0 : \mathbb{P}\{\xi = i\} > 0\}.$ 

#### Examples.

 $\varphi(t) = \mathbb{E} t^{\xi} = (1+t)^2/4 = \frac{1}{4} + \frac{t}{2} + \frac{t^2}{4}$ : **binary trees** with *n* (internal) nodes, where each binary tree (of size *n*) has equal probability.

 $\varphi(t) = \mathbb{E} t^{\xi} = 1/(2-t) = \frac{1}{2} + \frac{t}{4} + \frac{t^2}{8} + \cdots$ : planted plane trees, every rooted planar tree is equally likely.

 $\varphi(t) = e^{t-1}$ : Cayley trees, every rooted labeled tree is equally likely.

**Equivalent description:** Simply generated trees (introduced by Meir and Moon)

Cayley Trees: labeled, rooted, non-planar


Cayley Trees: labeled, rooted, non-planar



**Theorem** (de Bruijn, Knuth and Rice, Flajolet, Gao, Odlyzko and Richmond, Aldous)

Suppose that the second moment  $\mathbb{E}\xi^2$  is finite. Then

$$\boxed{\frac{1}{\sqrt{n}}H_n \stackrel{\mathsf{d}}{\longrightarrow} \frac{2}{\sigma} \max_{0 \le t \le 1} e(t),}$$

where  $(e(t), 0 \le t \le 1)$  denotes Brownian excursion of duration 1. Furthermore, if  $\varphi(t) = \mathbb{E} t^{\xi}$  exists for some t > 1 then we also have convergence of all moments. For every  $r \ge 0$  we have, as  $n \to \infty$ ,

$$\mathbb{E}(H_n^r) = 2^{r/2} \sigma^{-r} r(r-1) \Gamma(r/2) \zeta(r) \cdot n^{r/2} \left( 1 + O(n^{-\frac{1}{4} + \eta}) \right)$$

where  $\zeta(s)$  denotes the Riemann Zeta-function and  $(r-1)\zeta(r) = 1$  for r = 1 and  $\eta$  is any positive number.

Depth-first search.



$$\frac{T_n(\lfloor 2nt \rfloor)}{c_2\sqrt{n}} \to e(t) \quad \dots \quad \text{Brownian excursion}$$

 $h_T$  ... height of T

$$y_k(x) = \sum_{n \ge 1} \mathbb{P}\{|T| = n, h_T \le k\} x^n,$$
  
$$\implies y_{k+1}(x) = x\varphi(y_k(x)), \qquad y_0(x) = \varphi_0 x.$$
  
$$H(x) := \sum_{n \ge 1} \mathbb{E} H_n \cdot y_n \cdot x^n = \sum_{k \ge 0} (y(x) - y_k(x)).$$

A subtle analysis of the above recurrence yiels

$$H(x) = \frac{1}{\sigma^2} \log \frac{1}{1-x} + K + O\left(|1-x|^{\frac{1}{4}-\eta}\right)$$

for some constant K and every (fixed)  $\eta > 0$ 

$$\implies \mathbb{E} H_n = \frac{\sqrt{2\pi}}{\sigma} \cdot \sqrt{n} + O\left(n^{\frac{1}{4} + \eta}\right).$$

Theorem (Meir and Moon, Carr, Goh and Schmutz)

(i) Suppose that  $\varphi_i = \Pr\{\xi = i\} > 0$  for sufficiently large  $i \ge i_0$  and that  $\varphi_{i+1}/\varphi_i \to 0$  as  $i \to \infty$ . Then

 $\mathbb{P}\{|\Delta(T_n) - \delta(n)| \le 1\} = 1 + o(1),$ 

where  $\delta(n) = \max\{k \ge 0 : \mathbb{P}\{\xi \ge k\} \ge 1/n\}.$ 

(ii) If  $\varphi(t) = e^{t-1}$  then there exists a sequence  $\delta'(n)$  that is asymptotically equivalent to  $\delta'(n) \sim \frac{\log n}{\log \log n}$  such that  $\mathbb{P}\{\delta'(n) \leq \Delta(T_n) \leq \delta'(n) + 1\} = 1 + o(1).$ 

