



Analytic Combinatorics— A Calculus of Discrete Structures

Philippe Flajolet

INRIA Rocquencourt, France

Carleton University, Ottawa, 26 March 2008 Text in ACM-SIAM SODA'07; http://algo.inria.fr/flajolet/

Analysis of algorithms: What is the cost of a computational task?



Babbage (1837):

number of turns of the crank



On a data ensemble, as a function of size *n*?

- in the worst case
- typically: on average; in probability in distribution.

Also vital for randomized algorithms.

SURPRISE (1960-1970s): A large body of classical maths is adequate for many average-case analyses.

- Von Neuman 1946+Knuth 1978: adders=carry riples.
- Hoare 1960: Quicksort and Quickselect
- Knuth 1968–1973⁺: The Art of Computer Programming.
- Sedgewick: median of three, halting on small subfiles, etc



"The Unreasonable Effectiveness of Mathematics" [E. Wigner]

....BUT ...:

In the 1970s and 1980s, culmination of recurrences and real analysis ($\Sigma \rightarrow \int$) techniques.

- Limitations for richer data structures and algorithms

- analyses become more and more technical.

No clear relationship

Algorithmic structures — Complexity structures.

+ Explosion in difficulty: average-case \rightsquigarrow variance \rightsquigarrow distribution

ADVANCES (1990-2008)

Synthetic approaches emerge based on generating functions.

- A. Combinatorial enumeration: Symbolic methods. Joyal's theory of species [Bergeron-Labelle-Leroux 1998]; Rota-Stanley [books]; Goulden & Jackson's formal methods; Bender-Goldman's theory of "prefabs"; Russian school.
- B. Asymptotic analysis: Complex methods.
 Bender *et al.*. F-Odlyzko, 1990+: singularity analysis; Odlyzko's survey 1995; uses of saddle points and Mellin transform.
- **C.** Distributional properties: Perturbation theory. Bender, F-Soria; H.K. Hwang's Quasipowers, 1998; Drmota-Lalley-Woods....

AofA Books: Hofri (1995), Mahmoud (1993); Szpankowski (2001). + Analytic Combinatorics, by F. & Sedgewick (2008).

Part B. Complex asymptotics Part C. Distributions Part D. Frontiers

PART A. SYMBOLIC METHODS

How to enumerate a combinatorial class C?

- $C_n = \#$ objects of size n
- \heartsuit Generating function: $C(z) := \sum z^n$.

∑∑Ƴ∫ Symbolic approach

• An object of size *n* is viewed as composed of *n* atoms (with additional structure): words, trees, graphs, permutations, etc.

- Replace each atom by symbolic weight z:
- Class: \sum objects. Object: $\gamma \rightsquigarrow z^{|\gamma|}$.

Gives the Ordinary Generating Function (OGF):

$$\mathcal{C} \qquad \rightsquigarrow \qquad \mathcal{C}(z) := \sum_{\gamma \in \mathcal{C}} z^{\gamma} \equiv \sum_{n} C_{n} z^{n}.$$

Mathematician: "To count sheep, count legs and divide by 4."

E.g.: a class of graphs enumerated by # vertices



Principle (Symbolic method)

The OGF of a class: (i) encodes the counting sequence; (ii) is nothing but a reduced form of the class itself.

Several set-theoretic constructions translate into GFs.

disjoint union $\sum_{\mathcal{A} \oplus \mathcal{B}} =$ cartesian product $\sum_{\mathcal{A} \oplus \mathcal{B}} =$

$$\sum_{\mathcal{A} \oplus \mathcal{B}} = \sum_{\mathcal{A}} + \sum_{\mathcal{B}}$$
$$\sum_{\mathcal{A} \times \mathcal{B}} = \sum_{\mathcal{A}} \cdot \sum_{\mathcal{B}}$$

There is a micro-dictionary:

