

Hawkes Process と最適区間幅推定を用いた Web Data の解析

Web data analysis with hawkes process and histogram bin-width optimization

三宅 雅矩 *1
Masanori Miyake池上 高志 *1
Takashi Ikegami岡 瑞起 *2
Mizuki Oka橋本 康弘 *2
Yasuhiro Hashimoto

*1 東京大学大学院総合文化研究科

Graduate School of Arts and Science, The University of Tokyo

*2 筑波大学システム情報系

Department of Computer Science, University of Tsukuba

Many studies on statistical properties of hawkes process, which is the self-excitable extension of poisson point process, have been conducted these days, and there are a lot of applications to timeseries-analyses in fields such as seismology, neuroscience, financial trading, and web science. There are also theoretical studies about the relationship between the excitability of hawkes process and its burstiness, and the existence of burst transition has been proven with simulations. But there has been no demonstration of the burst transition using real data. We conducted timeseries-analyses of social web data with hawkes process, and we found there is a critical point concerning its excitability, and burst transitions occur around the point.

1. はじめに

自己励起型の点過程である Hawkes Process を用いた地震や神経細胞の活動, 金融取引パターンなどの解析および予測が従来より盛んに行われてきた [Filimonov 12]. 近年では非常に多数の要素が互いに影響を及ぼし合うような系の内部のダイナミクスを上手く捉えられるという特性を活かし, Social Networking Service(SNS) における人々の活動を Hawkes Process によってモデル化するなどといった研究も行われている [Oka 15].

一方で, 単純な Poisson 点過程とは全く異なる挙動を示すこの Hawkes Process の性質について, 理論的に解明しようとする試みもなされてきた. 特に, Hawkes Process におけるイベントの励起率について, ある値を境に Burst の生じない静的な定常状態から, Burst が頻繁に生じる非定常状態へと移行 (burst transition) するような critical point が存在することが理論的に示され, シミュレーションによる実験でもそれに従った結果が得られている [Onaga 14].

本研究では, SNS から得られた実データについて Hawkes Process を用いた fitting を行うとともに, データ時系列からヒストグラムを作成する上での最適区間幅推定を通して非定常状態の程度を測定した. これらの結果より, 実データにおいても Hawkes Process の励起率に関して critical point が存在し, その点の周辺を境に burst transition が生じることを実証することができた.

2. 解析手法と対象

2.1 Hawkes Process

Hawkes Process は正のフィードバックを考慮して Poisson Process を拡張した自己励起型の点過程である [Hawkes 71]. 1次元 Hawkes Process において, イベントの発生率 (強度関数) $\lambda(t)$ は過去に生じたイベントとの時間的關係によって変化し, 以下の式で表される.

$$\lambda(t) = \rho + \alpha \sum_{t_k < t} f(t - t_k) \quad (1)$$

ここで, ρ は常に保たれているベースラインの発生率であり, α はイベント発生率のジャンプ幅を示す.

連絡先: 三宅雅矩, miyake@sacral.c.u-tokyo.ac.jp

指数関数カーネル $f(x) = e^{-\beta x}$ を用いると,

$$\lambda(t) = \rho + \alpha \sum_{t_k < t} e^{-\beta(t-t_k)} \quad (2)$$

のように表現され, この形の Hawkes Process は最尤推定における利点などから多くの実データの解析において用いられている.

上記の指数関数カーネルを用いた Hawkes Process の内部でのダイナミクスを定量化する指標として, 以下に示される branching ratio が頻繁に用いられる.

$$\text{branching ratio} = \int_0^{\infty} \alpha e^{-\beta t} dt = \frac{\alpha}{\beta} \quad (3)$$

この branching ratio はある一つのイベントから内部的に生成される別のイベントの平均個数に対応し, この値が 1 より大きくなると非定常状態になることは, 従来より様々なデータを用いて確認されている.

