

計数問題における複合モンテカルロ法について

中村 将人 (金沢大学大学院 自然科学研究科)

小川 重義 (立命館大学理工学部 数理科学科)

1 はじめに

これから述べる複合モンテカルロ法は、バックトラック法とモンテカルロ法を組み合わせたアルゴリズムで、計数問題の近似解を得るために [1] により考案された。その手法について、[1],[2] において、精度と計算時間の両面において有効な探索が行えることが実験的に確かめられている。しかし、なぜそのような探索が可能になるかについての理論的考察は行われていない。

そこで、本稿ではバックトラック法による全探索法、モンテカルロ法、そして複合モンテカルロ法の計算量を比較することで、より定量的な考察を行うことを目的とする。

キーワード：計数問題、バックトラック法、モンテカルロ法

2 計数問題

ある配置の集合の中から、解となる要素をすべて数え上げる問題である。

計数問題の定式化

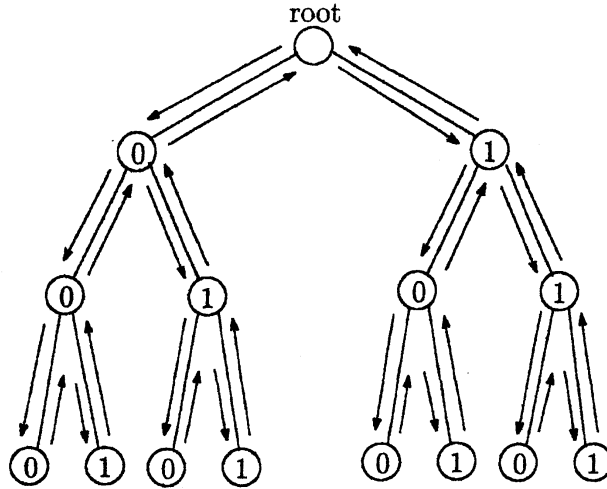
$$\begin{aligned} \Omega &= \{0,1\}^N \text{ 配置の集合} \\ S &\subset \Omega \text{ 解となる配置の集合} \\ x_i &\in \{0,1\} \quad i=1,2,\dots,N \\ \mathbf{x} &= (x_1, x_2, \dots, x_N) \quad (\in \Omega) \\ 1_S(\mathbf{x}) &= \begin{cases} 0 & , \mathbf{x} \notin S \\ 1 & , \mathbf{x} \in S \end{cases} \\ \text{このとき} \\ |S| &= \sum_{\mathbf{x} \in \Omega} 1_S(\mathbf{x}) \\ &\text{を推定したい。} \end{aligned}$$

したがって、正確な解の総数を求めるには、深さ N の 2 分木の全探索が必要となる。このような木を探索するアルゴリズムにバックトラック法がある。

(注) ここで、配置の集合 Ω として N 個の 0 と 1 からなる集合としているのは、あらゆる計数問題に対して適用できる手法を考えたいがためである。

3 バックトラック法

バックトラック法を用いると、効率よく配置を作ることができる。
 $N = 3$ の場合におけるバックトラック法で配置の作り方を図示したものを以下に示す。



今、どの枝も移動するのに要する時間を一定とすると、バックトラック法の計算量は木の枝の本数の2倍であるといえる。

$$2 \cdot \sum_{i=1}^N 2^i = 4(2^N - 1) \quad (1)$$

この式が示す通り、 N が大きくなることで計算量が増大し、実行不可能となる。そこで、ある程度の誤差を許容し、解の密度を推定することを考える。

4 近似値の精度について

単純に誤差 $|\hat{p} - p|$ をとったのでは、解の密度が小さい場合、ほぼ0になってしまい都合が悪い。そこで、相対誤差が ϵ 以下になる確率が $1 - \alpha$ 以上となる解の密度の推定値を得ることを目標と定める。

$$\text{Prob}\left\{\frac{|\hat{p} - p|}{p} < \epsilon\right\} > 1 - \alpha \quad (2)$$

