

非定常ポアソン過程における隣接時間区間のデータを流用した 平均事象発生数の擬似推定法

01506240 (株)日立製作所 生産技術研究所 *船木謙一 FUNAKI Kenichi
 (株)日立製作所 生産技術研究所 的場秀彰 MATOBA Hideaki

1. 緒言

非定常ポアソン過程(NHPP)は、店への客の到着数や台風の発生数など多くの時系列確率過程のモデルとして用いられている。NHPPの性質は、時刻 t の関数として表される平均値関数 $\Lambda(t)$ によって一意に決定される。したがって、シミュレーションモデルなどにおいてNHPPを用いるためには、対象事象の発生データを時系列に並べた標本データ系列を観測し、 $\Lambda(t)$ の形を推定することが必要となる。 $\Lambda(t)$ を表す関数式の形(パラメータ構成)が分かっている場合には、尤度関数を導き、観測した標本データ系列に対して最大尤度を与えるパラメータ値を求めれば良い[1]。標本データ系列の表現方法には、事象の発生時刻を観測して時点列として表す方法と単位時間区間当たりの事象発生数を観測して件数列として表す方法があるが、どちらの場合にも尤度関数は簡単に導くことができ、推定は容易である。しかし、実際には $\Lambda(t)$ の形を予め特定することが困難で、尤度関数も導けないことも多い。そのような場合には、時系列上を適当な時間区間に区切り、複数の標本データ系列を観測して各時間区間における事象発生数の区間平均値を求める方法がある[2]。この方法は、標本データ系列が多数得られる場合には有効であるが、標本データ系列が少ない場合には推定精度が悪くなるという問題がある。

本発表では、 $\Lambda(t)$ の形が分からず、かつ標本データ系列が一つまたは少数しか観測できない場合に、上記の事象発生数の区間平均値をとる方法において、隣接時間区間のデータを流用した擬似推定値を用いて推定精度を向上する方法を提案する。

2. 隣接時間区間のデータを流用した推定とその有効性

2.1 用語、記号の定義

まず、時系列上のある時間区間中の平均事象発生件数を推定する場合を考える。推定の対象となる時間区間を推定区間、推定区間に隣接する時間区間を併せた区間を参照区間と呼ぶことにし(図1)、以下のように記号を定義する。

- N : 観測した標本データ系列数
- H_a : 推定区間中の真の平均事象発生件数
- H_b : 参照区間中の真の平均事象発生件数
- τ_a, τ_b : 推定区間、参照区間の長さ

X_a : 推定区間中の観測事象発生件数の標本平均
 X_b : 参照区間中の観測事象発生件数の標本平均
 Y_a : 擬似推定値($\equiv X_b \cdot \tau_a / \tau_b$)
 但し、観測事象発生件数の標本平均とは、各標本データ系列の当該時間区間中に観測した事象発生件数を合計して、標本データ系列数Nで割った値。

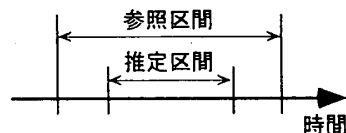


図1 推定区間と参照区間

2.2 擬似推定値の有効性

単純に標本データからの区間平均値を用いるならば、 H_a の推定値は X_a である。しかし、ここでは X_a の代わりに擬似推定値 Y_a を用いた方が統計的に精度が良くなる場合があることを示す。精度とは、真の平均値 H_a の周辺の一定領域に入る確率が高い方が良いという尺度で考える。

まず単純平均値 X_a について、任意の正定数 ϵ_a に対して、チェビシェフの不等式より、

$$Prob\{|X_a - H_a| \geq \epsilon_a\} \leq \frac{H_a}{N \cdot \epsilon_a^2} \quad \text{式1}$$

が成り立つ。これを变形して

$$Prob\{|X_a - H_a| < \epsilon_a\} \geq 1 - \frac{H_a}{N \cdot \epsilon_a^2} = \alpha \quad \text{式2}$$

とおくと

$$\epsilon_a = \sqrt{\frac{H_a}{(1-\alpha)N}} \quad \text{式3}$$

であるから、式2に戻して、 X_a が H_a の周辺で確率 α 以上で入る領域 D_x は

$$D_x = \{X_a | X_a \in (H_a - \sqrt{\frac{H_a}{(1-\alpha)N}}, H_a + \sqrt{\frac{H_a}{(1-\alpha)N}})\} \quad \text{式4}$$

と表される。次に、擬似推定値 Y_a についても同様にチェビシェフの不等式から、

$$Prob\{|Y_a - H_a| \geq \epsilon_a\} \leq \frac{1}{\epsilon_a^2} \left[\left\{ H_a - \frac{H_b \cdot \tau_a}{\tau_b} \right\}^2 + \frac{H_b \cdot \tau_a^2}{N \cdot \tau_b^2} \right] \quad \text{式5}$$

であり、これを变形して

$$Prob\{|Y_a - H_a| < \epsilon_a\} \geq 1 - \frac{1}{\epsilon_a^2} \left[\left\{ H_a - \frac{H_b \cdot \tau_a}{\tau_b} \right\}^2 + \frac{H_b \cdot \tau_a^2}{N \cdot \tau_b^2} \right] = \alpha \quad \text{式6}$$

とおくと、

$$\epsilon_a = \sqrt{\frac{\left\{ H_a - \frac{H_b \cdot \tau_a}{\tau_b} \right\}^2 + \frac{H_b \cdot \tau_a^2}{N \cdot \tau_b^2}}{1-\alpha}} \quad \text{式7}$$

