# ブログ空間における異常訪問行動の分析

山 本 和 紀<sup>†1</sup> 小 野 景 子<sup>†2</sup> 熊 野 雅 仁<sup>†2</sup> 木 村 昌 弘<sup>†2</sup>

本論文では,プログ空間における異常行動について調べる.まず,異常行動に関して,アクセス頻度異常」と「趣向比率異常」の 2 種類の概念を定義する.そして,実際のプログデータを用いて,これらの異常行動を特定し,それらの性質を分析する.特に,趣向比率異常のユーザと人気ブログを発見するイノベータには相関関係があることを示す.

# Analyzing anomalous behavior in Blogosphere

KAZUNORI YAMAMOTO,<sup>†1</sup> KEIKO ONO,<sup>†2</sup>
MASAHITO KUMANO<sup>†2</sup> and MASAHIRO KIMURA<sup>†2</sup>

We investigate anomalous behaviors in Blogosphere. First, we define two concepts for anomalous behavior, the access frequency anomaly and the share anomaly. Next, using real blog data, we identify these two anomalous behaviors, and analyze their properties. We show in particular that there is a correlation between the share anomaly users and the innovators of finding popular blogs.

# 1. はじめに

インターネットの普及により,手軽な情報発信ツールとしてプログが普及しており,ユーザのアクセスに着目したプログユーザの行動予測に関する研究が数多く成されている.例えば,アクセスの確率分布と Web ユーザビリティの関係を用いた研究<sup>1)</sup>,プログへのアクセ

ス数と社会現象に関する研究 $^{2)}$ ,ブログのアクセス履歴を可視化により分析する研究 $^{3)}$ などが報告されている。我々は,アクセスデータを用いて,イノベータ理論に基づいたブログユーザの行動予測の研究 $^{4)}$ ,急増する訪問行動の予測の研究 $^{5)}$ を行ってきた.

ところで、推薦システム、e コマース、およびサーバ負荷分散などにおけるコアテクノロジーの一つとして、プログユーザの大きな行動変化を予測する有効な手法の構築が挙げられる。したがって、本論文では、プログ空間における異常行動について調べる。まず、異常行動に関して、アクセス頻度異常」と「趣向比率異常」の 2 種類の概念を定義する。そして、実際のプログデータを用いて、これらの異常行動を特定し、それらの性質を分析する。特に、趣向比率異常のユーザと人気プログを発見するイノベータの相関関係や、現在において趣向比率異常を行うユーザが将来においても趣向比率異常を行う確率と、現在においてアクセス頻度異常を行うユーザが将来においてもアクセス頻度異常を行う確率との違いなどを調べる。

# 2. 分析データ

# 2.1 Doblog データ

本研究では  $Doblog^{*1}$ データセット $^{*2}$ のアクセスデータを使用する.アクセスデータは,データ ID,visitor ID,owner ID,アクセス時間で構成されるデータであり,2003 年 10 月から 2005 年 6 月までのデータが存在する.Doblog のデータセットを表 1 に示す.

表 1 Doblog のデータセット Table 1 Dataset of Doblog

項目	件数
記事数	1,540,077
ユーザ数	$53,\!525$
コメント数	$2,\!220,\!727$
トラックバック数	133,177
アクセス数	$12,\!542,\!581$

<sup>†1</sup> 龍谷大学 大学院 理工学研究科 電子情報学専攻 Division of Electronics and Informatics, Ryukoku University

<sup>†2</sup> 龍谷大学 理工学部 電子情報学科
Department of Electronics and Informatics, Ryukoku University

<sup>\*1 (</sup>株)NTT データ. http://www.doblog.com/

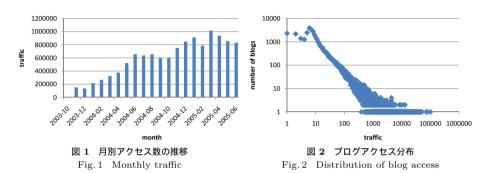
 $<sup>\</sup>star 2$  (株) ホットリンクと (株) NTT データの共同事業契約に基づき、(株) ホットリンクより提供。2003 年 10 月から 2005 年 6 月のデータを利用 .

IPSJ SIG Technical Report

# 2.2 訪問行動の統計解析

## 2.2.1 月別アクセス数の推移

図 1 は, Doblog における月別のアクセス数を示したグラフである.徐々にアクセス数が 増加し, Doblog が活発になっていることが分かる.



# 2.2.2 ブログアクセス分布

図 2 は , 全期間中 ( 2003 年 10 月から 2005 年 6 月まで ) に , ブログがアクセスされた回数の分布である . 一般的なロングテールの分布であることが分かる .

