

Distributional Analysis
of
Recursive Algorithms
by the
Contraction Method

Ralph Neininger
Universität Freiburg

INRIA–Rocquencourt, November 1999

Contents

Partial match query (pp. 3–40)

Internal path length (pp. 41–49)

FIND (pp. 50–54)

Partial Match Query

retrieve multiattribute records belonging to some k -dimensional domain

$$D = D_1 \times \dots \times D_k$$

Given:

file $F \subset D$

query $q = (q_1, \dots, q_k)$,

$$q \in (D_1 \cup \{*\}) \times \dots \times (D_k \cup \{*\})$$

Problem:

Find all records in the file F which *satisfy* the query q , i.e. find all $r = (r_1, \dots, r_k) \in F$ with

$$r_j = q_j \quad \text{for all } 1 \leq j \leq k \quad \text{with } q_j \neq *.$$

The *specification pattern* of a query q is the word $u = (u_1, \dots, u_k) \in \{S, *\}^k$ with

$$\begin{aligned} u_j = S & \quad \text{if } q_j \in D_j \quad (\text{specified}) \\ u_j = * & \quad \text{if } q_j = * \quad (\text{unspecified}) \end{aligned}$$

The model

For the attributes' domains assume

$$D_i = [0, 1] \quad \text{for all } 1 \leq i \leq k.$$

The **uniform probabilistic model** assumes all attributes (=components, coordinates) in the files and queries to be independent and uniformly distributed over the unit interval.

Appropriate for comparison based algorithms.

Comparison based structures

Quadtree (Finkel, Bentley '74)

K-d trees

- “classical” *K*-d tree (Bentley '75)

- *K*-d-*t* tree (locally balanced version)

- . (Cunto, Lau, Flajolet '89)

- random relaxed *K*-d tree (randomized version)

- . (Duch, Estivill-Castro, Martínez '98)

- squarish *K*-d tree

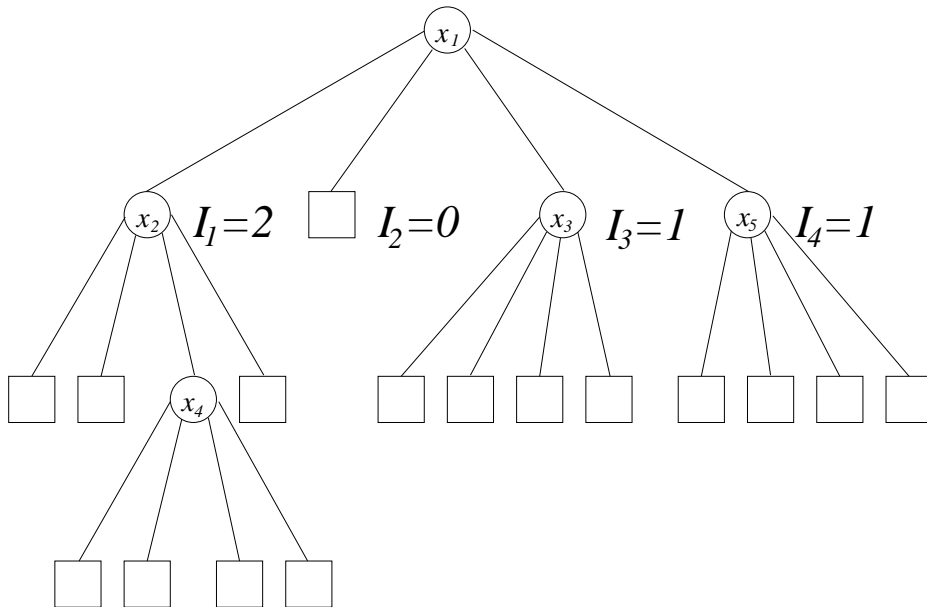
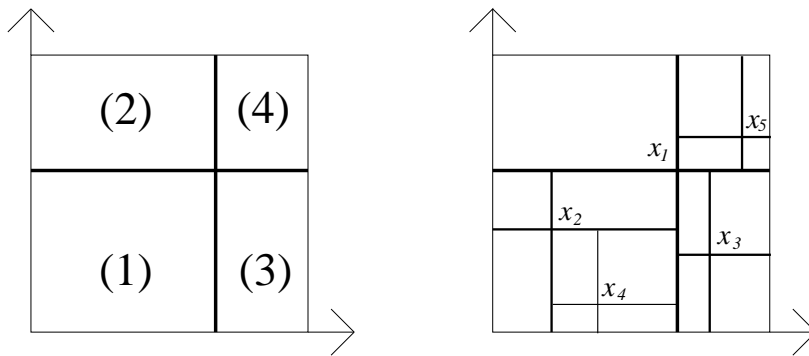
- . (Devroye, Jabbour, Zamora-Cura '99)

A digital structure: The 2-d trie

- . (Schachinger '99)

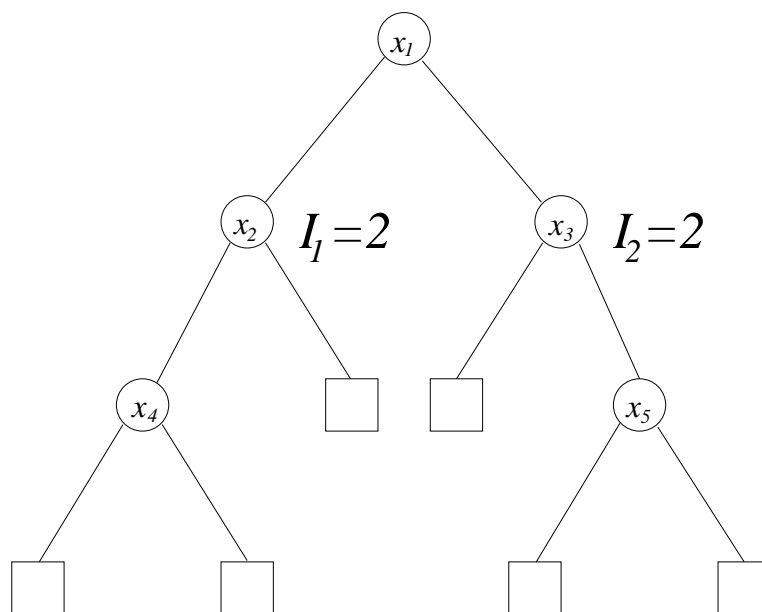
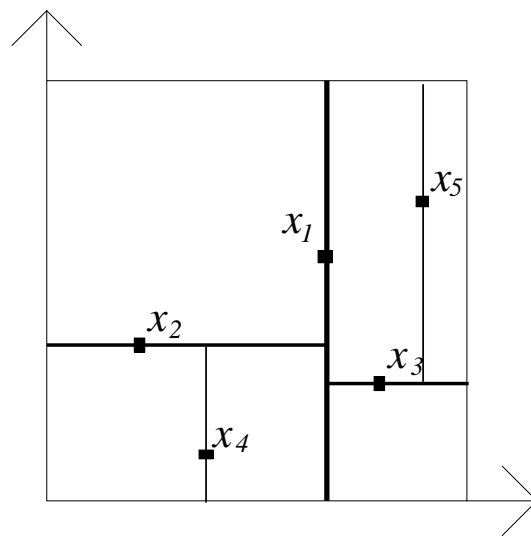
The Quadtree

The data partition the unit-cube recursively into quadrants. The quadtree corresponds to this partitioning.



The K -d-Tree

The components of the data are used cyclically as discriminators.