(iii) If  $\varphi(t) = 1/(2-t)$  then we have uniformly for all  $k \ge 0$  $\mathbb{P}\{\Delta(T_n) \le k\} = \exp\left(-2^{-(k-\log_2 n+1)}\right) + o(1)$ 

$$y_d(x) = \sum_{n \ge 1} \mathbb{P}\{|T| = n, \ \Delta(T_n) \le d\} x^n.$$

$$\implies y_d(x) = x \varphi_d(y_d(x)) \quad \text{with} \quad \varphi_d(t) = \sum_{i \le d} \varphi_i t^i.$$

$$\implies \mathbb{P}\{|T| = n, \ \Delta(T_n) \le d\} \sim C_d(\varphi'_d(\tau_d))^n n^{-3/2},$$
where  $\tau_d > 0$  is determined by  $\tau_d \varphi'_d(\tau_d) = \varphi_d(\tau_d).$ 

$$\implies \mathbb{P}\{\Delta(T_n) \le d\} = \frac{\Pr\{|T| = n, \, \Delta(T_n) \le d\}}{\Pr\{|T| = n\}} \sim \sqrt{2\pi} \sigma C_d \, (\varphi_d'(\tau_d))^n.$$

### Pólya Trees

 $t_n$  ... number of rooted unlabeled (non-planar) trees

$$t(x) := \sum_{n \ge 1} t_n x^n.$$

$$\implies t(x) = x \exp\left(t(x) + \frac{1}{2}t(x^2) + \frac{1}{3}t(x^3) + \cdots\right).$$

# Polya Trees



$$t(x) = \sum_{n \ge 1} t_n x^n \qquad t(x) = x e^{t(x) + \frac{1}{2}t(x^2) + \frac{1}{3}t(x^3) + \cdots}$$

# Pólya Trees

The **height** of Pólya Trees has the same properties as Galton-Watson trees.

**Theorem** (Goh and Schmutz)

Let  $\Delta(T_n)$  denote the maximum out-degree of Pólya trees of size n. Then

$$\mathbb{P}\{\Delta(T_n) \le k\} = \exp\left(-c_0 \eta^{k-\mu_n}\right) + o(1)$$

with  $c_0 = 3.262..., \eta = 0.3383..., \text{ and } \mu_n = 0.9227... \cdot \log n$ .

.

.

Storing Data

.

.

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

Storing Data

.

.



-

Storing Data

.

.



-

Storing Data

.

.



.

Storing Data

.

.



.

Storing Data

.

.



.

Storing Data

.

.



.

Storing Data

.

.



.

-

Storing Data

.

.



-

#### Storing Data

.



 $\Box$  ... free position

**Probabilistic Model:** 

Every permutation of  $\{1, 2, \ldots, n\}$  is equally likely.

 $\longrightarrow$  probability distribution on binary trees of size n

Height of median of (2t + 1)-variant:

(t = 0: usual binary search trees)

$$\Pr\{H_n \le k+1\} = \sum_{\substack{n_1+n_2=n-1}} \frac{\binom{n_1}{t}\binom{n_2}{t}}{\binom{n_1}{2t+1}} \Pr\{H_{n_1} \le k\} \cdot \Pr\{H_{n_2} \le k\}$$



Generating functions for the median of 2t + 1-variant:

$$y_k(x) = \sum_{n \ge 0} \Pr\{H_n \le k\} \cdot x^n$$

$$y_{k+1}(x)^{(2t+1)} = \frac{(2t+1)!}{(t!)^2} \left( y_k(x)^{(t)} \right)^2$$

with initial conditions  $y_0(x) = 1$ ,  $y_k(0) = 1$ .

Similarly defined as binary search trees (m = 2).

Here every node can store up to m - 1 items and has (at most) m subtrees.

There is also a *fringe balanced version* that corresponds to the *median* of 2t + 1-variant in the binary case.

 $\mathbf{V} = (V_1, V_2, \dots, V_m) \dots \text{ random vector supported on simplex}$  $\Delta = \{(s_1, \dots, s_m) : s_j \ge 0, s_1 + \dots + s_m = 1\} \text{ with density}$ 

$$f(s_1,\ldots,s_m) = \frac{((t+1)m-1)!}{(t!)^m} (s_1\cdots s_m)^t.$$

Lemma The functional equation

$$F(x/\rho_1) = \mathbb{E} (F(xV_1) \cdots F(xV_m))$$

has (up to scaling) a unique solution  $F^{(m,t)}(x)$  with the properties

$$1 - F^{(m,t)}(x) \sim d_1 x^{\beta_1} \log x \quad (x \to 0+)$$

and

$$\lim_{x \to \infty} F^{(m,t)}(x) = 0.$$

Let  $\beta_1 > 0$  be the positive solution of the equation

$$\sum_{j=0}^{(m-1)(t+1)-1} \log(\beta + t + 1 + j) - \log\left(\frac{(m(t+1))!}{(t+1)!}\right)$$
$$= \sum_{j=0}^{(m-1)(t+1)-1} \frac{\beta}{\beta + t + 1 + j}$$

and set

$$\rho_1 = \exp\left(\sum_{j=0}^{(m-1)(t+1)-1} \frac{1}{\beta_1 + t + 1 + j}\right).$$