2.2 最適区間幅推定

神経細胞のスパイク時系列などといったデータに関して, そのイベントの生成頻度 (レート) を推定するため, 適当な時間幅をもつ区間に分割した時系列内でのイベント生成率からヒストグラム (Peristimulus Time Histogram, PSTH) を作成することがよくある. この時間幅に関しては, 平均二乗誤差最小化の観点から最適区間幅を決定する手法が考案されており, 以下の手順に従って求められる [Shimazaki 07].

1. n 回の試行より得られたイベント時系列について, その観測期間 T を幅 Δ の N 個の bin に区切り, i 番目の bin に入るイベントの数 k_i を求める.
2. イベント数の平均および分散

$$\bar{k} = \frac{1}{N} \sum_{i=1}^N k_i, \quad v = \frac{1}{N} \sum_{i=1}^N (k_i - \bar{k})^2 \quad (4)$$

を計算する.

おける各イベントの間隔

$$\delta\tau_i = \tau_i - \tau_{i-1} \quad (13)$$

によって構成される時間列のパーセンタイル値は指数分布のそれと一致する。

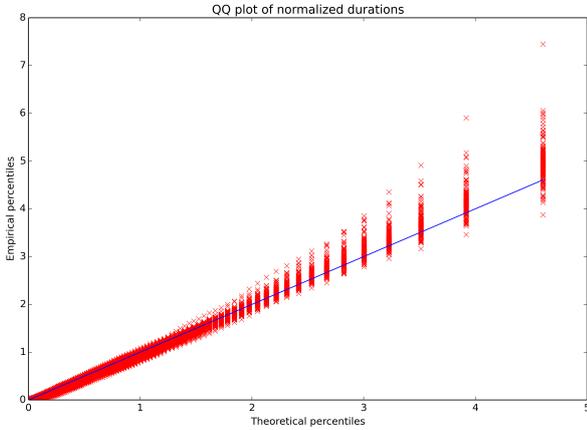


図 3: QQ-plot. 横軸は理論上のパーセンタイル, 縦軸は実データを用いて得られた各 tag についてのパーセンタイル値である。

これを各 tag についてプロットしたものが図 3 に示す QQ-plot である。fitting の当てはまりが良ければ, 理論値と実験によって得られた結果の値が一致するので, $y = x$ の直線に沿ってプロットされる。この図から, 今回の fitting が概ね良好であることが確認された。

3.2 branching ratio と最適区間幅との関係

先述の PSTH における最適区間幅の推定手法を各 tag の投稿日時の時系列データに対して適用し, 求めた最適区間幅の逆数と Hawkes Process による fitting の結果得られた branching ratio(先述の α_c に対応)との関係を, 図 4 に示す。

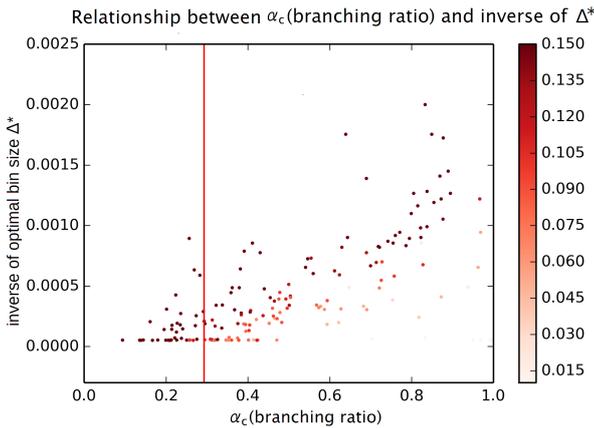


図 4: α_c (branching ratio) と最適な bin size Δ^* の逆数との関係。各点は各 tag の時系列データから得られた Hawkes Process における branching ratio (α_c に対応する) および最適区間幅の関係を示している。また, 各点の色は Hawkes Process における β の値に対応する。赤線は $\alpha_c \approx 0.2929$ を示している。

