

シミュレーションによる待ち行列モデルの最適化について

白川 浩

1. はじめに

たとえば生産組み立てライン、コンピュータシステム、情報通信ネットワーク等のように、世の中に存在するシステムには、待ち行列がネットワーク状に構成されたシステムとしてモデル化できるものが多い。これらのシステムの効率的な設計および運用を検討するときには、適当な待ち行列ネットワークモデルを考察することになる。しかし、一般に待ち行列ネットワークの性能評価を解析的に行なうことは困難であり、シミュレーションによりシステムの解析が行なわれることが多い。

原理的にシミュレーションによる解析はコストと時間さえかければ、所望の誤差水準内で行なうことができる。最も単純な方法としては、同様なシミュレーション・ランを独立に複数回行なうことによって、各種性能指標の推定および信頼区間を求めればよい。しかしシステムのパラメトリックな分析を行なう場合にはこのようなやり方は余りに多くの時間とコストがかかる。そのため多くの場合、シミュレーション結果を得るためになんらかの工夫を行なう。たとえば再生点を利用したシミュレーション[19]、各種の分散減少法[12]等が挙げられよう。

本稿ではそのような工夫のうち、とくに連続的パラメータのパラメトリック分析に有効な Infinitesimal Perturbation Analysis (IPA) について紹介する。IPA とは Y. C. Ho, X.R. Cao らによって提案され、現在急速に普及しつつあるパラメトリックなシミュレーション技法である [8, 9, 11, 15, 17, 25, 26]。IPA では、システム・パラメータに対する性能評価基準の感度分析を、その値の推定と同時に1回のランで行なってしまう。一般にこのような感度分析を単に力づくで行なうと、パラメータの数を N として合計 $N+1$ 回のランが必要となる!

またIPAで構成された感度情報(推定量)は、単純な推定量よりはるかに分散が少ない。このためIPAを利用すると、システム設計作業等を大幅に効率化できる。

以下、本稿では次のような項目について解説する。2. ではIPAの手続きについて、具体例をとおして説明する。3. ではIPAで得られた推定量の性質について検討する。4. では現在までに開発された各種のIPAアルゴリズムを概説する。5. ではIPAの改良形、ならびにそのほかの効率的な感度分析法について紹介する。最後に6. で本稿で解説した内容について要約する。

2. IPA とは

はじめに巡回型待ち行列ネットワークを例にとり、IPAについて説明しよう[9]。図1に示すような N 個の $GI/1/K$ 型待ち行列から構成される巡回型待ち行列ネットワークを想定する。なお各ステーションでの待ち室は有限であり、ステーション j でのサービスが終了してもステーション $j+1$ の待ち室が一杯であれば、そのジョブはステーション j のサーバーを占有し続ける (プロダクション型ブロッキング)。このときステーション1でのスループット(単位時間あたりのジョブ処理数) J を最大化するよう、各サーバーの平均サービス時間 $\bar{s}_j (1 \leq j \leq N)$ を設定するのが目的である。すなわち、 $T_{1n}(\bar{s})$ をこのシステムが適当な初期状態 (=再生点、存在を仮定) から出発してステーション1で n 個のジョブを終了するまでに要する時間とすれば、

$$\begin{aligned} \max J(\bar{s}) &= E \{ \lim_{n \rightarrow \infty} n / T_{1n}(\bar{s}) \} \\ (2.1) \quad &= \lim_{n \rightarrow \infty} n / T_{1n}(\bar{s}, \omega), \text{ w. p. 1.} \\ &= E \{ M_1(\bar{s}) \} / E \{ T_1(\bar{s}) \}. \end{aligned}$$



図1 巡回型待ち行列ネットワーク

しらかわ ひろし 東京工業大学 人文社会群

〒152 目黒区大岡山2-12-1

ただし $\bar{s}=(\bar{s}_1, \dots, \bar{s}_N)$ とし, $T_1(\bar{s})(M_1(\bar{s}))$ はこのシステムが再びその再生点に行き着くまでに要する時間(ステーション1でのジョブ処理数)を表わす. また全体のサービス処理能力には, 次のような資源制約があるものとする.

$$(2.2) \quad \sum_{i=1}^N n_i \bar{s}_i = C, \quad \bar{s}_i \geq 0, \quad 1 \leq j \leq N.$$

一般に $J(\bar{s})$ の値は各ステーションでのサービス時間分布, 待ち室容量等に複雑に依存しており, 陽にその形を示すことは困難である. したがって(2.1)を(2.2)のもとで最適化するには, $J(\bar{s})$ のシステム・パラメータに対する偏微係数 $\partial J(\bar{s})/\partial \bar{s}_j$ を何らかの方法で推定する必要がある.

