

論文 / 著書情報  
Article / Book Information

論題(和文)	アクセスログに基づくWebページ推薦におけるLCSの利用とその解析
Title(English)	Analyses of the Effects of Utilizing Web Access Log LCS for Web Page Recommendation
著者(和文)	山元理絵, 小林大, 吉原朋宏, 小林隆志, 横田治夫
Authors(English)	RIE YAMAMOTO, DAI KOBAYASHI, TOMOHIRO YOSHIHARA, TAKASHI KOBAYASHI, HARUO YOKOTA
出典(和文)	Proc. of IPSJ DBWeb2006, Vol. , No. , pp. 43-50
Citation(English)	Proc. of IPSJ DBWeb2006, Vol. , No. , pp. 43-50
発行日 / Pub. date	2006, 11
権利情報 / Copyright	<p>ここに掲載した著作物の利用に関する注意: 本著作物の著作権は(社)情報処理学会に帰属します。本著作物は著作権者である情報処理学会の許可のもとに掲載するものです。ご利用に当たっては「著作権法」ならびに「情報処理学会倫理綱領」に従うことをお願いいたします。</p> <p>The copyright of this material is retained by the Information Processing Society of Japan (IPSJ). This material is published on this web site with the agreement of the author (s) and the IPSJ. Please be complied with Copyright Law of Japan and the Code of Ethics of the IPSJ if any users wish to reproduce, make derivative work, distribute or make available to the public any part or whole thereof.</p>

# アクセスログに基づく Web ページ推薦における LCS の利用とその解析

山元 理 絵<sup>†</sup> 小林 大<sup>†,††</sup> 吉原 朋 宏<sup>†</sup>  
小林 隆 志<sup>†††</sup> 横田 治 夫<sup>†††,†</sup>

近年, Web サイトによる情報発信の重要性から, ユーザのニーズに適したサイト構築や情報提供の要求が高まってきている. Web アクセスログを Web ページ推薦に用いる方法は, クライアント側に手を加える必要がなく有用であるが, これまで提案されている手法では, 頻出アクセスパターンと僅かでも外れると適切な推薦ができない, あるいは順序を考慮できないといった問題点があった. 我々は, それらの問題を解決するために, Web アクセスログから LCS (Longest Common Subsequences) を抽出してページ推薦に利用する手法である WRAPL を提案している. 本稿では, 実際の Web アクセスログを用いた実験を通して WRAPL の効果を詳細に解析し, その実験結果から得られた知見を基に優先順位付け手法に対して改良を行い, その有効性を示す.

## Analyses of the Effects of Utilizing Web Access Log LCS for Web Page Recommendation

RIE YAMAMOTO,<sup>†</sup> DAI KOBAYASHI,<sup>†,††</sup> TOMOHIRO YOSHIHARA,<sup>†</sup> TAKASHI KOBAYASHI<sup>†††</sup>  
and HARUO YOKOTA<sup>†††,†</sup>

Sophisticated websites satisfying users' requirement becomes much more important to propagate information via websites, nowadays. Web page recommendation methods using web access logs are useful for them because they need no modification in client-side applications to meet the requirement. However, traditional methods have problems of insufficient recommendation precision caused by strict matching of access patterns or neglect of access sequences. To solve the problems, we are proposing WRAPL as a method of extracting LCSs (Longest Common Subsequences) from web access logs and using them to recommend web pages for an active session. In this paper, we analyze the effects of WRAPL using actual web access logs and propose an enhanced weighting method for it to improve the precision based on the analyses.

### 1. はじめに

近年, ビジネスの場としての Web の役割と情報量の増大から, Web personalization が注目され<sup>1)</sup>, 特にユーザの嗜好に合った Web ページをシステムがユーザに推薦し提示する Web ページ推薦が盛んに研究されている. その情報収集の方法としては, 閲覧した情報に対する各々のユーザの興味の有無を何らかの方法で収集し分析する方法や, Web サイトのアクセスログを分析する方法などが取られる.

前者の例としては, ユーザによるなぞり読みやリンククリック等の特徴的なマウス操作を利用している

TextExtractor<sup>2)</sup> がある. ユーザの嗜好を評価するためのフィードバックとして, ページ内のテキスト部分を, 文や行の単位で興味情報として抽出できるが, クライアント側にプログラムを埋め込まなければならない等の問題点がある.

一方, アクセスログを利用する方法は, クライアント側の変更が不要であることから, 様々な利用方法が研究されている. 当初は, 相関ルールを用いた頻出共起ページの組の抽出<sup>3),4)</sup> や, 利用頻度に基づくリンクの接続性の評価<sup>5)</sup>, バックトラックポイントの発見<sup>6)</sup> 等のアクセス解析の研究が中心であったが, 最近では, ユーザビリティの向上を目的とし, ユーザ行動の予測手法<sup>7)</sup> や Web ページ推薦に関する手法<sup>8)</sup> が提案されている.