Partial match query in a quadtree

Dimension $d=2$

$W = (U, V)$ the first key to be inserted (the root)

$q = (Y, *)$ the query

Y, U, V independent, uniformly distr. on $[0, 1]$

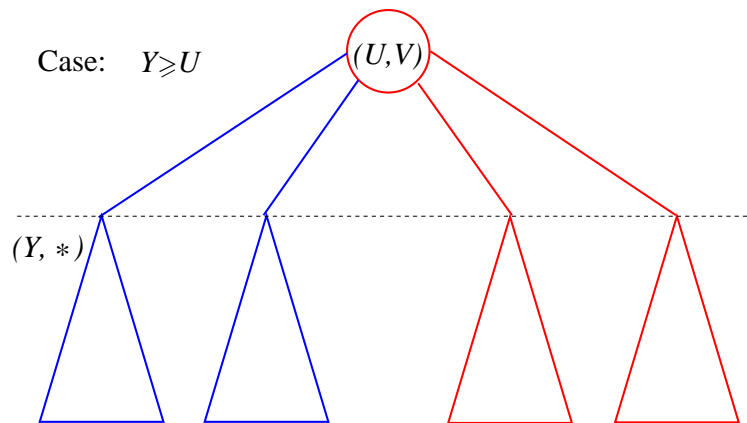
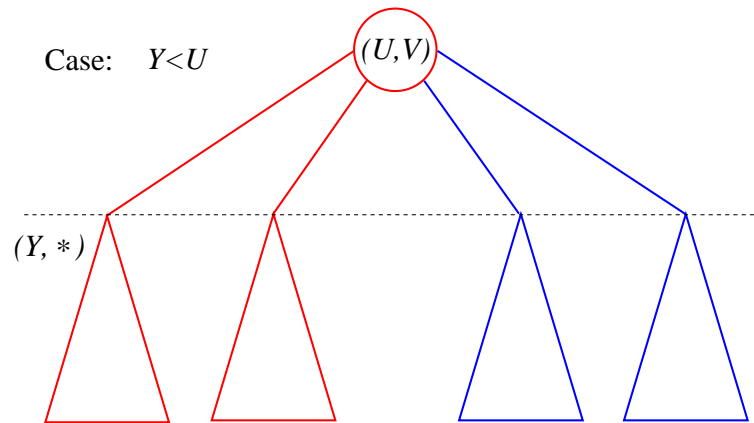
$I^{(n)}$ vector of the cardinalities of the subtrees of the root

$\langle W \rangle = (UV, U(1-V), (1-U)V, (1-U)(1-V))$ volumes of the generated quadrants

C_n cost of a partial match query in a tree of size n , measured as the number of nodes traversed during the search

Search for query $(Y, *)$

Search for $(Y, *)$ with root (U, V)



Distributional recursive equation for the cost C_n

$$C_n \stackrel{\mathcal{D}}{=} \mathbf{1}_{\{Y < U\}} \left(C_{I_1^{(n)}}^{(1)} + C_{I_2^{(n)}}^{(2)} \right) \\ + \mathbf{1}_{\{Y \geq U\}} \left(C_{I_3^{(n)}}^{(3)} + C_{I_4^{(n)}}^{(4)} \right) + 1$$

with

$Y, U, V, (C_i^{(1)})_{i \in \mathbb{N}}, \dots, (C_i^{(4)})_{i \in \mathbb{N}}$ independent

Y, U, V uniformly distributed on $[0, 1]$

$$C_i^{(k)} \stackrel{\mathcal{D}}{=} C_i, \quad 1 \leq k \leq 4$$

$I^{(n)}$ given $(U, V) = w$ multinomial $M(n - 1, \langle w \rangle)$
distributed

Moments of the cost C_n

$$\mathbb{E}C_n \sim \gamma n^{\alpha-1},$$

$$\alpha = \frac{\sqrt{17} - 1}{2}, \quad \gamma = \frac{\Gamma(2\alpha)}{2\Gamma^3(\alpha)}$$

(Flajolet, Gonnet, Puech, Robson '93)

$$\text{Var}(C_n) \sim \beta n^{2\alpha-2},$$

$$\beta = \frac{(2\alpha - 1)\Gamma(2\alpha)}{3\alpha(\alpha - 1)\Gamma^4(\alpha)} - \frac{\Gamma^2(2\alpha)}{4\Gamma^6(\alpha)}$$

(Martínez, Panholzer, Prodinger '98)

Normalization:

$$X_n = \frac{C_n - \mathbb{E}C_n}{n^{\alpha-1}}$$

The Limiting equation

$$\begin{aligned}
 X_n \stackrel{\mathcal{D}}{=} & \mathbf{1}_{\{Y < U\}} \sum_{k=1}^2 \left(\frac{I_k^{(n)}}{n} \right)^{\alpha-1} \left(X_{I_k^{(n)}}^{(k)} + \gamma \right) \\
 & + \mathbf{1}_{\{Y \geq U\}} \sum_{k=3}^4 \left(\frac{I_k^{(n)}}{n} \right)^{\alpha-1} \left(X_{I_k^{(n)}}^{(k)} + \gamma \right) \\
 & - \gamma + o(1)
 \end{aligned}$$

Recall

$$\begin{aligned}
 \frac{I^{(n)}}{n} & \xrightarrow{\mathbb{P}} \langle W \rangle \\
 & = (UV, U(1-V), (1-U)V, (1-U)(1-V))
 \end{aligned}$$

Limiting equation:

$$\begin{aligned}
 X \stackrel{\mathcal{D}}{=} & \mathbf{1}_{\{Y < U\}} \sum_{k=1}^2 \langle W \rangle_k^{\alpha-1} (X^{(k)} + \gamma) \\
 & + \mathbf{1}_{\{Y \geq U\}} \sum_{k=3}^4 \langle W \rangle_k^{\alpha-1} (X^{(k)} + \gamma) - \gamma
 \end{aligned}$$

The limiting operator

$$T : M^1(\mathbb{R}^1, \mathcal{B}^1) \longrightarrow M^1(\mathbb{R}^1, \mathcal{B}^1)$$

$$T(\mu) \stackrel{\mathcal{D}}{=} \mathbf{1}_{\{Y < U\}} \sum_{k=1}^2 \langle W \rangle_k^{\alpha-1} (Z^{(k)} + \gamma) \\ + \mathbf{1}_{\{Y \geq U\}} \sum_{k=3}^4 \langle W \rangle_k^{\alpha-1} (Z^{(k)} + \gamma) - \gamma$$

$Y, U, V, Z^{(1)}, \dots, Z^{(4)}$ are independent

$$W = (U, V)$$

Y, U, V uniformly distr. on $[0, 1]$

$$Z^{(1)}, \dots, Z^{(4)} \stackrel{\mathcal{D}}{=} \mu$$

The Wasserstein-metric

$$l_2(\mu, \nu) := \inf\{(\mathbb{E}(X - Y)^2)^{1/2} : X \stackrel{\mathcal{D}}{=} \mu, Y \stackrel{\mathcal{D}}{=} \nu\}$$

$$M_{0,2} := \{\mu \in M^1(\mathbb{R}^1, \mathcal{B}^1) : \mathbb{E}\mu = 0, \text{Var}(\mu) < \infty\}$$

The infimum is attained for “optimal couplings”.

$(M_{0,2}, l_2)$ is a complete metric space.

$$l_2(\mu_n, \mu) \rightarrow 0 \quad \iff \quad \mu_n \stackrel{\mathcal{D}}{\rightarrow} \mu \text{ plus} \\ \text{conv. of seconds moments}$$

Contraction property

Lemma: *The limiting operator $T : M_{0,2} \rightarrow M_{0,2}$ is a contraction w.r.t. l_2 :*

$$l_2(T(\mu), T(\nu)) \leq \xi l_2(\mu, \nu) \quad \forall \mu, \nu \in M_{0,2},$$

$$\xi = \frac{2}{\sqrt{19 - 3\sqrt{17}}} = 0.776\dots$$

Banach's fixed point theorem:

T has a unique fixed point ρ in $M_{0,2}$ and

$$l_2(T^n(\mu), \rho) \rightarrow 0$$

exponentially fast for any $\mu \in M_{0,2}$.

