

ウェブのアクセスパターンとリンク構造の同時可視化の一手法

川本 真規子[†] 伊藤 貴之[†]

[†]お茶の水女子大学理学部情報科学科 〒112-8610 東京都文京区大塚 2-1-1

E-mail: [†] {makiko, itot}@itolab.is.ocha.ac.jp

あらまし ウェブの可視化に関する研究は、ウェブサイト本体に関する情報（リンク構造や文書内容など）の可視化と、閲覧者のアクセス動向の可視化に大別される。この2種類の可視化を統合することで、ウェブサイト運営に関する有益な知見が得られると考えられる。そこで我々は、ウェブのアクセスパターンとリンク構造を同時に一画面上に可視化するという研究を進めている。本報告では、可視化手法「FRUITS Net」を用いたウェブの可視化手法を提案する。「FRUITS Net」はカテゴリ情報を同時に可視化できるネットワーク可視化手法である。本手法では、クローラによりリンク構造を、アクセスログファイルによりアクセスパターンを構築し、これらのデータを「FRUITS Net」を用いて可視化する。本手法により、アクセスパターンとリンク構造の対応関係を視覚的にとらえることができるようになり、リンク構造の再構築やページ構成の変更など、ウェブサイトのデザインを再検討する手がかりになると考えられる。

キーワード 可視化, アクセスパターン, リンク構造, アクセスログ

A Visualization Method for Web Access Pattern and Link Structure

Makiko Kawamoto[†] Takayuki Itoh[†]

[†] Department of Information Sciences, Ochanomizu University 2-1-1 Otsuka, Bunkyo-ku, Tokyo, 112-8610 Japan

E-mail: [†] {makiko, itot}@itolab.is.ocha.ac.jp

Abstract There have been two types of Web visualization techniques: visualization of Web sites themselves based on link structures or lexical contents, and visualization of browsers' behaviors. We think that integration of such two visualization techniques is very useful for Web site management, and therefore we are currently studying on visualization of access pattern and link structure on a single screen. This paper presents a Web visualization technique using 'FRUITS Net', which is a visualization technique for multiple-category-embedded network data. The presented technique constructs link structures using crawler software, and access patterns from access log files. It then integrates them and visualizes by 'FRUITS Net'. We expect that users can visually understand the relationship between access patterns and link structures, and utilize the knowledge for design and management of Web sites.

Keyword Visualization, Access Pattern, Link Structure, Access Log

1. はじめに

ウェブに関する情報可視化の研究は、1990年代中盤から非常に多く発表されている。ウェブ可視化の対象は、リンク構造や文書内容などウェブサイト本体に関する情報と、アクセス統計をはじめとする閲覧者情報に大別される。これら2種類の情報を一画面上に同時に可視化することで、ウェブサイトの構築や管理に関する有用な知見が得られることが期待される。

本研究では、「FRUITS Net」[1]という可視化手法を用いて、アクセスパターンとリンク構造の同時可視化を試みる。ここで本研究ではアクセスパターンを、複

数の閲覧者からアクセスされる同一ウェブページ群と定義する。これをリンク構造と同時可視化することにより、アクセスパターンとリンク構造が適切に対応しているか、人がリンクを辿ってページにアクセスしているか、いう知見を視覚的に得られると考えられる。そしてこの知見を、閲覧者の経路を反映した適切なリンク構造の再構築、ウェブサイトのページ構成の変更、などに活用できると考えられる。

2. 関連研究

アクセスログからのアクセスパターン抽出手法と

して、文献[2]では、閲覧者の1セッション内においてアクセスされた一連のURLの推移について、Longest Common Subsequence(LCS)アルゴリズムを用いて頻出のアクセスパターンを抽出するという方法を提案している。また文献[3]では、各ページに閲覧時間の長さに応じて重みを付け、グラフマイニングによりアクセスパターンを抽出する方法を提案している。閲覧者の興味遷移の抽出手法として、文献[4]では、ウェブアクセスログデータを解析し、閲覧者の興味やアクセスしている情報が時間と共にどのように変化しているのかを抽出して可視化している。また文献[5]では、ウェブサイトで提供されるサービスの関連性の分析を行う方法を提案している。

ウェブサイトのリンク構造を可視化する手法として、ウェブサイトを階層型グラフデータとして表現し、力学モデルを用いたグラフデータの画面配置手法により可視化する手法が挙げられる[6]。また、ウェブサイトのアクセス分布の可視化手法[7][8]もいくつか提案されている。文献[8]ではアクセス統計とリンク構造の同時可視化を試みているが、この手法でのリンク構造は1ページを根とした木構造に限定されている。

