Philippe Flajolet and Analytic Combinatorics: Tree Rewriting Systems by Jean-Marc Steyaert INRIA-team AMIB, Laboratoire d'Informatique Ecole Polytechnique

Principle:

- Trees are recursive structures
- Algorithms on trees are naturally recursive
- Translate recursive definitions into equations over o.g.f.
- Conclude by means of analytic combinatorics

Bibliography

- [15] PF: Analyse en moyenne de la détection des arbres partiels.(1978)
- [28] PF, JMS: On the analysis of tree-matching algorithms. (1980)
- [31] PF,JMS: A complexity calculus for classses of recursive programs over tree structures. (1981)
- [43] PF, JMS: Patterns and pattern-matching in trees. (1983)
- [66] PF,H.Prodinger: Level number sequences for trees.(1987)
- [67] PF, JMS: A complexity calculus for recursive tree algorithms. (1987)
- [89] PF,P.Sipala,JMS: Analytic variations on the dcommon subexpression problem.(1990)
- [177] PF,B.Chauvin,D.Gardy,B.Gittenberger: And/Or Trees Revisited.(2004)

Basics on trees (cf. B. Sedgewick's talk)

Simple families of trees (Meir & Moon): $T = \sum_{\omega \in \Omega} \omega(T, ..., T) = \Omega(T)$, where Ω is a finite set of operators ω of various arities $\delta(\omega)$.

Let $\Phi(u)$ be the enumerative polynomial of Ω : $[u^n]\Phi(u) = \operatorname{card}\{\omega \in \Omega | \delta(\omega) = n\}$; then the o.g.f. t(z) of \mathcal{T} satisfies: $t(z) = z\Phi(t(z))$.

Theorem(M&M): Under very general conditions the number of trees of size *n* satisfies, with $\rho = 1/\Phi(\tau)$ and $\tau \Phi'(\tau) = \Phi(\tau)$:

$$t_n = [z^n]t(z) = \sqrt{\frac{\Phi(\tau)}{2\pi\Phi''(\tau)}}\rho^{-n}n^{-3/2}(1+O(1/n)).$$

Proof: Track the first singular point on the real positive axis, which is a branching point of order 2. Then apply the Darboux method.

N.B.: All this process can be automated (cf. B. Salvy's talk)

PL-trees: Programming on trees

Procedures, Functions Conditionals: *if* cond *then* ... *else* ... Conditional iterations: *for* i=range *while* cond *do* ... *od*

```
Primitive operations on trees:

root(X) \rightarrow label of the root of tree X

deg(X) \rightarrow arity of the root of tree X

for 1 \le i \le deg(X), X[i] \rightarrow i-th root subtree of X:

X = root(X)(X[1], ..., X[deg(X)])

tests on label, arity
```

```
Example:
function equal(X,Y:T)
if root(X)<>root(Y) then assign(false) else assign(true);
for i:=1 to deg(X) while assign(equal(X[i],Y[i])) do nil od fi
```

EXAMPLE PATTERN. MATCHING IN TREES · Binary trees B = 0 + R = binary symbol (Cons) 1 BB leaf (2) A 2 . Pattern : A binary tree . Root - occurrence: 'T is obtained from P. by grafting binary trees to leaves of 2 T : . Iccurrence

. Algorithus root? (2,T) = if P=x then true else if T=x then false duif root? (?_ left, T_ left) then root? (I_right, T_right) of occ?(P,T) = root?(P,T);if $P \neq z$ then occ?(P, T left)oce? (I, Triflet) fi. . Problem : I being fixed, determine the average cost of nunning root? and oce? over all trees of size n.

From the program to its complexity

Idea: Associate to instructions, blocs, constructs, o.g.f. to capture the running time complexity!

Extension of combinatorial constructs (cf. B.Sedgewick)

 $au A(X_1, X_2, ..., X_m)$: cost of running A on $X_1, X_2, ..., X_m$ of sizes $n_1, n_2, ..., n_m$ consider $au A(z_1, z_2, ..., z_m) = \sum_{|X_1|=n_1, ..., |X_m|=n_m} au A(X_1, ..., X_m) z_1^{n_1} ... z_m^{n_m}$ which represents the cumulative costs of running A over all possible inputs, indexed by their sizes

In the same spirit for programs (functions) with boolean result, consider characteristic o.g.f. $\chi A(z_1, z_2, ..., z_m) = \sum_{|X_1|=n_1, ..., |X_m|=n_m} A_{true}(X_1, ..., X_m) z_1^{n_1} ... z_m^{n_m}$

Challenge: translate the constructors and basic operations $A = C(B_1, B_2, ..., B_k) \rightarrow \tau A = \Gamma(\tau B_1, \tau B_2, ..., \tau B_k, \chi B_1, \chi B_2, ..., \chi B_k)$

$$A = B_1; B_2 \to \tau A(z) = \tau B_1(z) + \tau B_2(z)$$

 $A = B(X[i]) \rightarrow \tau A(z) = z\tau B(z)\Psi(t(z))$, with $\Psi(u) = \Phi(u)/u$ and variants in case of equality between subtrees...

$$A = if \ Q \ then \ B \ else \ C \longrightarrow \tau A(z) = \tau Q(z) + \tau B(z|Q) + \tau C(z|\neg Q)$$

$$A = for \ i \ do \ B(X[i]) \ od \rightarrow \tau A(z) = z\tau B(z)\Phi'(t(z))$$

copying a tree $- > \tau \operatorname{copy}(z) = zt'(z)$

and also : conditional iteration, characteristic functions, etc.