ユーザ行動の予測手法<sup>7)</sup> では, ログ中のユーザのアクセスパスの統計を取ってユーザ行動をモデル化し, 全てのページに対し, 現在のアクセスパスに引き続いてアクセスされるための条件付確率を算出することで, 続くアクセスページを予測する. この中で, 予測モデ

<sup>†</sup> 東京工業大学大学院情報理工学研究科計算工学専攻  
Department of Computer Science, Graduate School of Information Science and Engineering, Tokyo Institute of Technology

<sup>††</sup> 日本学術振興会特別研究員 DC  
Research Fellow (DC), Japan Society for the Promotion of Science

<sup>†††</sup> 東京工業大学 学術国際情報センター  
Global Scientific Information and Computing Center, Tokyo Institute of Technology

ルの適用例の一つとして Web ページ推薦が挙げられているが、このような手法では推薦されるのは直後のページのみに限られてしまう。さらに、サイト内のナビゲーションが不適切な時やサイト規模が大きい時には、目的のページに到達するまでに後戻りや遠回りを含んだり、目標が複数存在することで、多種多様なパスを取りうるための確な推薦が難しい。

一方、関連ルールを用いて Web ページ推薦を行う手法<sup>8)</sup>では、アクセスログに apriori アルゴリズムを適用して 1 セッション中で頻繁に共起するページの組の集合及びページアクセスにおける関連ルールを作成してページを予測する。アクセスパスそのものを扱うのではなく、共起頻度の高いページ同士の関係を見て推薦を行うため、前述の直後のページの推薦に限られるという問題や、後戻りや遠回りを含む場合の問題を解消することができる。しかし、この手法ではページアクセスの順序情報を含まないため、ページ参照の順序に特徴的な傾向がある場合など、すでにアクセスしたページを推薦したり、ユーザにとってもはや不要となったページを推薦することで、推薦精度を低下させてしまう可能性がある。

我々は、順序情報を考慮しないという上記の問題を解決するため、アクセスログ中のシーケンスの LCS (Longest Common Subsequence) を用いることで、アクセスパターンのぶれを吸収した概括的なアクセス順序を利用して推薦精度を向上させる手法を提案している。

これまでに、LCS を用いたアクセスログの解析手法とその解析に基づくサイト構成の改善手法<sup>9)</sup>、その場合の LCS 抽出の効率化手法<sup>10)</sup>を提案し、その有効性を示してきた。Web サイト内における全てのセッションの URL の推移をシーケンスとして抽出し、そのそれぞれに対して他の全てのシーケンスとの LCS を求め、頻出アクセスパターンを発見するものである。

さらに、本稿の事前検討として、上述のアクセスログ解析によって抽出された LCS を用いて現在までのアクセス履歴から次にアクセスされるページを予測して推薦する手法に関して、研究会で報告を行った<sup>11)</sup>。LCS を用いることで、アクセスパスが完全に一致しない場合でも全体のアクセスの傾向の表現が可能になるとともに、順序情報を保持することができるため、実際の Web アクセスログを用いた実験において、関連ルールを用いる手法と比較して推薦精度が向上した実験結果が得られている。しかし、事前検討であったため実験内容とその解析が必ずしも十分でなく、予想に反する結果も含まれていた。そこで、本稿では、アク

セスログに基づく Web ページ推薦における LCS の利用効果をより詳細に解析するために、条件及び設定を変更したいくつかの実験を行うと同時に、それらの実験結果から得られた知見を基に手法の改良を行い、その効果を調べる。

## 2. Web アクセスログ解析への LCS の適用

本節では、我々が以前に提案してきた、アクセスログから LCS を抽出する方法<sup>9),10)</sup>について述べる。LCS を抽出するためには、まずアクセスログからユーザセッションを切り出し、データを精練する必要がある。その後、各セッションを比較することで LCS を抽出し、その頻度を集計する。

### 2.1 LCS

リスト  $x$  の部分列とリスト  $y$  の部分列の中で両方のリストに含まれるものを共通部分列という。共通部分列の中で最も長いものを最長共通部分列 (Longest Common Subsequences) と呼び、LCS と略す。

二つのリストの中に同じ要素が同じ順序で出現したものが共通部分列なので、LCS が長いということは二つのリストの類似性が高いことを表す。

これを URL シーケンスに適用することで、アクセスシーケンスが完全に一致しない場合でも、寄り道等の余分な情報を取り除くことにより各々のシーケンス間の類似性を発見することができる。と考える。

### 2.2 LCS を用いた Web アクセスログ解析

#### 2.2.1 ユーザセッションの抽出

アクセスシーケンス解析を行う際、蓄積されている未加工のアクセスログを精練して、マイニングに必要なデータのみを取りだし、ユーザのセッションを抽出する必要がある。

サイト内におけるユーザごとの移動情報を得るため、IP アドレスや Cookie を用いて各セッションに一意的なセッション ID を割り当てる。特に、ユーザのナビゲーションに合わせてリアルタイムに適切なページ推薦を行うことを目的とする本研究においては、Cookie を用いることが好ましいが、これはユーザの了承を必要とするため、Cookie 情報が利用できない場合は、複数のユーザを含む可能性があるため正確性は下がるものの、IP アドレスの情報などで代用する。