**Theorem** (Chauvin and D.)

Let  $m \geq 2$  and  $t \geq 0$  be integers. There exist sequences  $c_k$  with

$$\lim_{k \to \infty} \frac{c_{k+1}}{c_k} = \rho_1$$

such that

$$\mathbb{P}\{H_n^{(m,t)} \le k\} = F^{(m,t)}(n/c_k) + o(1).$$

Futher there exists  $\eta > 0$  with

$$\mathbb{P}\{|H_n^{(m,t)} - \mathbb{E}H_n^{(m,t)}| \ge y\} = O(e^{-\eta y}).$$

In particular we have, as  $n \to \infty$ ,

$$\mathbb{V}H_n^{(m,t)} = O(1).$$

**Combinatorial Description:** 

- labeled rooted tree
- labels are strictly increasing (starting at the root)
- no left-to-right order (non-planar)

.



.

.

(1)

.

.

•

1

.

-

.

.

1 2 3 .

.

.



.

.

.



.

.

.



.

.



Remark


Number of Recursive Trees:

$$y_n$$
 = number of recusive trees of size  $n$   
=  $(n-1)!$ 

The node with label j has exactly j - 1 possibilities to be inserted  $\implies y_n = 1 \cdot 2 \cdots (n - 1).$ 

**Generating Functions:** 

$$y(x) = \sum_{n \ge 1} y_n \frac{x^n}{n!} = \sum_{n \ge 1} \frac{x^n}{n} = \log \frac{1}{1-x}$$

$$y'(x) = 1 + y(x) + \frac{y(x)^2}{2!} + \frac{y(x)^3}{3!} + \dots = e^{y(x)}$$

$$R = \bigcirc + \bigcirc R + \bigcirc R + \bigcirc R + \cdots$$

A recursive tree can be interpreted as a root followed by an unordered sequence of recursive trees.  $(y'(x) = \sum_{n \ge 0} y_{n+1}x^n/n!)$ 

**Probability Model:** 

Process of growing trees

- The process starts with the root that is labeled with 1.
- At step j a new node (with label j) is attached to any previous node with probability 1/(j-1).

After n steps every tree (of size n) has equal probability 1/(n-1)!.

.

.

(1)

.

.

.

*p* = 1

.

.

.

1 *p* = 1 (2)

.

.

. 
$$p = 1/2$$
 (2)  $p = 1/2$ 

.



#### Theorem

Let  $H_n$  denote the height of random recursive trees of size n.

$$\mathbb{E} H_n = e \log n + O\left(\sqrt{\log n} \left(\log \log n\right)\right).$$

Furthermore we have (uniformly for all  $k \ge 0$  as  $n \to \infty$ )

$$\mathbb{P}\{H_n \leq k\} = F(n/y'_k(1)) + o(1),$$

where F(y) satisfies the integral equation

$$y F(y/e^{1/e}) = \int_0^y F(z/e^{1/e}) F(y-z) \, dz. \tag{1}$$

Moreover, as  $n \to \infty$ ,

$$\mathbb{V}H_n = O(1)$$

and there exist  $\eta > 0$  and c > 0 such that

$$\mathbb{P}\{|H_n - \mathbb{E} H_n| \ge y\} \le c e^{-\eta y}$$

for all  $y \ge 0$ .

**Theorem** (Szymanski, Devroye and Lu, Goh and Schmutz)

Let  $\Delta(T_n)$  denote the maximum out-degree or random recursive trees. Then we have  $\mathbb{E}\Delta(T_n) \sim \log_2 n$  and the distribution is given by

$$\mathbb{P}\{\Delta(T_n) \leq k\} = \exp\left(-2^{-(k-\log_2 n+1)}\right) + o(1).$$

.

.

(1)

.

.

.

(1) p = 1

.

.

.

1 *p* = 1 (2)

.

.

$$p = \frac{1}{3} \qquad \begin{array}{c} 1 \\ p = \frac{1}{3} \\ 2 \\ p = \frac{1}{3} \end{array}$$

.

.



