

平均計算時間に基づく計算複雑さの研究について

(解説)

東工大・工学部 渡辺 治 (Osamu Watanabe)

概要: NP 問題の平均計算時間の理論的研究について解説する.

1. 問題の背景

計算の複雑さの理論では, “時間計算量” はおもに最悪計算時間に基づいて測られている. (以下では, これを簡単に最悪時間計算量と呼ぼう.) たとえば, 「プログラムの計算時間が  $3n^2$  である」といえば, 普通, 「最悪の入力を考えても (その入力サイズ  $n$  にたいし)  $3n^2$  の計算時間でプログラムが答えを出す」ということを意味する. もちろん最悪時にも速いのはよいことだが, プログラム A の最悪時間計算量が  $10n^8$  で, プログラム B が  $n^2$  のとき, 簡単に B のほうが A より優れているとは言い切れない. ほとんどの場合で, A のほうが B より速い, つまり A が B よりも平均的に速いこともありうるからだ. それなら平均の計算時間に基づいて計算時間を測ればよいではないか, それが “平均時間計算量” である.

NP 問題, とくにその中でも NP 完全問題は, 手に負えない問題と予想されている. つまり, それらを解く効率的なアルゴリズムがなく, どんな速いコンピュータをもってしても天文学的 (いやそれ以上, 計算科学的! [Mi89]) 計算時間がかかると予想されている. ところがこれは “最悪時” の話である. もしかすれば平均的には速く解けるのかもしれない. たとえばつぎのような例を考えてみよう (この話は [Jo84] より).

ハミルトン閉路問題 (与えられたグラフにハミルトン閉路があるか否かを判定する問題) は, NP 完全問題としてよく知られている. ところが, ある妥当

な入力分布を仮定すると、この問題は平均時間  $O(n^2 \log n)$  で解くことができる。

定理： つぎのようなアルゴリズム H が存在する。

- (1) すべてのグラフにたいし、H はそれがハミルトン閉路を持つか否かを正しく判定する、そして
- (2) A の計算時間は、

$$\exists c, \forall n, \sum_{G \in \mathcal{G}_n} p_n(G) \times \text{time}_H(G) \leq c n^2 \log n + c.$$

ただし、 $\mathcal{G}_n$  は  $n$  頂点のグラフの全体。また、各  $n$  ごとに  $p_n$  は、 $n$  頂点のグラフ  $G$  が入力として与えられる確率で、つぎのように定義される：

$$p_n(G) = \left(\frac{1}{2} + \varepsilon\right)^e \times \left(\frac{1}{2} - \varepsilon\right)^{N-e}$$

ただし、 $N = n(n-1)/2$  ,  $e = \text{辺の数}$  .

(簡単にいえば、各辺を  $\frac{1}{2} + \varepsilon$  の確率でランダムに発生させたときに生成される確率。)

すなわち、確率分布  $p_H (= \{p_n\})$  に従って入力を与えられるとすると、アルゴリズム H は  $O(n^2 \log n)$  の平均時間計算量でハミルトン閉路問題を解いてしまうのである。この定理は、ハミルトン閉路問題が（最悪時間計算量的にみて）とうてい多項式時間内に解けないと予想されていることを考えると、驚異的な結果のように思える。しかし、そうでもない。実は、この確率分布で生成したグラフは、ほとんどの場合ハミルトン閉路を持ってしまうのだ。もう少し正確にいえば「 $p_H$  でグラフを生成すると、（ある方法で簡単に見つかる）ハミルトン閉路を持つグラフができる確率が  $1 - o(2^{-n})$  になる」ことが Chvátal により証明されている。したがって、ほとんどの入力ですぐにハミルトン閉路が見つかる。アルゴリズム H は、非常に運の悪い場合だけ、まじめにハミルトン閉路を捜せばよいのである。だから  $O(n^2 \log n)$  の平均計算時間で済むのだ。なんだかペテンにかけられたような気持ちだろう；アルゴリズム H の速さはハミルトン閉路問題の“易しさ”とはあまり関係ないからだ。もちろん、Chvátal の定理は非常におもしろい。しかし、この定理はハミルトン閉路問題の本質的な“難しさ／易しさ”を解析しているわけではない。では、平均