最も単純な方法としては \bar{s}_j と $\bar{s}_j + \Delta \bar{s}_j$ の2つのパラメータ値に対しシミュレーションを行ない, その差分により偏微係数を推定する方法である. すなわち

$$\begin{aligned} \Delta J(\bar{s})/\Delta \bar{s}_j &= [J(\bar{s} + \Delta \bar{s}_j) - J(\bar{s})]/\Delta \bar{s}_j \\ (2.3) \quad &= \lim_{n \rightarrow \infty} n/T_{1n}(\bar{s} + \Delta \bar{s}_j, \omega) \\ &\quad - \lim_{n \rightarrow \infty} n/T_{1n}(\bar{s}, \omega), \quad \text{w. p. 1.} \end{aligned}$$

により, 偏微係数 $\partial J(\bar{s})/\partial \bar{s}_j$ を推定する. この場合

$$(2.4) \quad \partial J(\bar{s})/\partial \bar{s}_j = \lim_{\Delta \bar{s}_j \rightarrow 0} \Delta J(\bar{s})/\Delta \bar{s}_j,$$

が成立するので, 漸近的に推定量の一致性が満たされる. しかしこのような方法により感度分析を行なう場合, 合計 $N+1$ 回のランが必要となり効率的な推定法とはいえない.

そこで(2.1)から, 次のようなサンプルパスごとの偏微係数について考えよう.

$$\begin{aligned} \partial \hat{J}(\bar{s}, \omega)/\partial \bar{s}_j &= \lim_{n \rightarrow \infty} \partial [n/T_{1n}(\bar{s}, \omega)]/\partial \bar{s}_j \\ (2.5) \quad &= -\lim_{n \rightarrow \infty} [n/T_{1n}^2(\bar{s}, \omega)] \cdot [\partial T_{1n}(\bar{s}, \omega)/\partial \bar{s}_j]. \end{aligned}$$

ここで仮に

$$(2.6) \quad \partial J(\bar{s})/\partial \bar{s}_j = \partial \hat{J}(\bar{s}, \omega)/\partial \bar{s}_j, \quad \text{w. p. 1.}$$

が成立すれば, (2.5)によって推定される偏微係数をもとに(2.1)の最適化を考えることができる. (2.6)の成立条件は後に3.で検討することとして, 以降では(2.5)の計算法について述べる.

まず $T_{1n}(\bar{s}, \omega)$ が, どのように構成されているか考えてみよう. ステーション i での k 番目のジョブのサービス時間の実現値を, $S_{ik}(\bar{s}_i, \omega)$ (ただし平均サービス時間は \bar{s}_i) とする. このとき $T_{1n}(\bar{s}, \omega)$ は, あるサンプルパスの構成手続き ϕ_n によって

$$(2.7) \quad T_{1n}(\bar{s}, \omega) = \phi_n(\{S_{ik}(\bar{s}_i, \omega) : 1 \leq i \leq N, k \geq 1\})$$

と表わせ,

$$(2.8) \quad \begin{aligned} \partial T_{1n}(\bar{s}, \omega)/\partial \bar{s}_j &= \sum_{k \geq 1} [\partial \phi_n/\partial S_{jk}] \cdot [S_{jk}(\bar{s}_j, \omega)/\partial \bar{s}_j] \end{aligned}$$

となる. したがってシステム・パラメータ \bar{s} に対する \hat{J} の偏微係数を計算するには, 次の2点が問題となる.

i) システム・パラメータの変化が, どのようにサービス時間の実現値に影響を与えるか?

ii) サービス時間の実現値の変化が, 最終的にどの程度 $T_{1n}(\bar{s}, \omega)$ に影響を与えるか?

i) の $S_{jk}(\bar{s}_j, \omega)/\partial \bar{s}_j$ の計算法がパタベーション生成ルールであり, ii) の $\partial T_{1n}(\bar{s}, \omega)/\partial \bar{s}_j$ ($1 \leq m \leq n, 1 \leq j \leq N$) の帰納的な計算法がパタベーション伝達ルールと呼ばれる.

はじめにパタベーション生成ルールについて説明する[24, 26]. ステーション j でのサービス時間分布を, $F_j(x; \bar{s}_j)$ (ただし平均サービス時間は \bar{s}_j) とする. このとき S_{jk} は, $[0, 1]$ 上の一様分布にしたがう確率変数 U_{jk} によって

$$(2.9) \quad S_{jk} = F_j^{-1}(U_{jk}; \bar{s}_j),$$

と与えられる. ただし $F_j^{-1}(U; \bar{s}_j) = \inf\{x; F_j(x; \bar{s}_j) \geq U\}$ とする. よって

$$\begin{aligned} \partial S_{jk}/\partial \bar{s}_j &= \partial F_j^{-1}(U_{jk}; \bar{s}_j)/\partial \bar{s}_j \\ (2.10) \quad &= -[\partial F_j(S_{jk}; \bar{s}_j)/\partial \bar{s}_j]/[\partial F_j(S_{jk}; \bar{s}_j)/\partial S_{jk}], \end{aligned}$$

により, サービス時間の実現値に対するシステム・パラメータの影響を計算できる.