であるから、 Y_a が H_a の周辺で確率 α 以上で入る領域 D_y は

$$D_y = \{Y_a | Y_a \in (H_a - \sqrt{\frac{\{H_a - \frac{H_b \cdot \tau_a}{\tau_b}\}^2 + \frac{H_b \cdot \tau_a^2}{N \cdot \tau_b^2}}{1 - \alpha}}, H_a + \sqrt{\frac{\{H_a - \frac{H_b \cdot \tau_a}{\tau_b}\}^2 + \frac{H_b \cdot \tau_a^2}{N \cdot \tau_b^2}}{1 - \alpha}})\} \quad \text{式 8}$$

と表される。 $D_y \subseteq D_x$ であれば推定値として Y_a を用いた方が真値 H_a により近い値を得る確率が高いといえるので、式 4 と式 8 より、

$$\sqrt{\frac{\{H_a - \frac{H_b \cdot \tau_a}{\tau_b}\}^2 + \frac{H_b \cdot \tau_a^2}{N \cdot \tau_b^2}}{1 - \alpha}} \leq \sqrt{\frac{H_a}{(1 - \alpha)N}} \quad \text{式 9}$$

を満たすことが、 Y_a による推定の方が X_a による推定よりも精度が良くなるための十分条件である。すなわち、ある推定区間と参照区間に対して、式 9 を満たす関係があるとき、単純区間平均値 X_a の代わりに Y_a を推定値とした方が良いといえる。

しかし、真の値 H_a 、 H_b は事前に知りえないので、式 9 の条件を吟味する際には、代わりにそれぞれ X_a および X_b を用いて計算することを提案する。

3. 推定アルゴリズムと実験結果

3.1 推定アルゴリズム

本発表による NHPP の推定では、時系列上を推定する時間区間(推定区間)に区切って、各推定区間ごとに平均事象発生件数を逐次推定していく。推定アルゴリズムを次のようにする。

```

For 全ての推定区間
  当該推定区間の単純区間平均値算出
  推定値 ← 単純区間平均値
  Do While 参照区間長さが限界値以下
    参照区間設定(逐次広げる)
    If 式 9 を満たす Then
      推定値 ← 擬似推定値
    Exit Do
  End If
Loop
Next
  
```

上記において、各推定区間に対する参照区間の設定には、参照区間長さの限界値を決めておき、推定区間の両側から限界値に達するまで逐次単位時間ずつ広げていく方法をとる。そして、各参照区間について式 9 の判定を行い、式 9 を満たしたら、そのときの参照区間による擬似推定値を採用する。

3.2 実験結果

提案する推定法の有効性を見るため、次の $\Lambda(t)$ を持つ NHPP を実験的に発生させ、その発生結果を標本データ系列と見立てて上記アルゴリズムを用いて推定した(図 2、表 1)。

$$\Lambda(t) = 6t^5 - 45t^4 + 130t^3 - 180t^2 + 210t + \frac{15}{2} - \frac{15}{2} \cos 6t + \frac{45}{8} \sin 8t$$

この $\Lambda(t)$ の形は不規則で、標本データから事前にその形を予想することが困難な場合をうまく表している。標本データには、0.05 時間ごとの事象発生件数を 1 系列観測したものをを用いた。また、推定区間を各 0.05 時間ごとにとり、参照区間長さの限界値は 0.15、0.25、0.35 時間の 3 通りを考えた。図 2 は、参照区間長さの限界値が 0.15 時間のときの単純区間平均値と提案した方法による推定値とを各推定区間ごとにプロットした例である。提案した方法による推定では、隣接時間区間の変動が式 9 を満たす範囲ならば、単純区間平均値よりも滑らかな値をとるように修正していることが分かる。これは、隣接時間区間のデータを似ていると判断して、これらの区間で平均化しているからである。表 1 は、各推定区間の真値との相対 2 乗誤差の全区間の平均値を、単純区間平均値を用いて推定した場合と上記アルゴリズムを用いて推定した場合とを比較したものである。但し、実験は 50 回繰り返し、表中の値はその平均を示している。この結果から、上記アルゴリズムによる推定値の方が全体として精度を上げていること、および参照区間長さの限界値を大きくとれば精度が上がるのが分かり、提案した方法の有効性が実験的に証明された。

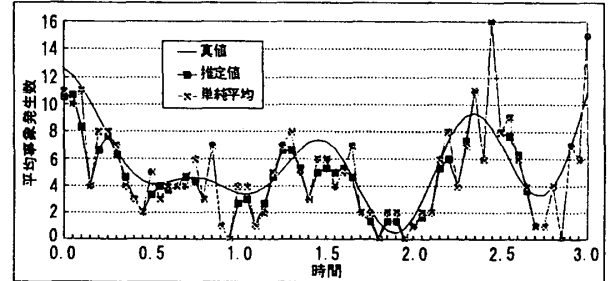


図 2 推定結果

表 1 各方法における誤差の比較

単純区間平均	提案した方法		
	限界値=0.15	0.25	0.35
0.22518	0.17909	0.17026	0.16711

4. 結言

本発表では、 $\Lambda(t)$ の形が分からず、かつ標本データ系列が一つまたは少数しか観測できない場合でも隣接時間区間のデータを流用して平均事象発生数の推定精度を向上する方法を提案した。また、その有効性を実験的に検証し、確かめることができた。

参考文献

- [1] I. Bar-David, "Communication under the Poisson regime," *IEEE Trans. Information Theory*, vol. IT-15, no. 1, pp. 31-37 (1969)
- [2] A. M. Law, et al, *Simulation modeling & analysis*, pp. 406-408, McGraw-Hill (1991)