 $\begin{array}{lll} \text{disjoint union} & \mathcal{C} = \mathcal{A} \cup \mathcal{B} & \Longrightarrow & \mathcal{C}(z) = \mathcal{A}(z) + \mathcal{B}(z) \\ \text{cartesian product} & \mathcal{C} = \mathcal{A} \times \mathcal{B} & \Longrightarrow & \mathcal{C}(z) = \mathcal{A}(z) \cdot \mathcal{B}(z) \end{array}$

Part B. Complex asymptotics Part C. Distributions Part D. Frontiers

Theorem (Symbolic method)

A dictionary translates constructions into generating functions:



♣ C = SEQ(A) ≡ {ε} + A + (A × A) + · · · .
Thus C = 1 + A + A² + A³ =
$$\frac{1}{1 - A}$$
.
♣ C = MSET(A) ≡ $\prod_{\alpha \in A}$ SEQ(α) \rightsquigarrow C = Exp[A],
with Exp[f] := $e^{f(z) + \frac{1}{2}f(z^2) + \cdots}$

Part B. Complex asymptotics Part C. Distributions Part D. Frontiers

More generating functions ...

• Labelled classes: via exponential GF (EGF) $\sum C_n \frac{z^n}{r!}$.

11	1 3 2 4	1 2 3 4	1 3 2 4	1 2 3 4	1 2 4 3	2 1 3 4	2 1 4 3	3 1 4 2	::
•••	1 2 4 3	1 2 4 3	1 - 3 4 - 2	1 = 2 4 = 3	1 = 2 3 = 4	2 1 4 3	2 = 1 3 = 4	3 1 2 4	

• Parameters: via *multivariate* GFs.



• Additional constructions: substitution, pointing, order constraints:

$$f \circ g, \qquad \partial f, \qquad \int f.$$

Part B. Complex asymptotics Part C. Distributions Part D. Frontiers

Linear probing hashing: From Knuth's original derivation (rec.):

We have
$$M(N,k) = \sum_{M=1}^{N} m P(h,k,n) = \frac{N(N+1)}{2} \left(1 - \frac{k-1}{N}\right) - \frac{N-k}{N^2} \sum_{X=1}^{N} \left(\frac{h(J+1)}{N} - \frac{k(X+1)}{2}\right) K(h)$$

= $\frac{1}{2} \left\{ N(N+1) \left(1 - \frac{k-1}{N}\right) - \left(1 - \frac{k}{N}\right) \sum_{X=1}^{N} \left[(2N+1) \left(1 - \frac{K}{N}\right) - N\left(1 - \frac{K}{N}\right)^2 \right] R(h,k,k).$



Get nonempty island by joining two islands by means of a gluing element.

 \rightsquigarrow wide encompassing extensions of original analyses [F-Poblete-Viola, Pittel, Knuth 1998, Janson, Chassaing-Marckert, ...].

Some constructible families

- Regular languages, FA, paths in graphs
- Unambiguous context-free languages



- Increasing trees
- Mappings

$$\begin{cases} m = Set(K) \\ K = Gycle(T) \\ T = Z * Set(T) \end{cases}$$

Some constructible families and generating fuctions

- Regular languages, FA, paths in graphs: \rightsquigarrow
- Unambiguous context-free languages \rightsquigarrow
- Terms trees → [+Pólya operators]



rational fns

algebraic functions .

implicit functions

• Increasing trees $\rightsquigarrow Y = \int \Phi(Y)$

differential equation

• Mappings \rightsquigarrow

$$\begin{cases} m = Set l \\ K = Gycle \\ T = Z * S$$

 $\mathsf{exp} \, \circ \, \mathsf{log} \, \circ \, \mathsf{implicit}$

$$M = \exp(K)$$

$$K = \log(1 - T)^{-1}$$

$$T = z \exp(T)$$

PART B. COMPLEX ASYMPTOTICS

- The continuous [=analysis] helps understand the discrete.
- The complex domain has powerful properties.

"The shortest path between two truths on the real line goes through the complex plane." — Jacques Hadamard

Erdös' proofs from the Book [cf Aigner-Ziegler]

Why are there infinitely many primes?