Remark



Number of Plane Oriented Trees:

$$y_n = \text{number of plane oriented trees of size } n$$
$$= 1 \cdot 3 \cdot 5 \cdots (2n - 3) = (2n - 3)!!$$
$$= \frac{(2n - 2)!}{2^{n-1}(n-1)!}$$

The node with label j has exactly 2j - 3 possibilities to be inserted  $\implies y_n = 1 \cdot 3 \cdots (2n - 3).$ 

**Generating Functions:** 

$$y(x) = \sum_{n \ge 1} y_n \frac{x^n}{n!} = \sum_{n \ge 1} \frac{1}{2^{n-1}} \binom{2(n-1)}{n-1} \frac{x^n}{n} = 1 - \sqrt{1-2x}$$

$$y'(x) = 1 + y(x) + y(x)^2 + y(x)^3 + \dots = \frac{1}{1 - y(x)}$$

$$R = \bigcirc + \bigcirc + \bigcirc + \bigcirc + \bigcirc + \bigcirc + \cdots$$

A plane oriented tree can be interpreted as a root followed by an **ordered** sequence of plane oriented trees.  $(y'(x) = \sum_{n \ge 0} y_{n+1}x^n/n!)$ 

**Probability Model:** 

Process of growing trees

- The process starts with the root that is labeled with 1.
- At step j a new node (with label j) is attached to any previous node of outdegree d with probability (d+1)/(2j-3).

After n steps every tree (of size n) has equal probability 1/(2n-3)!!.

**Probability Model:** 

Process of growing trees

- The process starts with the root that is labeled with 1.
- At step j a new node (with label j) is attached to any previous node of outdegree d with probability proportial to d + r (for some r > 0).

For d = 1 we get plane oriented trees.

.

.

(1)

.

.

.

*p* = 1

.

.

.

1 *p* = 1 (2)

.

.

.

$$p = c/(2+r)$$
 1  $p = c/(2+r)$   
2  $p = c/(1+r)$ 

.

.



#### **Generating Functions**

 $y_n \ldots$  weighted sum of plane oriented trees (according to probability distribution)

$$y(x) = \sum_{n \ge 1} y_n \frac{x^n}{n!} \dots$$
 generating function

$$y'(x) = \frac{1}{(1-y(x))^r}$$

$$\implies \qquad y(x) = 1 - (1 - (r+1)x)^{\frac{1}{r+1}}$$

#### **Degree distribution**

Set

 $\lambda_d = \lim_{n \to \infty} \mathbf{P} \text{ (a random node in a tree of size } n \text{ has out-degree } d)$  $= \lim_{n \to \infty} \frac{\text{expected number of nodes with out-degree } d}{n}$ 

Then

$$\lambda_d = \frac{(r+1)\Gamma(2r+1)\Gamma(r+d)}{\Gamma(r)\Gamma(2r+d+2)}$$

We have a scalefree distribution

$$\lambda_d \sim \frac{(r+1)\Gamma(2r+1)}{\Gamma(r)} \cdot d^{-2-r}.$$

#### Theorem

ullet

Suppose that  $r = \frac{A}{B} > 0$  is rational. Then

$$\mathbf{P}\{H_n \le k\} = F(c_k - d_r \log n) + o(1),$$

where  $c_k = k + O(\log k)$  and Set  $d_r = 1/((r+1)s)$  with  $r s e^{s+1} = 1$ .

Further,  $F(x) = \Psi(e^{-x})$ , where  $\Psi(y)$  is calculated by the following procedure.

Let  $\Phi(y)$  be the solution of

$$y^{\frac{1}{A+B}} \Phi(ye^{-1/d_r}) = \frac{\Gamma\left(1 + \frac{1}{A+B}\right)}{\Gamma\left(\frac{1}{A+B}\right)^{A+B+1}} \times \\ \times \int_{\substack{y_1 + \dots + y_{A+B+1} = y, y_j \ge 0 \\ y_1 + \dots + y_{A+B+1} = y, y_j \ge 0 }} \prod_{j=1}^{B+1} \left(\Phi(y_j e^{-1/d_r}) y_j^{\frac{1}{A+B}-1}\right) \\ \times \prod_{\ell=B+2}^{A+B+1} \left(\Phi(y_\ell) y_\ell^{\frac{1}{A+B}-1}\right) dy$$

Then

$$\Psi(y) = \frac{\Gamma\left(\frac{A}{A+B}\right)}{\Gamma\left(\frac{1}{A+B}\right)^{A}} \int_{z_{1}+\dots+z_{A}=1, z_{j}\geq 0} \prod_{j=1}^{A} \left(\Phi(yz_{j})z_{j}^{\frac{1}{A+B}-1}\right) d\mathbf{z}$$

Further, for all r > 0 there exist  $\eta > 0$  and c > 0 such that

$$\mathbb{P}\{|H_n - \mathbb{E} H_n| \ge y\} \le c e^{-\eta y}$$

for all  $y \ge 0$ .