Exponential rate \longrightarrow quick approximation of the fixed point by iterating the limiting operator

The Limit Law

Theorem: (Limit law for PMQ in 2-dim. quadrees)
Let X_n be the normalized number of traversed nodes and X the fixed point of the limiting operator (i.e. $T(\mathcal{L}(X)) = \mathcal{L}(X)$), then:

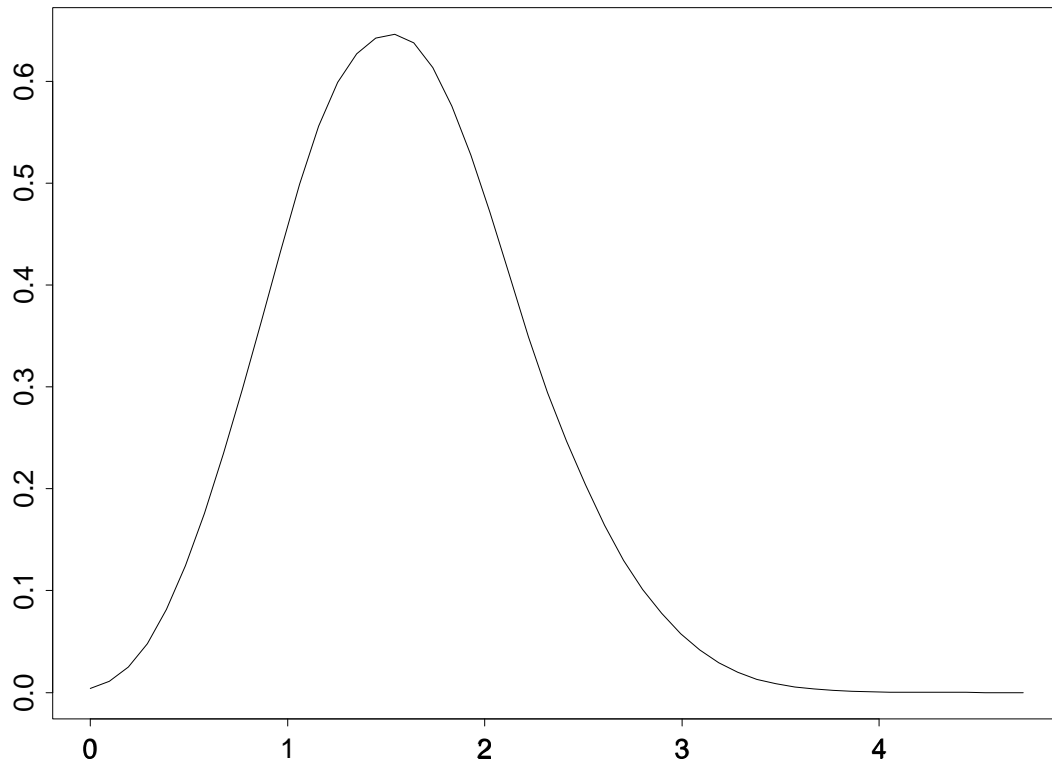
$$l_2(X_n, X) \rightarrow 0.$$

The translated limit $\tilde{X} := X + \gamma$ is the unique solution in $M_{\gamma,2}$ of the equation

$$Z \stackrel{\mathcal{D}}{=} U^{\frac{\alpha-1}{2}} \left(V^{\alpha-1} Z^{(1)} + (1-V)^{\alpha-1} Z^{(2)} \right).$$

Approximation of the limiting distribution

estimated density of the translated limiting distribution



Iterate 10 times

$$Z \stackrel{\mathcal{D}}{=} U^{\frac{\alpha-1}{2}} \left(V^{\alpha-1} Z^{(1)} + (1 - V)^{\alpha-1} Z^{(2)} \right).$$

starting with δ_γ . From this distribution 15 000 samples were produced. Then a S-Plus smoothing-routine was applied to the histogram of the data.

Multidimensional Quadtree

Dimension: d

$U = (U_1, \dots, U_d)$ the root

$q = (Y_1, \dots, Y_s, *, \dots, *)$ the query

$1 \leq s \leq d - 1$ number of specified components

$\langle U \rangle = (\langle U \rangle_0, \dots, \langle U \rangle_{2^d - 1})$
volumes of generated quadrants

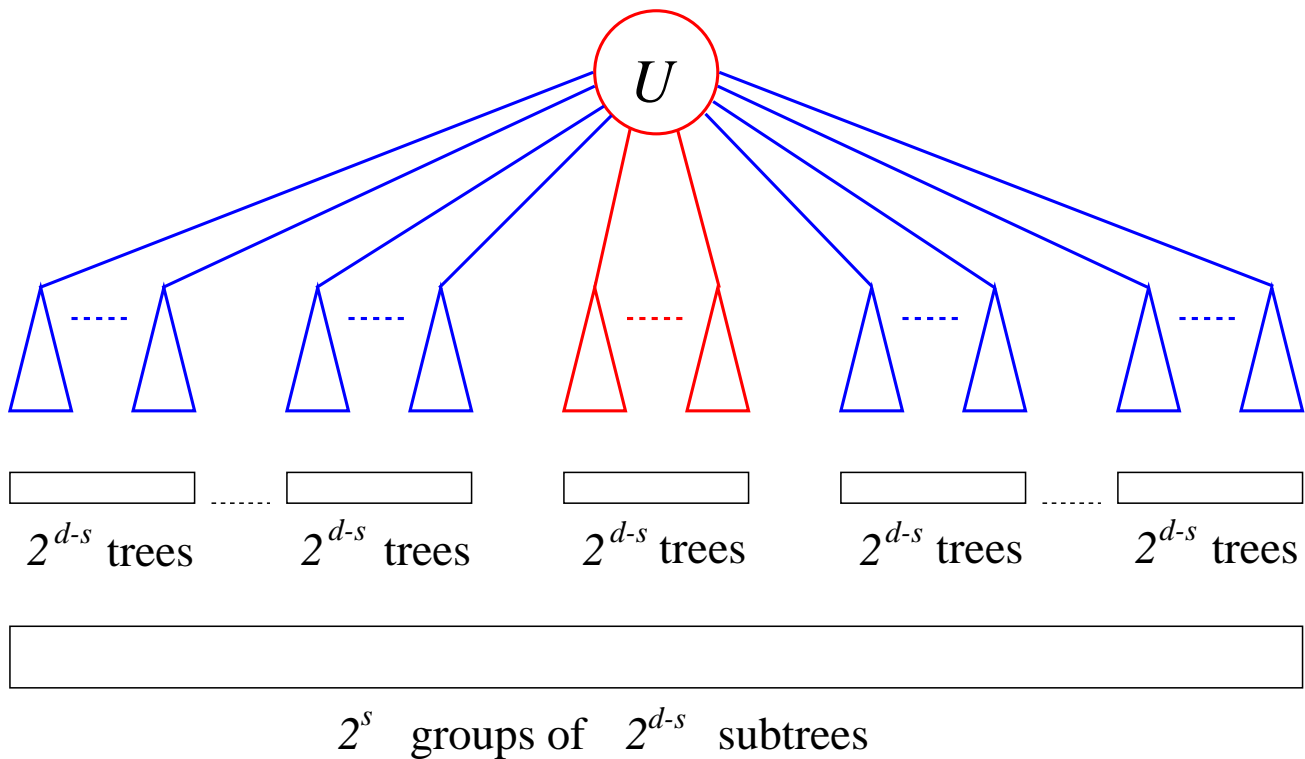
$I^{(n)}$ vector of the subtree-sizes

PMQ in the multidim. Quadtree

$U = (U_1, \dots, U_d)$ the root

$q = (Y_1, \dots, Y_s, *, \dots, *)$ the query

Decision in s components



The cost C_n

$$\mathbf{1}_{j_1, \dots, j_s}(Y, U) := \prod_{\substack{1 \leq i \leq s \\ j_i = 0}} \mathbf{1}_{\{Y_i < U_i\}} \prod_{\substack{1 \leq i \leq s \\ j_i = 1}} \mathbf{1}_{\{Y_i \geq U_i\}}$$

$(j_1, \dots, j_s \in \{0, 1\})$

Distributional equation for C_n :

$$C_n \stackrel{\mathcal{D}}{=} \sum_{j_1, \dots, j_s = 0, 1} \mathbf{1}_{j_1, \dots, j_s}(Y, U) \sum_{j_{s+1}, \dots, j_d = 0, 1} C_{I_j^{(n)}}^{(j)} + 1$$

The mean: (Flajolet, Gonnet, Puech, Robson '93)

$$\mathbb{E}C_n \sim \gamma_{s,d} n^{\alpha-1}$$

$\gamma_{s,d} > 0$ (unknown)

$\alpha \in (1, 2)$, satisfying the **indicial equation**

$$\alpha^{d-s} (\alpha + 1)^s = 2^d$$

Scaling

$\text{Var}(C_n)$: unknown

Assume: $\text{Var}(C_n) \sim \beta_{s,d} n^{2\alpha-2}$ (later proved)