セッション ID ごとに整理された URL の集合を時系列順に並べることで、各セッションで訪問者が行ったアクセスの URL シーケンスを得ることができる。

アクセスログの中には、画像ファイルへのアクセス等の解析処理を行う際に不必要な情報や目的に合わない情報が多く含まれる。そのため、セッション抽出時に、

前処理<sup>12)</sup>によりそれらの情報を取り除く必要がある。

### 2.2.2 LCS の抽出

二つのシーケンスから LCS を求めるには、動的計画法が用いられる。これにより長さ  $M$  の文字列  $X$  と長さ  $N$  の文字列  $Y$  の LCS を求める場合、 $O(MN)$  時間で LCS を求めることができる<sup>13)</sup>。

また、LCS を求める問題と等価である、SED (Shortest Edit Distance) を求める問題に関して、効率化された手法が提案されている<sup>14)</sup>。この手法では、比較する二つの文字列の差異が小さいほど必要とする時間計算量が小さくなるため、実際のデータに適用すると、多くの場合で  $O(MN)$  よりも大幅に小さい計算量で LCS の計算が可能になる。

### 2.2.3 LCS の頻度集計

Web アクセスログから得られた各セッションを、URL を各要素に持つシーケンスとみなし、各シーケンスについて総当たりに LCS を抽出する。各 LCS の出現頻度の集計を行い、高頻度で出現する LCS パターンを発見する。

LCS 抽出のための計算においては、全てのアクセスシーケンスに対して総当たりで求めるため、セッション数が大きくなるとその二乗に比例して時間計算量が増加してしまい、計算コストが大きいという問題があるが<sup>9)</sup>、本研究では、ハッシュを用いたアクセスシーケンスのフィルタリング手法やインクリメンタルな LCS 抽出手法、並列計算のためのアルゴリズム<sup>10)</sup>を用いることで、LCS 抽出にかかわる計算量をさらに抑えることとする。

## 3. LCS を用いた Web ページ推薦

本節では、2 節で説明した手順でアクセスログから抽出される LCS を用いたユーザにページを推薦する手法<sup>11)</sup>について説明する。

抽出された LCS とアクティブセッションのマッチングを行って推薦候補ページを選出する。次に、それぞれのページに得点を付加することで、推薦のための優先順位を決定し、その得点が上位のページから順に推薦する。

### 3.1 推薦候補ページの選出と得点付け

LCS を用いた Web ページ推薦手法では、アクセスログから抽出した LCS のそれぞれと、現在までのユーザのセッション (アクティブセッション) とのマッチングを行い、頻出 LCS の中で、ユーザの現在位置以降に現れているページを推薦する。

抽出された LCS の内、全セッション中において数え上げられた回数が閾値  $min.Count$  以上であり、かつ長

さが  $min.Length$  以上である LCS の集合を large LCS 集合と呼び、 $LL = \{lcs_1, lcs_2, \dots, lcs_k\}$  で表す。また、 $LL$  内の  $i$  番目の要素が全セッション中で数え上げられた回数を  $c_i$  と表す。ここで、 $min.Count$  と  $min.Length$  は、Web サイトの持つ特性に合わせて設定するパラメータである。このとき、長さ  $n$  のアクティブセッション  $act_n$  からそれに続くユーザのページアクセスを予測する。

$lcs_i$  と  $act_n$  の間で共通しているページを調べ、 $lcs_i$  の後半部分の中でまだアクセスされていないページがあれば、そのページはその後にアクセスされる可能性が高いと我々は考える。なぜなら、 $lcs_i$  と  $act_n$  を比較する際、それぞれが完全に一致する必要はなく、共通要素が多く存在する場合に  $lcs_i$  はそのアクセス傾向の特徴を表現していると捉えることができるからである。

また、推薦順を決定するための推薦候補のページへの得点付けに際しては、以下のような点を考慮する必要がある。まず、高い  $c_i$  を持つ  $lcs_i$  は、多くのセッションにおいて頻繁にナビゲートされた部分シーケンスであり、重視すべきである。また上述のように、 $lcs_i$  と一致の度合いが高い  $act_n$  は同じ傾向を示す可能性が高いと考え、高い得点を付加する。以上を考慮した上で、優先順位を付けて推薦を行うための手法として、WRAPL-FL (Web page Recommendation by Access Pattern Lcs with Frequency and matched Length based weighting) 法を定義する。large LCS 集合  $LL$  と長さ  $n$  のアクティブセッション  $act_n$  に対し WRAPL-FL 法では、あるページ  $p$  の得点の算出方法として次の式を用いる。ただし、ページ  $p$  は候補ページの集合に含まれるものとする。

$$point(p) = \sum_{lcs_i \in LL} |lcs_i \cap act_n| \cdot c_i^\alpha \quad (1)$$

$\alpha$  は  $c_i$  の重みであり、Web サイトの特徴からその影響度合いを考慮して適切に調節する。

### 3.2 ページ推薦手法 WRAPL-FL 法

以下の手順で推薦ページの決定を行う。

- (1)  $lcs_i$  と  $act_n$  の間で共通するページを抜き出す。
- (2)  $lcs_i$  より、一番目から共通部分の最後まで要素全てを除去する。
- (3) 残ったページを推薦ページの候補とし、それぞれの  $point$  に  $|lcs_i \cap act_n| \cdot c_i^\alpha$  を加える。
- (4)  $LL$  の中の全ての要素に対して (1)~(3) を行い、候補ページの中で得点の総和が上位のページを推薦する。