**Remark.** The concentration property of the height will be proved in the last part of the talk.

Theorem (Mori)

Let  $\Delta(T_n)$  denote the maximum degree of scale-free trees with parameter r > 0. Then

$$\frac{\Delta(T_n)}{n^{\frac{1}{1+r}}} \to \mu \quad (a.s.)$$

for some random variable  $\mu$  (that is related to the degree distribution of  $T_n$ ). Further

$$\frac{\Delta(T_n) - \mu n^{\frac{1}{1+r}}}{\sqrt{\mu n^{\frac{1}{1+r}}}} \xrightarrow{\mathsf{d}} N(0,1).$$

Tries are rooted trees which are used to store data which are labeled with a (possibly) infinite string of symbols from a binary (or generally finite) alphabet.

Each string  $\mathbf{x}_j$  defines an infinite path in the (infinite) binary tree. Let  $u_j$  denote the node where the suffix part starts. Then we can *trim* the tree by cutting away everything below node  $u_j$ . The node  $u_j$  becomes now a leaf representing  $\mathbf{x}_j$ . If we repeat this procedure for all labels  $\mathbf{x}_1, \ldots, \mathbf{x}_n$  then we obtain a finite binary tree with n nodes, called the **trie**.

It is usual to assume that the strings follow the Bernoulli model with probability  $p \in (0, 1)$ , or equivalently they are obtained from a memoryless source on 2 symbols.

#### **Theorem** (Flajolet)

Let  $H_n$  denote the height of tries in the symmetric case  $p = \frac{1}{2}$ . Then

$$\mathbb{P}\{H_n \le k\} = \exp\left(-2^{-(k-2\log_2 n)-1}\right) + o(1).$$

Furthermore, as  $n \to \infty$ ,

$$\mathbb{E} H_n = 2 \log_2 n + O(1), \quad \mathbb{V} H_n = O(1)$$

and there exist  $\eta > 0$  and c > 0 such that

$$\left|\mathbb{P}\{|H_n - \mathbb{E} H_n| \ge y\} \ll e^{-\eta y}\right|$$

for all  $y \ge 0$ .

Remark. 
$$\mathbb{P}{H_n \leq k} = \frac{n!}{2^{kn}} \binom{2^k}{n}$$

#### Theorem (Devroye)

Let  $H_n$  denote the height of tries of generated by iid labels with f(x) as a density on [0,1]. If  $C := \int_0^1 f(x)^2 dx < \infty$  then

$$\mathbb{P}\{H_n \le k\} = \exp\left(-C \, 2^{-(k-2\log_2 n)-1}\right) + o(1).$$

Furthermore,

$$\mathbb{E} H_n = 2 \log_2 n + O(1).$$

If  $\int_0^1 f(x)^2 dx = \infty$  then  $\mathbb{E} H_n = \infty$  for all  $n \ge 2$ .

#### Theorem (Pittel)

Let  $H_n$  denote height of tries of generated by iid labels on an *m*-ary alphabet with probability distribution  $p_1, p_2, \ldots, p_m$  and set

$$b = \left(\sum_{i=1}^{m} p_i^2\right)^{-\frac{1}{2}}$$

•

Then

$$\mathbb{P}\{H_n \leq k\} = \exp\left(-\frac{1}{2}b^{-2(k-\log_b n)}\right) + o(1).$$

# **Digital Search Trees**

Digital search trees are again rooted trees which are used to store data which are labeled with a (possibly) infinite string of symbols from a binary (or generally finite) alphabet.

The empty string is stored in the root, while the first item occupies the right or left child of the root depending whether its first symbol is "1" or "0". The remaining items are always stored in the next available node according to the rule that we move to the right if the next symbol is "1" and we move to the left if the next symbol is "0". In the same way we can also search for a specific item.

We assume that the 0-1-strings follow the Bernoulli model with probability  $p \in (0, 1)$
# **Digital Search Trees**

#### Theorem

Let  $H_n$  denote the height of digital search trees in the symmetric case  $p = \frac{1}{2}$ . Then there exists a sequence  $k_n$  that is asymptotically given by

$$k_n = \log_2 n + \sqrt{2\log_2 n} - \log_2 \left(\sqrt{2\log_2 n}\right) + O(1)$$

such that

$$\mathbb{P}\{k_n \le H_n \le k_n + 1\} = 1 + o(1).$$

Furthermore, if  $0 \le y \le c_1 \sqrt{\log n}$  (for some constant  $c_1 > 0$ ) then there exist constants  $c_2, c_3 > 0$  with

$$\mathbb{P}\{|H_n - k_n| \ge y\} \le c_2 e^{-yc_3\sqrt{\log_2 n}}.$$

 $y_n \ldots$  weighted sum of plane oriented trees (according to probability distribution)