# 2.3 対象ユーザ

2004 年 1 月 1 日から 2004 年 12 月 31 日までのアクセスデータをもとに足切りを行い,対象となるユーザを同定した.各ユーザについて,ユーザが任意のプログにアクセスした回数を算出し,平均値以上の値を持つユーザを対象ユーザとした.なお,期間中に任意のプログに 1 回でもアクセスを行ったユーザは 21,628 人,平均アクセス回数 105.21 で足切りをした結果,対象ユーザは 2.884 人となった.

#### 2.4 対象ブログ

2004 年 1 月 1 日から 2004 年 12 月 31 日までのアクセスデータをもとに足切りを行い,対象となるプログを同定した.各プログについて,任意のユーザにアクセスされた回数を算出し,平均値以上の値を持つプログを対象プログとした.なお,期間中に任意のユーザから 1 回でもアクセスされたプログは 25,193 個,アクセスされた回数の平均値 95.56 で足切りをした結果,対象プログは 3994 個となった.

## 3. 異常訪問行動の分析法の提案

これまでの研究により、ブログへのアクセス数をもとにブログへのアクセスを増加させる ユーザの存在について明らかにした.しかしながら、どの程度のアクセス数の増加が起こる のか、その持続性については明らかでない。

図3は,ユーザの1日のプログ訪問回数によって色を変え,プロットしたグラフである. ユーザの1日のプログ訪問回数が少ないほど白,多いほど黒で表現している.図3より, ユーザがプログを訪問する頻度は一様ではないことが分かる.これより,毎日決まったプロ グにしか訪問しない「通常行動ユーザ」に対し,突然アクセス回数が急増する「異常行動 ユーザ」が存在することが分かる.

本研究では,このユーザに着目し,アクセス数の急増を異常行動と捉え定式化を行い,4 つの観点から異常行動の分析を行う.

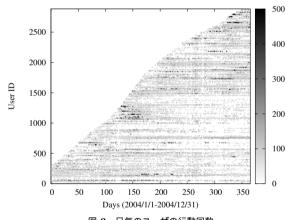


図 3 日毎のユーザの行動回数 Fig. 3 Daily user's activity

# 3.1 異常行動の定義

プログユーザの異常行動を「アクセス頻度」と「趣向比率」の 2 つの観点から分析を行う「アクセス頻度」は、ユーザ $u_n$  のプログ $b_m$  へのアクセス回数に着目し,前半期間におけるユーザ $u_n$  のプログ $b_m$  へのアクセス回数  $f_{n,m}(t-1)$  と,後半期間におけるユーザ $u_n$ 

IPSJ SIG Technical Report

のブログ  $b_m$  へのアクセス回数  $f_{n,m}(t)$  から,以下の式 1 より異常度  $a_{n,m}(t)$  を算出する.

$$a_{n,m}(t) = \frac{f_{n,m}(t) - f_{n,m}(t-1)}{f_{n,m}(t-1)} \tag{1}$$

異常度  $a_{n,m}(t)$  が 1 を超えている事例を異常な行動事例とする、また、異常な行動を行っ たユーザを異常行動ユーザと定義する「アクセス頻度」における異常行動ユーザは、ブロ グ bm への活動量が急増したユーザである.

「趣向比率」については,アクセス回数  $f_{n,m}(t)$  の代わりに,以下の値を用いる.

$$f'_{n,m}(t) = \frac{f_{n,m}(t)}{\sum_{m} f_{n,m}(t)}$$
 (2)

 $f'_{n,m}(t)$  を「趣向比率」と考え,その値を 式 1 に代入し,異常度  $a_{n,m}(t)$  を求める.異常 度  $a_{n,m}(t)$  が 1 を超えている事例を異常な行動事例とし,異常な行動を行ったユーザを異常 行動ユーザと定義する「趣向比率」における異常行動ユーザは、訪問するブログの傾向が 変化したユーザである.

- 3.2 異常訪問行動の分析方法
- 3.2.1 指標1: 行動異常度分布

ある期間における行動の異常度を「アクセス頻度」と「趣向比率」の観点からどの程度存 在するのか,分布を示す、

3.2.2 指標2:異常行動ユーザの恒常性

異常行動ユーザの恒常性を分析を行う、ある期間で異常行動したユーザのうち、次の期間 にも異常な行動をするようなユーザ数とその割合を調べことで、アクセス頻度に関する異常 行動ユーザと趣向比率に関する異常行動ユーザの恒常性の違いを比較する.