時間という立場から見て、ハミルトン閉路問題の本質的な“難しさ／易しさ”を議論するにはどうすればよいのだろうか？

平均時間計算量の方が最悪時間計算量より、より現実的であるにもかかわらず、その理論的研究があまりされていない。そのもっとも大きな理由は「妥当な入力（確率）分布を決めることができない」という点だろう。入力分布は、その問題が考えられている状況に応じて違ふのだから、ひとつの入力分布にたいして理論を作り上げてもただの机上の空論になりかねない。一方、すべての入力分布にたいして計算時間を議論することは、最悪時間計算量を解析することに他ならない。また、哲学的（？）に言えば、ひとつの入力分布のもとでの平均時間計算量の解析が、その問題の平均的難しさの本質的解析になるかどうかは疑問である。たとえば、ハミルトン閉路問題は、入力分布  $\rho_H$  では平均的に多項式時間で解けるけれども、その事実から「他の入力分布でもハミルトン閉路問題が易しい」ことは示せない。こういった問題点に決着をつけないことには平均時間計算量の理論的研究を始めることができないのである。

ところが最近（といっても1989年頃）、Levin がこういった問題に決着をつけるひとつの糸口を示した。つまり平均時間計算量を議論する枠組みを提案し、ある問題にたいしては“妥当な入力分布”を示したのである。また、それに関していくつかの研究が報告されている。以下ではその概要を述べたい。なお、以下の議論は、計算の複雑さの理論で普段行われている議論の方法にのっとり行い、よく知られている用語・記号は定義せずに用いる（[HU79]などを参照されたい）。ここでのアルファベット  $\Sigma$  は  $\{0,1\}$  とする。

## 2. Levin の枠組み

Levin の枠組みは“多項式時間計算可能性”を平均時間計算量の立場から議論するためのものである。平均時間計算量を一般的に考えるため、Levin の枠組みでは、問題と入力分布の組を扱う。入力分布は、 $\Sigma^*$  上の分布関数で与える。ただし、つぎの条件を満たす関数  $\mu$  が  $\Sigma^*$  上の分布関数と呼ばれる：

- (i)  $\forall x_1, x_2 \in \Sigma^*, x_1 \leq x_2 \rightarrow \mu(x_1) \leq \mu(x_2),$
- (ii)  $\lim_{x \rightarrow \infty} \mu(x) = 1.$

(注:  $\Sigma^*$  上の順序としては辞書式順序を用いる)

すなわち, 分布関数  $\mu$  のもとでは, 入力  $x$  は  $\hat{\mu}(x) = \mu(x) - \mu(\hat{x})$  という確率分布で与えられると考える. (ただし,  $\hat{x}$  は  $\Sigma^*$  上で  $x$  のひとつ手前の元. つまり,  $\hat{\mu}$  は  $\mu$  の微分.)

入力分布のこの決め方には少し違和感を感じられる読者も多いだろう. たとえばグラフの入力分布を考えると, 普通は  $\rho_H$  のように入力のサイズ (頂点数) ごとの確率分布を考える. それに対し, 上の方法は, すべてのグラフに対する確率分布を考えているのだ. たしかに  $\rho_H$  のような入力分布の与え方は直感に合うが, その方法だとサイズの決め方やコーディングの方法に依存した定義になってしまう. 一方, 上の方法はそれらに影響されないので, 一般論に向いている. また, 入力のサイズごとの議論も簡単に上の枠組みにはめ込むことができる. たとえば, 頂点  $n$  のグラフがすべて長さ  $m$  の  $\Sigma^*$  の文字列でコーディングされるとすると,  $\rho_H$  に対応する分布関数  $\mu_H$  には,

$$\forall n, \forall G \in \mathcal{J}_n, \hat{\mu}_H(G \text{ のコード }) = \frac{6}{\pi^2 m^2} \rho_n(G).$$

$$\left( \text{注: } \sum_{m=1}^{\infty} \frac{6}{\pi^2 m^2} = 1 \right)$$

を用いればよい.

さらに, Levin の枠組みでは, つぎのような条件を満たす分布関数 だけを考える:

(iii)  $\mu$  は多項式時間内に計算可能.

“自然な” 入力分布の多くはこの条件を満たしている. 以下では, 単に分布関数といったならば (i), (ii), (iii) の条件を満たす関数を指すことにしよう. Levin 流議論では, 集合  $D$  と分布関数  $\mu$  の組  $(D, \mu)$  を対象とするのである. この組を分布付き問題と呼ぶことにしよう.