$$y(x) = \sum_{n \ge 1} y_n \frac{x^n}{n!} \dots$$
 generating function

$$y'(x) = \frac{1}{(1-y(x))^r}$$

$$\implies y(x) = 1 - (1 - (r+1)x)^{\frac{1}{r+1}}$$

$$y_n = n! (-1)^{n-1} (r+1)^n {\binom{1/(r+1)}{n}}.$$

$$y_k(z) = \sum_{n \ge 0} y_n \mathbb{P}\{H_n \le k\} \frac{z^n}{n!}$$
  
$$\implies y'_{k+1}(z) = \frac{1}{(1 - y_k(z))^r} \qquad (y_0(z) = 0, \ y_{k+1}(0) = 0).$$
  
$$Y_k(z) := y'_k(z)$$

$$\implies Y'_{k+1}(z) = r Y_{k+1}(z)^{1+\frac{1}{r}} Y_k(z) \qquad (Y_1(z) = 1, Y_{k+1}(0) = 1).$$

$$Y_k(z) = y'_k(z) = \sum_{n \ge 0} y_{n+1} \mathbb{P}\{H_{n+1} \le k\} \frac{z^n}{n!}.$$

### Lemma 1

 $Y_1(z), Y_2(z), \overline{Y}_1(z), \overline{Y}_2(z) \dots$  non-negative, continuous functions, defined for  $z \ge 0$ 

$$\begin{split} Y_1(0) < \overline{Y}_1(0), \quad Y_2(0) < \overline{Y}_2(0), \\ Y'_2(z) = r Y_2(z)^{1+\frac{1}{r}} Y_1(z), \quad \overline{Y}'_2(z) = r \overline{Y}_2(z)^{1+\frac{1}{r}} \overline{Y}_1(z). \\ \\ \overline{Y}_1(z) - Y_1(z) \quad \text{has exactly one positive zero} \end{split}$$

 $\implies \overline{Y}_2(z) - Y_2(z)$  has at most one positive zero.

Proof of Lemma 1.

$$y_j(z) := \int_0^z Y_j(t) dt \quad \text{and} \quad \overline{y}_j(z) = \int_0^z \overline{Y}_j(t) dt \qquad (j = 1, 2)$$
  

$$\overline{y}'_1(z) - y'_1(z) = \overline{Y}_1(z) - Y_1(z)$$
  

$$\implies \quad \overline{y}'_1(z) - y'_1(z) \text{ has exactly one positive zero } \zeta.$$
  

$$\implies \quad \overline{y}_1(z) - y_1(z) \text{ increasing for } 0 \le z \le \zeta \text{ and decreasing for } z \ge \zeta.$$
  

$$\implies \quad \overline{y}_1(z) - y_1(z) \text{ has at most one positive zero.}$$
  

$$\implies \quad \overline{Y}_2(z) - Y_2(z) = \overline{y}'_2(z) - y'_2(z)$$
  

$$= \frac{(1 - \overline{y}_1(z))^{-r} - (1 - y_1(z))^{-r}}{\overline{y}_1(z) - y_1(z)} (\overline{y}_1(z) - y_1(z))$$

has at most one positive zero, too.

### Lemma 2

The sequence  $Y_k(1/(r+1))$  is log-concave, that is,

$$\frac{Y_{k+2}(1/(r+1))}{Y_{k+1}(1/(r+1))} \le \frac{Y_{k+1}(1/(r+1))}{Y_k(1/(r+1))}.$$

Corollary

The limit

$$\lim_{k \to \infty} \frac{Y_{k+1}(1/(r+1))}{Y_k(1/(r+1))} =: \rho \ge 1$$

exists.

#### Proof of Lemma 2

For  $0 \leq \gamma < 1$  set

$$V_k(z,\gamma) = \begin{cases} (1 - (r+1)z)^{-r/(r+1)} & \text{for } 0 \le z \le \frac{1}{r+1}(1-\gamma), \\ \gamma^{-r/(r+1)}Y_k\left(\frac{z - \frac{1}{r+1}(1-\gamma)}{\gamma}\right) & \text{for } \frac{1}{r+1}(1-\gamma) \le z \le \frac{1}{r+1}. \end{cases}$$

$$\implies V_{k+1}'(z,\gamma) = r V_{k+1}(z,\gamma)^{1+\frac{1}{r}} V_k(z,\gamma),$$

 $(V_k(0) = 1 \text{ and } V_k(1/(r+1), \gamma) = \gamma^{-r/(r+1)}Y_k(1/(r+1)).)$ 

Lemma 1  $\implies$   $Y_{k+1}(z) - V_k(z, \gamma)$  has (at most) one positive zero.

$$\gamma_k := \left(\frac{Y_k(1/(r+1))}{Y_{k+1}(1/(r+1))}\right)^{1+\frac{1}{r}}$$