本研究で用いる「FRUITS Net」[1]は、ノード配置とノード着色の工夫により、リンク構造とカテゴリ情報（本報告ではアクセスパターン）の同時可視化を目指す手法である。FRUITS Netでは、力学モデルに基づく画面配置アルゴリズムにより、

条件1 共通のカテゴリを有するノードが画面上で近くに配置される

条件2 リンク長の総計とリンク間交差を減らすの2条件を満たすような配置を実現する。さらに、テンプレートを用いた空間充填モデルに基づく画面配置アルゴリズムによって配置結果を修正することにより、

条件3 ノードが画面上で重ならない

条件4 配置結果の画面占有面積を減らすという2条件も同時に満たすような配置を実現する。

3. 提案手法

本手法では前処理として、

- アクセスログからのアクセスパターン構築
 - クローラを用いたリンク構造構築
- により入力データを生成する。そして、これらのデータを「FRUITS Net」で同時可視化する。

3.1 入力情報の定義

本報告では、アクセスログファイルの入手可能な1ドメインを対象として、そのウェブサイトのトップページからクローラによってリンクを辿ることにより、リンク構造を構築する。

また本報告では、標準的なアクセスログとして、閲

覧者IPアドレス、アクセス日時、アクセスされたファイル名、リンク元ページのURL、使用しているOS名やブラウザ名、などが記録されているアクセスログファイルの使用を前提とする。

ただし、このアクセスログファイルから、閲覧者の全てのアクセス履歴を抽出できるとは限らない。例えばウェブブラウザの「戻る」ボタンを押した場合などには、キャッシュされたウェブページを再表示するためにサーバへのアクセスが発生せず、結果として閲覧履歴がアクセスログファイルに記録されないことがある。そのため本研究では、以下の2種類の立場を想定するものとする。

立場 a 閲覧者がどのページからどのページへ辿ったという経路情報を一切参照しない

立場 b 徹底的に閲覧者の経路を記録する

[立場 a]では、アクセスパターンとリンク構造を同時に可視化することで、画面上で閲覧者が辿ったリンクの軌跡を想像できるが、その正当性は保証されない。[立場 b]では、正当性のある形で閲覧者の軌跡を可視化できる。しかし、そのためにはウェブサイトの各ページにトラッキングコードを埋め込む、あるいは閲覧者側のパソコンに特定のプログラムをインストールしてデータを採るなど、特殊な方法で閲覧者の全軌跡を記録する必要がある。だが、これらの方法を使用すると、他のウェブサイトでの応用が困難になる、あるいは限られた閲覧者の軌跡しか採れない、といった制限が生じる。そのため現時点では、我々は[立場 a]を前提として研究を進めている。しかし原理的には、提案手法は[立場 b]を前提とすることも可能である。

3.2 アクセスパターン構築

アクセスパターン構築における我々の実装は、以下のとおりである。本処理ではまず、アクセスログファイルを読み込み、閲覧者とURLの一覧を作成する。ただし我々の実装では、画像や音楽などのコンテンツファイルのURLを削除し、それ以外のURLだけを対象とする。続いて本処理では、閲覧者のIPアドレスの数を n 、アクセスされたURLの数を m として、 $n \times m$ の表を作成する。表の各欄には、各閲覧者から各URLへのアクセス回数の集計結果を記録する。

続いて本処理では、閲覧者のデンドログラムを構築する。このとき1閲覧者のアクセス回数を m 次元ベクトルとして、すべての閲覧者ペアについてベクトル間余弦を算出し、これが最大であるペアを併合する。この処理を再帰的に反復することで、デンドログラムを構築する。そして、このデンドログラムを用いて閲覧者を階層的にクラスタリングする。続いて各クラスタに対して、所定人数を超える閲覧者がアクセスしたページを抽出することで、アクセスパターンのデータを

構築する。現時点での我々の実装では、3 ページ以上のページが抽出されたアクセスパターンのみを可視化の対象としている。

なお我々の実装では、抽出されたアクセスパターン全てを可視化するのではなく、抽出されたアクセスパターンの中から手動で 10~15 個を選んで可視化している。この理由は「FRUITS Net」で可視化する際に、アクセスパターン 1 つに対して 1 色を割り当てていることにある。抽出されたアクセスパターン全てを可視化するには、使用される色が多すぎて、人間の目による識別は難しい。そのため現時点では、可視化対象となるアクセスパターンを 10~15 種類に限定している。将来的には、抽出されたアクセスパターンに優先度をつけ、アクセスパターンを自動選択できるようにしたいと考えている。