• Generating series for costs can be
computed recursively on the pattern structure
(e.g. counting cost 1 for a generalized tot)

$$\tau \operatorname{toot} [n](2) = B(2) = \sum 1.2^{|t|}$$

 $\operatorname{troot} [n](2) = B(2) + 2 \operatorname{troot} [n](2) \cdot B(2)$
 $+ 2 \cdot B(2) \cdot \operatorname{troot} [n](2) \cdot B(2)$
 $+ 2 \cdot B(2) \cdot \operatorname{troot} [n](2)$
 $= 3 B(2) - 2$
 $\operatorname{troot} [n](2) = 2B(2) - 1$
 $\operatorname{troot} [n](2) = B(2) - 2$
 $\operatorname{troot} [n](2) = B(2) - 2$
 $\operatorname{troot} [n](2) = (n + 2 \cdot \operatorname{root} [n](2) \cdot B(2)$
 $+ 2 \operatorname{troot} [n](2) \cdot \operatorname{troot} [n](2)$
 $= (5 - 32)B - 4$
 $\operatorname{troot} [n](2) = (5 - 22)B - 4$
 $\operatorname{troot} [n](2) = (7 - 72 - 2^2)B - 6 + 22 = 53/16$
Accerace # of comparisons

Complexity of iteration over subtrees

A(X) = B(X); for i = 1..deg(X) do A(X[i]) od (the mapboth of Lisp)

Theorem:

Under very large analytic conditions, if the coefficients of the complexity descriptor of B,

 $b_n \sim c.n^{\alpha}.\rho^{-n}$ for some constants c and $\alpha \leq 1/2$, then the average complexity of A $\overline{\tau a_n} \sim c.\theta \frac{\Gamma(\alpha-1/2)}{\Gamma(\alpha)} n^{\alpha+1/2}$

Proof: from the true nature of the singularity and the basic asymptotics of the o.g.f. t(z) for trees; translating the program, we get the equation:

 $\tau a(z) = \tau b(z) + z \tau a(z) \Phi'(t(z))$ from which $\tau a(z) = \tau b(z) \cdot zt'(z)/t(z)$

TOP-DOWN RECURSION
(2. Flajolet, JTIS)
Formal derivation (generalized)

$$S = set of symbols$$
 constants operators variables...
 $D = set of derivation rules$
with $T = \int_{T.1}^{T.2} T.v(s)$

U

$$D(T) = \frac{|\text{header}| = \gamma(4)}{|\text{header}| = \gamma(4)}$$

$$T.i_{1} T.i_{2} T.i_{k} \# \text{ copies} = \alpha(4)$$

$$D(T.i_{1}) D(T.i_{3}) D(T.i_{4})$$

$$\# \text{ derived subtress} = \beta(4)$$

Equations for the family of trues

$$\begin{aligned}
\mathcal{F} &= \sum_{\substack{\delta \in S' \\ \downarrow}} \delta \left(\underbrace{\mathcal{F}, \mathcal{F}, \dots \\ \nu(\delta)} \right) \\
&= \sum_{\substack{\delta \in S' \\ \downarrow}} 3 \cdot \left(\underbrace{\mathcal{F}(3)} \right)^{\nu(\delta)}
\end{aligned}$$

Complexity of Differentiation algorithms

Algebraic expressions with $+, -, x, \sqrt{,}/$

Thm : The average running time of standard differentiation is $O(n^{3/2})$ compared to a $O(n^2)$ worst case behaviour

If sharing of subexpressions is allowed then the average running keeps linear

The complexity calculus applies more generally to generalized tree transducers, specified along the same lines;

4 types of average behaviours are possible:

- linear
- sesquilinear
- quadratic
- exponential

Size of the derived trees a cost of derivation
trans
$$(\Delta(T.1,...,T.m)) = \gamma(\Delta)$$

 $+ \Sigma' \alpha(\Delta); |T.i|$
 $+ \Sigma \beta(\Delta); trans(T.i)$
trans $(z) = \Sigma$ trans(T) $z^{|T||}$
 $T \in \mathbb{P}$
Using "complexity calculus "* rules
trans $(z | T.root = \Delta) =$
 $\gamma(\Delta) \cdot z \cdot (f(z))$
 $+ \alpha(\Delta) \cdot z^2 \cdot f'(z) \cdot (f(z))^{V(\Delta)-1}$
 $+ \beta(\Delta) \cdot z \cdot trans(z) \cdot (f(z))^{V(\Delta)-4}$
trans $(z) = \frac{z \cdot E(f(z)) + z \cdot f'(z) \cdot A(f(z))}{1 - 2 \cdot B(f(z))}$
 $E(\omega) = \Sigma \gamma(\Delta) \omega^{V(S)}$
 $A(\omega) = \Sigma \alpha(\Delta) \omega^{V(S)-1}$