次にパタベーション伝達ルールについて考える. さきに述べた巡回型待ち行列ネットワークのサンプルパスの構成手続き ϕ_n は, 次のような3つの基本的な処理の繰り返しに帰着できる.

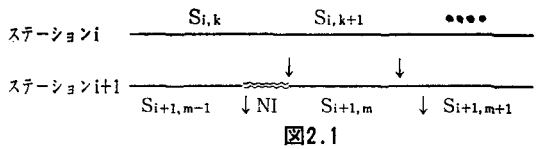


図 2.1

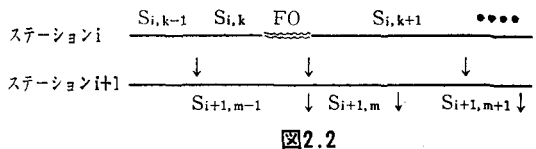


図 2.2

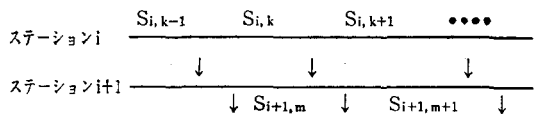


図 2.3

i)あるステーションでのサービス終了後に、次のステーションが空である場合(No Input;NI):そのジョブを次のステーションに移動させ、(待ち室にジョブがあれば)次のジョブのサービスを行なう。また同時に次のステーションでのサービスを開始させる。(図2.1)

ii)あるステーションでのサービス終了時に、次のステーションの待ち室が一杯の場合(Full Output;FO):次のステーションでサービス中のジョブがサービスを終了するまで、そのジョブをサーバーで待避させる。そしてサービス終了時に、直ちにジョブを移動させる。(図2.2)

iii)あるステーションでのサービス終了後に次のステーションの待ち室に空きがあり、かつそこでサービス中のジョブがある場合:そのジョブを次のステーションに移動させ(待ち室にジョブがあれば)次のジョブのサービスを行なう。(図2.3)

この基本処理において、サービス時間の実現値が微小変動した場合に、サンプルパス上の各ジョブのサービス開始時点がどのように変化し伝達されてゆくか考えてみよう。うまずi)の場合には、ステーション $i+1$ での $S_{i+1,m}$ 終了までの微小変動は、引き続きNIの存在によって消去されてしまう。そして $S_{i+1,m}$ のサービス開始以降の変動は、ステーション i での $S_{i,k}$ 終了までの微小変動を引き継ぐことになる。すなわち

$$(2.11) \quad \partial A_{i,k+1}/\partial \bar{s}_j \leftarrow \partial A_{i,k}/\partial \bar{s}_j \quad (j \neq i), \\ \partial A_{i,k+1}/\partial \bar{s}_i \leftarrow \partial A_{i,k}/\partial \bar{s}_i + \partial S_{i,k}/\partial \bar{s}_i,$$

および

$$(2.12) \quad \partial A_{i+1,m}/\partial \bar{s}_j \leftarrow \partial A_{i,k+1}/\partial \bar{s}_j.$$

ただし $A_{i,k}$ は、ステーション i での k 番目のジョブのサービス開始時点を表すものとする。次にii)の場合には、ステーション i での $S_{i,k}$ 終了までの微小変動が、引き続きFO存在によって消去されてしまう。そして $S_{i,k}$ でのサービス開始以降の変動は、ステーション $i+1$ での $S_{i+1,m}$ のサービス開始時点の微小変動を引き継ぐことになる。すなわち

$$(2.13) \quad \partial A_{i,k+1}/\partial \bar{s}_j \leftarrow \partial A_{i+1,m}/\partial \bar{s}_j.$$

最後にiii)の場合には、ステーション i での $S_{i,k}$ 終了までの微小変動が、そのまま $S_{i,k+1}$ のサービス開始時点の微小変動へと引き継がれる。よって(2.11)と同様な処理を考えればよい。

以上の手続きから明らかなように、パターン生成ルールおよびパターン伝達ルールを実行すれば、1回のランにより各システム・パラメータに対する偏微係数を推定できる。なおIPAのための追加的処理は

