

⑯ 20 の扉について

高野金作

戦後アメリカから流行してきた 20 の扉 (*twenty questions*) は数学的には次の様に考えられる。

有限集合 A とそれに属する未知の確定要素 x とがあって、 A の任意の部分集合に対し、それが x を含むか含まないかを確めることができる時に、未知要素 x を求めること。そのためには A を互に共通要素のない二つの空でない部分集合 A_0 と A_1 とに分けて、 x を含むのがどちらであるかを確める(第1問)。

その結果例えば x は A_0 に含まれることが分ったとすれば、 A_0 を再び前の様に二つの部分集合 A_{00} と A_{01} とに分けて、 x を含むのがどちらであるかを確める(第2問)。その結果例えば x は A_{01} に含まれることが分ったとすれば、 A_{01} を再び前の様に二つの部分集合 A_{010} と A_{011} とに分けて、 x を含むのがどちらであるかを確める(第3問)。この様にしてゆくと、 A は有限集合であるから、いつかは唯一つの要素を含む部分集合に到達して x を求めることができる。

実際の場合には x を含む方の部分集合を次々に二つに分けていくのがよいのであるが、理論的には起り得る部分集合の系列、

(1) $A_0, A_1, A_{00}, A_{01}, A_{10}, A_{11}, A_{000}, A_{001}, A_{010}, A_{011}, \dots$

をきめてかゝることもできる。すると A のどの要素をとっても

それだけからなる集合 A_{v_1, v_2, \dots, v_i} が存在する。

このときその要素を 0 と 1 との有限系列 v_1, v_2, \dots, v_i 又は二進小数 $0, v_1, v_2, \dots, v_i$ で表はすことができる。即ち、集合列 (1) をきめておけば、 A の任意の要素を 0 と 1 との有限系列又は二進小数で表はすことができる。この二進小数表示は次の性質をもつ。

C1 一つの要素がある二進小数で表はされたとすれば、その小数点下の桁数を ℓ とするとき、小数点下 ℓ 桁まで x と一致する様な二進小数は、他の異なる要素を表はすのに用いられない。

逆に、有限小数 A の各要素に二進小数を対応させ、条件 C1 が成立する様にしておくと、(1) の集合の系列を作ることができる。

A の要素をその小数第 1 位が 0 又は 1 であるに従つて A_0 又は A_1 に入れ、 A_0 の要素をその小数第 2 位が 0 又は 1 であるに従つて A_{00} 又は A_{01} に入れ、以下同様にやればよい。

二進小数の小数点下の桁数（系列で言えば項数）をその長さということにする。20 の席では、解答に到達するまでの項間数はあるべく少い方がよいから、二進小数の長さはなるべく短い方がよい。解答者には本当の解は分つていないわけであるが、解に対する確率分布は分つているとしよう。前に未知の確定要素として考えた X を、 A の要素の値をとる確率変数の実現値として考えるわけである。さうすると

C2 X の値を表はす二進小数の長さの平均値を最小にする。
様子表はし方が望ましい。そのためには A の各要素をどんな方法で二進小数で表はしたらよいであろうか。

問題を整理して述べると次の様になる。有限集合 $A = \{a_1, a_2, \dots, a_n\}$ を値域とする確率変数 X があって、 X は a_1, a_2, \dots, a_n の値をそれぞれ確率 p_1, p_2, \dots, p_n でとること（ $p_1 + p_2 + \dots + p_n = 1$ ）， A の要素 a_1, a_2, \dots, a_n をそれぞれ二進小数 $x_1, x_2, \dots,$

\cdots, x_n で表はし、條件 C_1 及び C_2 を成立させるには、 x_1, x_2, \cdots, x_n をどの様にとつたらよいか。換言すれば、(1) の集合列をどの様にとつてゆけばよいか。そのためには、 A の要素を確率の減少する順に並べ、適当に番号をつけかえて、それを a_1, a_2, \cdots, a_n ($p_1 \geq p_2 \geq \cdots \geq p_n$) とし、確率がほぼ半分になる様に集合を分けてゆけばよいことが知られている。

詳言すれば、 $p_1 + \cdots + p_k = 1/2$, $p_{k+1} + \cdots + p_n = 1/2$ なる点をとつて、 $A_0 = \{a_1, \cdots, a_k\}$ と $A_1 = \{a_{k+1}, \cdots, a_n\}$ とに分け、 A_0 は $p_1 + \cdots + p_i = 1/4$, $p_{i+1} + \cdots + p_k = 1/4$ なる i をとつて、 $A_{00} = \{a_1, \cdots, a_i\}$ と $A_{01} = \{a_{i+1}, \cdots, a_k\}$ とに分け、以下同様にしてやってゆけばよい。

問題の動物であるとき、それが人間である確率が $1/2$ に近い様であれば、動物を人間とそれ以外の動物とに分ければよいし、不めるべきものが人間の部分であるとき、特別な理由がなければ、上半身と下半身とに分けて割けばよい。これらのこととは 20 の扇の解答者達がいつもやっていることである。