ここまで述べた、LCS を用いた推薦手法の概要を図 1 に示す。例えば、ページ推薦のステップにおいて、図のように  $act_3 = (A, B, C)$  が与えられた時、

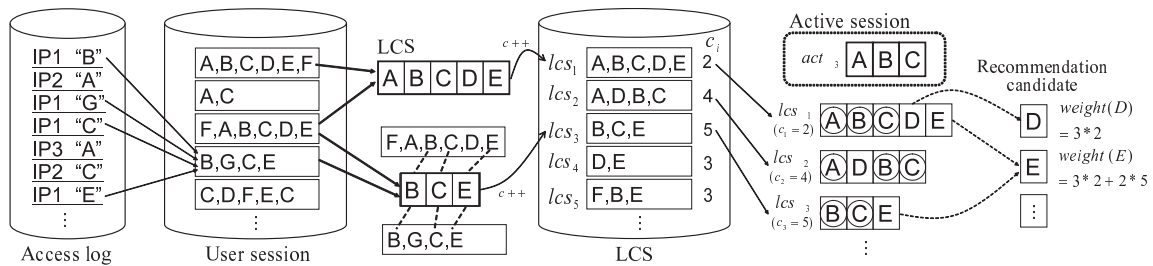


図 1 LCS を用いた推薦

$lcs_1 = (A, B, C, D, E)$  と一致する部分は  $(A, B, C)$  であり、それに続くページ  $D, E$  が推薦の候補ページに加えられる。また  $lcs_2 = (A, D, B, C)$  については、同様に  $(A, B, C)$  が一致するものの、共通部分の最後のページ  $C$  以降に続くページはないため、ここから推薦候補に加えられるページはない。さらに、 $lcs_3 = (B, C, E)$  では  $E$  となる。したがって、この例で  $lcs_1 \sim lcs_3$  から推薦されるページの候補は  $\{D, E\}$  となる。

ここでは、ページ  $E$  が 2 つの LCS で推薦候補となっている。したがって、各 LCS から算出された得点の和がページ  $E$  の得点となる。

#### 4. 推薦精度評価のための準備

我々は、LCS を用いた推薦手法 WRAPL-FL 法の有効性を確認するために、実際のアクセスログを使って実験を行い、その効果を測定すると共に、1 節で紹介した Mobasher らが提案している apriori アルゴリズムによって生成される相関ルールを用いた推薦手法を単純化した手法を実装し、その推薦精度を比較した<sup>11)</sup>。本節では、実験対象として用いたデータの説明と評価指標の定義、さらに比較手法の詳細について説明する。

##### 4.1 実験対象データと評価指標

実験対象として、“The Internet Traffic Archive” (<http://ita.ee.lbl.gov/>) で配布されているいくつかの Web サイトのアクセスログの内、NASA の Web サイトでの 1995 年 8 月 1 日から 8 月 31 日までの Web サーバへのリクエストに対するアクセスログを用いた。このデータにはセッション情報が含まれていなかったため、同一 IP アドレスからのアクセスを同一ユーザからのアクセスとみなし、ページアクセスの間隔が 1,200 秒以上のときにセッションを分割した。ログ全体の中に出現した固有 URL 数は 1,276 で、総セッション数は 39,900 であった。ここで、各 URL のログへの出現頻度には大きな偏りがあったため、各セッション中の出現割合が 0.5% に満たない URL を取り除き、さらに推薦の評価に利用できないため長さが 3 以下のセッショ

ンも除外した結果、URL 数は 174、総セッション数は 23,663 となった。

全セッションの内、時期が早い方の 75% を学習セットとして、そこから各手法ごとに頻出パターンを抽出した。残りの 25% の新しいセッションの集合をテストセットとみなしてページ推薦を行い、その評価を行った。

評価に用いる指標として、以下に定義される *precision* と *coverage* を用いる。

$$precision(Recom) = \frac{|Recom \cap eval|}{|Recom|} \quad (2)$$

$$coverage(Recom) = \frac{|Recom \cap eval|}{|eval|} \quad (3)$$

ここで、*Recom*、*eval* はそれぞれ対象アクティブセッションから導かれた推薦ページの組、対象アクティブセッションに引き続いて実際にアクセスされたページの組 (評価セット) を表す。*precision* は推薦の正確性の指標であり、推薦されるページ数に対する正解ページ数の割合で表現される。また、*coverage* は評価セットをどれだけ網羅しているかの指標であり、評価セットのページ数に対する正解ページ数の割合で表現される。

ページ推薦を行うために、テストセットの各セッション (テストセッション) のはじめの  $n$  ページをアクティブセッション  $act_n$  とみなして large LCS 集合  $LL$  中の全要素とマッチングを行う。そこから推薦ページのランク付けを行って上位のページを推薦し、そのセッションの残りのページ *eval* と比較することで *precision* と *coverage* を求めた。

