

エントロピーと同時分布推定

堀 部 安 一

1. はじめに

周辺分布から同時分布は一意には定まらない—これはたいへん明らかなことなので、確率統計の教科書にも取り立てては出てこないほどである。たとえば二値0, 1をとる変数 x_1, x_2, x_3 に関する三通りの2次元分布 $p_{(1,2)}, p_{(1,3)}, p_{(2,3)}$ ($p_{(i,j)}$ は (x_i, x_j) の分布をあらわす) がつぎのように与えられているとしよう。

$x_1 x_2$	$p_{(1,2)}(x_1, x_2)$	$x_1 x_3$	$p_{(1,3)}(x_1, x_3)$	$x_2 x_3$	$p_{(2,3)}(x_2, x_3)$
0 0	3/8	0 0	1/8	0 0	1/8
0 1	1/8	0 1	3/8	0 1	3/8
1 0	1/8	1 0	3/8	1 0	3/8
1 1	3/8	1 1	1/8	1 1	1/8

これらは、つぎの分布 $p = \{p(x_1, x_2, x_3)\}$ の周辺分布になっていることがわかる。

$x_1 x_2 x_3$	$p(x_1, x_2, x_3)$
0 0 0	t
0 0 1	$3/8-t$
0 1 0	$1/8-t$
0 1 1	t
1 0 0	$1/8-t$
1 0 1	t
1 1 0	$2/8+t$
1 1 1	$1/8-t$

$$(0 \leq t \leq 1/8)$$

すなわち、 p の (x_i, x_j) —周辺分布(これを $p_{(i,j)}$ であらわす)をとると、

$$p_{(i,j)} = p_{(i,j)}, 1 \leq i < j \leq 3, t \in [0, \frac{1}{8}].$$

多数の変数に関する同時分布を伴う数理的モデルでは、いくつかの低次元周辺分布だけから、同時分布を推定しなければならないことがおこりうる—限られたデータから直接信頼度高く得られるのは低次元分布だけだからである。

以下で述べる最大エントロピー的同時分布推定は、この問題に対する1つの興味ある方向と思われる。その際に、divergence とよばれる情報理論的な量 $D(p, q) = \sum p(x) \log(p(x)/q(x))$ の演ずる役割が基本的であることがわかる。

2. 記号など

N 個の変数 x_1, \dots, x_N があり、簡単のためどれも二値0, 1をとるとし、われわれには $\binom{N}{2}$ 通りの2次元分布 $p_E, E = \{i, j\}, 1 \leq i < j \leq N$, のみがあたえられているとしよう。ここで p_E は、 $E = \{i, j\}$ なら (x_i, x_j) の分布をあらわしている。添字 $1, \dots, N$ を頂点とし、 $E = \{i, j\}$ を頂点 i, j 間の辺とする。 N 頂点完全グラフ K_N において、各辺 E に分布 p_E が対応していると考えてもよい。1節の例では図1のようである。

N 次元分布 $p = (p(x)), x = (x_1, \dots, x_N)$, は 2^N 個の成分をもつ(ベクトル)点とみられる。あらゆる N 次元分布の全体を Π とする。 $p \in \Pi$ の (x_i, x_j) —周辺分布を $p_{(i,j)}, E = \{i, j\}$, であらわし、 $p_{(i,j)} = p_E$ ならば p は p_E を満たすと言おう。すべての