シミュレーション全体の手間から考えればそれほど問題にはならない量である。

3. IPAにもとづく推定量の性質

ここではIPAにもとづく偏微係数の推定量が、不偏性、一致性、最小分散性といった基本的な性質を有しているかどうか考察する。システム評価基準としては次のクラスを想定する。

$$(3.1) \quad J(\theta) = E\{L(\theta, \xi)\}.$$

ただし $\theta = (\theta_1, \dots, \theta_N)$ は対象とする確率システムの操作可能なパラメータ、 $\xi = (\xi_1, \dots, \xi_n) \in R^n (n \leq \infty)$ は再生点の間のシステムの不確実な事象を表す n 個の独立な確率変数の組である。なお2.で考察した評価関数 $J(\bar{s})$ はこのようなクラスの評価関数をもとに構成されており、一般にこのクラスの評価関数について考察しておけば十分といえる。

2.で示したように、パターン分析では J に対する偏微係数の推定量を、サンプルごとの偏微係数により構成する。ここで

$$(3.2) \quad \partial J(\theta)/\partial \theta_j = \partial E\{L(\theta, \xi)\}/\partial \theta_j \\ = E\{\partial L(\theta, \xi)/\partial \theta_j\},$$

が成立するものと仮定しよう。このとき m 再生点でのパターン平均

$$(3.3) \quad \partial \hat{J}_m(\theta)/\partial \theta_j = [\sum_{k=1, m} \partial L(\theta, \xi_k)/\partial \theta_j]/m,$$

は、 $\partial J(\theta)/\partial \theta_j$ の不偏推定量になっている。すなわち

$$(3.4) \quad E\{\partial \hat{J}_m(\theta)/\partial \theta_j\} = \partial J(\theta)/\partial \theta_j, m \geq 1.$$

さらに $\{\partial L(\theta, \xi_k)/\partial \theta_j : k \geq 1\}$ は独立で同一の分布にしたがう確率変数列 (*i. i. d.*) となるので、大数の弱法則から

$$(3.5) \quad \partial \hat{J}_m(\theta)/\partial \theta_j \xrightarrow{P} \partial J(\theta)/\partial \theta_j, m \rightarrow \infty.$$

が成立し、一致性も満たされる。

次に最小分散性について考える。 $L(\theta, \xi)$ は、任意の θ に対して ξ の各要素ごとに単調な関数であるとしよう。このときGal[10]らによって示された共通乱数(Common Random Number)の性質から、

$$\text{Var}\{[L(\theta + \Delta \theta_j, \xi) - L(\theta, \xi)]/\Delta \theta_j\} \\ (3.6) \leq \text{Var}\{[L(\theta + \Delta \theta_j, \eta) - L(\theta, \xi)]/\Delta \theta_j\} \\ \Delta \theta_j > 0,$$

が成立する。ただし $\eta = (\eta_1, \dots, \eta_n) \in R^n$ は、 ξ と同一の分布にしたがう任意の n 個の確率変数の組である。この結果

$$(3.7) \quad \text{Var}\{\partial L(\theta, \xi)/\partial \theta_j\} \\ = \text{Var}\{\lim_{\Delta \theta_j \rightarrow 0} [L(\theta + \Delta \theta_j, \xi) - L(\theta, \xi)]/\Delta \theta_j\}$$

$= \lim_{\Delta\theta_j \rightarrow 0} \text{Var}\{[L(\theta + \Delta\theta_j, \xi) - L(\theta, \xi)]/\Delta\theta_j\} < \infty$,
 が成立するならば, 任意の $\partial L(\theta)/\partial\theta_j$ に対する不偏推定量の中で, $\theta L(\theta, \xi)/\partial\theta_j$ にもとづく偏微係数の推定量が最小分散性を持つことがわかる[4].

以上により, パタベーションアナリシスによる感度情報の妥当性は, 仮定(3.2), (3.7)が成立するか否かにあるといえよう. 以降では(3.2), (3.7)の成立条件について検討する.

$L(\theta + \Delta\theta_j, \xi)$ を θ_j について1次までテーラー展開すると,

$$(3.8) \quad L(\theta + \Delta\theta_j, \xi) = L(\theta, \xi) + [\partial L(\theta, \xi)/\partial\theta_j] \cdot \Delta\theta_j + o(\theta, \Delta\theta_j, \xi)$$

と近似できる. ただし $L(\theta, \xi)$ は連続で, θ について2階微分可能と仮定する. よって

$$(3.9) \quad \lim_{\Delta\theta_j \rightarrow 0} E\{0(\theta, \Delta\theta_j, \xi)/\Delta\theta_j\} = 0,$$

$$(3.10) \quad \lim_{\Delta\theta_j \rightarrow 0} E\{0(\theta, \Delta\theta_j, \xi)/\Delta\theta_j^2\} = 0,$$

が, (3.2)および(3.7)が成立するための必要十分条件となる. ここで $0(\cdot)$ の定義より,

$$(3.11) \quad \lim_{\Delta\theta_j \rightarrow 0} 0(\theta, \Delta\theta_j, \xi)/\Delta\theta_j = 0, \quad \text{w. p. 1.}$$