#### 3.2.3 指標3:ブログの被異常度

ユーザが異常行動したブログについて,ブログ視点で分析を行う.以下の式3で,ブログ に訪問した異常行動ユーザの異常度の合計をブログ毎に求め,その値をブログの被異常度  $A_m$  とする.

$$A_m(t) = \sum_{n} a(t) \{ a(t); a_{n,m}(t) > 1 \}$$
(3)

プログ $b_m$  に関する被異常度 $A_m$  のすべての期間における平均値と分散を求める.

3.2.4 指標4:ブログに訪問した異常行動ユーザ数に対する全体訪問ユーザ数 ブログ  $b_m$  に訪問した異常行動ユーザ数に対して,それ以降の期間に訪問した全ユーザ数 を求める、異常行動ユーザ数と全ユーザ数に相関があれば、異常行動ユーザが相関のあった 期間に対し影響を及ぼしていることがいえ、異常行動ユーザに着目することで、将来のブロ グへのアクセスに関して予測が可能となる.

# 4. 分析結果

本研究では,図4のように,3か月の期間(period)を5つ設定し,前半期間・後半期間 の組み合わせ (term)を 4 つとして,分析を行った.

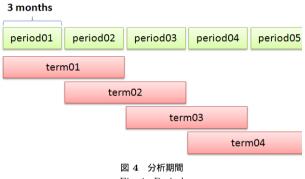


Fig. 4 Period

#### 4.1 指標1:行動異常度分布

図 5 にアクセス頻度,図 6 に趣向比率について,ともに term02 の期間における異常度の 分布を示す. 横軸に異常度の下限値,縦軸にアクセス数を示す. 横軸 10 における値は異常 度が 10 以上の場合のアクセス数を表す . ユーザ $u_n$  がブログ $b_m$  への行動した事例ごとに異 常度  $a_{n,m}(t)$  を算出し,異常度  $a_{n,m}(t) > x$  となる事例がいくつ存在するかを示している. アクセス頻度 (図 5) について,  $a_{n,m}(t) > 0$  であった事例 (アクセスが増加した事例) は 6572 件存在し,そのうち, $a_{n,m}(t) > 1$  であった事例(異常行動事例)は3055 件存在した. また,趣向比率(6)について, $a_{n,m}(t) > 0$ であった事例(アクセスが増加した事例)が 7280 件存在し,そのうち, $a_{n,m}(t) > 1$  であった事例(異常行動事例)は 2860 件であった. アクセス頻度,趣向頻度,双方について,異常度 $a_{n,m}(t) > 10$ を超えるような異常行動 事例が約100件存在していることが分かる.

IPSJ SIG Technical Report

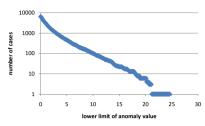


図 5 アクセス頻度に関する, 行動異常度分布 (term2)

Fig. 5 Distribution of anomaly action(access frequency)

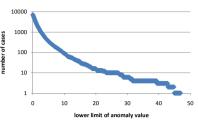


図 6 趣向比率に関する,行動異常度分布 (term2)

Fig. 6 Distribution of anomaly action(share)

## 4.2 指標2:異常行動ユーザの恒常性

表 2 は , ある期間において異常行動したユーザが , 次の期間以降にも異常行動するユーザ数とその割合を算出した表である .

アクセス頻度,趣向比率ともに,どの期間においても,次の期間にも異常行動するユーザの割合は約6割以上であることが分かった.

アクセス頻度において異常行動したユーザよりも,趣向比率において異常行動したユーザのほうが,異常行動を継続する割合が高いといえる.

表 2 異常行動ユーザの恒常性

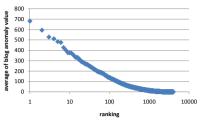
Table 2 Anomaly user's constancy

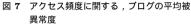
	term01	term02	term03	term04
term01(Access frequency)	626(100%)	380(60.7%)	263(42.0%)	200(31.9%)
term01(Share)	664(100%)	478(72.0%)	346(52.1%)	270(40.7%)
term02(Access frequency)	-	1118(100%)	669(59.8%)	482(43.1%)
term02(Share)	-	1291(100%)	892(69.1%)	628(48.6%)
term03(Access frequency)	-	-	1183(100%)	767(64.8%)
term03(Share)	-	-	1517(100%)	1053(69.4%)

### 4.3 指標3:ブログの被異常度に関する結果

図 7 はアクセス頻度に関する異常行動ユーザに着目 , 図 8 は趣向比率に関する異常行動ユーザに着目した時の結果である.ブログ  $b_m$  の被異常度  $A_m(t)$  を term01 から term04 の期間で算出し , その平均値が高いブログから順に並べたグラフである.