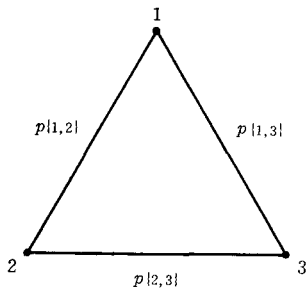


図 1 K_3

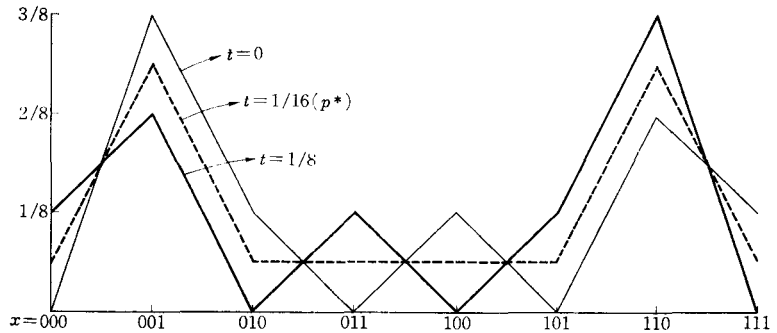


図 2

E について $p_{/E} = p_E$ となっている p の全体を Π_0 とする ($\Pi_0 \subset \Pi$). Π_0 が空だと意味がないので $\Pi_0 \neq \emptyset$ とする. 1 節の例では t を $[0, 1/8]$ 内で動かして得られる分布 $(t, 3/8-t, \dots, 1/8-t)$ の全体が Π_0 である.

$E = \{i, j\}$ のとき $x = (x_1, \dots, x_N)$ の部分ベクトル (x_i, x_j) を x_E であらわす. また, すべての辺 E , すべての x_E について $p_E(x_E) > 0$ としておこう——もしある $E = \{i, j\}$ と $(x_i, x_j) = (a, b)$ に対して $p_E(a, b) = 0$ なら, $(x_i, x_j) = (a, b)$ なるすべての x に対して $p(x) = 0$ とする $p \in \Pi$ のみに話を限定すればよい.

問題は, Π_0 に属すどの分布を推定分布とし, それをどのように見出せばよいかということになる.

3. エントロピー最大の分布

分布 p のエントロピー $H(p)$ は, その分布が《支配する》事象生起の不確定さの程度を測る量であって,

$$H(p) = -\sum_x p(x) \log p(x)$$

で与えられる. 対数の底は以後 2 とする.

『われわれは 2 次元分布族 $\{p_E\}$ しかわかっていない. したがって $\{p_E\}$ を満たす分布のうち(すなわち Π_0 内の分布のうち)最も不確定さの大きい分布を推定分布とせざるをえない』という立場をとるならば,

$$H(p^*) = \max_{p \in \Pi_0} H(p)$$

なる $p^* \in \Pi_0$ を推定分布とするのが自然であろう.

p^* は, $H(p)$ が狭義凸関数 $[-u \log u$ の狭義凸性より] であり, Π_0 が凸集合をなすことより, 一意に決まる. [Π_0 の凸性は, $p, p' \in \Pi_0, \lambda + \lambda' = 1, (\lambda p + \lambda' p')_{/E} = \lambda p_{/E} + \lambda' p'_{/E} = \lambda p_E + \lambda' p'_E = p_E$ より]

1 節の例での p のエントロピーは,

$$H(p) = h(t) = -3t \log t - (3/8-t) \log(3/8-t) - 3(1/8-t) \log(1/8-t) - (2/8+t) \log(2/8+t).$$

$$\max_{p \in \Pi_0} H(p) = \max_{0 \leq t \leq 1/8} h(t) = h(1/16) \text{ より,}$$

$$p^* = (1/16, 5/16, 1/16, 1/16, 1/16, 1/16, 5/16, 1/16).$$

図 2 では, この例について, 分布が折線であらわされていて, 折線群が Π_0 を代表的に示している.

4. $\{p_E\}$ の不完全な利用

Π_0 に属さない分布でも p^* に《近い》ものが, $\{p_E\}$ を使って簡単に得られる場合がある. 分布 q が分布 p からみてどの程度遠まっているか, 離れているかを, divergence とよばれる量

$$D(p, q) = \sum_x p(x) \log \frac{p(x)}{q(x)}$$

で測ることにしよう. つねに $D(p, q) \geq 0$ で, $= 0$ となるのは $p = q$ のときに限るという性質がある. $D(p, q)$ については付記にて 2, 3 述べることとし, ここでは, $D(p, q)$ が小さくなればなるほど q は p によく《似た》分布になるということだけを記しておこう.