よってルベグの有界収束定理から,

$$(3.12) \quad \exists K > 0, \exists \Delta\theta_j^0 > 0, \\ \text{s.t. } |0(\theta, \Delta\theta_j, \xi)/\Delta\theta_j| < K, \quad \text{w. p. 1,} \\ \text{for } \forall |\Delta\theta_j| \in (0, \Delta\theta_j^0),$$

が満たされれば(3.9)および(3.10)が成立する. Cao[1]はこの十分条件を満たす関数 $L(\theta, \xi)$ を, 一様微分可能と定義した. しかし(3.12)は単に(3.9) ((3.10))が成立するには強すぎる条件であり, 実際には次の条件を満たせば十分といえる.

$$\exists \text{ 確率変数 } K, \text{ s.t. } E\{K\} < \infty (E\{K^2\} < \infty), \exists \Delta\theta_j^0 > 0, \\ (3.12) \quad \text{s.t. } |0(\theta, \Delta\theta_j, \xi(\omega))/\Delta\theta_j| < K(\omega), \quad \text{w. p. 1,} \\ \text{for } \forall |\Delta\theta_j| \in (0, \Delta\theta_j^0).$$

一般に条件(3.13)を, ある確率システムについて検証することは容易ではない. 特に $\xi = (\xi_k)_{k \geq 1}$ のような無限点列によって L が構成されている場合, この点は大きな問題となる[6, 26]. なおたとえ $L(\theta, \xi)$ が連続でなくても, 不連続点のジャンプサイズが(3.9)のような条件を満たせば, やはり(3.2)が成立する[1, 23].

IPA アルゴリズムについて

IPAの適用法(アルゴリズム)は, 具体的にどのような確率システムを考えるかに大きく依存している. しかもこれらのアルゴリズムで構成した推定量が条件(3.2)を満たしているか否かは, 各モデルごとに考察する必要がある.

ある. このためIPAを適用するアルゴリズムは, 各待ち行列システムごとに開発されているのが現状である. そこで本節では, 現在までにどのようなタイプのシステムとその評価基準並びに操作可能なパラメーターに対し, IPAアルゴリズムが開発されているかを簡単に紹介する.

4.1 GI/GI/1 型待ち行列[26]

Suri and Zazanis[26]は, GI/GI/1 型待ち行列の平均系内滞在時間に対するIPAアルゴリズムを与えた. また到着分布およびサービス分布の操作パラメーターが満たすべき条件を示している. 特に M/G/1 型待ち行列においては, IPAによって得られた偏微係数の推定量が, 不偏性・一致性を満たすことが証明されている.

4.2 巡回型待ち行列ネットワーク[9]

Cao and Ho[9]は, ブロッキングを伴う巡回型行列ネットワークのスループットに対するIPAアルゴリズムを与えた. 操作パラメーターとしては, 各サーバーの平均サービス時間を想定している. なお, ある固定されたジョブ数 $n (< \infty)$ にもとづくスループット ($= E\{n/T_{in}(\delta)\}$), ただし T_{in} は2. で定義したとおり) については, IPAによって得られた偏微係数の推定量が不偏性を満たすことが証明されている. ただしこのようなスループットの定義では, $n/T_{in}(\delta, \omega)$ を得るための1回のランで再生点を利用することができず, 多少の問題が残ろう.

4.3 閉鎖型待ち行列ネットワーク [2, 3, 5, 6, 7, 8, 15, 17, 18]

IPAにおいて現在までに最もよく研究されたモデルが, 閉鎖型待ち行列ネットワークであろう. モデルとしては複数のジョブ種があり, 一般的な $\cdot/GI/1/K/FCFS$ 待ち行列から構成され, サーバーのブロッキングによる閉塞を認めた一般的なルーティング・ルールの待ち行列ネットワークまで研究されている[17]. またシステム評価基準としては, (ジョブ種ごとの) スループット[3, 15, 17], 各ステーションでの平均滞在時間[8]等が取り扱われている. さらに操作可能なパラメーターとしては, 平均サービス時間のほかルーティング確率[18]についても検討されている. このようにIPAアルゴリズムに関してはかなり一般的なものが開発されているといえよう.