Normalisation:

$$X_n := \frac{C_n - \mathbb{E}C_n}{n^{\alpha-1}}$$

→ modified recursion + $I^{(n)}/n \rightarrow \langle U \rangle$

→ guess for limiting equation

→ limiting operator

$$T : M_1(\mathbb{R}, \mathcal{B}) \rightarrow M_1(\mathbb{R}, \mathcal{B})$$

$$\begin{aligned} T(\mu) &\stackrel{\mathcal{D}}{=} \sum_{j_1, \dots, j_s=0,1} \mathbf{1}_{j_1, \dots, j_s}(Y, U) \\ &\quad \times \sum_{j_{s+1}, \dots, j_d=0,1} \langle U \rangle_{j_1, \dots, j_d}^{\alpha-1} (Z^{(j)} - \gamma_{s,d}) \\ &\quad \quad \quad + \gamma_{s,d} \end{aligned}$$

Contraction property

Lemma:

$T : M_{0,2} \rightarrow M_{0,2}$ is a contraction w.r.t. l_2 :

$$l_2(T(\mu), T(\nu)) \leq \xi l_2(\mu, \nu) \quad \forall \mu, \nu \in M_{0,2},$$

$$\xi = \frac{1}{\sqrt{\alpha^s (\alpha - 1/2)^{d-s}}} < 1.$$

Remark to the proof:

First check that T is well-defined:

$$\mathbb{E}T(\mu) = \dots = \gamma_{s,d} \left(\frac{2^d}{(\alpha + 1)^s \alpha^{d-s}} - 1 \right) \stackrel{!}{=} 0$$

→ indicial equation

Limit law

Theorem: (Limit Law for PMQ in Quadtrees)

X_n the normalized cost for PMQ

X the fixed point of the limiting operator, then

$$l_2(X_n, X) \rightarrow 0.$$

The translated limit $\tilde{X} := X + \gamma_{s,d}$ is the unique solution in $M_{\gamma_{s,d},2}$ of the equation

$$\tilde{X} \stackrel{\mathcal{D}}{=} \prod_{i=1}^s U_i^{\frac{\alpha-1}{2}} \sum_{j_{s+1}, \dots, j_d=0,1} \bar{U}_{j_{s+1}, \dots, j_d}^{\alpha-1} \tilde{X}^{(j_{s+1}, \dots, j_d)}$$

with

$$\bar{U}_{j_{s+1}, \dots, j_d} := \prod_{\substack{s+1 \leq i \leq d \\ j_i=0}} U_i \prod_{\substack{s+1 \leq i \leq d \\ j_i=1}} (1 - U_i),$$

the volumes of generated quadrants in the subspace of the last $d - s$ components.

Asymptotic for the variance

Corollary: (Variance of PMQ in Quadrees)

The variance of the cost C_n for a partial match query in a d -dimensional quadtree with $1 \leq s \leq d - 1$ components specified satisfies:

$$\text{Var}(C_n) \sim \beta_{s,d} n^{2\alpha-2},$$

$$\beta_{s,d} = \left[\frac{((2\alpha - 1)B(\alpha, \alpha) + 1)^{d-s} - 1}{\alpha^s(\alpha - 1/2)^{d-s} - 1} - 1 \right] \gamma_{s,d}.$$

Proof now is easy:

$$\begin{aligned} \text{Var}(C_n) &= \text{Var}(n^{\alpha-1} X_n) = (\text{Var}(X) + o(1)) n^{2\alpha-2} \\ &\sim \beta_{s,d} n^{2\alpha-2}. \end{aligned}$$

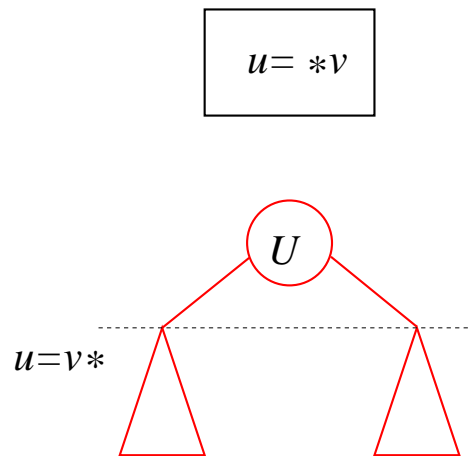
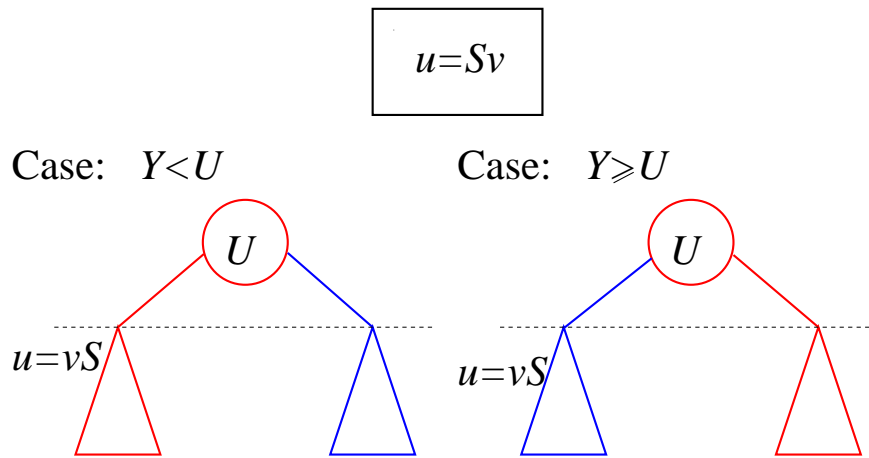
Calculate $\text{Var}(X) = \beta_{s,d}$: Take square and expectation in

$$\tilde{X} \stackrel{\mathcal{D}}{=} \prod_{i=1}^s U_i^{\frac{\alpha-1}{2}} \sum_{j_{s+1}, \dots, j_d=0,1} \bar{U}_{j_{s+1}, \dots, j_d}^{\alpha-1} \tilde{X}^{(j_{s+1}, \dots, j_d)}.$$

PMQ in K -d-trees

Specification pattern specified in the components

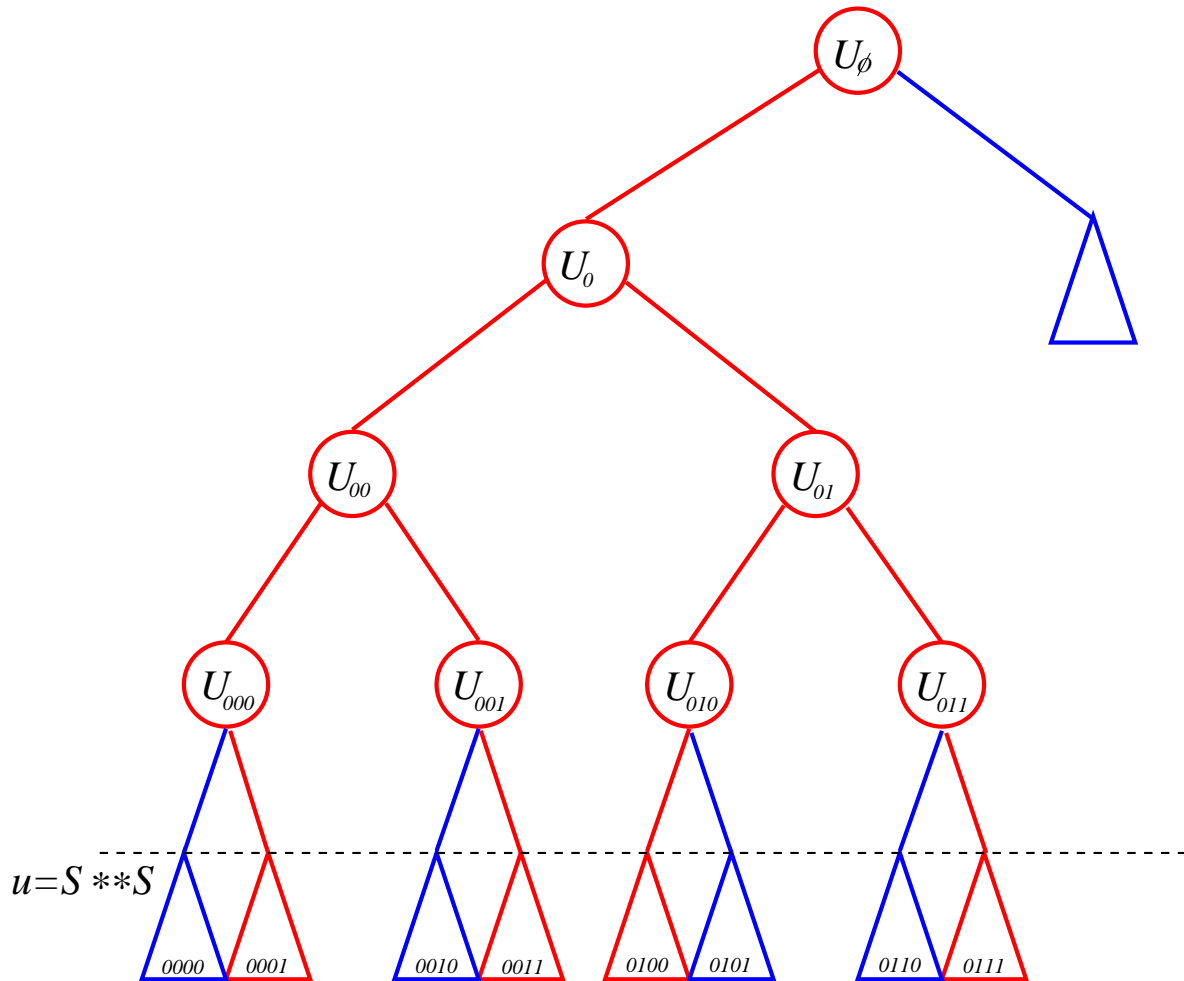
$$1 \leq r_1 < r_2 < \dots < r_s \leq K, \quad 1 \leq s \leq K - 1$$