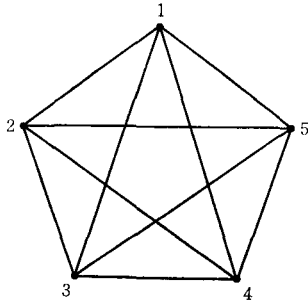


図 3 K_5

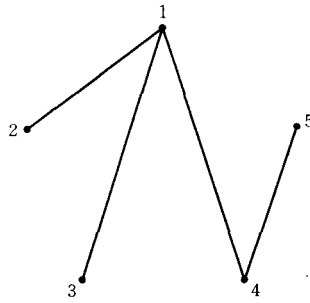


図 4 τ

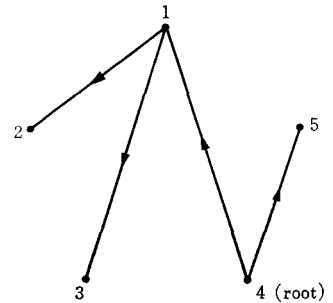


図 5 rooted tree

《情報》 $\{p_E\}$ から同時分布 $p \in \Pi$ を構成する場合、 $\{p_E\}$ の《利用程度》が大きくなれば、 $D(p^*, p)$ は小さくなることを例でもって示してみたい。

K_5 , $N=5$ を考えよう (図 3)。

0) $\{p_E\}$ をまったく使わずに一様分布 $p_0, p_0(x) = 2^{-5}$, を作ると、 $D(p^*, p_0) = 5 - H(p^*)$ 。

1) $p(x_i, x_j), 1 \leq i < j \leq 5$ (混乱はおきないので、 p_E の E など省略) より、その 1 次元周辺分布 $\{p(x_i)\}, 1 \leq i \leq 5$ をとり、 $p_1(x) = p(x_1) \cdots p(x_5)$ を作る。 p^* が $\{p(x_i)\}$ を満たすことより $D(p^*, p_1) = H_1 + \cdots + H_5 - H(p^*)$ 。ここに、 H_i は $\{p(x_i)\}$ のエントロピー。 $H_i \leq \log 2 = 1$ だから明らかに $D(p^*, p_1) \leq D(p^*, p_0)$ 。

2) たとえば $p_2(x) = p(x_1, x_2)p(x_3, x_4)p(x_5)$ を作る。 p^* がこれら成分分布を満たすことより、

$$D(p^*, p_2) = H_{12} + H_{34} + H_5 - H(p^*)$$

ここに、 H_{ij} は $\{p(x_i, x_j)\}$ のエントロピー。

$H_{ij} \leq H_i + H_j$ であるから $D(p^*, p_2) \leq D(p^*, p_1)$ 。

3) K_5 の完全木 (spanning tree) τ をとる (図 4)。 どれかの頂点たとえ 4 を root とする τ の rooted tree を考える (図 5)。

root には $\{p(x_4)\}$ を、有向辺 (i, j) には $\{p(x_i, x_j)\}$ より得た $\{p(x_j|x_i)\}$ を対応させこれらの積によって分布 p_r を作る:

$$p_r(x) = p(x_4)p(x_5|x_4) \\ p(x_1|x_4)p(x_2|x_1) \\ p(x_3|x_1)$$

これはたしかに確率分布であり

root を変えると表現は変わっても分布としては同じものが得られることを容易に示すことができる。 p^* が p_r の成分分布を満たすことより、

$$D(p^*, p_r) = H_4 + H_{5|4} + H_{1|4} + H_{2|1} + H_{3|1} \\ - H(p^*)$$

ここに、 $H_{j|i}$ は $\{p(x_j|x_i)\}$ に対する条件付エントロピー。 さて x_i と x_j との相互情報量を I_{ij} とあらわすと、 $I_{ij} = H_j - H_{j|i} = H_i - H_{i|j}$ だから、

$$D(p^*, p_r) = -(I_{12} + I_{13} + I_{14} + I_{45}) \\ + (H_1 + \cdots + H_5) - H(p^*)$$