$$p = \frac{|S|}{|\Omega|} \text{ 解の密度, } \hat{p} : \text{ 解の密度の推定量}$$

こうすることで、配置の総数 $|\Omega| (= m \text{ とする})$ よりも小さいサンプル数 n で、計算が実行できることが期待できる。

(※) 中心極限定理から次のことがわかっている。

$$\text{Prob}\left\{\frac{|\hat{p} - p|}{\sqrt{\text{Var}(\hat{p})}} < z_{\frac{\alpha}{2}}\right\} \xrightarrow{n \rightarrow \infty} 1 - \alpha \quad (3)$$

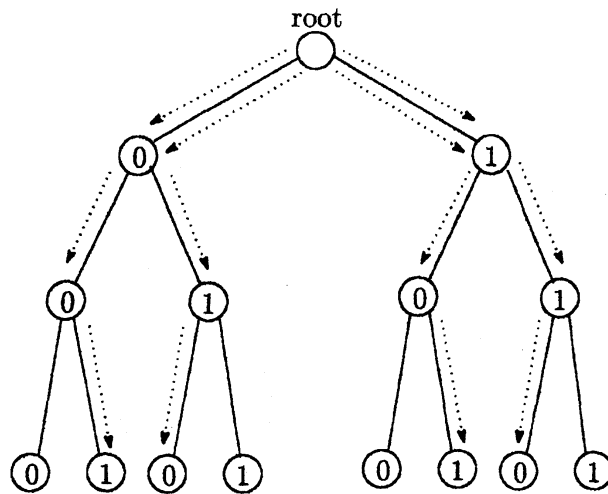
ここで、 $z_{\frac{\alpha}{2}}$ は $N(0, 1)$ の $\frac{\alpha}{2}$ 点

したがって、式(2)を満たすには密度の推定量の分散に課す条件は次式で表せる。

$$z_{\frac{\alpha}{2}} \sqrt{\text{Var}(\hat{p})} < \epsilon p \quad (4)$$

5 モンテカルロ法

具体的な推定量としてすぐに思いつのがモンテカルロ法である。ここでも、 $N = 3$ の場合の配置の作り方を以下に示す。

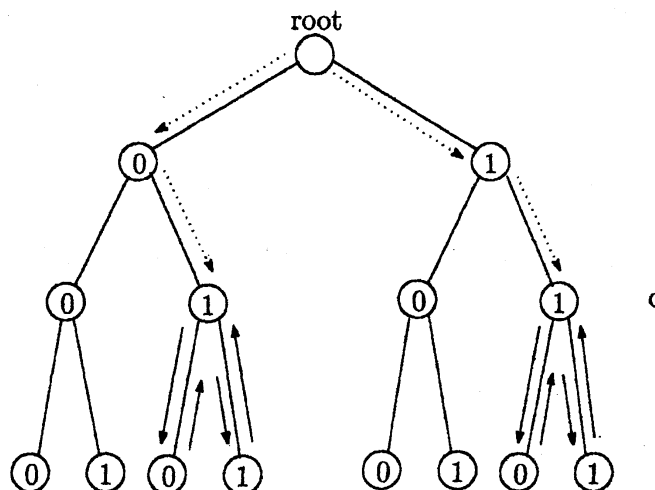


どの葉を調べる確率も等しくなるように探索を行い、見つかった解の頻度から推定する。しかしこの方法は、常に根からスタートするため、探索の効率においてバックトラック法に劣るという欠点がある。

そこで、バックトラック法とモンテカルロ法の長所を組み合わせたアルゴリズムである複合モンテカルロ法が考案された。

6 複合モンテカルロ法

$N = 3$ の場合、複合モンテカルロ法による配置の作り方を図示すると以下のようになる。



上図において、モンテカルロ法で降下する枝は点線の矢印、バックトラック法で移動する枝は実線の矢印で表現されている。こうして、各グループごとの解の頻度を調べ、それらの平均値を推定量とする。