$$C_n^{Sv} \stackrel{D}{=} \mathbf{1}_{\{Y < U\}} C_{I_1^{(n)}}^{vS} + \mathbf{1}_{\{Y \geq U\}} \bar{C}_{I_2^{(n)}}^{vS} + 1$$

PMQ in K -d trees

Example: $K = 4$, $u = (S ** S)$



$$D_n := \{0, 1\}^n, D_0 := \{\emptyset\}$$

$$D(K) := \bigcup_{n=0}^{K-1} D_n$$

$$\sigma|_j := (\sigma_1, \dots, \sigma_j) \in D_j \text{ for } \sigma \in D_n, 0 \leq j \leq n$$

The recursion for the cost C_n

Distributional recursive equation for C_n :

$$C_n \stackrel{\mathcal{D}}{=} \sum_{\sigma \in D_K} \mathbf{1}_{\sigma}(\mathcal{U}, Y) C_{I_{\sigma}^{(n)}}^{(\sigma)} + N_n,$$

$Y, \mathcal{U}, (C_i^{(\sigma)})_{i \in \mathbb{N}}$ ($\sigma \in D_K$) independent

$$(C_i^{(\sigma)})_{i \in \mathbb{N}} \sim (C_i)_{i \in \mathbb{N}}, \quad \sigma \in D_K$$

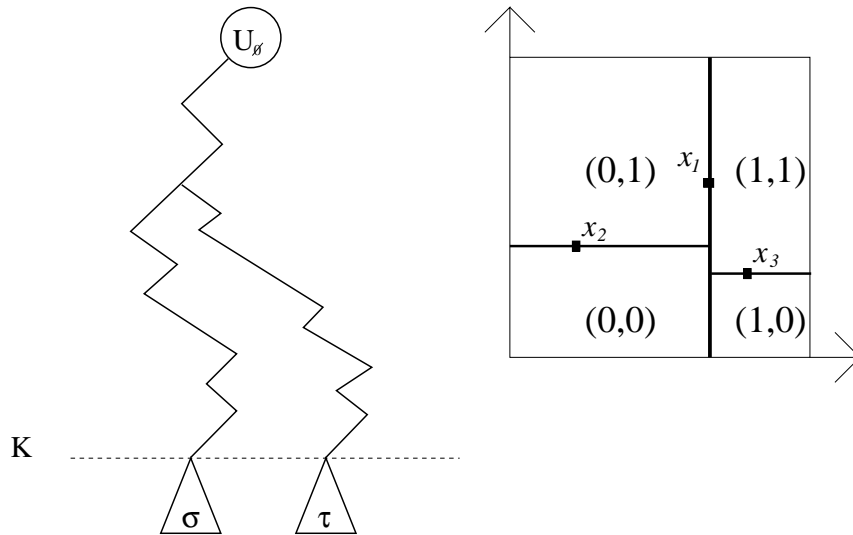
N_n number of nodes traversed on levels $0, \dots, K - 1$

$$0 \leq N_n \leq 2^K$$

$I^{(n)} = (I_{\sigma}^{(n)})_{\sigma \in D_K}$ vector of the cardinalities of the subtrees on level K

$$\mathbf{1}_{\sigma}(\mathcal{U}, Y) := \prod_{\substack{1 \leq j \leq s \\ \sigma_{r_j} = 0}} \mathbf{1}_{\{Y_j < U_{\sigma|_{r_j-1}}\}} \prod_{\substack{1 \leq j \leq s \\ \sigma_{r_j} = 1}} \mathbf{1}_{\{Y_j \geq U_{\sigma|_{r_j-1}}\}}$$

The subtree-sizes $I^{(n)}$



Volumes generated by the discriminators:

$$\langle \mathcal{U} \rangle_{\sigma} := \prod_{\substack{1 \leq j \leq K \\ \sigma_j = 0}} U_{\sigma|j-1} \prod_{\substack{1 \leq j \leq K \\ \sigma_j = 1}} (1 - U_{\sigma|j-1})$$

Given the levels $0, \dots, K - 1$ are full and given the discriminators \mathcal{U} , $I^{(n)}$ is multinomial distributed.

$$\frac{I^{(n)}}{n} \xrightarrow{\mathbb{P}} \langle \mathcal{U} \rangle = (\langle \mathcal{U} \rangle_{\sigma})_{\sigma \in D_K}$$

Scaling

The mean: (Flajolet, Puech '86)

$$\mathbb{E}C_n \sim \gamma_u n^{\alpha-1}$$

$\gamma_u > 0$ depends on the specification pattern u ,
 $\alpha = \alpha(s, K) \in (1, 2)$ satisfying the indicial equation

$$\alpha^{K-s}(\alpha + 1)^s = 2^K$$

Assume: $\text{Var}(C_n) \sim \beta_u n^{2\alpha-2}$

Normalized version:

$$X_n := \frac{C_n - \mathbb{E}C_n}{n^{\alpha-1}}$$

Distributional recursion for the normalized cost:

$$X_n \stackrel{\mathcal{D}}{=} \sum_{\sigma \in D_k} \mathbf{1}_\sigma(\mathcal{U}, Y) \left(\frac{I_\sigma^{(n)}}{n} \right)^{\alpha-1} \left(X_{I_\sigma^{(n)}}^{(\sigma)} + \gamma_u \right) - \gamma_u + o(1)$$

The limiting operator

Limiting equation:

$$X \stackrel{\mathcal{D}}{=} \sum_{\sigma \in D_k} \mathbf{1}_{\sigma}(\mathcal{U}, Y) \langle \mathcal{U} \rangle_{\sigma}^{\alpha-1} (X^{(\sigma)} + \gamma u) - \gamma u$$

Limiting operator:

$$T : M_1(\mathbb{R}, \mathcal{B}) \rightarrow M_1(\mathbb{R}, \mathcal{B})$$

$$T(\mu) \stackrel{\mathcal{D}}{=} \sum_{\sigma \in D_K} \mathbf{1}_{\sigma}(\mathcal{U}, Y) \langle \mathcal{U} \rangle_{\sigma}^{\alpha-1} (Z^{(\sigma)} + \gamma u) - \gamma u$$

$\mathcal{U}, Y, Z^{(\sigma)}, \sigma \in D_K$ independent, $Z^{(\sigma)} \sim \mu$.