ここで右辺の I_{ij} の和の部分だけが完全木 τ のとり方に依存していて、 τ の辺に対応して I_{ij} があらわれている。 したがって、 K_5 の各辺 $E = \{i, j\}$ について I_{ij} を計算しておき、これを E の《重み》とすれば、辺の重み和最大の完全木 τ_0 が簡単に見つかる ([4]) ので、 $p_{r_0} = p_3$ とすれば任意の τ に対して $D(p^*, p_3) \leq D(p^*, p_r)$ となる。 したがって

$$D(p^*, p_2) = H_{12} + H_{34} + H_5 - H(p^*) \\ = -(I_{12} + I_{34}) + (H_1 + \cdots + H_5) \\ - H(p^*) \geq D(p^*, p_3) \text{ となる。}$$

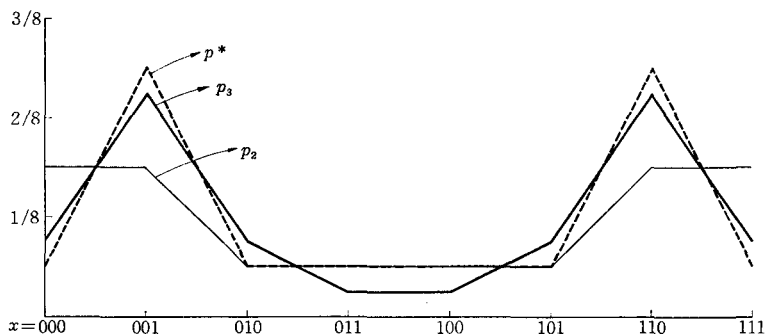


図 6

p_3 を求めるには10通りの p_E をすべて用いたが、 p_3 自身には $N-1=4$ 個の p_E しか取り入れていないことに注意されたい。1節の例で、 p_2 として $p(x_1, x_2)p(x_3)$ をとり、 p_3 として、($I_{12}=I_{13}=I_{23}$ より)たとえば、 $p(x_1)p(x_2|x_1)p(x_3|x_2)$ をとると、図2と対比して図6のようになる。

以上の計算では、 p^* に関しては、それが $\{p_E\}$ を満たしているということしか使っていないので、 p^* を、任意の $p \in \Pi_0$ でおきかえてもよい。したがって、divergenceの意味で、 p_0, p_1, p_2, p_3 の順に Π_0 に近くなっていると考えられる。

5. 反復計算

p^* を求めることは、凸集合 Π_0 上での凸関数 $H(p)$ の最大化問題にほかならないが、 2^N 個もの変数 $p(x)$ を扱わなければならない。ここでは別の方法により、あるクラスに属する分布 p^0 を初期分布として、反復的に $p^0 \rightarrow p^1 \rightarrow \dots \rightarrow p^*$ に至らしめることを考えよう。

前節の p_0, p_1, p_2, p_3 のように、積の形 $cq_1 \dots q_r$ をしていて、 c は定数、各因子 q_k は x_1, \dots, x_N のうち高々2個の変数にしか依存しないようになっている分布を、初期分布 p^0 に選ぼう。 p^* に近いという意味から、 p_3 を p^0 にとるのもよい。このような p^0 に対しては、前節の具体的な計算からもわかるように、 $D(p, p^0) = -\sum_x p(x) \log p^0(x) - H(p)$ において、 $-\sum_x p(x) \log p^0(x)$ は、 $p \in \Pi_0$ によらずに一定 c_0 であることに注意したい。

さて、 p^n から p^{n+1} を求めるのに、 p^n が(divergenceの意味で)《最も満たしていない》 p_E を p_{E_n} とする(greedy!):

$$D(p_{E_n}, p^{n/E_n}) = \max_E D(p_E, p^{n/E}).$$

そこで p_{E_n} が満たされるように p^n に《押し込む》:

$$(1) \quad p^{n+1}(x) = \frac{p^n(x)}{p^{n/E_n}(x_{E_n})} \cdot p_{E_n}(x_{E_n})$$