##### 4.2 比較手法

相関ルールを用いた Web ページ推薦手法<sup>8)</sup> では、長さ  $n$  のアクティブセッション  $act_n$  が与えられると要素数  $n+1$  の頻出アイテムセットを探索し、アクセスされた  $n$  個のページを全て含むアイテムセットから、差分のページを推薦するという手法が提案されている。これは、ページの組  $\{A, B\}$ 、 $\{A, B, C\}$  がそれぞれ頻出アイテムセットに含まれる時、 $\{A, B, C\}$  の共起頻度が高

ければ、{A,B} へのページアクセスに引き続いて C へのアクセスが起こる確率も高くなるという仮定に基づいている。また、アクティブセッションの長さである  $n$  の値を、推薦ページが見つかるまで下げていくという方法で coverage を上げている。

この手法では、推薦ページの順位付けを行っていないため、confidence 値で順位付けを行うこととする。

## 5. LCS を用いたページ推薦手法に対する解析

我々は本稿の事前検討として、前節の設定で実験を行ったが、比較手法の相関ルールを用いた Web ページ推薦手法の文献内で提案されている全ての改善内容を実装していなかったため、今回機能を追加して改めて比較を行った (5.1 節)。その結果、WRAPL-FL 法の優位性は実証できたが、アクティブセッション長が短い場合で長い場合に比べて良いという、予想に反する結果が得られた。そこで、本稿ではその原因について解析を行うため、以下の 3 点に着目した。

- 今回対象とした Web サイトでは、セッション長が短いものが大部分を占めていた。
- アクティブセッションと LCS とのマッチングにおいて、一致要素数が 1 の LCS も判断材料として用いている。
- 推薦順位の決定において、LCS の出現回数とアクティブセッションとの一致要素数のみを考慮している。

これらが推薦精度に与える影響を調べるために、5.2 節~5.4 節では実験の条件を変更して詳細な解析を行い、LCS を用いた推薦手法の特徴について調査する。さらに、5.5 節ではそれらから得た知見を利用し、推薦精度の改善を目指す。

### 5.1 LCS を用いた手法と比較手法による推薦精度の比較

実験に用いたパラメータは、比較手法では一定の最小サポート値 0.008 を用い、要素数が 1 から 5 の頻出アイテムセットを生成した。LCS を用いた手法では、 $min.Count$  を 150、 $min.Length$  を 3 として  $LL$  を作成し、推薦ページの得点付けの式 (1) は  $\alpha = \frac{1}{2}$  の場合を用いた。それぞれについて推薦精度を測定した結果を図 2 に示す。図中の LCS, AR はそれぞれ、LCS を用いた手法、比較手法を用いた推薦に対応しており、 $n$  はアクティブセッションの長さを表す。また、図における各点は、推薦する上位ページの数  $|Recom|$  を 1 から 14 の間で変化させた場合に対応しており、右の点ほど  $|Recom|$  が大きい場合を表している。

図 2 から、同じ長さのアクティブセッションに対す

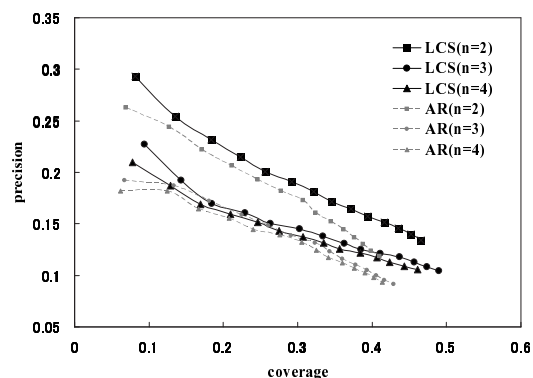


図 2 LCS を用いた手法と比較手法による推薦の評価

るページ推薦の精度について、比較手法に比べ、LCS を用いた手法で良い結果が得られていることが確認できる。

しかし、一般にアクティブセッション長  $n$  の値を大きくするほど履歴から抽出されたパターンとの一致の度合いが大きくなり、precision は改善するはずであるが、図 2 から、 $n$  が大きいものに比べ小さいもので良い結果が得られたことがわかる。

### 5.2 テストセッション長を固定した場合の比較

今回対象としたサイトでは、セッション長ごとのセッション数に大きな偏りがあった。セッション長が 3~9 のテストセッションの数は順に、3189, 1834, 1139, 717, 500, 354, 217 という分布になっており、短いアクセスでセッションを終えるユーザが多く、長いセッションの数が相対的に少なかった。

長さ  $n$  のアクティブセッションから推薦を行う際には、テストセッション長が  $n+1$  以上である必要があることから、調査の対象となるテストセッションに差が生じ、その結果等しい条件で比較が行えなかった可能性がある。そこで、ここではまず、テストセッション長を固定して実験を行うことで、可能な限り近い条件で推薦精度の比較を行うことを考える。

4.1 節と同様のデータに対し、テストセッション長を 6 に固定した場合の結果を図 3 に示す。凡例中の  $n$  は、アクティブセッション長を表している。上の 3 つ (fixed) が固定長 6 のテストセッションを用いた場合、下の 3 つがテストセッション長を固定しない場合の結果である。

図から、 $n = 2, 3$  の場合には大きな差は見られないが、 $n = 4$  の場合で精度が低下しているため、テストセッションの差は精度に影響を与えていることがわかる。しかし、近い条件下で精度を比較した場合にも、 $n$  が小さいアクティブセッションからの推薦でより良

