

2 × n 分割表の Perfect Sampling

東京大学 *来嶋秀治 KIJIMA Shuji
01605000 東京大学 松井知己 MATSUI Tomomi

1. はじめに

2元分割表は正の整数からなる行和と列和を持ち、表中の値として非負整数をとる表(行列)である。2元分割表は医療統計の分野などで統計データを扱うのに用いられる。分割表に対する統計学的興味の一つに行と列の相関関係があげられる。行と列の独立性の検定に χ^2 検定があるが、周辺和の中に小さな値があると、 χ^2 検定の当てはまりが悪くなるため、正確検定が行われる。

正確検定に用いる p 値の計算には、周辺和を満たす分割表を全列挙する方法があるが、分割表の全列挙は一般に困難なため、MCMC(マルコフ連鎖モンテカルロ)法がよく用いられる。MCMC法はモンテカルロ法を実行する際、マルコフ連鎖の定常分布からのサンプルを用いる手法である。しかし、定常分布は無限回の推移の後に得られるため、定常分布からのサンプリングは一般に困難である。本報告では $2 \times n$ 分割表の一樣標本抽出法にサンドイッチングの技法を応用した CFTP(Coupling From The Past)理論を用い、Perfect Sampling(厳密に一樣な分布からの標本抽出)を行う手法を提案する。

2. マルコフ連鎖と標本抽出アルゴリズム

整数(非負整数、正整数)全体の集合をそれぞれ \mathbb{Z} (\mathbb{Z}_+ , \mathbb{Z}_{++}) で表すことにする。ベクトル $r = (r_1, r_2) \in \mathbb{Z}_{++}^2$ と $s = (s_1, \dots, s_n) \in \mathbb{Z}_{++}^n$ は正整数 $N \in \mathbb{Z}_{++}$ に対して、 $\sum_{i=1}^2 r_i = \sum_{j=1}^n s_j = N$ を満たすとする。行和 r および列和 s をもち、非負整数を表値にとる 2 行 n 列の 2 元分割表全体の集合 $\Sigma_{r,s}$ を、

$$\Sigma_{r,s} \stackrel{\text{def}}{=} \left\{ X \in \mathbb{Z}_+^{2 \times n} \mid \sum_{j=1}^n X(i,j) = r_i \quad (1 \leq i \leq 2), \quad \sum_{i=1}^2 X(i,j) = s_j \quad (1 \leq j \leq n) \right\}$$

で定義する。但し、 $X(i,j)$ は分割表 X の i 行 j 列の値を表す。集合 $\Sigma_{r,s}$ を状態空間にもつマルコフ連鎖 \mathcal{M} を定義する。マルコフ連鎖 \mathcal{M} の推移 $X^t \mapsto X^{t+1}$ ($t \in \mathbb{Z}$) は一樣実数乱数 $\lambda \in [1, n]$ が与えられた時、以下のように定義される。

Step 1: i 列と $i+1$ 列 ($i = 1, \dots, n-1$) について、 $a = X^t(1, i) + X^t(1, i+1)$, $b = X^t(2, i) + X^t(2, i+1)$ として、 $\theta_{X^t}(i) \stackrel{\text{def}}{=} \min\{a, b, s_i, s_{i+1}\} + 1$ を定義する。

Step 2: 乱数 $\lambda \in [1, n]$ に対して、 $[\lambda]$ 列と $[\lambda] + 1$ 列を選び、 $\phi_{X^t}(\lambda) \stackrel{\text{def}}{=} [\theta_{X^t}([\lambda]) (\lambda - [\lambda])]$ として、 $X^{t+1} \in \Sigma_{r,s}$ を

$$\begin{aligned} X^{t+1}(1, j) &= \begin{cases} \phi_{X^t}(\lambda) + \min\{a, s_i\} - \min\{a, b, s_i, s_{i+1}\} & (j = [\lambda]), \\ a - X^{t+1}(1, j-1) & (j = [\lambda] + 1), \\ X^t(1, j) & (\text{otherwise}), \end{cases} \\ X^{t+1}(2, j) &= s_j - X^{t+1}(1, j), \end{aligned}$$

と置く。

このマルコフ連鎖 \mathcal{M} は、明らかにエルゴード性を満たす。また \mathcal{M} は detailed balance equations を満たし、定常分布として一樣分布をもつ。

今、分割表 $X_U, X_L \in \Sigma_{r,s}$ をそれぞれ、北西隅のルールおよび南西隅のルールに従ってできる表とする。すなわち、 X_U は $X(1, 1) = s_1$ から順に $X(1, 2) = s_2, X(1, 3) = s_3, \dots$ と、 $\sum_{j=1}^k X(1, j) = r_1$ となるまで 1 行目を決め、残りの表値は周辺和を満たすように一意に決める。同様に X_L は 2 行目の 1 列目から表値を決める。すな