こうして、 p^{n+1} が p_{E_n} を満たす分布になったことは明らかだが、そのつぎの段階でもこの p_{E_n} が満たされているという保証はない。

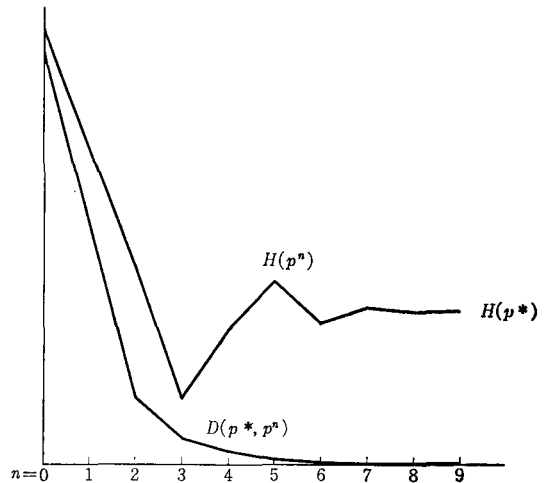


図 7

(1)によると、計算時間を喰うことはよしとするなら、この反復計算では常時 $O(N^2)$ の記憶量でよいことがわかる： $p^0(x) = 2^{-N}$ の場合、各辺 E には2つの表 $\{p_E(x_E)\}$ と $\{r_E(x_E)\}$ が対応していると考え、初期では $r_E(x_E)$ はすべて1にセットされている。 n 段階にきたとき、辺 E_n においては、

$$r_{E_n}(x_{E_n}) \leftarrow r_{E_n}(x_{E_n}) \cdot \frac{p_{E_n}(x_{E_n})}{p^{n/E_n}(x_{E_n})}$$

とし、他の辺 E では r_E は変更しないとすればよい。表 $\{p_E(x_E)\}, \{r_E(x_E)\}$ は4項目を含むにすぎず、 E は $\binom{N}{2}$ 個あるため記憶量は $O(N^2)$ である。

(1)より、ただちにつぎの式を得る:

任意の $p \in \Pi_0$ について

(2) $D(p, p^n) - D(p, p^{n+1}) = D(p_{E_n}, p^{n/E_n})$.
したがって右辺分 (≥ 0) だけ、 p^n は Π_0 に近づき p^{n+1} となることがわかり、 $D(p, p^n)$ は単調減少し、ある値(c_p とおく)に収束する。

1節の例で、 $p^0(x) = 1/8$ としたときの、 $H(p^n)$ と $D(p^n, p^*)$ の変化の様子が図7に示してある。

6. 収束性の証明

$p^n \rightarrow p^*$ ($n \rightarrow \infty$)を示そう。 Π は閉じた単体を形成しているから、 $\{p^n\}$ から収束部分列 $\{p^{n_\nu}\}$ がとれ、 $p^{n_\nu} \rightarrow \bar{p}$ ($\nu \rightarrow \infty$)とする。(2)で、 n を n_ν でおきかえ、 $\nu \rightarrow \infty$ とすると、左辺は $\rightarrow 0$ 。ところ

が任意の E について $D(p_E, p^{n\nu/E}) \geq D(p_E, p^{n\nu})$. したがって $\bar{p}_{/E} = p_E$. ゆえに $\bar{p} \in \Pi_0$.

(2) は任意の $p \in \Pi_0$ について成立し, その右辺は p に依存しないので,

$$\begin{aligned} D(p^*, p^n) - D(p^*, p^{n+1}) \\ = D(\bar{p}, p^n) - D(\bar{p}, p^{n+1}). \end{aligned}$$

したがって,

$$\begin{aligned} D(p^*, p^0) - D(p^*, p^{n\nu}) \\ = D(\bar{p}, p^0) - D(\bar{p}, p^{n\nu}). \end{aligned}$$

ここで $\nu \rightarrow \infty$ とすると, $D(p^*, p^0) - c_{p^*} = D(\bar{p}, p^0)$ 前節の p^0 に関する注意より,

$D(p^*, p^0) = c_0 - H(p^*)$, $D(\bar{p}, p^0) = c_0 - H(\bar{p})$ であるから,

$$H(\bar{p}) - H(p^*) = c_{p^*}$$

ところが, $H(p^*)$ は最大エントロピーであるから $H(p^*) \geq H(\bar{p})$. ゆえに $c_{p^*} = 0$.

