

Web ページの更新時間間隔の測定と統計的性質*

4X-05

松岡高治†

森川 綾‡

水谷正大§

東京情報大学情報システム学科¶

1 Web ページ空間の調査

Web ページ空間は膨大で複雑なリンク関係を持ち、巨大な空間構造の存在が予感されている。例えば、Web ページのリンク数とその頻度にはベキ法則が成立するとか、任意の Web ページ間の平均リンク数が Small World 性 [D.J. Watts and S.H. Strogatz(1998)] を有するなどの報告 [A.-L. Barabasi and R. Albert(1999)] がある。

本稿ではこのような Web の空間的性質ではなく Web の動的性質に着目し、その研究方法の 1 つとして Web ページ群の更新時間間隔の測定実験を行い、Web 空間の時間構造を探るための統計分析法を提案する。

2 URL の更新時間列

URL url で指定される Web ページに対して次のような HTTP HEAD 要求

```
HEAD url HTTP/1.1
If-Modified-Since: Thu, 01 Jan 1970 00:00:00 GMT
```

を送り、その応答を記録するという測定を継続的に行うことを考える。この HEAD 要求に対して Last-Modified 時間を返す URL 集合を $\mathcal{U} = \{u_i\}$ 、その要素数を $N = |\mathcal{U}|$ とする。ある URL u_i を一定の測定間隔 Δt で Last-Modified 時間を L 回測定し続けて得られる Last-Modified 時間の列を

$$(\ell_1(i), \ell_2(i), \dots, \ell_j(i), \dots, \ell_L(i))$$

とする。このとき、差分列 $(\ell_2(i) - \ell_1(i), \ell_3(i) - \ell_2(i), \dots, \ell_j(i) - \ell_{j-1}(i), \dots, \ell_L - \ell_{L-1})$ において零でない項はページ更新があったことを示すが、これを順番に並べた列

$$I(u_i) = (I_1(i), I_2(i), \dots, I_{n_i}(i))$$

* Measurement of Renewal Intervals of Webpages and their Statistical Properties

† Takaharu Matsuoka

‡ Aya Morikawa

§ Masahiro Mizutani

¶ Dept. of Information Systems, Tokyo University of Information Sciences

を測定間隔 Δt によって得られた URL u_i の更新時間間隔列と定義する。ただし、周期間隔 Δt で測定しているために (さらに実際の測定では正確に周期的とはならない)、ここで得られる更新時間間隔は真の値ではないことを注意しておく。こうして、対象としている URL 集合全体に対して同様の測定を行って更新時間間隔列の集合

$$I(\mathcal{U}) = \{I(u_1), \dots, I(u_i), \dots, I(u_N)\}$$

が得られる。

3 測定の方法

ページ更新時間間隔測定の対象である URL 集合を Web 空間から抽出するための合理的方法を確立することは容易ではない。そこで、本稿では

- Google(<http://www.google.co.jp/>)
- Yahoo(<http://www.yahoo.co.jp/>)
- Infoseek(<http://www.infoseek.co.jp/>)

の Web 検索システムのディレクトリサービスから提供されている URL 群 (それぞれ、約 19,000、27,000、19,000 個) を測定対象の候補とした。実際の測定に際しては、CGI や SSI ページなどのために HEAD 要求に対して Last-Modified を返さなかったり、Web サーバ側の時間設定不良などのために更新時間に矛盾をきたす URL 群が存在する。これらを取り除き、測定期間にわたって常に矛盾のない Last-Modified 時刻を取得できた URL 群を測定対象としての URL 集合 \mathcal{U} とした。測定間隔 $\Delta t = 1$ 日として、2001 年 11 月 21 日から測定実験を行った。

4 着目する統計量

いま、十分に長期間に渡って測定を繰り返した結果、十分に長い更新時間列が得られたとする。このとき、次の統計量に着目する。

4.1 URL ごとの更新時間間隔の平均分布

URL u_i の更新時間間隔列 $I(u_i) = (I_1(i), I_2(i), \dots, I_{n_i}(i))$ において、時間区間 $[T_k, T_{k+1})$ に属する更新時間間隔の個数 $m_k(u_i)$

$$m_k(u_i) = \#\{I_j(i) \in [T_k, T_{k+1}) \mid I_j(i) \in I(u_i)\}$$

を求め、時間区間 $\{[T_k, T_{k+1})\}_k$ 上の頻度 $\frac{m_k(u_i)}{n_i}$ を考える。URL 集合 U 内のすべての URL についてこれを求めて、各時間区間 $[T_k, T_{k+1})$ 上での URL についての平均値