3.3 リンク構造のデータ構築

リンク構造のデータ構築にはクローラを使用する。本処理では、アクセスログファイルを入手したサイトのトップページを指定し、そこからリンクで繋がっているページを抽出してリストを作り、得られた URL をもとにリンクのグラフを構築する。我々の実装では、オープンソースとして提供されている「JSpider」[7] というクローラを採用している。

3.4 可視化

本手法では 3.2 節および 3.3 節で構築されたデータを統合し、「FRUITS Net」を用いて可視化する。本処理ではまず、3.2 節および 3.3 節で抽出された URL を統合し、ディレクトリ構造に基づいて階層的に格納することで、各 URL を葉ノードとする木構造を生成する。そして各ノードにリンク構造を付加することで、階層型ネットワークを形成する。さらに、各 URL にアクセスパターン情報を付加することで、可視化のための入力データを構築する。

本手法において「FRUITS Net」を適用する利点は、以下のとおりである。まず[条件 1][条件 2]により、同一アクセスパターンに属する URL や、リンクされた URL が、画面上の近い位置に配置されるので、アクセスパターンとリンク構造の関係を視覚的に理解しやすくなる。また[条件 3][条件 4]により、できるだけ多くの URL を一画面上で一覧できるようになる。

4. 適用事例

本章では、我々の所属研究室のウェブサイトの本手法を適用した事例を報告する。

本手法による可視化結果において、各ノードはウェブページ、カテゴリ情報を表す色はアクセスパターン、リンクはウェブページ間のハイパーリンクを表す。ノードの大きさはリンク数に比例し、たくさんリンクを

持つノードほど大きく描かれる。図 1 は可視化結果画面の一例である。右端は GUI の操作部分であり、色一覧がアクセスパターンに使われている色を表している。色一覧から色を選択すると、選択された色で色付けされているノードがハイライトされる。

図 1 は、7 月のアクセスログファイルから抽出されたアクセスパターンを用いて可視化を行った結果である。この結果では、教員の担当講義の資料ページへのアクセスパターンに該当する部分を拡大表示している。この可視化結果から、教員のトップページから担当科目のページへアクセスし、そこから専門科目の資料のページへアクセスした、という閲覧者の軌跡を想像できる。この結果から、目的のページにアクセスするために、トップページから順にリンクを辿ったアクセスが多数あったことがわかる。

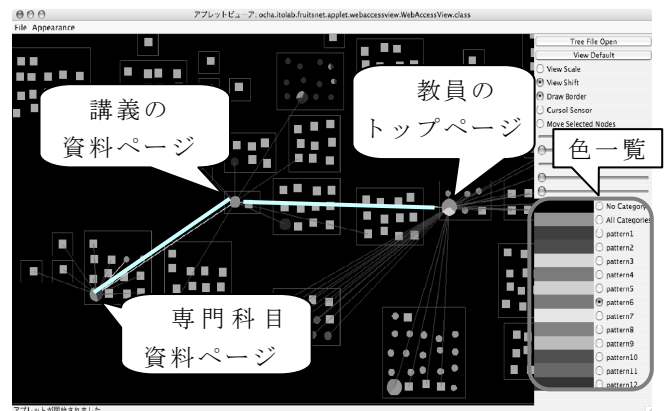


図 1: 可視化結果 1 (教員の講義科目の資料ページに関するアクセスパターン)

別の事例として、研究室のメンバーのページに関連するアクセスパターンに関する可視化結果を示す。図 2 は、メンバー一覧のページから、同学年全員のホームページにアクセスする閲覧者が複数いることを示している。図 3 は、同クラス内での大半のページにアクセスされ、特定の個人のホームページだけをアクセスする閲覧者が複数いることを示している。

本事例を通して我々が発見したアクセスパターンは、主として以下の 3 種類の分布に分類できる。

[分布 1] 直線型のアクセスパターン (図 1 参照)

[分布 2] 1 つのノードを中心とした放射型のアクセスパターン (図 2 参照)

[分布 3] 同クラス内だけにアクセスが集中している同クラス型のアクセスパターン (図 3 参照)

直線型のアクセスパターンをとる閲覧者は、目的の 1 ページにアクセスするためにリンクを辿ってアクセスしているということが想像される。放射型のアクセスパターンをとる閲覧者は、あるページを中心として紹介されている複数のページに関心があるということが