- Combinatorial proof ©Euclid: n! + 1 is divisible by a prime > n.
- Analytic proof ©Euler: consider a (Dirichlet) generating function

$$\begin{aligned} \zeta(s) &= \sum_{\substack{n\geq 1 \\ p \text{ Prime}}} \frac{1}{n^s} \\ &= \prod_{\substack{p \text{ Prime}}} \frac{1}{1-1/p^s}. \end{aligned}$$

We have $\zeta(1^+) = +\infty$ while the finiteness of primes would imply $\zeta(1^+) < \infty$, a contradiction.

~ Riemann, Hadamard, de la Vallée-Poussin: Prime Number Theorem.

Complex asymptotics and GFs

- formal z yields formal generating function as "power series";
- real z gives us a real function with convergence interval;



• complex z gives us a function of a complex variable.

Surface

(here: modulus of OGF of balanced trees)



in \mathbb{R}^{4} with $\langle \Re, \Im \rangle$.

Analytic function := *smooth* transformation of the complex plane.



 \implies Analytic functions satisfy rich closure properties.



Definition

f(z) has singularity at boundary point ζ if it cannot be made analytic around ζ .

E.g.: f discontinuous, infinite, oscillating, derivative blows up, etc.



 \heartsuit Analytic properties of GF provide coefficients' asymptotics.

Principle (Singularity Analysis)

Singularities determine asymptotics of coefficients. A singularity at ζ of f(z) implies a contribution to f_n like $\zeta^{-n}\vartheta(n)$, where $\vartheta(n)$ is subexponential.

Theorem: $R_{conv} = \rho_{sing}$



LOCATION of **SINGULARITY**: by rescaling, $f(z/\zeta)$ is singular at 1. A factor of ζ^{-n} corresponds to a singularity at ζ .

NATURE of SINGULARITY: examine simple functions singular at 1:

Function	\longrightarrow	Coefficient	
$\frac{1}{(1-z)^2}$		n+1	\sim n
$\frac{1}{1-z}\log\frac{1}{1-z}$		$H_n\equiv 1+\tfrac{1}{2}+\cdots$	$\sim \log n$
$\frac{1}{1-z}$		1	~ 1
$\frac{1}{\sqrt{1-z}}$		$4^{-n}\binom{2n}{n}$	$\sim rac{1}{\sqrt{\pi n}}.$



Let *L* be a slowly varying function, like $\log^{\beta} \log \log^{\delta}$.

Theorem (Singularity analysis)

Under a Camembert condition, the following implication is valid

$$f(z) \approx \frac{1}{(1-z)^{\alpha}} L\left(\frac{1}{1-z}\right) \longrightarrow [z^n] f(z) \approx n^{\alpha-1} L(n).$$

Works for equality (=) with full asymptotic expansions; for $\mathcal{O}(.)$, o(.), hence \sim .

[F., Odlyzko 1990]; closures ∂, \int, \odot [Fill, F., Kapur 2005]; [F., Sedgewick 2008]

Proof of Singularity Analysis Theorems:

Cauchy's coefficient formula: $[z^n]f(z) = \frac{1}{2i\pi} \oint f(z) \frac{dz}{z^{n+1}}$.



 $z \to 1 + rac{t}{n} \quad z^{-n} \to e^{-t}; \qquad dz \to rac{dt}{n}; \qquad (1-z)^{-\alpha} \to (-t/n)^{-\alpha}.$

Singularity analysis works *automatically* for wide classes of generating functions.



- $\begin{array}{ll} & & \text{Rational [Perron-Frobenius]} & \rightarrow \zeta^{-n} n^k \\ & & \text{Implicit} & \rightarrow \zeta^{-n} n^{-3/2} \\ & & \text{Algebraic [Newton-Puiseux]} & \rightarrow \zeta^{-n} n^{p/q} \end{array}$
- Holonomic [linear ODEs] $\rightarrow \zeta^{-n} n^{\alpha} (\log n)^k$

Universality in trees and maps



<u>TREES:</u> $Y = z\Phi(Y)$ universality of $\sqrt{-}$ -singularity. Counting is universally $C \cdot A^n n^{-3/2}$. Height and width are $\approx \sqrt{n}$. Path length is $\approx n\sqrt{n}$, &c.