図 9 はアクセス頻度に関する異常行動ユーザに着目,図 10 は趣向比率に関する異常行動





異常度 Fig. 7 Average of anomaly values(access



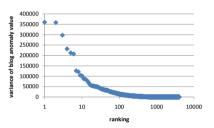
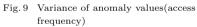


図 9 アクセス頻度に関する, プログの被異常 度の分散



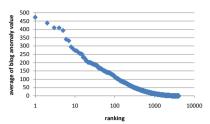


図 8 趣向比率に関する, ブログの平均被異常度

Fig. 8 Average of anomaly values(share)

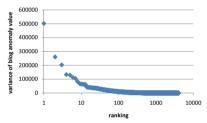


図 10 趣向比率に関する, プログの被異常度の 分散

Fig. 10 Variance of anomaly values(share)

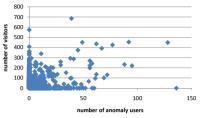
ユーザに着目した時の結果である.ブログ  $b_m$  の被異常度  $A_m(t)$  を term01 から term04 の期間で算出し,その分散が大きいブログから順に並べたグラフである.

アクセス頻度,趣向比率の双方において,分散値が高いブログが10個程度存在している. これより,期間によって被異常度が急激に変化するようなブログが少ないことがわかった. 多くのブログは異常行動され続ける傾向にあることが示唆された.

4.4 指標4:プログに訪問した異常行動ユーザ数に対する全体訪問ユーザ数 term1 においてプログに訪問した異常ユーザが他の期間におけるアクセス数に及ぼす影響を考察する.

図 11,図 12 は,term01 においてブログ  $b_m$  に訪問した異常行動ユーザ数と,period01 においてブログ  $b_m$  に訪問した全ユーザ数との関係を表している.図 11 はアクセス頻度,図 12 は趣向比率に関するグラフであるが,どちらも相関がない.

IPSJ SIG Technical Report



800 700 600 5 500 0 20 40 60 80 100 120 number of anomaly users

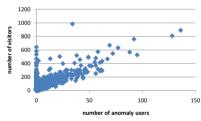
図 11 アクセス頻度に関する , プログに訪問した 異常行動ユーザ数に対する全体訪問ユー ザ数 (period01)

図 12 趣向比率に関する, ブログに訪問した異常行動ユーザ数に対する全体訪問ユーザ数(period01)

Fig. 11 Relation between the number of anomaly users and the number of visitors in blog(access frequency)(period01)

Fig. 12 Relation between the number of anomaly users and the number of visitors in blog(share)(period01)

図 13,図 14 は,term01 においてプログ  $b_m$  に訪問した異常行動ユーザ数と,period02 においてプログ  $b_m$  に訪問した全ユーザ数との関係を表している.term01 においてプログ  $b_m$  に訪問した異常行動ユーザ数と,period02 においてプログ  $b_m$  に訪問するユーザ数には相関がある.



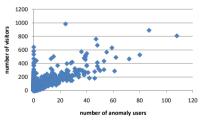
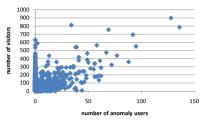


図 13 アクセス頻度に関する結果 (period02) Fig. 13 Access frequency(period02)

図 14 趣向比率に関する結果 (period02) Fig. 14 Share(period02)

図 15, 図 16, 図 17, 図 18 は , term01 においてプログ  $b_m$  に訪問した異常行動ユーザ数と , period03 または period4 においてプログ  $b_m$  に訪問した全ユーザ数との関係を表している . term01 においてプログ  $b_m$  に訪問した異常行動ユーザ数と , period03 においてプログ  $b_m$  に訪問するユーザ数には中程度の相関があり , period4 では相関がないことが分かる . 異常行動ユーザ数とプログに訪問するユーザ数には相関関係があり , 相関は期間によって



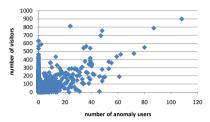
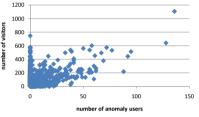


図 15 アクセス頻度に関する結果 (period03) Fig. 15 Access frequency(period03)

図 16 趣向比率に関する結果 (period03) Fig. 16 Share (period03)



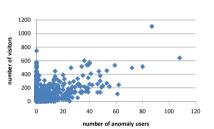


図 17 アクセス頻度に関する結果 (period04) Fig. 17 Access frequency(period04)

図 18 趣向比率に関する結果 (period04) Fig. 18 Share (period04)

異なることが分かった.term01 においてブログ  $b_m$  に訪問した異常行動ユーザ数に対して最も相関があるのは,period02 におけるブログ  $b_m$  に訪問した全ユーザ数である.それ以降の期間になると相関が薄れていく傾向が見られた.異常行動ユーザに着目することで,直後の訪問ユーザ数を予測できるといえる.