最適区間幅が無限大に近い(最適区間幅の逆数が 0 に近い)部分では, データ時系列は Burst が生じていない定常状態で

あり, 最適区間幅が小さくなる(逆数が大きくなるにつれ)につれ動的な非定常状態となるが, 実験の結果より branching ratio が大きくなるにつれ最適区間幅の逆数の値も上昇する傾向にあることが確認された。

そして, 若干のばらつきはあるものの, branching ratio の値が 0.2~0.4 の部分で定常状態と非定常状態との transition が生じており, 0.4 を超えると最適区間幅は有限の値となり, 必ず Burst が観測されるような状態となることが示された。これは, 実データでも理論値 $\alpha_c \approx 0.2929$ の付近で transition が生じることを実証するものであると考えられる。

一方で, $\alpha_c < 0.2929$ であるにも関わらず最適区間幅が小さな値を取るようなデータも存在しており, このようなゆらぎの原因については今後細かく検証していく必要がある。

4. 議論

今回の実データを用いた Hawkes Process による fitting およびその結果と Burst 状態との関係性の検証より, 理論上およびシミュレーションによって確認された critical excitability 付近において burst transition が実際に生じていることが確認された。この結果は現実のデータの統計的性質を新たに解明するものであると同時に, 今回 Web データを扱ったことによる, Web における人々のインタラクションのダイナミクスの解析や Web 広告への応用にも繋がる可能性があると考えられる。

一方で, branching ratio と最適区間幅との関係においてゆらぎが大きいなどといった理論とは異なる結果も得られたことについて, 2 週間ごとの総投稿数に応じた rescaling が適当であるか等を含め, 実データを扱う上での処理の方法についても一度検討する必要があると考えられる。また, 社会活動の昼夜差に起因するイベント発生率のゆらぎ等については今回の rescaling では捉えきれない可能性が高い。このような複数の要因を考慮に入れたモデルについての検証も引き続き進めたい。

最後に, 今回のモデルでは時間の経過および Web サービスのユーザーの増加に伴う総投稿数の増加を外的な要因として捉え, その影響を取り除くための rescaling を行ったが, 時間の経過およびユーザーの増加に伴って Web サービスの内的構造が変化している可能性も考えられる。こうした Web の進化に関する考察についても議論していきたい。

謝辞

貴重な解析用データを提供いただいた, Tunnel 株式会社の皆様に感謝いたします。

参考文献

- [Daley 07] Daley, D. J., and Vere-Jones, D.: An introduction to the theory of point processes: volume II: general theory and structure. Springer Science and Business Media (2007)
- [Filimonov 12] Filimonov, V., and Sornette, D.: Quantifying reflexivity in financial markets: Toward a prediction of flash crashes. Physical Review E, 85(5), 056108 (2012)
- [Hawkes 71] Hawkes, A. G.: Point spectra of some mutually exciting point processes. Journal of the Royal

Statistical Society. Series B (Methodological), 438-443
(1971)

[Ogata 81] Ogata, Y.: ‘ On Lewis ’ simulation method for
point processes ’, IEEE Transactions on Information
Theory 27(1), 2331 (1981)

[Oka 15] Oka, M., Hashimoto, Y., and Ikegami, T.: Open-
ended evolution in a web system. In Late breaking pa-
pers at the European Conference on Artificial Life (p.
17) (2015)

[Onaga 14] Onaga, T., and Shinomoto, S.: Bursting tran-
sition in a linear self-exciting point process. Physical
Review E, 89(4), 042817 (2014)

[Shimazaki 07] Shimazaki, H., and Shinomoto, S.: A
method for selecting the bin size of a time histogram.
Neural computation, 19(6), 1503-1527 (2007)

[Taylor 16] Taylor, T., Bedau, M., Channon, A., Ackley,
D., Banzhaf, W., Beslon, G., ... and McMullin, B.:
Open-ended evolution: perspectives from the OEE
workshop in York. Artificial Life (2016)