[Tutte⁺]: universality of $C \cdot A^n n^{-5/2}$ for Rooted maps.

[Bender-Gao-Wormald 2002] \rightsquigarrow Gimenez–Noy [2005⁺]: Planar graphs $n! C \cdot A^n n^{-7/2}$. Fusy: random generation is $O(n^2)$.

Trees, walks, and hashing: moment pumping ~> Airy distribution.



Louchard, Takacs, F.-Poblete-Viola.



Quadtrees and the holonomic framework



$$F(z,u) = 1 + 2^{3}u \int_{0}^{z} \frac{dx_{1}}{x_{1}(1-x_{1})} \int_{0}^{x_{1}} \frac{dx_{2}}{1-x_{2}} \int_{0}^{x_{2}} F(x_{3},u) \frac{dx_{3}}{1-x_{3}}.$$
Partial Match Query (1/2) : $PMQ_{n}^{(1/2)} \approx n^{(\sqrt{17}-3)/2}.$

Stanley-Lipshitz-Zeilberger-Gessel: Holonomic framework = linear ODEs with rational coefficients.

A theory of special functions. Equality is decidable; asymptotics are "essentially" decidable.

PART C. DISTRIBUTIONS

Runs in permutations: $\Phi(x) := \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{x} e^{-t^{2}/2} dt$

For combinatorial class \mathcal{F} with parameter χ , get bivariate GF F(z, u) which is deformation of F(z, 1) = F(z). $[z^n]F(z, u)$ is proportional to the *probability generating function* of χ on \mathcal{F}_n .

* For functions F(z) with finite singularities, usually,

 $[z^n]F(z)=\rho^{-n}n^\delta,$

 ρ given by location and n^{δ} by nature of sings.

* For F(z, u), expect to get uniform & analytic

 $[z^n]F(z,u) = \rho(u)^{-n}n^{\delta}$ or $\rho^{-n}n^{\delta(u)} \equiv \rho^{-n}e^{\delta(u)\log n}$,

via perturbation analysis.

Quasi-Powers approximation:= $PGF_n(u) \approx B(u)^{large(n)}$.

Theorem (H-K. Hwang's Quasi-Powers Theorem)

In the Quasi-Powers situation, $PGF_n(u) \approx B(u)^{large(n)}$, one has: (i) convergence to a Gaussian law

$$\mathbb{P}_n\left[\frac{\chi-\mathbb{E}[\chi]}{\sqrt{\mathbb{V}[\chi]}}\leq x\right]\to \frac{1}{\sqrt{2\pi}}\int_{-\infty}^x e^{-t^2/2}\,dt$$

(*ii*) speed of convergence; (*iii*) moment estimates; (*iv*) a large deviation principle.

Works for movable singularity & variable exponent!

[Bender, Richmond, F., Soria, Hwang] Based on: continuity theorem for characteristic functions; Berry-Essen inequalities; differentiability properties of holomorphic functions; basic large deviation theory.

- A "conceptual" proof: polynomials over finite fields.
- * Polynomials are sequences of coeffs $\implies P(z)$ has pole.

* Polynomials are multisets of irreducibles $\implies P \approx \exp(I)$, so that I(z) is logarithmic.

The density of irreducibles is $\sim q^n/n$.

* Bivariate relation $P(z, u) \approx e^{ul(z)}$ implies movable exponent implies Gaussian law.

The **number** of irreducibles is asymptotically **normal**, with log n scaling.

Cf Prime Number Theorem and Erdös–Kac. Analysis of polynomial factorization [F-Gourdon-Panario]. The **exp-log** schema [F-Soria]. Cf [Arratia-Barbour-Tavare].