Limit law

Theorem: (Limit law for PMQ in K -d trees)

X_n the normalized cost for PMQ

X the fixed point of the limiting operator, then

$$l_2(X_n, X) \rightarrow 0.$$

Corollary: (Variance of PMQ in K -d trees)

$$\text{Var}(C_n) \sim \beta_u n^{2\alpha-2}$$

with

$$\beta_u = \left[c_\alpha \sum_{l \in \mathbb{U}} \left(\frac{2(2\alpha - 1)}{\alpha^2} \right)^{K-l} \eta_\alpha^{s-s_l} - 1 \right] \gamma_u^2,$$

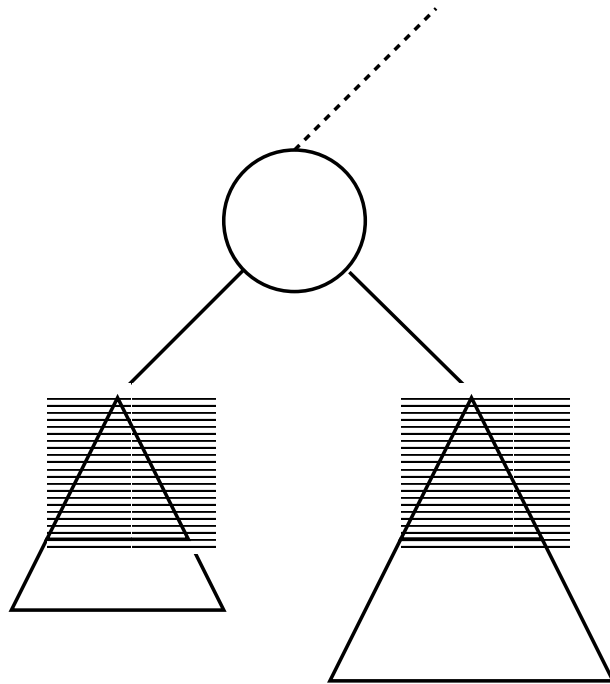
$$c_\alpha = \frac{(2\alpha - 1)B(\alpha, \alpha)}{\alpha^s (\alpha - 1/2)^{K-s} - 1},$$

$$\eta_\alpha = \frac{\alpha(8\alpha^2 - 2\alpha - 2 - \alpha(\alpha + 1)B(\alpha, \alpha))}{2(\alpha + 1)(2\alpha - 1)(2\alpha + 1)}.$$

K -d-t trees

Locally balanced version of K -d trees introduced in Cunto, Lau and Flajolet ('89).

Subtrees of size greater than $2t$ have at least t nodes on each side. ($t = 0$: k -d trees)



$$\mathbb{E}C_n^{(t)} \sim \gamma_{t,u} n^{1-s/K+\theta(s/K,t)}$$

$$\theta(s/K, t) \longrightarrow 0 \quad \text{for} \quad t \rightarrow \infty$$

Random relaxed K-d tree

Each internal node contains a key plus an associated discriminant $j \in \{1, \dots, K\}$.

The discriminant determines, which component is used for the comparison.

randomized version: The discriminants are independent r.v., uniformly distributed on $\{1, \dots, K\}$.

Duch, Estivill-Castro, Martínez ('98)

Martínez, Panholzer, Prodingler ('98)

$1 \leq s \leq K - 1$ components specified

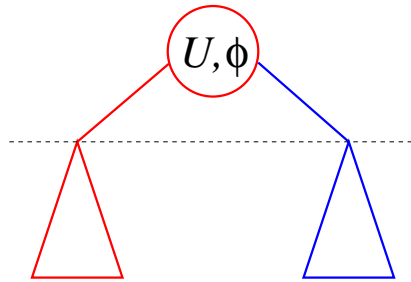
$$\rho := \frac{s}{K}$$

PMQ in random-r-K-d-trees

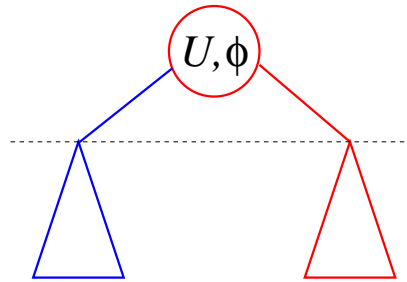
ϕ the random discriminant associated to the root

ϕ -th component specified

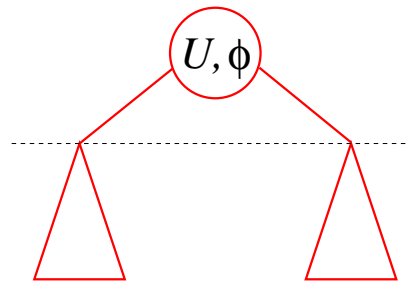
Case: $Y < U$



Case: $Y \geq U$



ϕ -th component unspecified



$$C_n \stackrel{\mathcal{D}}{=} B_\rho \left(\mathbf{1}_{\{Y < U\}} C_{I_1^{(n)}}^{(1)} + \mathbf{1}_{\{Y \geq U\}} C_{I_2^{(n)}}^{(2)} \right) \\ (1 - B_\rho) \left(C_{I_1^{(n)}}^{(1)} + C_{I_2^{(n)}}^{(2)} \right) + 1$$

The moments

$$\mathbb{E}C_n = \gamma n^{\alpha-1} + \mathcal{O}(1)$$

$$\alpha = \frac{1 + \sqrt{9 - 8\rho}}{2}, \quad \gamma = \frac{\Gamma(2\alpha - 1)}{(1 - \rho)\alpha\Gamma^3(\alpha)}$$

$$\text{Var}(C_n) \sim \left(\frac{8\Gamma(2\alpha)}{\alpha^2(\alpha - 1)^2(2\alpha - 1)(3\alpha - 2)\Gamma^4(\alpha)} - \frac{4\Gamma^2(2\alpha)}{\alpha^4(\alpha - 1)^2(2\alpha - 1)^2\Gamma^6(\alpha)} \right) n^{2\alpha-2}$$

$$X_n = \frac{C_n - \mathbb{E}C_n}{n^{\alpha-1}}$$

Limiting operator

$$T : M_1(\mathbb{R}, \mathcal{B}) \rightarrow M_1(\mathbb{R}, \mathcal{B})$$

$$\begin{aligned} T(\mu) \stackrel{\mathcal{D}}{=} & B_\rho \left(\mathbf{1}_{\{Y < U\}} U^{\alpha-1} (X^{(1)} + \gamma) \right. \\ & \left. + \mathbf{1}_{\{Y \geq U\}} (1 - U)^{\alpha-1} (X^{(2)} + \gamma) \right) \\ & + (1 - B_\rho) \left(U^{\alpha-1} (X^{(1)} + \gamma) \right. \\ & \left. + (1 - U)^{\alpha-1} (X^{(2)} + \gamma) \right) - \gamma \end{aligned}$$

Lemma:

$T : M_{0,2} \rightarrow M_{0,2}$ is a contraction w.r.t. l_2

$$l_2(T(\mu), T(\nu)) \leq \xi l_2(\mu, \nu) \quad \forall \mu, \nu \in M_{0,2},$$

$$\xi = \left[2 \frac{1 - \rho + \sqrt{9 - 8\rho}}{9 - 8\rho + \sqrt{9 - 8\rho}} \right]^{1/2} < 1$$

Limit law

Theorem: (Limit law for PMQ in random-r-K-d-trees)