Analytic study

$$f \text{ has its dominant singularity in } 3= \rho \in \mathbb{R}^{+}$$
and locally

$$f(3) = \tau + \Im(1 - \frac{\pi}{\rho})^{1/2} + \text{ s.o.t.} \quad \tau = f(\rho)$$

$$f'(3) = \tau' + \Im'(1 - \frac{\pi}{\rho})^{1/2} + \text{ s.o.t.}$$
numerator = $\rho^{2}\alpha \Im'(1 - \frac{\pi}{\rho})^{1/2} + \text{ s.o.t.}$
denominator } depends on $B(\tau)$!? $\Phi'(\tau)$
and trans }
$$B \neq \text{ constant}$$

$$- B(\tau) < \Phi'(\tau) = \frac{\rho^{2}\alpha \Im'}{1 - \rho \Theta(\tau)} (1 - \frac{\pi}{\rho})^{-1/2} + \text{ s.o.t.}$$

$$- B(\tau) = \Phi'(\tau) = \frac{\rho^{2}\alpha \Im'}{1 - \rho \Theta(\tau)} (1 - \frac{\pi}{\rho})^{-1} + \text{ s.o.t.}$$

$$- B(\tau) > \Phi'(\tau) = \frac{\rho \otimes \beta'(\tau)}{\tau \otimes \beta'(\tau)} (1 - \frac{\pi}{\rho})^{-1} + \text{ s.o.t.}$$

$$- B(\tau) > \Phi'(\tau) = \frac{\rho \otimes \beta'(\tau)}{\tau \otimes \beta'(\tau)} (1 - \frac{\pi}{\rho})^{-1} + \text{ s.o.t.}$$

$$- B(\tau) > \Phi'(\tau) = \frac{\rho \otimes \beta'(\tau)}{\tau \otimes \beta'(\tau)} (1 - \frac{\pi}{\rho})^{-1} + \text{ s.o.t.}$$

$$- B(\tau) > \Phi'(\tau) = \frac{\rho \otimes \beta'(\tau)}{\tau \otimes \beta'(\tau)} (1 - \frac{\pi}{\rho})^{-1} + \text{ s.o.t.}$$

$$- B(\tau) > \Phi'(\tau) = \frac{\rho \otimes \beta'(\tau)}{\tau \otimes \beta'(\tau)} + \frac{\rho \otimes \beta'(\tau)}{\tau \otimes \beta'(\tau)} = \frac{\rho \otimes \beta'(\tau)}{\tau \otimes \beta'(\tau)} + \frac{\rho \otimes \beta'(\tau)}{\tau \otimes \beta'(\tau)} = \frac{\rho^2 \otimes \beta'(\tau)}{\tau \otimes \beta'(\tau)} + \frac{\rho$$

More algorithms

Tree compatibility: returns the greatest common root part of X and Y; average cost is O(1) (compare to naive string matching)

Tree matching: searching a motif in a tree, at all positions; average cost is O(n) the constant depending on the motif's shape

Tree simplification: applying t - t = 0 syntactically; average cost is linear in the input size

Further developments (by the AofA community): recursive tree simplification, boolean expressions, unification, higher order differentiation, rewriting systems, etc.

Limiting laws can also be derived in a systematic way

Tree compaction

How much space can be saved by sharing systematically the identical subtrees of a given binary tree?

Complete binary tree of height h has size $n = 2^{h+1} - 1$ and after compaction has size h + 1Right comb of height h has size n = h + 1 and is not changed by compaction

Theorem: The average size of a tree of size n after compaction is $O(n/\sqrt{\log \log n})$

Comment: such a result cannot be obtained directly by the complexity calculus! a weaker version along these lines gives o(n)

Automating the process

by P. Flajolet, B. Salvy, P. Zimmermann



Dérivation formelle

function diff(e:expression):expression;

case e of

plus(e1,e2)	:	plus(diff(e1),diff(e2));
times(e1,e2)	:	plus(times(diff(e1),copy(e2)),
		times(copy(e1),diff(e2)));
expo(e1)	:	times(diff(e1),copy(e));
zero	:	zero;
one	:	zero;
x	•	one

end;

function copy (e:expression):expression;

```
measure plus,times,expo,zero,one,x : 1;
```

#analyze"diff";;

Average cost for diff on random inputs of size n is:

1/2 1/2 3/2 (126 + 6) Pi n 1/2 ----- + O(n) 1/2 3/2 3/4 1/2 (-1 + 2 6) 6 23 0

Work in progress...