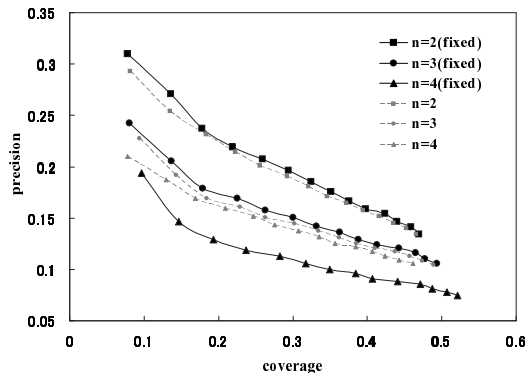


図3 テストセッション長を6に固定した場合の比較

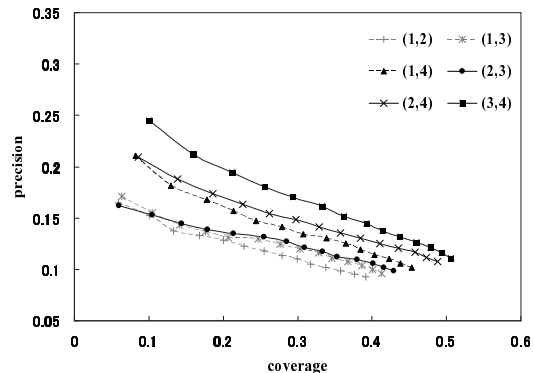


図5 アクティブセッション中のマッチ位置による比較

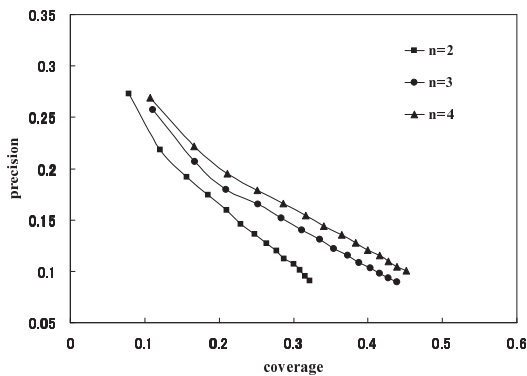


図4 一致要素数が1のものを除いた場合の比較

い結果が得られるという傾向は変わらなかった。また、テストセッション長を5,7に固定した場合でも、同様の傾向が得られた。

### 5.3 一致要素数が1のLCSを除いた場合の比較

WRAPL-FL法では、図1中の $lcs_3$ と $act_3$ とのマッチングの例からもわかるように、一致要素数が1の場合(つまり $|lcs_i \cap act_n| = 1$ の場合)にもそれに続く推薦候補を選出し、得点付けを行っている。

しかし、一致要素数が1の場合、アクティブセッションと過去のパターンとのマッチ度が低いため、推薦精度に悪影響を与えてしまう可能性がある。そこで、一致要素数が1の場合には得点は付加しないものとした上で、同様に精度を測定する実験を行った。

結果を図4に示す。図から、一致要素数が1の場合を判断材料から取り除くと、アクティブセッションを長く取った場合で、短い場合に比べて良い結果が得られることが確認できる。

精度の順番が逆転する要因として、 $n=2$ のときに、一致要素数1も含めて推薦を行った場合に比べて精度が大きく低下していることが挙げられるであろう。こ

れは、 $act_2$ とLCSのマッチングにおいて、 $act_2$ との一致要素数が1であるLCSが全体の約25%を占めているのに対し、一致要素数が2であるLCSは2%弱にとどまっていることに起因すると考える。一方、 $n=3,4$ の結果ではほぼ全体的に、一致要素数1を含める場合に比べ、精度の向上が確認できた。したがって、 $n$ が大きい場合には、一般的な傾向に当てはまり、マッチの度合いが高いLCSが精度向上に影響を与えていることがわかる。

### 5.4 アクティブセッション中のマッチ位置による比較

5.3節のおわりで触れたように、今回対象としたWebサイトは規模が大きく、ユーザのアクセス行動も非常に多様であるため、アクティブセッションとの一致要素数が多いLCSは非常に少ない。そこで本節では、その中で、アクティブセッション中のどの部分がLCSと一致するとき、より精度の高い推薦が行えるかについて調査することによって、推薦精度に与えるマッチ位置の影響を検討する。

これまででは、テストセッションのはじめの $n$ ページを $act_n$ とみなしてページ推薦を行い、そのセッションの残りの部分 $eval$ と比較することでprecisionとcoverageを求めてきた。しかし今回は、セッション前半のどの部分が後半ページとの関連が深いかを調査するために、テストセッション中のはじめの4ページの中から任意の2ページの全ての組み合わせを順序を変えずに取り、それらに対し5ページ目以降の部分の評価セット( $eval$ )として推薦精度を測ることとする。