6.1 カットオフレベルの調節

複合モンテカルロ法において最も重要なパラメータであるカットオフレベル c について説明する。カットオフレベルとはモンテカルロ法で根からランダムに降下する枝の本数のことである。上の図で言えばカットオフレベルは2である。なぜ重要かといえば、カットオフレベルにより推定量の分散が変化するからである。例えば、 $c = N$ とすればモンテカルロ法に完全に一致する。

また、カットオフレベルを決めることは集合 Ω を交わりの無い部分集合 (グループ) に分けることに相当する。

6.2 複合モンテカルロ法による推定量

これまでのことを定式化すると以下ようになる。

複合モンテカルロ法の定式化

$a := |S|$ 解の総数

s : m の約数 ($= 2^{N-c}$) グループ内の要素数

$b := \frac{m}{s}$ グループの数

$A(j)$: j 番目のグループの解の個数, $1 \leq j \leq b$

$X_k = A(r_k)$, (ただし, $r_k \sim 1$ から b までの離散一様分布 *i.i.d.*) (5)

$\hat{p} = \frac{1}{s} \frac{1}{n'} \sum_{k=1}^{n'} X_k$, (ただし, $sn' = n$) (6)

$E(\hat{p}) = p$, $Var(\hat{p}) = \frac{\sum_{j=1}^b \{A(j)\}^2 - \frac{a^2}{b}}{nm}$ (7)

式 (7) と、相対誤差に関する条件式 (4) から必要なサンプル数が求まる。

$$n > \frac{z_{\frac{\alpha}{2}}^2 \left(\sum_{j=1}^b \{A(j)\}^2 - \frac{a^2}{b} \right)}{m\epsilon^2 p^2} \quad (8)$$

6.3 複合モンテカルロ法の計算量

複合モンテカルロ法の計算量も、これまで同様移動する木の枝の本数で見積もることにする。したがって、まず c 本の枝をランダムに降下し、その後の枝をバックトラック法で移動する、という試行を n' 回繰り返すので以下の式のようになる。

$$\begin{aligned} & (c + 2 \cdot \sum_{i=1}^{N-c} 2^i) n' \\ &= (c + 4(2^{N-c} - 1)) \frac{n}{2^{N-c}} \\ &= (c2^{c-N} + 4(1 - 2^{c-N})) \frac{z_{\frac{\alpha}{2}}^2 \left(\sum_{j=1}^b \{A(j)\}^2 - \frac{a^2}{b} \right)}{m\epsilon^2 p^2} \\ &= ((c-4)2^c + 2^{N+2}) \frac{z_{\frac{\alpha}{2}}^2 \left(\sum_{j=1}^b \{A(j)\}^2 - \frac{a^2}{b} \right)}{\epsilon^2 a^2} \quad (9) \end{aligned}$$

6.4 計算量比

ここで、全探索を行う場合と比較がしやすいように計算量比を定義しておく。

$$\begin{aligned} \text{計算量比}(c) &:= \frac{\text{あるカットオフレベルにおける計算量(式(9))}}{\text{バックトラック法における計算量(式(1))}} \\ &= \frac{((c-4)2^c + 2^{N+2}) \frac{z_{\frac{a}{b}}^2 (\sum_{j=1}^b \{A(j)\}^2 - \frac{a^2}{b})}{\epsilon^2 a^2}}{4(2^N - 1)} \end{aligned}$$

こうすると、各カットオフレベルの計算量がバックトラック法による全探索を行う場合の計算量の何倍になるかが分かる。この式を最小にするような c は、精度と計算量の両面で最適なカットオフレベルといえることができる。

しかし、解の分布を仮定しない限り 未知のパラメータがあるため、この計算量比を具体的に求めることができない。

そこで、次節からの例題で解の分布を仮定してみる。なお、近似値の精度に関するパラメータは、 $\epsilon = 0.1, z_{\frac{a}{b}} = 1.96$ としている。