一方, IPAアルゴリズムによって得られる推定量の理論的な性質については, 次の結果が得られている (ただしいずれの場合も, 平均サービス時間をパラメーターとする). まず1ジョブ種の $\cdot/M/1/K/FCFS$ 型待ち行列から成る閉鎖型ネットワークのスループットに関しては, た

とえブロッキングが存在しても一致性が満たされることが証明されている[2,7,15]。しかし複数ジョブ種の場合には、ある条件が成立しない限り一致性は保証されない[3,5]。また1ジョブ種で $\cdot/GI/1/\infty/FCFS$ 型待ち行列から成る閉鎖型ネットワーク(ブロッキングは生じない!)のスループットに関しては、ある固定されたジョブ数 $n(<\infty)$ にもとづくスループット(4.2と同じもの)について、IPAによる推定量が不偏性を満たすことが証明されている[6](なお[6]ではサービス時間分布が指数分布に限定されているが、この点に関しては容易に一般化できる)。

4.4 開放型待ち行列ネットワーク[8,19]

開放型の待ち行列ネットワークについては、なぜかIPAの研究成果はほとんど発表されていない。なおIPAアルゴリズムとしてはCao and Ho[8]によって、1ジョブ種の $\cdot/GI/1/K/FCFS$ 型待ち行列から成る開放型ネットワークの平均ネットワーク内滞在時間に関するものが与えられている。またCaoらと同様なアルゴリズムにもとづくシミュレーション実験の結果が、倉本[19]らによって報告されている。

4.5 出生死滅過程(M/M/m/k型待ち行列)[11]

M/M/m/k型待ち行列については解析的に評価でき、IPAを適用する必要はないと考えられよう。しかし理論的見地から、このモデルについて非常に興味深い結果が得られている。M/M/m/k型待ち行列について、その空き率 $\pi_0(=1-\text{利用率})$ をシステム評価基準としよう。このとき到着率 λ_j (ただし j は系内ジョブ数)についての感度 $\partial\pi_0/\partial\lambda_j$ を、単純にIPAを適用して推定すると不偏推定量が得られない。しかしGlasserman[11]は、元々のM/M/m/k型待ち行列をそれと等価な巡回型待ち行列ネットワークとみなしてIPAを適用すれば、この問題が回避できることを示した。これはIPAを適用する場合、どのような確率モデルの表現を想定するかによってIPAの有効性が左右されることを示している。

4.6 Discrete Event Dynamic System(DEDS,[25])

IPAの考え方は、待ち行列システムに限らずより一般的な離散型確率システムに適用できるはずである。この観点からSuri[25]は、IPAが適用可能となるDEDSとあるクラスの評価関数ならびにパタベーションを定義し、一般的なIPAアルゴリズムを与えている。

5. IPAの改良形ならびにその他の感度分析法

IPAでは、サンプルバスごとのシステム評価基準が、パラメータの変化に応じ連続的に摂動可能と想定している。しかしブロッキングを伴うモデルや複数ジョブ種を想定した場合、サンプルごとのシステム評価基準はパラメータ変化に対し不連続となる[3,8,13,15,16,17]。したがってこのようなモデルの感度情報を推定するには、この不連続な変化をなんらかの方法により考慮しなくてはならない。ここではこのような試みのうち、代表的な方法を3つ紹介する。

5.1 Finite Perturbation Analysis(FPA,[3,8,16])

FPA[16]では、IPAで想定されていたサンプルバスの微小変動が、ある有限の大きさをとるものとした場合のパタベーション伝達ルールを考える。ただしサンプルバスの摂動としてはあくまで局所的な1次近似を考え、1回のランをもとにして偏微係数を推定する。

具体的には、IPAで想定したNI(次のステーションが空の場合)、FO(次のステーションの待ち室が一杯の場合)に加え、PNI(Potential N.I.; 次のステーションのジョブ数が1)およびPFO(Potential F.O.; 次のステーションの待ち室が1つしか空いていない)も考慮して、サンプルが有限大の摂動をした場合のパタベーション伝達を行なっていく。なおパタベーション伝達ルールの詳細に関しては、[16]を参照されたい。

一般にサンプルごとのシステム評価基準がパラメータの微小変化に対して不連続に変化する場合、IPAによる偏微係数の推定値よりもFPAによるものの方が経験的によい近似値を与える[3,8,16]。しかしFPAによる推定量が一致性等の性質を有するか否かは、ごく限られたモデル[3]を除き未だ証明されていない。

5.2 Smoothed Infinitesimal Perturbation Analysis(SIPA,[13])

IPAではサンプルごとの性能評価基準 $L(\theta, \xi)$ の平均化をまったく考えずに、直接サンプルごとの偏微係数 $\partial L(\theta, \xi)/\partial\theta_j$ から性能評価基準 $J(\theta)$ についての偏微係数 $\partial J(\theta)/\partial\theta_j$ を推定した。しかし $L(\theta, \xi)$ が不連続に変化する場合、平均ジャンプサイズ如何によってはこの推定量では偏りが生じてしまう。そこでSIPA[13]では、 $L(\theta, \xi)$ がスムーズな連続関数となるまで条件付き期待値をとり、条件付き期待値ごとの偏微係数 $\partial E\{L(\theta, \xi)|\eta\}/\partial\theta_j$