結果を図5に示す。凡例中の(a,b)はそれぞれ、はじめの4ページ中のa,b番目のページを順に2ページ取り、 $act_2$ とみなした場合の結果に対応している。

図より、推薦精度は(3,4)の場合で最も良く、以下は順に(2,4),(1,4),(2,3),(1,3),(1,2)となっている

ことが確認できる。この結果から明らかに、アクティブセッションの後方のページがユーザの今後のアクセス行動と深いかかわりを持つことがわかる。したがって、推薦を行うための履歴として用いるページの中では、前方のページに比べて後方のページを重視すべきである。

### 5.5 アクティブセッション中の LCS とのマッチ位置の考慮

5.3, 5.4 節で得られた特性を利用して、アクティブセッションと LCS とのマッチングを行う際に、それらのマッチ位置を推薦ページの優先順位付けに反映させる方法を提案する。

#### 5.5.1 マッチ位置を考慮した得点付け

$lcs_i$  と  $act_n$  とのマッチングの際に、アクティブセッション中におけるマッチ位置を考慮するために、次のようにマッチ位置重み  $l_i$  を定義する。前述のように、アクティブセッションとのマッチ位置が後方にある LCS の方がより重要であるため、 $l_i$  は後方ページが重視されるように設定する必要がある。

そこで、 $act_n$  と  $lcs_i$  を比較し、 $act_n$  の  $m$  ページ目が  $lcs_i$  とマッチした場合、 $l_i$  に  $m$  を加算する。例えば、 $act_4 = (A, B, C, D)$ 、 $lcs_i = (B, D, E)$  のとき、 $act_4$  中の 2, 4 番目のページ  $B$  と  $D$  が  $lcs_i$  と一致するため、 $l_i = 2 + 4 = 6$  となる。

このようにして得られた  $l_i$  を、WRAPL-FL 法による推薦ページの優先順位付けの式(1)に掛け合わせることで新たな得点付けの式を以下で定義し、これを用いた推薦手法を WRAPL-FLP (Web page Recommendation by Access Pattern Lcs with Frequency, matched Length and Position based weighting) 法と呼ぶ。

$$point(p) = \sum_{lcs_i \in LL} |lcs_i \cap act_n| \cdot c_i^\alpha \cdot l_i^\beta \quad (4)$$

ただし、 $\beta$  は  $l_i$  の重みである。

#### 5.5.2 WRAPL-FLP 法による順位付けの評価

WRAPL-FLP 法による改良の有効性を確認するため、LCS を用いたページ推薦手法において、推薦ページの優先順位決定に、WRAPL-FL 法、WRAPL-FLP 法を用いた場合の推薦精度の比較を行った。

図 6~8 はそれぞれ、長さ  $n$  のアクティブセッションからのページ推薦において、WRAPL-FLP 法を採用した場合の結果を表している。凡例の  $\beta$  は、式(4)における一致位置に応じた得点  $l_i$  の重みを表す。すなわち、 $\beta = 0$  は WRAPL-FL 法に対応している。

グラフより、アクティブセッションにおける LCS とのマッチ位置を考慮することで、結果が改善されることが確認できる。

### 5.5.3 考察

図 6~8 のように、得点付けの式において、 $l_i$  の重み ( $= \beta$ ) をかなり大きくしても精度は向上する。従って、今回対象としたサイトでは、アクティブセッション中で後方に現れるページはその後のアクセスに対して大きな関連性を持つことがわかる。

また、図 2, 6~8 より、同じ長さのアクティブセッションから推薦を行った場合、関連ルールによる推薦に比べ提案手法が優れていることがわかる。このことから、関連ルールを用いて共起頻度のみを考慮するよりも、順序情報を用いた方が、より正確な推薦を行うことができると考える。

WRAPL-FL 法を用いて推薦ページの優先順位付けを行った場合に比べ、WRAPL-FLP 法を用いた場合にはその差は狭まるものの、やはりアクティブセッション長を短く取った方が推薦精度が良い。一方で、アクティブセッションとの一致要素数が 1 の LCS を除いた場合にアクティブセッション長が長いもので良い結果が得られた。このことから、今回対象とした Web サイトでは、直前のアクセスページ(実際にリンクが張られているページ)からの推薦が最も精度が良いという特徴があると考えられる。現在この Web サイトは存在しないためリンク構造等を解析することはできないが、ページ配置が細分化、階層化されており URL 数が非常に多く、また短いセッションが非常に多いことから考えて、目的のページまで迷わずにナビゲートをするユーザが多いためにこのような傾向が現れると推測する。

## 6. おわりに

本稿では、アクセスログ解析によって抽出した LCS を用い、ユーザの過去のアクセス行動からそれに続くアクセスページを推薦する手法である WRAPL について解析を行った。

実際のアクセスログを用いて、その一部から LCS を抽出し、残りのログに対して WRAPL-FL 法を用いて推薦するシミュレーションを行った。その際、条件を限定して実験を行うことにより解析し、手法の問題点や改善点について検討した。さらに、それらから得た知見を活かして推薦ページの優先順位付け方法を改良した WRAPL-FLP 法を提案し、それによって推薦精度が向上されることを確認した。

本研究の今後の課題について述べる。まず、本研究の最終的な目的は、Web サイトの利用者に対し、リアルタイムに Web ページを推薦することであるため、今後は、計算量の見積もりや他手法との比較、さらに、