 $\implies V_k(1/(r+1), \gamma_k) = Y_{k+1}(1/(r+1)).$ 

 $\implies \zeta = 1/(r+1)$  is the only positive zero of  $Y_{k+1}(z) - V_k(z,\gamma)$ .

$$\implies Y_{k+1}(z) \le V_k(z, \gamma_k) \text{ for } 0 \le z \le \frac{1}{r+1}$$

+ integration

$$\implies Y_{k+2}(1/(r+1)) \leq V_{k+1}(1/(r+1),\gamma_k) \\ = \gamma_k^{-r/(r+1)} Y_{k+1}(1/(r+1)) \\ = \frac{Y_{k+1}(1/(r+1))^2}{Y_k(1/(r+1))}.$$

Lemma 3

For  $0 \le z < \frac{1}{r+1}$  and  $k \ge 1$  we have  $Y(z) - Y_k(z) \le \left(\frac{2r+1}{r+1}\right)^k \sum_{\ell \ge k} \frac{1}{\ell!} \left(\log \frac{1}{1-(r+1)z)}\right)^\ell.$ 

Proof of Lemma 3

By induction:

$$k = 1$$
:  $Y_1(z) = 1$  and  $(1 - (r+1)z)^{-r/(r+1)} \le (1 - (r+1)z)^{-1}$ .

 $k \to k + 1:$   $Y(z)' - Y_{k+1}(z)' = r \left( Y(z)^{2 + \frac{1}{r}} - Y_{k+1}(z)^{1 + \frac{1}{r}} Y_k(z) \right)$   $\leq r \left( Y(z)^{2 + \frac{1}{r}} - Y_k(z)^{2 + \frac{1}{r}} \right)$   $\leq r \left( 2 + \frac{1}{r} \right) Y(z)^{1 + \frac{1}{r}} \left( Y(z) - Y_k(z) \right)$  $= (2r + 1) \frac{1}{1 - (r + 1)z} \left( Y(z) - Y_k(z) \right)$ 

Lemma 4

Suppose that  $1 < C < e^{r/(r+1)}$  and  $C > \left(2 + \frac{1}{r}\right) e \log C$  $\implies \frac{Y_{k+1}(1/(r+1))}{Y_k(1/(r+1))} \ge C$ 

for all  $k \geq 1$ .

Corollary.  $\rho > 1$ .

**Proof of Lemma 4.** Set  $c = C^{-r/(r+1)}$  and  $z_0 = \frac{1}{r+1}(1-c^k)$ .

Lemma 3  $\implies$  (with some  $c_1 > 0$ )

$$\begin{split} Y(z_0) - Y_k(z_0) &\leq \left(\frac{2r+1}{r+1}\right)^k \sum_{\ell \geq k} \frac{1}{\ell!} \left(k \log \frac{1}{c}\right)^\ell \\ &\leq c_1 \left(\frac{2r+1}{r+1}\right)^k \frac{\left(k \log \frac{1}{c}\right)^k}{k!} \\ &\leq c_1 \left(\frac{2r+1}{r+1}e \log \frac{1}{c}\right)^k \\ &= c_1 \left(\left(2+\frac{1}{r}\right)e \log C\right)^k. \end{split}$$

 $Y(z_0) = C^k, \ C > \left(2 + \frac{1}{r}\right) e \log C$  $\implies Y_k(1/(r+1)) \ge Y_k(z_0) \ge C^k(1+o(1)).$ 

$$\implies \rho = \lim_{k \to \infty} \frac{Y_{k+1}(1/(r+1))}{Y_k(1/(r+1))} \ge C.$$

Lemma 2  $\implies$  the sequence  $\frac{Y_{k+1}(1/(r+1))}{Y_k(1/(r+1))}$  is decreasing.

$$\implies \frac{Y_{k+1}(1/(r+1))}{Y_k(1/(r+1))} \ge C.$$

### Lemma 5

If  $n \geq Y_k(1/(r+1))^{1+\frac{1}{r}}$  then

$$\mathbb{P}\{H_n \le k\} = O\left(Y_k(1/(r+1)) \cdot n^{-r/(r+1)}\right).$$

Conversely if  $n \leq Y_k(1/(r+1))^{1+\frac{1}{r}}$  then

$$\mathbb{P}\{H_n > k\} = O\left(Y_k(1/(r+1))^{-1-\frac{1}{r}} \cdot n\right).$$

### Proof of Lemma 5.

Define 
$$z_k$$
 by  $Y(z_k) = Y_k(1/(1+r))$ .  
 $\implies z_k = \frac{1}{r+1} \left( 1 - \eta_k^{-1 - \frac{1}{r}} \right)$  with  $\eta_k = Y_k(1/(1+r))$ .