から性能評価基準 $J(\theta)$ についての偏微係数 $\partial J(\theta)/\partial \theta_j$ を推定する。ただし $\eta = (\eta_1, \dots, \eta_m)$ はキャラクターゼーションと呼ばれ、 η についての条件付き期待値 $L(\theta, \eta) = E\{L(\theta, \xi) | \eta\}$ が計算でき、かつ $L(\theta, \eta)$ が θ について連続でかつ η のサンプルごとに微分可能となるようにうまく選ばれた不確実事象の部分集合である。すなわち

$$(5.1) \quad \begin{aligned} \partial J(\theta)/\partial \theta_j &= \partial E\{E\{L(\theta, \xi) | \eta\}\} / \partial \theta_j \\ &= E\{\partial E\{L(\theta, \xi) | \eta\} / \partial \theta_j\} \\ &= E\{\partial L(\theta, \eta) / \partial \theta_j\}, \end{aligned}$$

が成立するように、キャラクターゼーション η は選ばれる。このような η さえ構成できれば、あとは $L(\theta, \eta)$ を η のサンプルごとのシステム評価基準とみなして、通常のIPAを実行すればよいわけである。

Gong and Ho[13]には、いくつかの待ち行列システムについて、 η の選定法ならびにSIPAによる推定値の不偏性・一致性が成立する例が示されている。一般にこの η を見つけるのは、簡単なことではない。また η の構成法によっては、条件付き期待値 $L(\theta, \eta)$ の計算や η のサンプルの生成法が難しくなることがある。

5.3 Score Function Approach (SFA, [4, 21, 22])

IPAでは1回のランによって、各システム・パラメータについての感度情報を同時に得ることができた。このような1つのサンプルパスによる感度分析法は、じつはIPA以外にも様々な方法が考えられる。Rubinstein[21, 22]は同一のサンプルパスから異なるパラメータに対応するシステム評価基準を得るのに、サンプルパスではなく確率測度をパラメータの微小変動に応じ変化させることを提案している。すなわちパラメータ θ の下での $\xi \in R^n$ に対する確率密度関数を $f(\theta, \xi)$ とすると、

$$(5.2) \quad \begin{aligned} \partial J(\theta)/\partial \theta_j &= \partial \left[\int_{R^n} L(\xi) f(\theta, \xi) d\xi \right] / \partial \theta_j \\ &= \int_{R^n} [L(\xi) \partial f(\theta, \xi) / \partial \theta_j] d\xi \\ &= E\{L(\xi) \partial \ln(f(\theta, \xi)) / \partial \theta_j\}, \end{aligned}$$

が成立する。ただし $L(\cdot)$ は直接 θ には依存せず、 ξ によって評価されるものとする。また偏微分と積分は交換可能と仮定する。この結果シミュレーション時に $\partial \ln(f(\theta, \xi)) / \partial \theta_j$ を記録しておけば、IPAと同様、1回のランにより各システム・パラメータに対する感度情報が得られることになる。ここで偏微分と積分が交換可能となるためには、

$$(5.3) \quad \int_{R^n} [L(\xi) \partial^2 f^2(\theta, \xi) / \partial \theta_j^2] d\xi < \infty,$$

が成立すればよい。

一般に条件(5.3)は、IPAが不偏性を有するために要求

される条件(3.13)よりもはるかに弱い。これはIPAによる推定量に偏りがあるような場合でも、依然SFAでは不偏推定量となりうることを示している[4]。しかし不確実事象の次元数 n が大きくなるにつれて、SFAによる推定量の分散は発散する傾向にある。したがってIPAとSFAの選択は、考察する確率モデルに依存して決定されるものといえよう。

6. まとめ

本稿では、IPAを中心として1回のシミュレーション・ランにより各システムパラメータに対する偏微係数を推定する方法について解説した。条件(3.13)の下では、IPAによる推定量は、不偏性、一致性ならびに最小分散性といった推定量が持つべき性質を兼ね備えたものとなる。したがってこの条件が満たされる限り、IPAは非常に望ましい感度分析法といえよう。しかし複数ジョブ種もしくはブロッキングが存在する場合には、パラメータの微小変動に対し、必ずしもシステム評価基準が連続には変化しない。このような場合IPAによる推定量は偏りを持ち、IPAを利用してシステム設計等を行なうと誤った決定をしてしまう恐れがある。この対策としてFPA, SIPA, SFA といったいくつかの改良法が提案されている。

なお本稿ではIPAによっていかに偏微係数の推定量を構成するかを中心に説明し、求めた偏微係数を用いてどのように最適化を行なうかについては説明を割愛させていただいた。これらについては[19, 22]を参照されたい。

References

- [1] X. R. Cao(1985), "Convergence of Parameter Sensitivity Estimates in a Stochastic Experiment," IEEE Tr. Auto. Con., 30, 845-853.
- [2] X. R. Cao(1987), "Realization Probability in Closed Jackson Queueing Networks and Its Application," Ads. Appl. Prob., 19, 708-738.
- [3] X. R. Cao (1987), "First-Order Perturbation Analysis of a Simple Multi-Class Finite Source Queue," Perf. Eval., 7, 31-41.
- [4] X.R. Cao(1987), "Sensitivity Estimates Based on One Realization of a Stochastic System," J. Stat. Comp. Sim., 27, 211-232.
- [5] X. R. Cao(1988), "Realization Probability in Multi-Class Closed Queueing Networks," Euro. J. Ope. Res., 36, 393-401.