分布付き問題  $(D, \mu)$  をひとつ考えよう. 我々は, その問題が平均的にみて多項式時間で解けるか, という点を議論したい. そのために “平均計算時間が多項式以下” という概念を定義しよう. あるアルゴリズム  $A$  と  $k > 0$  が存在して,

$$\exists c, \forall n, \sum_{x \in \Sigma^n} \hat{\mu}(x) \times \text{time}_A(x) \leq c \cdot n^k,$$

となれば平均計算時間が多項式以下だ、というのが自然な考え方だろう。ところがこれはあまりうまくない。アルゴリズム A が上の条件を満たしていても、その高々 2 乗の計算時間を要するアルゴリズム B で、

$$\forall k, \forall c, \exists n, \sum_{x \in \Sigma^n} \hat{\mu}(x) \times \text{time}_B(x) > c \cdot n^k,$$

となることがある。これでは「どんな計算モデルでも多項式時間計算可能性という概念は一定である」という点が崩れてしまう。これに対し、もう少し弱い基準を考える。

$$\exists c, \forall n, \sum_{x \in \Sigma^n} \hat{\mu}(x) \times (\text{time}_A(x))^{1/k} \leq c \cdot n.$$

これでも A の平均計算時間が多項式以下とってよい（ような気がする）。しかも、もし A が上の条件を満たすならば、それと高々多項式程度でしか計算時間がちがわないアルゴリズムはすべて上の条件を満たすことが示せる。そこでつぎのような定義が妥当だろう。

定義： 分布付き問題  $(D, \mu)$  にたいし、つぎのようなアルゴリズム A が存在するとき、それを多項式平均時間計算可能という。

(1) A は  $D$  を認識する決定性アルゴリズム、そして

$$(2) \exists k, \sum_{x \in \Sigma^*} \hat{\mu}(x) \times (\text{time}_A(x))^{1/k} / |x| < +\infty.$$

(注：この式は上の式をいいかえたものに過ぎない)

たとえば、 $(\text{HAM}, \mu_H)$  はアルゴリズム H のおかげで、多項式平均時間計算可能となる。

P, NP に対応して、つぎの計算量のクラスが考えられている：

$$\text{ave-P} = \{(D, \mu) : (D, \mu) \text{ は、多項式平均時間計算可能}\},$$

$$\text{dist-NP} = \{(D, \mu) : D \in \text{NP}\}.$$

$P \subseteq \text{NP}$  とことなり、 $\text{ave-P} \subseteq \text{dist-NP}$  であるとは限らない点に注意。（なぜ、 $\text{ave-NP}$  を考えないのだろうか？）我々の最大の関心事はつぎの疑問だ：

疑問:  $\text{dist-NP} \subseteq \text{ave-P}$  ?

さて, NP 完全問題のように  $\text{dist-NP}$  完全問題というのを定義できないだろうか? そのためには“多項式時間還元性”なるものをここでも定義する必要がある. そのためにいくつかの概念を定義しよう.

定義:

(1) 分布関数  $\mu_1$  と  $\mu_2$  にたいし, もしある多項式  $p$  で,

$$\forall x \in \Sigma^*, \mu_1(x) \leq p(|x|) \cdot \mu_2(x)$$

ならば,  $\mu_1$  は  $\mu_2$  に支配されるといい,  $\mu_1 \preceq \mu_2$  と記述する.

(2)  $\Sigma^*$  から  $\Sigma^*$  への関数  $f$  と分布関数  $\mu$  にたいし,  $f(\mu)$  をつぎのような分布関数と定義する:

$$\forall y, f(\mu)(y) = \sum_{f(x)=y} \mu(x).$$

定義: 任意の分布付き問題  $(D_1, \mu_1)$  と  $(D_2, \mu_2)$  にたいし, つぎのような関数  $f$  が存在するとき,  $(D_1, \mu_1)$  は  $(D_2, \mu_2)$  へ還元可能といい,  $(D_1, \mu_1) \alpha (D_2, \mu_2)$  と記述する.

(1)  $f$  は多項式時間計算可能な  $\Sigma^*$  から  $\Sigma^*$  への関数,

(2)  $f$  により,  $D_1 \leq_m^p D_2$ ,

(3)  $f(\mu_1) \preceq \mu_2$ .

(注: この定義の条件は Levin [Le86] のそれより強い. 私は [Le86] の定義には誤りがあると思う.)

これらの定義の意味についてここでは詳しく述べないが, つぎのことが成り立つように定義したと考えてよいだろう.

定理:  $(D_1, \mu_1) \alpha (D_2, \mu_2) \wedge (D_2, \mu_2) \in \text{ave-P} \rightarrow (D_1, \mu_1) \in \text{ave-P}$ .

$\text{dist-NP}$  完全性の定義は従来通りつぎのようになる:

定義:  $(D, \mu)$  は, つぎの条件を満たすとき  $\text{dist-NP}$  完全である.

(1)  $(D, \mu) \in \text{dist-NP}$ ,

(2) すべての  $(D', \mu')$  が,  $(D', \mu') \leq (D, \mu)$ .