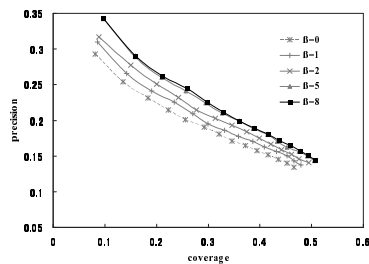


図 6 一致位置の考慮による改良 (n=2)

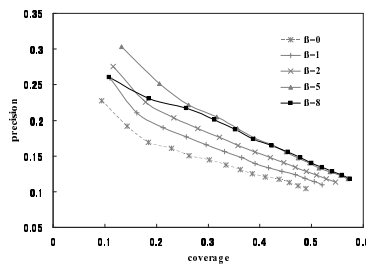


図 7 一致位置の考慮による改良 (n=3)

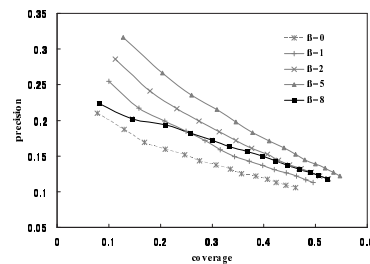


図 8 一致位置の考慮による改良 (n=4)

計算量を削減するための手法についても検討する必要がある。

今回のシミュレーションで採用した評価方法では、アクティブセッションより後にアクセスされたページ全てを正解としたため、サイト内で実際にリンクが張られており、直後にアクセスされやすいページを推薦し、それが正解とされるケースが多かった。しかし、階層化された商業サイト等でコンテンツページを優先的に推薦する場合においては、サイト構造における深さなどを考慮すべきであると考えられる。したがって、正解ページにも優先順位を付け、それに応じて評価する等、目的に合わせて評価方法も工夫すべきであろう。

さらに、Web サイトの規模や構造の違いによりユーザのアクセスパターンに異なる傾向がある場合、それに適した推薦ページの優先順位付けの方法について再考する必要がある。本手法を他の Web サイトに適用し、サイトの特徴と推薦精度の関連性を調査することも今後の課題である。

謝辞 本研究の一部は、文部科学省科学研究費補助金特定領域研究(18049026)、東京工業大学 21 世紀 COE プログラム「大規模知識資源の体系化と活用基盤構築」および独立行政法人科学技術振興機構戦略的創造研究推進事業 CREST の助成により行なわれた。

#### 参 考 文 献

- 1) Eirinaki, M. and Vazirgiannis, M.: Web mining for web personalization, *ACM TOIT*, pp.1-27 (2003).
- 2) 土方嘉徳: 情報推薦・情報フィルタリングのためのユーザプロファイリング技術, 人工知能学会誌, Vol.19, No.3, pp.365-372 (2004).
- 3) Agrawal, R. and Srikant, R.: Fast Algorithms for Mining Association Rules, *Proc. of 20th Intl. Conf. on Very Large Data Bases*, pp.487-499 (1994).
- 4) Sarawagi, S., Thomas, S. and Agrawal, R.: Integrating association rule mining with relational database systems: alternatives and implications, *Proc. of ACM SIGMOD Intl. Conf. on Management of data*, pp. 343-354 (1998).

- 5) Mobasher, B., Cooley, R. and Srivastava, J.: Creating Adaptive Web Site Through Usage-Based Clustering of URLs, *Proc. of the 99 Workshop on Knowledge and Data Engineering Exchange*, pp.19-25 (1999).
- 6) Srikant, R. and Yang, Y.: Mining web logs to improve website organization, *Proc. 10th Intl. Conf. on WWW*, pp.430-437 (2001).
- 7) Pitkow, J. and Pirollo, P.: Mining longest repeating subsequences to predict WWW, *Proc. of the 1999 USENIX Annual Technical Conf.* (1999).
- 8) Mobasher, B., Dai, H., Luo, T. and Nakagawa, M.: Effective personalization based on association rule discovery from web usage data, *Proc. 3rd Intl. Workshop on Web information and data management*, pp. 9-15 (2001).
- 9) 宇根田純治, 横田治夫: Web ログの共通シーケンス解析, 信学技報DE2002-2, 電子情報通信学会 (2002).
- 10) 戸田誠二, 横田治夫: LCS を用いたアクセスログ解析の並列処理による性能向上, 第 13 回データ工学ワークショップ論文集, DEWS2004 7-B-5 (2004).
- 11) 山元理絵, 小林 大, 小林隆志, 横田治夫: Web アクセスログの LCS を用いた Web ページの推薦手法, 信学技報DE2006-40, 電子情報通信学会 (2006).
- 12) Banerjee, A. and Ghosh, J.: Concept-based clustering of clickstream data, *Proc. 3rd Intl. Conf. on Information Technology*, pp.145-150 (2000).
- 13) Cormen, T.H., Leiserson, C.E., Rivest, R.L. and Stein, C.: *Introduction to Algorithms*, MIT Press (1990).
- 14) Wu, S., Manber, U., Myers, G. and Miller, W.: An O(NP) sequence comparison algorithm, *Information Processing Letters*, Vol.35, No.6, pp.317-323 (1990).