Applications to Analysis of Algorithms

 — search trees: binary, multiway, locally balanced, paged; quicksort and quickslect;

— ***multidimensional search**: *k*-d-trees, quadtrees; paged, relaxed.

— *digital structures: tries, ternary search tree hybrids, multidimensional trees; protocols, leader election; skip lists, ...

- *data compression: LZ algorithms, suffix trees.

— *hashing: random/uniform probing; LPH; paged; alternative displacements;

- priority trees, heaps, mergesort, sorting networks;

-- *symbolic manipulation: polynomial GCDs, factorization;
 symbolic differentiation and term-rewritings;

 quantitative data mining: probabilistic & approximate counting; Loglog counting; adaptive sampling.



Patterns in sequences



Many kinds of patterns are recognized by finite automata leading to rational functions whose poles move smoothly.

"In random sequence, the number of pattern occurrences is asymptotically normal, for a great variety of patterns and information sources."

[Guibas-Odlyzko; Régnier-Szpankowski; Nicodème-Salvy-F; Vallée]

"Borges' Theorem" for local patterns is known to hold in: words, trees, permutations, search trees, maps, etc. [Devroye, Martinez, F., Bender, Gao, Noy-Elizalde, ...]

A CONTRACTOR OF A

Digital structures and data compression

* <u>Digital trees</u> aka *"tries"* & variants are amenable to analytic combinatorics: GFs, singularity analysis, <u>Mellin transforms</u>, saddle point method \cong analytic depoissonization.

E.g., Jacquet-Szpankowski DST eqn $F(z, u) = \int F(pz, u) \cdot F(qz, u)$

Vallée's dynamical sources: "The cost of radix-sorting of n continued fractions depends on the Riemann hypothesis."

* <u>Suffix trees</u> too: combine with pattern analyses.

"Redundancy of Lempel-Ziv compression algorithms can be precisely quantified."

The trie saga. [De Bruijn-Knuth, F.-Sedgewick, Devroye, Pittel, Jacquet-Szpankowski-Louchard, Vallée-F.]. [Szpankowski's red book]. Cf [Devroye-Szpankowski@SODA'07] ...

PART D. FRONTIERS

* Organize the field into analytic-combinatorial schemas exhibiting universal properties. Towards a theory of combinatorial processes

* Expand the scope of analytic methods to hard computational problems.

* Determine decidable classes and work out decidability algorithms within symbolic manipulation systems like Maple, Mathematica.

Hard combinatorial problems (NP): What are the feasible/ unfeasible regions for random problem instances? E.g.: answer not known for 3-SAT [2-SAT: BoBoCh+]



Represent problem by n-dimensional Cauchy integral; estimate by SP. E.g. 3-regular graphs from general graphs.

$$RG_n^{(3)} = \frac{1}{(2i\pi)^n} \int \cdots \int \prod_{1 \le i < j \le n} (1+z_i z_j) \frac{dz_1 \cdots dz_n}{z_1^4 \cdots z_n^4}$$

B. Mc Kay has developed a specific calculus. (Gives access to exponentially sparse families and can "filter" according to many constraints.)

Computability within the calculus of analytic combinatorics.

Algorithms and programs for "automatic combinatorics"?

Theorem (Properties of specifications)

For the core language of constructions: (i) counting sequences are computable in $O(n^{1+\epsilon})$; (ii) GF equations are computable; (iii) partially decidable asymptotic properties; (iv) random generation by either recursive method or Boltzmann models is achievable in low polynomial time.

[F-Salvy-Zimmermann] [Duchon-F-Louchard-Schaeffer] [F-Fusy-Pivoteau]



$$Tr(z) = \frac{z^2 \left(1 + \log((1-z)^{-1})\right)}{\left(1 - z^2 \left(1 + \log((1-z)^{-1})\right)\right)} \left(1 - \frac{z^2 \left(1 + \log((1-z)^{-1})\right) e^{\left(\log((1-z)^{-1})\right)^2}}{1 - z^2 \left(1 + \log((1-z)^{-1})\right)}\right)^{-1}$$

noathar with the armoneion