Levin の仕事の大事な点は, このような dist-NP 完全な問題を実際に示した点にあると思う. ここでは, Gurevich [Gu87] が示した dist-NP 完全問題を考えよう.

$$K = \{ \langle M, x, 0^l, 0^t \rangle : \exists w \in \Sigma^l, M \text{ accepts } \langle x, w \rangle \text{ in } t \text{ steps} \},$$

$$\hat{\mu}_K(\langle M, x, 0^l, 0^t \rangle) = \frac{6}{\pi^2 (m+n+l+t)^2 \cdot 2^m \cdot 2^n}$$

ただし,  $M$  は決定性チューリング機械,  $m$  は  $M$  の記述の長さ, そして,  $n$  は  $x$  の長さ.

定理:  $(K, \mu_K)$  は dist-NP 完全.

つまり, ave-P  $\stackrel{?}{\supseteq}$  dist-NP? という問題は,  $(K, \mu_K) \in \text{ave-P?}$  という具体的な問題の難しさに帰着されたのだ. これは“妥当な入力分布”という面からみても重要である. 問題  $K$  では, その平均計算時間が入力分布  $\mu_K$  のもとで多項式以下ならば, どんな入力分布のもとでも多項式時間以下に解くことができるのである. この意味で,  $(K, \mu_K)$  は  $K$  の平均時間計算量を本質的に表しているといえよう. ハミルトン閉路問題にたいしても, このような入力分布を定義できる. (コメント: Levin らは, “自然な”入力分布を好んでいるようだ. ところが多くの場合, 自然な入力分布で完全問題をつくるのは難しい [Gu87]. これはコーディングの問題からくるのだと思う. しかし, なぜ自然な入力分布にこだわるのだろうか?)

従来の計算量のクラスとの関連について述べておこう. つぎの定理は簡単に証明できる.

定理:

$$(1) \quad P = NP \rightarrow \text{dist-NP} \subseteq \text{ave-P},$$

$$(2) \quad \text{DXT} \neq \text{NEXT} \rightarrow \text{dist-NP} - \text{ave-P} \neq \emptyset.$$

### 3. 暗号の安全性を調べるのに適した枠組み

最近では, ある問題の“計算の複雑さ”を用いて暗号系を作ることが行われて

いる。そうした暗号系では、生成した暗号鍵の平均的安全性が問題となる。具体例を示そう。つぎのようなアルゴリズム  $X$  を考える：

アルゴリズム  $X$  :

入力  $n$  にたいし、乱数を利用して約  $n$  ビットの素数の組  $(p, q)$  とその積  $m = p \cdot q$  を出力する。

この  $X$  の出力のうち  $m$  だけを見て、 $(p, q)$  を求める問題を考える。この種の問題が暗号鍵の安全性の問題と違ってよいだろう。さて、この問題（すなわち、ある種の整数にたいする素因数分解）の平均的難しさを解析したい。それには、入力分布を  $X$  が を出力する確率に基づいて決めるのが妥当だろう。ところが、そのような入力分布の中には、Levin の条件 (iii) を満たさないものがあり、Levin 流の枠組みではこの問題を議論できないのである。

Ben-David ら [BDGL89] は、条件 (iii) の代わりにつぎのような条件 (iv) を提案している：

(iv) 乱数を利用して  $\hat{\mu}$  の確率分布で入力例を生成する多項式時間アルゴリズムが存在する。

そして、NP 問題  $D$  と、(i), (ii), (iv) を満たす分布関数  $\mu$  の組  $(D, \mu)$  のクラスを  $\text{sample-NP}$  と定義した。このクラスは、我々が今考えた問題にピッタリのクラスである。このクラスについても完全問題の存在が示されている。

#### 4. 最後に

この分野はいまのところあまり注目されていない。しかし、ことによったら宝の山かもしれない、と私は思うのですが、...



## 参考文献

- [BCGL89] S. Ben-David, B. Chor, O. Goldreich, and M. Luby,  
On the theory of average case complexity,  
in "Proc. 21st STOC", ACM (1989), 204-216.
- [Gu87] Y. Gurevich, Complete and incomplete randomized NP problems,  
in "Proc. 28th FOCS", IEEE (1987), 111-117.
- [HU79] J. Hopcroft and J. Ullman, "Introduction to Automata Theory,  
Languages, and Computation", Addison-Wesley (1979).  
(邦訳) オートマトン 言語理論 計算論 II, サイエンス社.
- [Jo84] D. Johnson, The NP-completeness column: an ongoing guide,  
J. Algorithms 5 (1984), 284-299.
- [Le86] L. Levin, Average case complete problems, SIAM J. Comput. 15  
(1986), 285-286.
- [Mi89] S. Miyano, proposed in this meeting.