

WWW 上で文献検索プランニングを行なうソフトウェアエージェント NaviPlan

NaviPlan: a Software Agent for Planning Concept Understanding on WWW

大澤 幸生*¹ 山田 誠二*²
Yukio Ohsawa Seiji Yamada

- * 1 大阪大学大学院基礎工学系研究科
Graduate School of Engineering Science, Osaka University
- * 2 東京工業大学大学院総合理工学研究科知能システム科学専攻
CISS, Interdisciplinary Graduate School of Science and Engineering, Tokyo Institute of Technology

19YY 年 MM 月 DD 日 受理

Summary

We propose navigation planning to automatically generate a sequence of Web pages by which a user systematically understand a target concept. First, with a planning framework, we formalize the browsing task in the WWW. Action is defined as the understanding of a Web page, and an operator for a Web page consists of conditional/effect knowledge. Here, we developed and employed a method to generate an operator from a Web page by extracting condition/effect terms with keyword extraction techniques. The implemented the navigation planning system is tested by comparing with methods using a search engine and link-tracing (like a Web robot). As results, we found out navigation planning is a promising approach to assist the concept understanding in the WWW.

1. はじめに

WWW 上の Web ページを見ることにより、“ある概念(目標概念と呼ぶ)を理解する”という作業が、多くのユーザにより日常的に行われている。しかし、この概念理解に有用な Web ページを探す作業は、検索エンジンや Web ロボット [Cheong 96] を利用することにより部分的には効率化、自動化できるものの、結局検索されたページを見てチェックするところはユーザがやらねばならず、それによる時間のロスは大きい。

本研究では、概念理解のための有用な Web ページの探索を人工知能におけるプランニングとして捉え、その作業全体を自動化することを目的とする。ここで、プランとは、目標概念を理解するために、ユーザが見るべき Web ページの系列である。このように、目標概念を理解するための Web ページの系列を生成するプランニングを、ナビゲーションプランニングと呼ぶ。

ナビゲーションプランニングでは、一つの Web ページ

を見て理解することを、プランニングにおける一つの行為として捉え、その行為を U-オペレータとして定義する。本研究で構築したナビゲーションプランニング・システム *NaviPlan* はプランニング過程に必要に応じて、Web ページから U-オペレータを自動生成する。

従来の Softbot [Etzioni 94], Occam [Kwok 96] などは、情報収集の手続きプランとして、UNIX コマンドやデータベースの検索コマンドの系列を自動生成する。一方、*NaviPlan* はユーザの概念理解を誘導するようなプランを、ユーザが見るべき Web ページの系列で表現する。Web ロボット [Cheong 96] も検索エンジンの URL データベース構築のために URL 情報を自動収集するが、リンクの張られている Web ページ間を探索するに過ぎない。*NaviPlan* は、後述するようにリンクが明示的に張られていない Web ページも探索することが可能であり、目標概念の理解のために適切に制御される。又、Web Watcher [Armstrong 95], Letizia [Lieberman 95] のように次に見るべきページを提示するだけでなく、*NaviPlan* は Web ページ系列を

生成 / 提示する .

2. ナビゲーションプランニング

ユーザが目標概念を理解するために有用な Web ページをブラウズする行為は、以下のような対応によりプランニング [Russell 95] として定式化できる .

- 行為 : Web ページに記述されている概念を理解する .
- 状態 : ユーザの知識状態 . 既知の概念を表す単語の集合により記述 .
- 初期状態 : ユーザの初期の知識状態 .
- 目標状態 : ユーザが理解したい目標概念 . 単語の集合により記述 .
- オペレータ : Web ページを見て、知識を獲得する行為を表す U-オペレータ $U-Op(URL)$ は、以下の要素からなる .
 - ラベル : Web ページの URL でラベル付け .
 - 条件知識 : その Web ページを理解するために必要な知識 $C = \{c_1, \dots, c_i\}$. C の要素 (Web ページ中の単語である) を条件語という .
 - 効果知識 : その Web ページを理解することにより得られる知識 $E = \{e_1, \dots, e_j\}$. E の要素 (Web ページ中の単語) を効果語という .

3. Web ページからの U-オペレータ生成

NaviPlan は U-オペレータを自動生成する . というのは、世界中の膨大な数の Web ページすべてについて、U-オペレータを生成しておくのは、現実的には不可能だからである . 以下 U-オペレータをどのように抽出するか述べる .