7 例題

7.1 <例題1> 解が一様に散らばっている場合(バランスしたツリー構造)

$$|\Omega| = 2^{100} \quad \text{配置の総数} \quad |S| = 724 \quad \text{解の総数}$$



● 解となる配置

○ 解ではない配置

解のばら撒き方

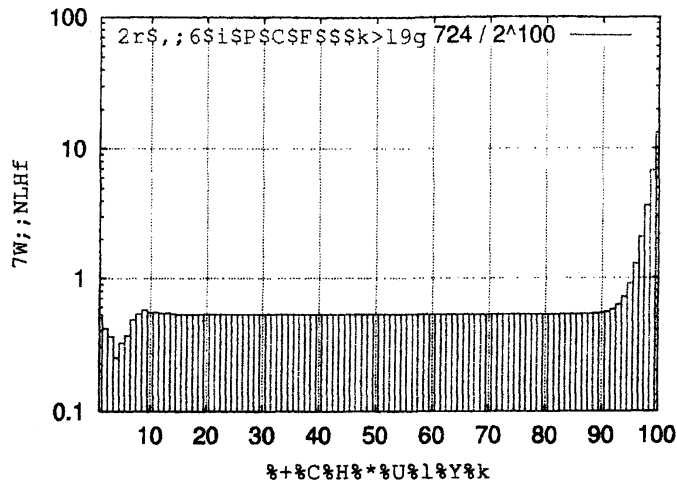
1. Y_1 を Ω の要素から等しい確率で選び、その配置を解とする。
2. Y_2 を $\Omega - \{Y_1\}$ の要素から等しい確率で選び、その配置を解とする。
3. 同様に解となる配置を決めていく。
4. Y_a を $\Omega - \{Y_1, Y_2, \dots, Y_{a-1}\}$ の要素から等しい確率で選び、その配置を解とする。

このように解をばら撒くことでどのような解の配置も等しい確率で起こる。

略証 a 個の 1 と $m-a$ 個の 0 からなる任意の配置が起こる確率は、

$$a! \cdot \frac{1}{m} \cdot \frac{1}{m-1} \cdots \frac{1}{m-a+1} = \frac{1}{\binom{m}{a}}$$

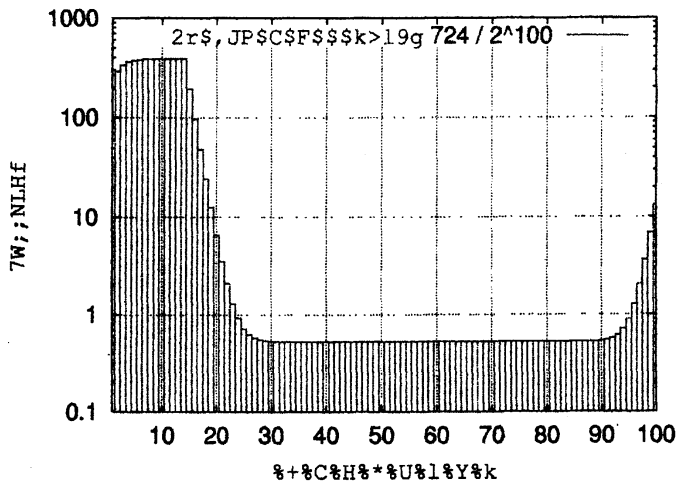
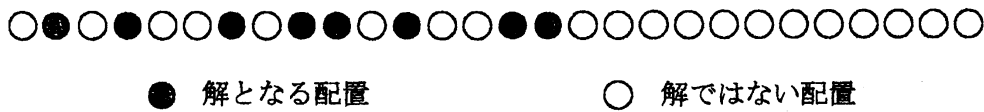
である。



この場合、カットオフレベルを深くしすぎなければ、どのカットオフレベルでも1以下となっているので全探索をする場合に比べ、有効な推定が行える。
 しかし、実際に解く問題ではこのように都合よく解が散らばってしてくれるとは限らない。

7.2 <例題2> 解が偏っている場合 (バランスしていないツリー構造)