ある病気の治療にある薬草が特効があることが分っているが、どの成分が主として効いているのか分らないときにも、同様に考えることができることである。全成分を二つの組に分けてどちらの組が効くかを実験によって確かめ、有効な方を又二つの組に分けて同様にしてやってゆけばよいが、その分け方が問題である。

薬学や医学等の知識又は経験によって、どの成分が効くか大凡その見当はつけ得るであろう。確定成分である等の未知の有効成分を確率度数の実現値とみて、その豫想確率によって計算したときの実験回数の平均値を最小にするには、上記の方法を適用することができる。

これらのこととは通信工学に於ける *encoding* の問題と全く同じであるので、それを紹介して証明もそれについて行うことにしておきたい。

る。

通信に於ては送ることになつて *message* を、具体的な通信方法で実際送ることのできる符号 (*Symbol*) の系列で、一对一の対応がつく様に表はさなければならない。

数学的に言えば、起り得る *message* の集合 M と符号の系列の集合 S とがあつて、 M から S の中元の一対一の対応 f を求めることが問題になる。これを *encoding* (符号化) という。

M に属する各 *message* に対する f の値を求める装置は *encoder* 又は送信機といはれる。

受信地に於てこの逆の対応 f^{-1} を求めるものを *decoder* 又は受信機という。簡単のため、通信路に起る雑音や歪みは考へないことにする。

こゝでは M が有限集合で符号の種類も有限個である場合についてだけ考える。符号が λ 種類あるとすれば、それらを $0, 1, 2, \dots, (\lambda - 1)$ で表はすことができる。これらの符号から成る λ 項系列は、 0 と 1 との間のある実数の表進法による小数表示に於ける最初の λ 個の数字の系列に対応する。

例えは $\lambda = 10$ ならば、 $14159 \dots$ なる系列が十進小数 $x = 0.14159 \dots$ を小数第5位までとつて6位以下を切り捨てたものと見ることができ。従つて

$$0.14159 \leq x < 0.14160.$$

各種類の符号を用いたときの *encoding* の問題は、可能な *message* の各々に 0 と 1 の間のある実数を対応させ、それを表進法によつてある桁まで展開したもので表はす問題である。このとき 1415 及び 14159 の様な二つの系列で二つの異なる *message* を表はすわけにはいかないことに注意しなければならない。何と云れば、 f' を求める受信機は意味のある系列 (ある *message* に対応する系列) を受け取る度に

それに対応する message を再現するからである。

故に若し M なる系列が意味をもっているときに、系列 M が送られると、 M に対応する message が作られ、 M は次の message に属することになる。

M に属する messages を実数 x_1, x_2, \dots, x_n で表はしたとき、それらの小数点以下の桁数 l_1, l_2, \dots, l_n の満たすべき必要充分条件は、任意の x_i, x_j ($i \neq j$) で、それらを $\ell = \min(l_i, l_j)$ 桁までとったとき、一致しないことである。

さて M に属する message の中には比較的頻繁に送られるものと稀にしか送られないものとがあるって、その確率分布を考えることができます。各符号を送るに要する時間は同一であるものとしておく。さうすると、もっと頻繁に起るもののは短い系列で表すべきであり、稀にしか起らないものは長い系列で表してもよいことになる。このために M の確率分布が與えられたときに、系列の長さの平均値を最小にする様な encoding が望ましい。それは次の様にして得られる。

M が n 個の message を含むものとし、それを確率の小さい方から並べ、 i 番目の message の確率を $p(i)$ ($i=1, 2, \dots, n$) とする。 $p(1) \leq p(2) \leq \dots \leq p(n)$, $\sum p(i) = 1$ 。
さて i 番目の message を実数

$$(2) \quad x_i = \sum_{n=1}^i p(n)$$

で表はす。

$p(i)$ の有効数字が小数第 l_i 位から始まるとすれば、 x_i は小数第 l_i 位まで (l_i 位未満切捨) 求めておけばよいことは明らかである。 l_i は次の不等式

$$k^{-l_i} \leq p(i) < k^{-l_i+1}$$

従つて

$$(3) -\log_k p(i) \leq l_i < -\log_k p(i) + 1$$

を満たす。 message を表はすに用いられる符号の平均数は、
 $\sum p(i)l_i$ で、これは (3) により次の不等式をみたす。

$$(4) -\sum_i p(i) \log_k p(i) \leq \sum_i p(i) l_i < -\sum_i p(i) \log_k p(i) + 1.$$

そこで

$$(5) H = -\sum p(i) \log_k p(i)$$

とおけば、(4) により

$$(6) H \leq \sum_i p(i) l_i < H + 1$$

が成立つ。(5)で定義されるHはMのエントロピーと呼ばれ、通信工学では重要な役割を演する。対数の底kを変えることは、Hに定数を掛けることになるので、エントロピーの単位を変えることに相当する。エントロピーの定義としては普通 $k=2$ の場合をとっている。そのときの単位を binary digit 略して bit という。