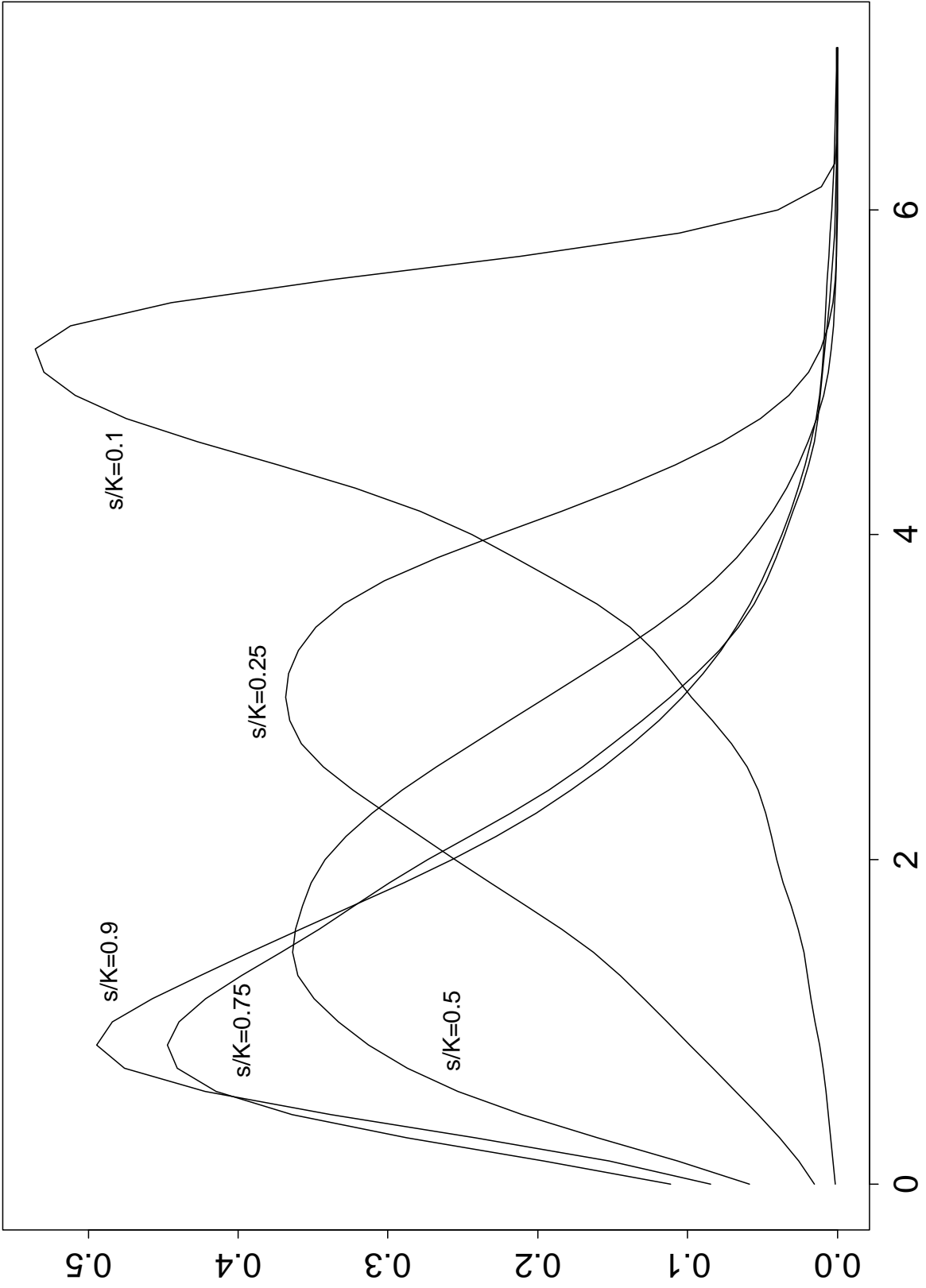
X_n the normalized cost for PMQ

X the fixed point of the limiting operator, then

$$l_2(X_n, X) \rightarrow 0.$$

The translated limit $\tilde{X} := X + \gamma$ is the unique solution in $M_{\gamma,2}$ of the equation

$$Z \stackrel{\mathcal{D}}{=} B_\rho U^{\frac{\alpha-1}{2}} Z^{(1)} + (1-B_\rho) \left(U^{\alpha-1} Z^{(1)} + (1-U)^{\alpha-1} Z^{(2)} \right)$$



Laplace Transform

Theorem: (Laplace transforms)

$$\mathbb{E} \exp(\lambda X) < \infty \quad \text{for all } \lambda \in (-\lambda_0, \lambda_0).$$

Assume

$$0 < \frac{s}{K} \leq \frac{\ln(4/3)}{\ln(5/3)} = 0.563\dots \quad \text{for the } K\text{-d tree,}$$

$$0 < \frac{s}{K} \leq \frac{\ln \left(\frac{4+2t}{3+2t} \right)^{t+1}}{\ln \left(\frac{5+2t}{3+2t} \right)^{t+1}} \quad \text{for the } K\text{-d-t tree,}$$

$$0 < \frac{s}{K} \leq 0.625 \quad \text{for the random relaxed } K\text{-d tree,}$$

$$0 < \frac{s}{d} \leq \frac{\ln(4/3)}{\ln(5/3)} = 0.563\dots \quad \text{for the quadtree.}$$

Then

$$\mathbb{E} \exp(\lambda X_n) \rightarrow \mathbb{E} \exp(\lambda X) < \infty \quad \text{for all } \lambda \in \mathbb{R}.$$

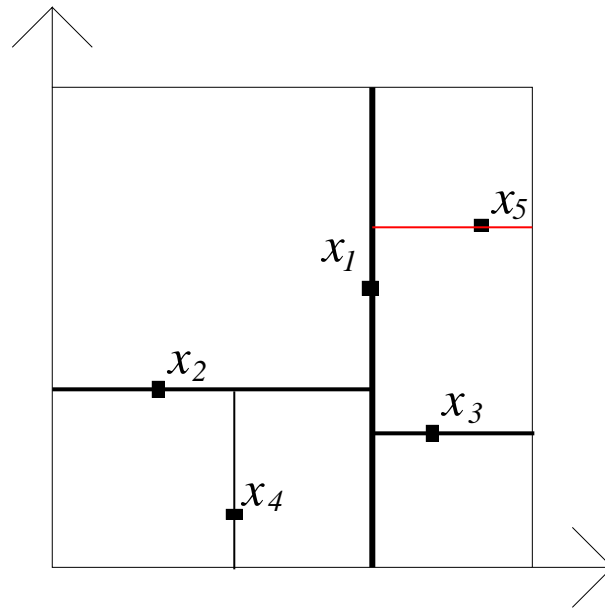
$$\mathbb{P}(C_n \geq a_n) \leq c_\lambda \exp \left(-\lambda \frac{a_n}{n^{\alpha-1}} \right) \quad \text{for all } \lambda > 0.$$

The squarish K -d tree

Devroye, Jabbour, Zamora-Cura '99

- relaxed K -d tree

- choose the discriminant at a node, so that the corresponding quadrant is cut perpendicular to the longest edge



$$\mathbb{E}C_n = \Theta(n^{1-s/K})$$

Problem: Prove

$$\mathbb{E}C_n \sim \gamma n^{1-s/K}.$$

Internal path length in random trees

Internal path length (IPL): Sum of the depths of the nodes in a tree.

- Cost (=key comparisons) for building up the tree.
- BST: Internal path length is distributed as the running time of the Quicksort algorithm.

The Quadtree

Y_n : IPL of the random d -dimensional quadtree,
 $Y_1 = 1, Y_2 = 3, \dots$

Distributional recursive equation for Y_n :

$$Y_n \stackrel{\mathcal{D}}{=} \sum_{k=0}^{2^d-1} Y_{I_k^{(n)}}^{(k)} + n$$

with

$I^{(n)}, (Y_i^{(0)}), \dots, (Y_i^{(2^d-1)})$ independent

$$Y_i^{(k)} \sim Y_i, \quad 0 \leq k \leq 2^d - 1$$

$I^{(n)}$ vector of the subtree sizes.

Moments of Y_n

$$\mathbb{E}Y_n = \frac{2}{d}n \ln n + \mu_d n + o(n)$$

(Flajolet, Labelle, Laforest, Salvy '95)

Var Y_n : unknown

Assume: Var $Y_n \sim v_d n^2$ (later proved)

Normalisation:

$$X_n := \frac{Y_n - \mathbb{E}Y_n}{n}.$$

The modified recursion

$$X_n \stackrel{D}{=} \sum_{k=0}^{2^d-1} \frac{I_k^{(n)}}{n} X_{I_k^{(n)}}^{(k)} + C_n(I^{(n)})$$

with

$$C_n(i) = 1 + \frac{1}{n} \left(\sum_{k=0}^{2^d-1} \mathbb{E}Y_{i_k} - \mathbb{E}Y_n \right)$$

for $i = (i_0, \dots, i_{2^d-1})$ with $\sum i_k = n - 1$.