$|\Omega| = 2^{100}$ 配置の総数 $|S| = 724$ 解の総数



このグラフから、解の分布がある程度偏っている場合においても、良いカットオフレベルを選べば、複合モンテカルロ法が有効に機能することがわかる。では、実際の計数問題に適用してみる。

7.3 <例題3> M-Queens Problem

M × M のチェス盤上に Queen を M 個置き、それらが互いに打ち合うことの無い配置の総数を求める問題。

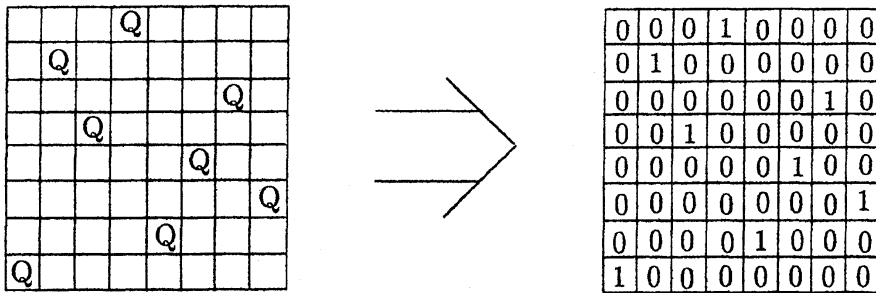


図 1: 8-Queens Problem の解の配置の例

このように配置を表現すると、配置の総数は

$$N = M \times M = 8 \times 8$$

$$|\Omega| = 2^N = 2^{64}$$

となる。

(注) M-Queens Problem の場合、どの行と列にも 1 個しか Queen を置けないことから、更に枝刈りが可能だがどのような計数問題に対しても適用できるようにしたいので、枝刈りなどの各問題固有の性質は考えないことにする。

10-Queens Problem

Queen の数が 10 のときの計算量比は以下のようになった。

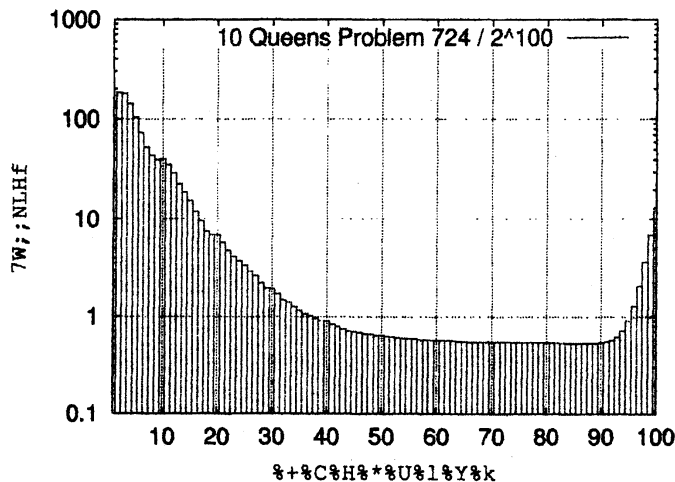


図 2: 10-Queens Problem $p = \frac{724}{2^{100}}$

以上の例題から、解の配置が偏っていると複合モンテカルロ法が有効に機能しないことが見てとれる。そこで、解の配置の偏り度を定義し、その偏り度が大きい値をとるような解の配置は非常に少ないことを次節で示す。

8 解の配置が偏る確率は少ないことの証明

まず、解の配置 Φ を以下のように表現する。

$$\begin{aligned}\Phi &= \{ \mathbf{y} = (y_1, y_2, \dots, y_m) \in \{0, 1\}^m \mid \sum_{i=1}^m y_i = a \} \\ \forall \mathbf{z} &= (z_1, z_2, \dots, z_m) \in \Phi \\ z_i &= \begin{cases} 0, & i \text{ 番目の配置が解でない} \\ 1, & i \text{ 番目の配置が解} \end{cases}\end{aligned}$$