わち、

$$X_U \stackrel{\text{def.}}{=} \left(X(i, j) \in \mathbb{Z}_+ \left| \begin{array}{l} \exists k, r_1 = \sum_{j=1}^k X(1, j) \leq \sum_{j=1}^k s_j, X(2, j) = 0 \ (j = 1, \dots, k-1) \end{array} \right. \right),$$

$$X_L \stackrel{\text{def.}}{=} \left(X(i, j) \in \mathbb{Z}_+ \left| \begin{array}{l} \exists k, r_2 = \sum_{j=1}^k X(2, j) \leq \sum_{j=1}^k s_j, X(1, j) = 0 \ (j = 1, \dots, k-1) \end{array} \right. \right),$$

となる。それぞれ X_U を NW(North-West) 表、 X_L を SW(South-West) 表と呼ぶことにする。
マルコフ連鎖 \mathcal{M} の定常分布からの標本抽出アルゴリズム P を以下のように定義する。

Step 1: シミュレーション開始時刻を $t = -1$ とし、 λ_t を空列とする。

Step 2: 一様乱数 $\lambda(t) \in [1, n]$ を生成して λ_t の先頭に挿入し、乱数列 $\lambda_t := (\lambda(t), \lambda(t+1), \dots, \lambda(-1))$ とする。

Step 3: 時刻 t における 2 つの状態 $Y_U, Y_L \in \Sigma_{r, s}$ を $Y_U^t := X_U, Y_L^t := X_L$ とする。

Step 4: Y_U, Y_L を初期状態とし、数列 $\lambda_t = (\lambda(t), \lambda(t+1), \dots, \lambda(-1))$ を用いて、時刻 t から 0 までそれぞれマルコフ連鎖 \mathcal{M} を推移させる。

a) もし、 Y_U^0, Y_L^0 が時刻 0 で一致して $Y_U^0 = Y_L^0$ となっていれば、 $X_0 := Y_U^0 = Y_L^0$ を標本とする。

b) 一致していなければ $t := t - 1$ として Step 2 にもどる。

この標本抽出法は厳密に一様分布に従う。すなわち以下の定理が成り立つ。

定理 1 アルゴリズム P からの標本は厳密に一様分布に従う。

3. CFTP とサンドイッチング

CFTP 理論は 1996 年 Propp and Willson によって提案された。CFTP の重要な主張は以下の定理である。

定理 2 有限離散空間 Ω 上のマルコフ連鎖 MC はエルゴード性を満たすとする。 Ω 中のすべての要素について、時刻 $t < 0$ から同一の一様乱数列 $\lambda \in [0, 1]^t$ を用いて MC にしたがって推移させる。この時、時刻 0 の状態がすべて一致 (coalescence) して、 $X \in \Omega$ となっていれば、 X はマルコフ連鎖 MC の定常分布に従って得られた標本である。

したがって CFTP 理論によって厳密に定常分布に従う標本を得ることが出来る。しかし、CFTP 理論では全状態を初期状態とするマルコフ連鎖の推移を追跡する必要があり、全状態からのシミュレーションは現実には不可能である。CFTP の技法の 1 つとして、サンドイッチングがある。これは複数の特定の状態からの推移が一致することが全状態からの推移が一致することの十分条件になっていれば、その特定の状態についてのみシミュレーションをすれば良いというアイデアである。したがってサンドイッチングが適用できるマルコフ連鎖に対して、その特定の状態からの CFTP によって、厳密に定常分布に従う標本を得ることが出来る。

我々は $2 \times n$ 分割表のマルコフ連鎖 \mathcal{M} に対してサンドイッチングが適用できることを示した。それが以下の定理である。

定理 3 マルコフ連鎖 \mathcal{M} について、時刻 $t < 0$ における状態 $X_U, X_L \in \Sigma_{r, s}$ から、乱数列 $\lambda \in [1, n]^t$ を用いて得られる推移が、時刻 0 において一致して $X_0 \in \Sigma_{r, s}$ となるならば、全状態からのマルコフ連鎖は時刻 0 において状態 X_0 を取る。

また離散化 Dirichlet 分布に対しても、同様の手法を用いて perfect sampling することができる。

参考文献

[1] O. Häggström, "Finite Markov Chains and Algorithmic Application," London Mathematical Society, Student Texts 52, Cambridge University Press, 2002.

[2] J. Propp and D. Wilson, "Exact sampling with coupled Markov chains and applications to statistical mechanics," *Random Structures and Algorithms*, 9 (1996), pp. 232-252.