想像される。また、同クラスタ型のアクセスパターンをとる閲覧者は、同じディレクトリ内のページにだけ関心があるということが想像される。このように、可視化結果に現れたアクセスパターンの形から、サイトを訪問した閲覧者がどのような意図を持ってアクセスしているのかということを目視的にとらえ、想像することができる。

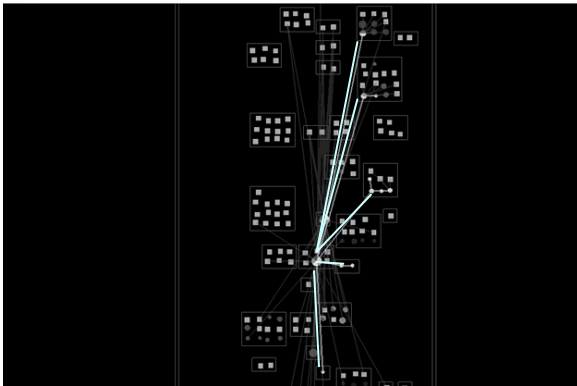


図 2: 可視化結果 2 (同学年の各メンバーのページに関するアクセスパターン)

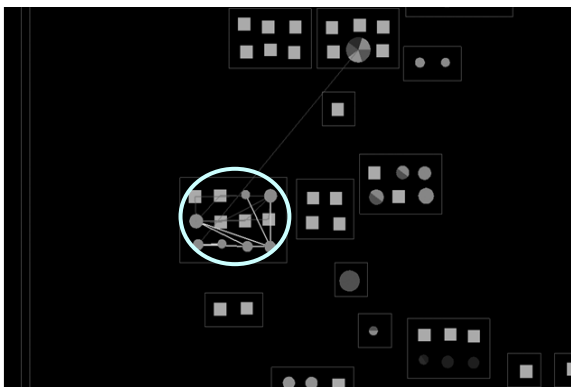


図 3: 可視化結果 3 (1人のメンバーのウェブページ群に関するアクセスパターン)

5. まとめ

本報告では、「FRUITS Net」を用いたアクセスパターンとリンク構造の同時可視化の一手法を提案した。本手法を用いて、アクセスパターンとリンク構造を対応させて同時可視化することで、アクセスパターンごとに異なる分布を画面上で発見できることがわかった。本報告で紹介した適用事例では、3種類の分布を発見することができた。

今後はさらに異なる分布を有するアクセスパターンの発見を目指すとともに、各アクセスパターンに対応する適切なサイト構成の提案を行いたい。また、各ページのアクセス数を高さで表示できるようにするなど「FRUITS Net」の機能の拡張を試みたい。また、現在実装しているアクセスパターン抽出手法についても

再検討したいと考えている。

文 献

- [1] T. Itoh, C. Muelder, K.-L. Ma, J. Sese, A Hybrid Space-Filling and Force-Directed Layout Method for Visualizing Multiple-Category Graphs, 2009 IEEE Pacific Visualization Symposium, pp. 121-128, 2009.
- [2] 宇根田純治, 横田治夫, Web ログの共通シーケンス解析, 電子情報通信学会信学技報 DE2002-2, 2002.
- [3] 三原宏一郎, 寺邊正大, 橋本和夫, ページ閲覧時間を考慮した Web ログマイニング手法の提案, 情報処理学会, pp. 39-44, 2007.
- [4] 山田和明, 中小路久美代, 上田完次, Web ユーザの行動履歴解析のためのデータマイニング, 電子情報通信学会 W12 研究会資料, pp. 59-64, 2005.
- [5] 末永高志, 岡田崇, 石打智美, Web アクセスログデータの系列情報を利用したサービスの関連性の分析, 電子情報通信学会信学技報 DE2005-17, 2005.
- [6] 土井淳, 伊藤貴之, 力学モデルを用いた階層型グラフデータ画面配置手法の改良手法とウェブサイト視覚化への応用, 芸術科学会論文誌, Vol. 3, No. 4, pp. 250-263, 2004.
- [7] 山口裕美, 伊藤貴之, 池端裕子, 梶永泰正, 階層型データ視覚化手法「データ宝石箱」とウェブサイトの視覚化, 画像電子学会誌, Vol. 32, No. 4, pp. 407-417, 2003.
- [8] 山縣修, 中村泰明, アクセス確率による Web サイトのリンク構造可視化ツール, 可視化情報学会論文集, Vol. 26, No. 6, pp. 43-50, 2006.
- [9] 「JSpider」
<http://j-spider.sourceforge.net/>