ただし、 i 番目の配置とは $i-1$ の 2 進展開の配置に対応することにする。

ここで、計算の便宜上 $\frac{a}{m} := p < \frac{1}{2}$ を仮定しておく。もし、 $p > \frac{1}{2}$ であるような解の配置に対してこれから定義する偏り度を適用したいのであれば、0 と 1 を反転すればよい。

8.1 局所密度 $\delta_k(\mathbf{z})$ の定義

次式で局所密度を定義する。

$$\delta_k(\mathbf{z}) := \frac{1}{a} \sum_{i=k}^{k+a-1} z_i \quad (10)$$

ただし、 $1 \leq k \leq m$, $z_{m+i} = z_i$

(※) ここで、 $z_{m+i} = z_i$ というようにループ状にしたのは、

$$\sum_{k=1}^m \delta_k(\mathbf{z}) = a \quad (11)$$

としたいためである。

8.2 平均局所密度 $\bar{\delta}(\mathbf{z})$

この局所密度の平均を計算する。

$$\begin{aligned}\bar{\delta}(\mathbf{z}) &:= \frac{1}{m} \sum_{k=1}^m \delta_k(\mathbf{z}) \\ &= p, \quad (\because (11))\end{aligned}$$

8.3 $\bar{\delta}(\mathbf{z})$ まわりの 2 次モーメント $M_2(\mathbf{z})$

$$\begin{aligned}M_2(\mathbf{z}) &:= \frac{1}{m} \sum_{k=1}^m (\delta_k(\mathbf{z}) - \bar{\delta}(\mathbf{z}))^2 \\ &= \frac{1}{m} \sum_{k=1}^m (\delta_k(\mathbf{z}))^2 - p^2\end{aligned} \quad (12)$$

この式 (12) を散らばり具合を表す指標としたのでは、0 と 1 の間の値にならない。そこで、(12) の最大値を求め、その値で割って規格化することにする。

8.4 $M_2(\mathbf{z})$ の最大値 M_2^{max}

$M_2(\mathbf{z})$ を最大にするような解の配置は $11 \dots 100 \dots 0$ となる場合で、容易に計算される。ただし、一番最初に仮定したとおり、1の数のほうが0の数よりも少ないとする。

$$\begin{aligned} M_2^{max} &= \frac{1}{a^2 m} (a^2 + 2 \sum_{k=1}^{a-1} k^2) - p^2 \\ &= \frac{1}{a^2 m} (a^2 + 2 \frac{(a-1)a(2a-1)}{6}) - \frac{a^2}{m^2} \\ &= \frac{m(2a^2 + 1) - 3a^3}{3am^2} \end{aligned} \quad (13)$$

8.5 偏り度 $D(\mathbf{z})$ の定義

これらを用いて解の配置の偏り度を次式で定義する。

$$D(\mathbf{z}) := \frac{M_2(\mathbf{z})}{M_2^{max}} \quad (14)$$

どのような解の配置も等確率と仮定したとき、ここで定義した偏り度の期待値は、

$$E(D(\mathbf{z})) = \frac{3am(1-p)(m-a)}{(m-1)(m(2a^2+1) - 3a^3)} \quad (15)$$

と求まる。

8.6 マルコフの不等式

$\forall \alpha > 0$ に対して、

$$\text{Prob}\{D(\mathbf{z}) \geq \alpha\} \leq \frac{E(D(\mathbf{z}))}{\alpha} \quad (16)$$

がいえる。

(確認) もしレアイベント ($p \ll 1$) ならば、

$$E(D(\mathbf{z})) \approx \frac{3am^2}{2a^2m^2} = \frac{3}{2a}$$

となり、解の数 a が大きければ確かに0に近い値となり、解が偏った配置となる確率は小さいといえる。したがって、複合モンテカルロ法が有効に機能しないような解の配置は非常に少ないことが言えた。 \square

参考文献

- [1] M.Terada. Estimating Number of Solutions with Hybrid Monte Carlo Method. (1997)
- [2] 紺谷真紀子 計数問題における複合モンテカルロ法とその有効性 (2003)