[タグ構造による抽出] 情報検索の分野において提案されたキーワード抽出法のほとんどは単語の出現頻度による . しかし多くの Web ページは HTML により構造化されているので、ここでは、そのタグ構造 ($\langle TITLE \rangle$, $\langle Hn \rangle$, $\langle A HREF = \dots \rangle$) を利用する . まず $\langle A HREF = URL \rangle$ と $\langle /A \rangle$ の間の単語を条件語の候補とする . というのは Web ページにリンクされている単語は、その Web ページを理解するために重要なことが多いからである . また Web ページのタイトルと見出しは、そのページの主張点を表していることが多いので、 $\langle TITLE \rangle$ と $\langle /TITLE \rangle$ および $\langle Hn \rangle$ と $\langle /Hn \rangle$ の間の単語を効果語の候補とする .

しかし、HTML タグだけでは条件 / 効果語を高い精度で求めることができない . そこで、次のキーワード

抽出法 *KeyGraph* を併用することにする .

[*KeyGraph* による抽出] *KeyGraph* [Osawa 98] では、まず文章の土台となる概念を表す単語を得るので、これらを *NaviPlan* では条件語の候補とする . 次に *KeyGraph* では、その文章の主張と見なされる単語を得るので *NaviPlan* ではこれらを効果語の候補とする . *KeyGraph* ではこれらのキーワードを条件 / 効果らしさを表わす実数値 (図 1) と共に得ることができる . *KeyGraph* の併用により、タグ構造のないプレーンテキストからの条件 / 効果知識の抽出も可能となる .

ナビゲーションプランニングでは以下のように、条件知識、効果知識を、この *KeyGraph* とタグ情報の組み合わせによって抽出する .

4. プランニング手続き

ナビゲーションプランニングの手続きを、図 2 に示す . 簡単にいうと、目標概念を起点として U-オペレータをつないで行くのである . 付加的な入力として、文脈語を用いている . 文脈語は、目標概念が属する領域を表し、同じ単語でも使われる領域により意味が異なる場合に付加することで、探索の精度を上げる効果がある .

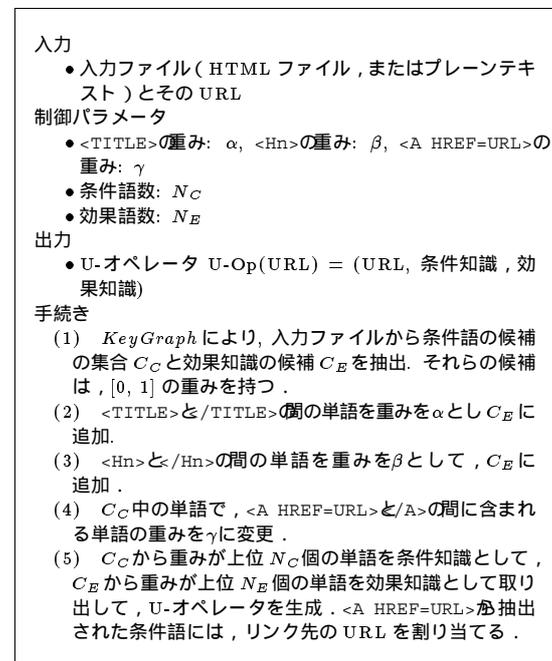


図 1 Ext-OP: U-オペレータ生成手続き

5. NaviPlanの実験による性能評価

NaviPlanは、Perlで記述されており、ユーザとのインタフェースとしてWebブラウザを用いる。図3に、目標概念が“concept formation”（概念形成）で、文脈御“AI”の場合のNaviPlanの出力プラン（深さ制限 $l = 4$ ）を示す。Webページ1は、目標概念に直接結びつくものであり、概念形成についてILP（Inductive Logic Programming）を基にして説明されている。しかし、ILP自体についての説明はページ1にはない。ILPについては、Webページ2~4のうち、3で詳しく述べられ、4には関連する機械学習の論文アブスト

入力

- 目標状態（= 目標概念） G_0 ：ユーザが理解したい概念を表す単語の集合。
- 文脈語 GC ：目標状態の背景領域を表す単語集合。

制御パラメータ

- 初期状態 IS ：ユーザの初期の知識状態。
- 深さ制限 l ：プランの最大長。
- Webページの制限数 w ：検索エンジンにより得られる関連Webページの最大数。
- ビーム探索の幅 b
- 状態ノードの評価関数 H