- [6] X. R. Cao (1988), "A Sample Performance Function of Closed Queueing Networks," *Ope. Res.*, 36, 128-136.
- [7] X.R. Cao(1989), "Calculation of Sensitivities and Realization Probabilities in Closed Queueing Networks with Finite Buffer Capacities," *Ads. Appl. Prob.*, 21, 181-206.
- [8] X. R. Cao and Y. C. Ho (1987), "Estimating the Sojourn Time Sensitivity in Queueing Networks Using Perturbation Analysis," *J.Opt. Th. Appl.*, 40, 559-582.
- [9] X.R. Cao and Y.C. Ho (1987). "Sensitivity Analysis and Optimization of Throughput in a Production Line with Blocking," *IEEE Tr. Auto Con.*, 32, 959-967.
- [10] S.Gal, R. Y. Rubinstein and A. Ziv (1984), "On the Optimality and Efficiency of Common Random Numbers," *Math. Com. Sim.*, 26, 502-512.
- [11] P. Glasserman (1988), "Infinitesimal Perturbation Analysis of a Birth and Death Process," *Ope. Res. Let.*, 7, 43-49.
- [12] P. W. Glynn and D. L. Iglehart (1988), "Simulation Methods for Queues:an Overview," *Que. Sys.*, 3, 221-256.
- [13] W. B. Gong and Y. C. Ho(1987), "Smoothed (Conditional) Perturbation Analysis of Discrete Event Dynamical Systems," *IEEE Tr. Auto. Con.*, 32, 858-866.
- [14] Y. C. Ho (1987), "Performance Evaluation and Perturbation Analysis of Discrete Event Dynamic Systems," *IEEE Tr. Auto. Con.*, 32., 563-572.
- [15] Y.C. Ho and X.R. Cao (1983), "Perturbation Analysis and Optimization of Queueing Networks," *J. Opt. Th. Appl.*, 40, 559-582.
- [16] Y.C. Ho, X.R. Cao and C.Cassandras(1983), "Infinitesimal and Finite Perturbation Analysis for Queueing Networks," *Automatica*, 19, 439-445.
- [17] Y. C. Ho, R. Suri, X. R. Cao, G. W. Diehl, J. W. Dille and M. Zazanis (1984), "Optimization of Large Multiclass (Non-Productform) Queueing Networks Using Perturbation Analysis," *Lar. Sca. Sys.*, 7, 165-180.
- [18] Y.C. Ho and X.R. Cao (1985), "Performance Sensitivity to Routing Changes in Queueing Networks and Flexible Manufacturing System Using Perturbation Analysis. *IEEE J. Rob. Auto.*, 1, 165-172.
- [19] 倉本剛, 森雅夫, 白川浩 (1989), "Perturbation Analysis を用いた待ち行列の最適化手法に関する実験," *Tokyo Institute of Technology, Tech. Rep. J-8.*
- [20] M. Meketon and P. Heidelberger (1982), "A Renewal Theoretic Approach to Bias Reduction in Regenerative Simulations," *Man. Sci.*, 28, 173-181.
- [21] R.Y. Rubinstein(1986), "The Score Function Approach for Sensitivity Analysis of Computer Simulation Models," *Math. Com. Sim.*, 28, 351-379.
- [22] R. Y. Rubinstein (1986), *Monte Carlo Optimization, Simulation and Sensitivity of Queueing Networks*, John Wiley & Sons.
- [23] R. Y. Rubinstein (1988), "Convergence of Perturbation Analysis Estimates for Discontinuous Sample Functions:A General Approach," *Ads. Appl. Prob.*, 20, 59-78.
- [24] R. Suri(1983), "Implementation of Sensitivity Calculations on a Monte Carlo Experiment," *J. Opt. Th. Appl.*, 625-630.
- [25] R. Suri (1987), "Infinitesimal Perturbation Analysis for General Discrete Event Systems," *J. A. C. M.*, 34, 686-717.
- [26] R. Suri and M. A. Zazanis(1988), "Perturbation Analysis Gives Consistent Sensitivity Estimates for the M/G/1 Queue," *Man. Sci.*, 34, 39-64.