$$\langle m_k \rangle_U = \frac{1}{|U|} \sum_{i=1}^{|U|} \frac{m_k(u_i)}{n_i}$$

を‘URL ごとの更新時間間隔の平均分布’と定義する。

4.2 全 URL に関する更新時間間隔分布

一方、URL 全体にわたる更新時間間隔列の集合 $I(U) = \bigcup_i I(u_i)$ において、時間区間 $[T_k, T_{k+1})$ 上の更新時間間隔の URL 集合 U についての個数 $M_k(U)$

$$M_k(U) = \#\{I_j \in [T_k, T_{k+1}) \mid I_j \in I(U)\}$$

を求め、時間区間 $[T_k, T_{k+1})$ 上の頻度 $\langle M_k \rangle_U = \frac{M_k(U)}{\sum_i n_i}$ を‘全 URL に関する更新時間間隔分布’と定義する。 $M_k(U) = \sum_{i=1}^{|U|} m_k(u_i)$ である。

4.3 平均更新時間間隔分布

URL $u_i \in U$ の更新時間間隔の平均 $\langle I(i) \rangle = \frac{1}{n_i} \sum_{k=1}^{n_i} I_k(i)$ が U において時間区間 $[T_k, T_{k+1})$ に属する個数 $m_k(U)$

$$m_k(U) = \#\{\langle I(i) \rangle \in [T_k, T_{k+1}) \mid u_i \in U\}$$

の頻度 $\frac{m_k(U)}{|U|}$ を‘平均更新時間間隔分布’と定義する。

4.4 更新時間間隔列のメンバ数分布

URL u_i の更新時間間隔列 $I(u_i) = (I_1(i), I_2(i), \dots, I_{n_i}(i))$ に関する時間区間 $[T_k, T_{k+1})$ 上のメンバ関数 $C_{[T_k, T_{k+1})}(u_i)$ を

$$C_{[T_k, T_{k+1})}(u_i) = \begin{cases} 1 & \#\{I_j \in [T_k, T_{k+1}) \mid I_j \in I(u_i)\} \geq 1 \\ 0 & \text{それ以外} \end{cases}$$

とする。このとき、時間区間 $[T_k, T_{k+1})$ 上の全 URL 集合にわたるメンバ数頻度

$$\frac{\sum_{i=1}^{|U|} C_{[T_k, T_{k+1})}(u_i)}{\sum_{k=0}^{\infty} \sum_{i=1}^{|U|} C_{[T_k, T_{k+1})}(u_i)}$$

を‘更新時間間隔列のメンバ数分布’と定義する。

4.5 更新時間間隔列のリターン写像の和

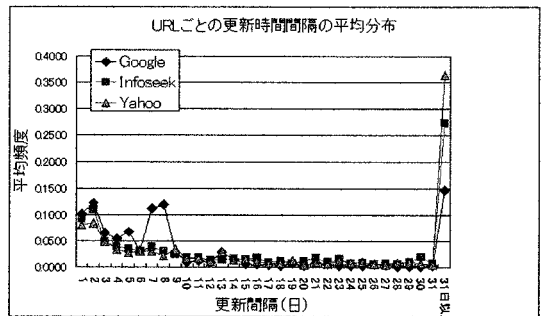
更新時間間隔列 $(I_1, \dots, I_{n-1}, I_n, I_{n+1}, \dots)$ について、時間間隔 I_n の次の更新時間間隔 I_{n+1} がリターン写像 $T: [0, \infty) \rightarrow [0, \infty)$ によって

$$I_{n+1} = T(I_n)$$

と与えられると仮定する。すなわち、座標点 $(I_1, I_2), (I_2, I_3), \dots, (I_{n-1}, I_n), (I_n, I_{n+1}), \dots$ という隣接対は写像 T の関数グラフの上に乗ることになる。URL u の更新時間間隔列 $I(u)$ における隣接対をプロットして得られるリターン写像 T_u のグラフを全 URL U にわたって重ね合わせたものを‘更新時間間隔列のリターン写像の和’と定義する。

5 測定結果

URL ごとの更新時間間隔の平均分布の例を次に示す。詳細は発表で述べる。



本研究は文部科学省助成の東京情報大学学術フロンティア推進事業から一部補助を受けた。