Set

$$\tilde{Y}(z) = (z_k(r+1))^{r/(r+1)} Y(z_k(r+1)z)$$
$$\implies \tilde{Y}'(z) = r\tilde{Y}(z)^{2+\frac{1}{r}}.$$

 $\tilde{Y}(0) < 1 + Lemma 1 \Longrightarrow$ 

• 
$$ilde{Y}(z) \leq Y_k(z)$$
 for  $0 \leq z \leq z_k$ ,

• 
$$\tilde{Y}(z) \ge Y_k(z)$$
 for  $z \ge z_k$ .

 $\mathbb{P}\{H_{n+1} \le k\} \le \mathbb{P}\{H_n \le k\}$ 

 $z \ge z_k \Longrightarrow$ 

$$\begin{split} \tilde{Y}(z) &\geq Y_k(z) \\ &\geq \sum_{\ell=0}^{n-1} y_{\ell+1} \mathbb{P}\{H_{\ell+1} \leq k\} \frac{z^{\ell}}{\ell!} \\ &\geq \mathbb{P}\{H_n \leq k\} \sum_{\ell=0}^{n-1} y_{\ell+1} \frac{z^{\ell}}{\ell!} \end{split}$$

 $n \ge \eta_k^{1+\frac{1}{r}}, \ z = \frac{1}{r+1}:$   $\tilde{Y}(1/(r+1)) \le c_2 \eta_k,$   $\sum_{\ell=0}^{n-1} y_{\ell+1} \frac{(r+1)^{\ell}}{\ell!} \ge c_2 n^{r/(r+1)}$   $\implies \qquad \mathbb{P}\{H_n \le k\} \le c_4 \eta_k n^{-r/(r+1)}.$ 

$$0 \leq z \leq z_k \Longrightarrow$$

$$Y(z) - \tilde{Y}(z) \geq Y(z) - Y_k(z)$$

$$\geq \sum_{\ell=n-1}^{\infty} y_{\ell+1} \mathbb{P}\{H_{\ell+1} > k\} \frac{z^{\ell}}{\ell!}$$

$$\geq \mathbb{P}\{H_n > k\} \sum_{\ell=n-1}^{\infty} y_{\ell+1} \frac{z^{\ell}}{\ell!}.$$

$$n \leq \eta_k^{1+\frac{1}{r}}, \ z' = \frac{1}{r+1} \left(1 - \frac{1}{n}\right) \leq z_k.$$

$$Y(z') - \tilde{Y}(z') \leq c_5 n^{1+\frac{r}{r+1}} \eta_k^{-1-\frac{1}{r}},$$

$$\sum_{\ell=n-1}^{\infty} y_{\ell+1} \frac{(z')^{\ell}}{\ell!} \geq c_6 n^{r/(r+1)}$$

$$\Longrightarrow \qquad \left[\mathbb{P}\{H_n > k\} \leq c_7 \eta_k^{-1-\frac{1}{r}} n\right]$$

#### Lemma 6

Let 
$$k(n) := \max\{\ell \ge 1 : Y_{\ell}(1/(r+1))^{1+\frac{1}{r}} \le n\}$$
. Then  
 $\mathbb{E} H_n = k(n) + O(1)$ 

and there exist  $\eta > 0$  and c > 0 such that

$$\mathbb{P}\{|H_n - \mathbb{E}H_n| > y\} \le ce^{-\eta y}$$

(2)

for all y > 0.

#### Proof of Lemma 6.

Lemma 4  $\Longrightarrow$ 

$$\frac{Y_{k(n)+\ell}(1/(r+1))}{Y_{k(n)}(1/(r+1))} \ge C^{\ell} \qquad (\ell \ge 0).$$

$$\implies \mathbb{P}\{H_n > k(n) + \ell\} \le c_8 Y_{k(n) + \ell} (1/(r+1))^{-1 - \frac{1}{r}} n \\ \le c_8 C^{\frac{r+1}{r}\ell} Y_{k(n)} (1/(r+1))^{-1 - \frac{1}{r}} n \\ \le c_9 C^{-\frac{r+1}{r}\ell}.$$

Similarly

$$\mathbb{P}\{H_n \le k(n) - \ell\} \le c_{10}C^{-\ell}.$$

#### Remark

The above proof provides stong concentration around the mean but it does not say where the mean value actually is.

It is related to the the actual growth of  $Y_k(1/(r+1))$ . Thus one has to analyze the recurrence  $Y'_{k+1}(z) = r(Y_{k+1}(z))^{1+\frac{1}{r}}Y_k(z)$  in more detail.

Thank You!