出力

- プラン P ：U-オペレータ $U-Op$ (URL) の系列。

手続き

- $N_0 = [n_0] = [(G_0, [])]$, $i = 0$ で初期化 (N_i は、深さ i での状態ノード系列。状態ノードは、 $n_i = (G_i, P_i)$ で記述される。ここで、 G_i は副目標を表す単語の集合、プラン P_i は、U-オペレータ系列。)
- もし $i = l$, または、 N_i が空なら、プラン P を出力して終了。
- 以下の手続きを、 N_i のすべてのノード n に対して適用して、 $N_{i+1} = []$ で初期化。
 - 差異 $D = G \cap IS$ を抽出。
 - 検索エンジンに $D \cup GC$ を入力して、 w 個の関連WebページのURLを得る。
 - 得られたURLのhttpサーバとのTCP/IP接続により、そのWebページ(HTMLファイル、またはテキストファイル)を獲得。
 - 得られたWebページから、手続き $Ext-OP$ を使って、U-オペレータを生成。そして、効果知識 E と n の副目標との交わりが空でないU-オペレータに対して、下の手続きを適用。
 - そのU-オペレータの E を n の G から取り除いたものに、U-オペレータの C を追加することにより、 G' を得る。
 - n のプラン P の先頭に、そのU-オペレータを加えることにより、 P' を得る。
 - 状態ノード (G', P') を N_{i+1} に追加。
- N_{i+1} の各状態ノードにつき、評価関数 H による評価を行い、その上位 b 個のノードで、 N_{i+1} を更新。
- $i = i + 1$ として、(2)へ。

図2 ナビゲーションプランニング手続き

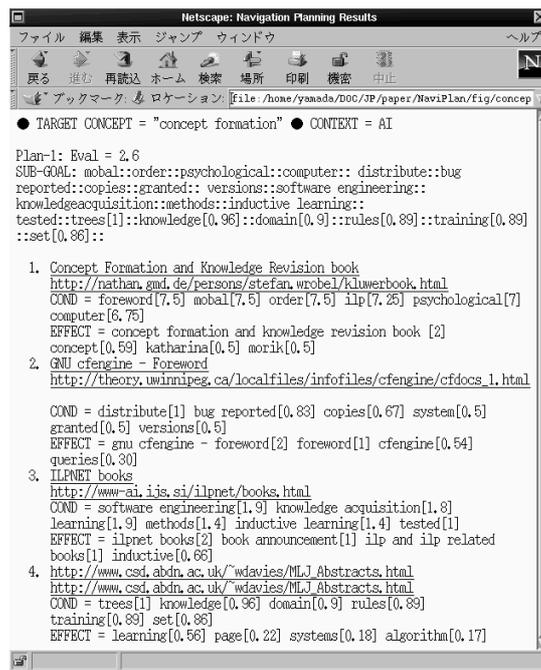


図3 “Concept formation” のプラン

ラクトが紹介されている。実際、ユーザが1, 3, 4のWebページを読むことによって概念形成はもとより、それとILP、機械学習との関係まで理解ができる。

以下、NaviPlanの性能を、実際にシステムを用いた結果から評価した結果を示す。NaviPlanの中で用いる検索エンジンと以下の実験で比較に用いる検索エンジンとして、広範囲に検索可能な検索エンジンであるMetaCrawler*1を用いた。

5.1 実験方法

われわれは、以下の5つの手法を比較した。

手法1: 検索エンジンのみ 目標概念と文脈語を検索キーとして入力する。ユーザは出力されたWebページを、ソートされた優先順に目標概念を理解するまで読み進める。

手法2: 検索エンジン+リンクの利用 手法1と同様の入力と出力ページをまず得る。次にユーザは、手法1と同様に優先順位に従ってWebページをゴールを理解するまで読み進めるが、途中でページ中のリンクに従って他のページに移っても良い。

手法3: NaviPlan ユーザは、NaviPlanに目標概念と文脈語を入力する。そして、読み進めるべきとNaviPlanが判断した順番に目標概念を理解するま

*1 <http://www.metacrawler.com/>

で読み進める。

手法 4: NaviPlan + リンクの利用 ユーザは、手法 3 と同様の入力を NaviPlan に与え、その結果として得られたページを、手法 3 に加えて手法 2 と同様のリンク利用を行いながら目標概念を理解するまで読み進める。

この実験において設定した NaviPlan のパラメータは、図 2 では、 $IS = []$, $l = 4$, $w = 8$, $b = 4$ で、図 1 では、 $\alpha = 2$, $\beta = 1$, $\gamma = 1$, $N_C = 6$, $N_E = 4$ とした。一度のプランニングに要した時間は 2 ~ 3 時間であったが、そのほとんどの時間は WWW サーバとの通信による HTML ファイルの獲得に費やされた。