この encoding の方法から分る様に、message の数nがkに比較して充分大きく、 $\max\{p(i)\}$ が充分小さいならば、message を表はす符号系列の初の符号（確率密度である）は、0, 1, ..., (k-1)なる符号を一様な確率(M内の確率即ち各 message の出現確率で考えて)でとる。第2の符号についても同様である。第1の符号と第2の符号とは独立である。

更に $\max\{p(i)\}$ が k^{-3}, k^{-4}, \dots に比較して小さければ、第3第4...の符号についても同様である。

勿論、符号の長さが 3, 4, ... 以上であるという條件の下での

語である。従って、message を表はす系列は長進法における random sequence を構成する。但し message の出現確率で考えたときのことと、 x_1, x_2, \dots, x_n に一様な確率を與えたときの議論ではない。

さてわれわれの encoding が実用上符号の平均数を最小ならしめることを証明しよう。先づ $n = k^s$ (s は正整数) ですべての $p(i)$ が k^{-s} に等しい場合を考える。 $-\log_k p(i) = s$ となるから、どの message も丁度 S 個の符号で表はされる。

従って各 message を表はす符号系列の平均の長さは S である。符号系列の平均の長さが S より小さい様な encoding は存在しない。これを証明しよう。任意の encoding を考え、符号系列の最小の長さを ℓ とする。 $\ell < S$ ならば、長さの平均を減少させて最小の長さを $\ell+1$ にすることができる。

何となれば、長さ ℓ の系列が m 個あったとすれば、これらの系列に更にもう一つの符号をつけ加えて出来る長さ $\ell+1$ の $m\ell$ 個の系列で、長さ ℓ の m 個の系列及び長い方の $m(k-1)$ 個をあきかえればよいからである。これを遂次行ってゆけば、長さの平均を減少させて最小の長さを S にすることができる。最小の長さが S の場合にはわれわれの encoding が平均の長さを最小にすることは明らかである。

次に一般の場合には長さの平均が $H - 1$ より小か $encoding$ は存在しないことを帰謬法によつて証明する。(M のエントロピー H は M の確率分布だけできまる量で、encoding とは無関係である。) 存在したとし、そのときの i 倍目の message に対応する系列の長さを ℓ_i (前のものとは異なる) とすれば、仮定により

$$(7) \quad \sum p(i) \ell_i < - \sum p(i) \log_k p(i) - 1$$

ここで $\ell_1, \ell_2, \dots, \ell_n$ を固定しておけば、両辺とも $p(1), p(2),$

----, $p(n)$ の連續函数であるから, $p(i)$ を適当な

$$(8) \quad p_i(i) = \frac{n_i}{k^s} \quad (S, n_i \text{ は整数}, \sum n_i = k^s)$$

でもさかえても, (7) は成立つ。即ち

$$(9) \quad \sum p_i(i) l_i < -\sum p_i(i) \log_k p_i(i) - 1 = -\sum p_i(i) \log_k n_i + S - 1.$$

故に

$$(10) \quad \sum p_i(i) (l_i + \log_k n_i + 1) < S$$

$[a]$ で a 以上の最小の整数を表はすことにはすれば ($a \leq [a] < a+1$)

$$(11) \quad \sum p_i(i) (l_i + [\log_k n_i]) < S.$$

第 i 番目の系列に更に $[\log_k n_i]$ 個の符号をつけ加えて長さを $l_i + [\log_k n_i]$ にすることにより, 少くとも n_i 個の系列を得るから, それらの中から適当に n_i 個とり, これをすべての i について考え, この様にしてできる $\sum n_i = k^s$ 個の各系列に改めて確率 k^{-s} を與えることにはすれば, 前に証明した結果と矛盾する。

以上により, われわれの encoding では符号系列の長さの平均値 $\sum p(i) l_i$ は H と $H+1$ との間 (H を含む) にあり, 任意の encoding におけるそれは $H-1$ 以上であることが分った, こゝに H は M のエントロピーである。

さてある $p(0 < p < 1)$ と整数 k とか存在し

$$\sum_{i=1}^{k-1} p(i) < p \leq \sum_{i=1}^k p(i) \quad \text{で且つ } p(k) \text{ は充分小さい}$$

ならば, H は充分大きくなり得るから, この場合には近似的に $\sum p(i) l_i \approx H \approx H-1$ としてよい。

従って, われわれの encoding は実際的見地からは最適のものである。

これで証明は終る。

前に述べた 20 の扉の場合の符号は yes と no との二つで、
 $k=2$ の場合である。

あとがき。

通信工学における encoding の話は G.A. Barnard, The theory of information, J.R.S.S. Vol. 13, No. 1, 1951 による。

但し、証明に不充分な点がある様に思はれるので、部分的な改良及び補足を加えた。

ここで述べた encoding の方法は $k=2$ の場合は, C.E. Shannon 及び R.M. Fano によって独立に発見されたものである。

(C. E. Shannon and W. Weaver, The mathematical theory of communication, 1949, pp. 29-30 参照)

なお、国沢清典、通信工学における coding の統計的改良、確率過程第1回シンポジウム講演集、1952年、を参照せよ。

(1953 年 2 月 21 日)