Using the expectation formula it follows

$$C_n(i) = 1 + \frac{2}{d} \sum_{k=0}^{2^d-1} \frac{i_k}{n} \ln \left(\frac{i_k}{n} \right) + o(1)$$

Recall $\frac{I^{(n)}}{n} \xrightarrow{\mathbb{P}} \langle U \rangle$.

($\langle U \rangle$ volumes generated by the root.)

The limiting equation

Entropy functional

$$C : T_{2^d-1} \rightarrow \mathbb{R}, \quad C(x) := 1 + \frac{2}{d} \sum_{k=0}^{2^d-1} x_k \ln x_k$$

defined on the simplex

$$T_{2^d-1} := \left\{ x \in [0, 1]^{2^d} : \sum_{k=0}^{2^d-1} x_k = 1 \right\}.$$

Limiting operator:

$$T : M^1(\mathbb{R}^1, \mathcal{B}^1) \longrightarrow M^1(\mathbb{R}^1, \mathcal{B}^1)$$

$$T(\mu) \stackrel{\mathcal{D}}{=} \sum_{k=0}^{2^d-1} \langle U \rangle_k Z^{(k)} + C(\langle U \rangle)$$

$U, Z^{(0)}, \dots, Z^{(2^d-1)}$ independent, U unif. on $[0, 1]^d$ and $Z^{(k)} \sim \mu$.

The limit law

Theorem: (Limit law for IPL in quadtree) *Let Y_n be the IPL in a d -dimensional quadtree, X_n be the normalized version and X the fixed point of the limiting operator, then:*

$$l_2(X_n, X) \rightarrow 0,$$

$$\text{Var}(Y_n) \sim v_d n^2, \quad v_d = \frac{21 - 2\pi^2}{9d(1 - (2/3)^d)},$$

$$\mathbb{E} \exp(\lambda X_n) \rightarrow \mathbb{E} \exp(\lambda X), \quad \lambda \in \mathbb{R},$$

$$\mathbb{P}(|Y_n - \mathbb{E}Y_n| \geq \epsilon \mathbb{E}Y_n) = O(n^{-k}) \quad \forall k \in \mathbb{N}.$$

General split trees

Other trees are of a similar type:

BST, m -ary search trees, median-of- $(2k + 1)$ search trees, . . .

→ random split tree model (Devroye '98)

→ analysis uniformly valid for the depth and height

IPL:

$$Y_n \stackrel{\mathcal{D}}{=} \sum_{k=1}^b Y_{I_k^{(n)}}^{(k)} + n$$

$I^{(n)}$ has conditionally given a **splitting vector** $\mathcal{V} = (V_1, \dots, V_b)$ a multinomial structure:

$$\mathbb{P}^{I^{(n)} | \mathcal{V}=v} = M(n - s, v),$$

$$\mathbb{P}^{I_1^{(n)}} = B(n - s, V_1)$$

$V \sim V_1$ is called the **splitter** .

Expectation of the IPL

$$\mathbb{E}Y_n = b \sum_{i=0}^{n-s} \mathbb{P}(I_1^{(n)} = i) \mathbb{E}Y_i + n$$

with $\mathbb{P}^{I_1^{(n)}} = B(n-s, V_1)$, i.e.

$$\mathbb{P}(I_1^{(n)} = i) = \int_0^1 \binom{n-s}{i} p^i (1-p)^{n-s-i} d\mathbb{P}^V(p)$$

BST: V uniformly on $[0, 1]$ distributed

m -ary tree: V beta($1, m-1$) dist.

. (min. of m unif. on $[0, 1]$ dist. r.v.)

quadtree: V product of d unif. dist. r.v.

median of $(2k+1)$ tree: V beta($k+1, k+1$) dist.

. (median of $2k+1$ unif. on $[0, 1]$ dist. r.v.)

Problem

$$\mathbb{E}Y_n = b \sum_{i=0}^{n-s} \mathbb{P}(I_1^{(n)} = i) \mathbb{E}Y_i + n,$$

$$\mathbb{P}(I_1^{(n)} = i) = \int_0^1 \binom{n-s}{i} p^i (1-p)^{n-s-i} d\mathbb{P}^V(p)$$

Find conditions on V which imply

$$\mathbb{E}Y_n = \mu^{-1} n \ln n + dn + o(n)$$

with

$$\mu = b\mathbb{E}[V \ln(1/V)]$$

The formula is valid for the examples mentioned above!

The Find-algorithm

Search for an order statistic in a file of n keys using the partitioning procedure of quicksort.

- pick the pivot element at random
- partition the set of elements into the elements smaller resp. greater than the pivot
- continue recursively in the set containing the desired statistic

Model: Assume

- the keys permuted uniformly at random
- the order statistic M_n uniformly over $\{1, \dots, n\}$
- the pivot Z_n chosen as the median of 3 independent, unif. over the keys distributed elements

Recursion for the cost C_n

C_n the number of key comparisons of Find applied to a set of n keys (without the comparisons for finding the median of three)

$$C_n \stackrel{D}{=} \mathbf{1}_{\{Z_n > M_n\}} C_{Z_n-1}^* + \mathbf{1}_{\{Z_n < M_n\}} C_{n-Z_n}^{**} + n - 1$$

$Z_n, M_n, (C_i^*), (C_i^{**})$ independent

$$C_i^*, C_i^{**} \sim C_i$$

M_n unif. dist. over $\{1, \dots, n\}$

Z_n dist. as the median of three uniformly on $\{1, \dots, n\}$
dist. r.v.

Scaling of C_n

$$\mathbb{E}C_n = \frac{5}{2}n + O(\ln n)$$

(Kirschenhofer, Martínez, Prodinger '97)

Assume: $\text{Var}C_n \sim m_2 n^2$

$$Y_n := \frac{C_n - \mathbb{E}C_n}{n}$$

Modified recursion

$$\begin{aligned} Y_n &\stackrel{\mathcal{D}}{=} \mathbf{1}_{\{Z_n > M_n\}} \frac{Z_n - 1}{n} \left(Y_{Z_n - 1}^* + \frac{5}{2} \right) \\ &+ \mathbf{1}_{\{Z_n < M_n\}} \frac{n - Z_n}{n} \left(Y_{n - Z_n}^{**} + \frac{5}{2} \right) \\ &- \frac{5}{2} + \frac{n - 1}{n} + o(1). \end{aligned}$$

Limiting equation/operator

$$Y \stackrel{\mathcal{D}}{=} \mathbf{1}_{\{T>U\}}T(Y^* + 5/2) \\ + \mathbf{1}_{\{T<U\}}(1 - T)(Y^{**} + 5/2) \\ - 3/2 \\ \stackrel{\mathcal{D}}{=} X(Y + 5/2) - 3/2$$

with X, Y independent, X has density $t \mapsto 12t^2(1 - t)$ for $t \in [0, 1]$.

The translation $\bar{Y} := Y + 5/2$ is the solution of

$$\bar{Y} = X\bar{Y} + 1 \quad \text{with} \quad \mathbb{E}\bar{Y} = 5/2$$

and

$$\frac{C_n}{n} \xrightarrow{\mathcal{D}} \bar{Y}.$$

Moments

Existence and convergence of the Laplace transform implies convergence of all moments, i.e.

$$\mathbb{E}[(C_n)^k] \sim m_k n^k \quad \text{for } n \rightarrow \infty$$

with $m_0 = 1, m_1 = 5/2$ and

$$m_k = \left(12 + \frac{144}{k^2 + 7k}\right) \sum_{j=0}^{k-1} \binom{k}{j} \frac{m_j}{(j+3)(j+4)}$$

for $k \geq 2$.

Fourier transform

$$\mathbb{E} \exp(it\bar{Y}) = e^{it} \phi(t)$$

where $\phi : \mathbb{R} \rightarrow \mathbb{C}$ is smooth and given by the linear homogeneous differential equation of second order

$$t^2 \phi''(t) + 8t \phi'(t) - 12(e^{it} - 1) \phi(t) = 0$$

with the initial conditions $\phi(0) = 1$ and $\phi'(0) = 3i/2$.