5.2 実験結果と評価

43 通りの検索キーについて上記の 5 つの手法を比較した結果、手法 1 と手法 3 については、ともに 26 通りの目標概念が正しく理解された。リンク利用を手法 1 と手法 3 と併用した手法 2 と手法 4 ではそれぞれ、33 通りと 32 通りに増加したので、リンク利用が有効であることがわかる。

次に、ユーザが目標概念理解までにどれだけ労力を使ったか、すなわちどれだけ多くの Web ページを読んだかを評価した。リンクの利用がユーザの理解を促進することは先述のとおりであるから、リンク利用を合わせた手法 2 と手法 4 との比較を行なった。図 4 は、そのデータを示すグラフである。横軸は、検索エンジン (MetaCrawler) で得られた出力ページの数 (検索語がどれだけ Web で普及しているかを表わしていると考えられる) であり、縦軸は、実際に目標概念を理解できるまでに読む必要があった Web ページ数である。図 4 から以下の 2 つの特徴が得られる。

- MetaCrawler (手法 2) による検索よりも NaviPlan (手法 4) によるプランニングがユーザの効率的な理解を助ける効果がある。
- Web 上でよく普及した目標概念では、手法 2 は多くのページをユーザに読ませてしまう。このことは、その概念が広く知れ渡っていて、改めて直接定義する (すなわち MetaCrawler で得られる) Web ページが少ないせいであろう。

6. まとめと課題

ユーザが理解したい目標概念の概念理解に有用な Web ページの系列を提示することができるナビゲーションプランニングを提案し、その実装システム NaviPlan を構築した。WWW 上でのユーザのブラウジングをプ

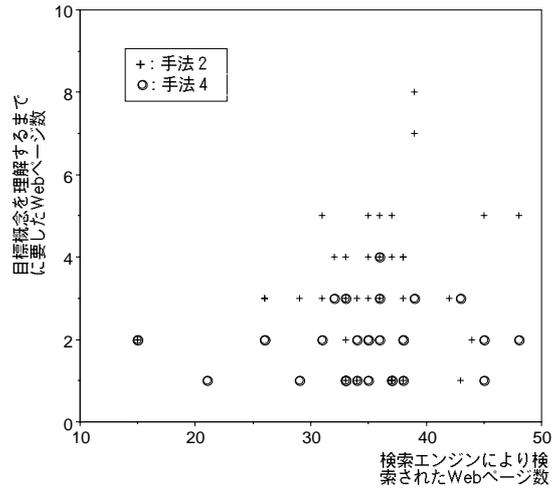


図 4 手法 2 と手法 4 の比較。

ランニングの枠組で定式化し、必要に応じて Web ページから U-オペレータを自動生成して、プランニングを行い、有用な Web ページの系列を得る。検索エンジン等との実験的比較を行い、NaviPlan の有用性を示した。

しかし、現状では U-オペレータの精度は十分に高くない。その原因は、U-オペレータで用いたヒューリスティックにある。例えば、見出しのタグ <Hn> を利用して効果知識の抽出を行っているが、見出し情報を用いて得られた効果知識が信頼できないことは [Osawa 98] でも指摘されている。このヒューリスティックの精度を高めることは、NaviPlan の第一の検討課題である。

参考文献

- [Armstrong 95] Armstrong, R., et al.: WabWatcher: A Learning Apprentice for the World Wide Web, *The AAAI Spring Symposium on Information Gathering from Heterogeneous, Distributed Environment*, (1995).
- [Cheong 96] Cheong, F.C.: *Internet Agents: Spiders, Wanderers, Brokers, and Bots*, New Riders (1996).
- [Etzioni 94] Etzioni, O. and Weld, D.: A SoftBot-based Interface to the Internet, *Communication of the ACM*, Vol.37 No.7, 72-76 (1994).
- [Kwok 96] Kwok C. T. and Weld, D. S.: Planning to Gather Information, *Proc. of AAAI-96*, 32-39 (1996).
- [Lieberman 95] Lieberman, H.: Letizia: An Agent that Assists Web Browsing, *Proc. of IJCAI-95*, 924-929 (1995).
- [Osawa 98] Ohsawa, Y. et al.: KeyGraph: Automatic Indexing by Co-occurrence Graph Based on Building Construction Metaphor, *Proc. of IEEE Advanced Digital Library Conference*, 12-18 (1998)
- [Russell 95] Russell, S. and Norvig, P.: *Artificial Intelligence - A Modern Approach*, Prentice-Hall (1995).

[担当編集委員: × × , 査読者: × ×]