異常行動ユーザと直後の訪問ユーザ数の相関が高いことが分かった.ブログにおいて異常行動ユーザが訪問を始める時間との関係性について調べた.

図 19 はアクセス頻度に関する異常行動ユーザが訪問し始めた時間を,図 20 は趣向比率に関する異常行動ユーザを,それぞれ,プログ  $b_k$  の日毎のアクセス数とともに示した.プログ  $b_k$  は,2004 年 6 月頃から急激にアクセス数が伸びている.これらの結果より,イノベータ理論におけるイノベータ的な特徴が現れていることがわかる.異常ユーザのユーザ数と訪問するユーザ数との間に見られた相関関係は異常ユーザがイノベータ的な役割を果たすことに起因していると考えられる.

アクセス頻度に関する異常行動ユーザ(図19)は2004年4月頃から訪問し始めている

IPSJ SIG Technical Report

のに対し,趣向比率に関する異常行動ユーザ(図 20)は 2004年2月頃から訪問し始めていた.これより,アクセス頻度に関する異常行動ユーザよりも,趣向頻度に関する異常行動ユーザのほうが,よりイノベータ的な要素を持っていることが分かった.

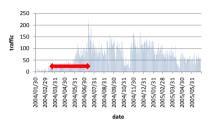


図 19 あるブログにおける日毎のアクセス数と, アクセス頻度に関する異常行動ユーザの 新規訪問

Fig. 19 Relation between the number of accesses and the first day of the anomaly user access(Access frequency)

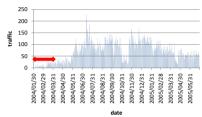


図 20 あるプログにおける日毎のアクセス数と 趣向比率に関する異常行動ユーザの新規 訪問

Fig. 20 Relation between the number of accesses and the first day of the anomaly user access(Share)

# 5. ま と め

本研究では、ブログ空間における異常訪問行動の分析を行った.まず、異常行動に関して、「アクセス頻度異常」と「趣向比率異常」の2種類の概念を定義した「アクセス頻度」と「趣向比率」というの2つの視点から、日毎のユーザの行動回数と行動異常度分布を求めることにより、アクセス頻度異常行動と趣向比率異常行動が多数存在することを観測した.また、趣向比率異常のユーザと人気ブログを発見するイノベータの間に相関関係があることや、現在において趣向比率異常を行うユーザが将来においても趣向比率異常を行う確率のほうが、現在においてアクセス頻度異常を行うユーザが将来においてもアクセス頻度異常を行う確率よりも高いことを観測した.

さらに,ブログの被異常度を算出し,異常行動ユーザがどのブログに異常行動をしたかを分析した.ブログの被異常度の分散において,期間による分散値が大きいブログは少なく,分散値が小さいブログがほとんどであったことから,ブログの被異常度は時期によってあまり変化しないことが観測された.すなわち,異常行動されやすいブログは継続して異常行動され,異常行動されにくいごとが示唆された.

# 参考文献

- 1) 山縣 修,柳下 孝義:アクセス確率分布と Web ユーザビリティの関連の評価, The bulletin of Health Science University (5), pp.33-44, (2009).
- 2) Heather A. Johnsona, Michael M. Wagnera: Analysis of Web Access Logs for Surveillance of Influenza, MEDINFO 2004 Amsterdam: IOS Press, pp.1202-1206(2004).
- 3) Takayuki Itoh, Member, IEEE Computer Society, Yumi Yamaguchi, Yuko Ikehata, and Yasumasa Kajinaga: Hierarchical Data Visualization Using a Fast Rectangle-Packing Algorithm, IEEE TRANSACTIONS ON VISUALIZATION AND COMPUTER GRAPHICS, VOL. 10, NO. 3, pp.302-313(2004).
- 4) 山本 和紀, 伊藤 政志, 熊野 雅仁, 木村 昌弘: イノベータ理論を用いたプログユーザ の行動予測, ネットワーク生態学シンポジウム (NetecoSymp 2009), P9-10 (2009).
- 5) 山本 和紀,熊野 雅仁,木村 昌弘:アクセス履歴を用いたプログ空間における急増する訪問行動の予測,情報処理学会50周年記念全国大会講演論文集,5A-3(2010).
- 6) 総務省情報通信政策研究所 (IICP) 調査研究部: プログの実態に関する調査研究の結果 (2008).