これより $D(p^*, p^n) \rightarrow 0$, したがって $p^n \rightarrow p^*$.

7. おわりに

以上において, x_i は 0, 1 しかとらないとしたが, 一般の多値としても本質的変更はいらない. また, E としてあらゆる $\{i, j\} \subset \{1, \dots, N\}$ を考え, 完全グラフをみてきたが, この仮定も必要なく一般の単純グラフを考えればよい. さらに E として $\{1, \dots, N\}$ の任意の部分集合をとっても (このとき既知の p_E は $|E|$ 次元分布となる), 5, 6 節は記憶量の点を除けばそのままよい. 構造的にはハイパーグラフ的扱いになると思われるが皮相の見を越えていない.

付記—divergence について

$D(p, q)$ は, 情報理論的議論にしばしば現われて特異な役割をもつ量であり, informational divergence, discrimination information などとよばれることがある. p を《基準》として, それから, q がどの程度 diverge しているか, 違いがあるか, 離れているか, どの程度判別しうるか, などをあらわす量である. この場を借りて, このような意味を裏づけるような事実を 2, 3 取り上げて略述しておくのが適当と思われる.

1) $D(p, q)$ は q について凸関数である:

$$D(p, (1-\lambda)q + \lambda q') \leq (1-\lambda)D(p, q) + \lambda D(p, q')$$

[これは $-\log u$ の凸性より]. とくに $q' = p$ ととると $D(p, \lambda p + (1-\lambda)q) \leq (1-\lambda)D(p, q)$. ここで, λ を 1 に近づけると分布 $\lambda p + (1-\lambda)q$ は p にだんだん《似てくる》が, このことが $D(p, \lambda p + (1-\lambda)q)$ が 0 に近くなることであらわされている.

2) variation $V(p, q) = \sum_x |p(x) - q(x)|$ との関係は

$$D(p, q) \geq \frac{1}{2} V^2(p, q)$$

(もっとシャープな評価については[5]).

3) 分布 p をもつ情報源を 2 元符号化する場合, 各 x に約 $-\log p(x)$ ビットを使えば, 平均はほぼ $H(p)$ ビットの《コンパクト》な符号化ができる. もし, p が不明でその推定分布 q に対して同様に符号化を行なうと, 平均 $\sum_x p(x)(-\log q(x))$ ビットになる. 両者の差は $\sum_x p(x)(-\log q(x)) - H(p) = D(p, q)$ である.

4) いま分布 p または分布 q にしたがって独立な標本が次々と得られるとする. 決定者はこの標本が, p にしたがって出現するのか, q にしたがうのかを, 標本の《出かた》を観測しながら決定しなければならない. 逐次決定理論によると, 決定を誤る確率を許容値に固定しておけば, p が真のとき, 決定に至るまでに要する標本の大きさの期待値と $D(p, q)$ との積は大体一定している.

参考文献

- [1] Brown, D. T.: A Note on Approximations of Discrete Probability Distributions, *Information and Control* 2(1959), 386-392.
- [2] Forney, G. D.: Information Theory, Dept. Electrical Eng., Stanford Univ. Course Notes. EE 376, Winter 1972.
- [3] 国沢清典: エントロピー・モデル, 日科技連出版社, 1975.
- [4] Kruskal, J. B. Jr.: On the Shortest Spanning Subtree of a Graph and the Traveling Salesman Problem, *Proc. Amer. Math. Soc.* 7(1956), 48-50.
- [5] Toussaint, G. T.: Sharper Lower Bounds for Discrimination in Terms of Variations, *IEEE Trans. Information Theory*, IT-21(1975), 99-100.

(ほりべ・やすいち 静岡大学)