

個人の感性モデルに基づく 対話型遺伝的アルゴリズムを用いた記事推薦システムの提案

Suggestion of the article recommendation system based on the user's Kansei model using interactive Genetic Algorithm

宮地 正大*¹ 廣安 知之*² 三木 光範*³ 横内 久猛*²
Masahiro MIYAJI Tomoyuki HIROYASU Mitsunori MIKI Hisatake YOKOUCHI

*¹同志社大学大学院工学研究科 *²同志社大学生命医科学部
Graduate School of Engineering Doshisha University Department of Life and Medical Sciences

*³同志社大学理工学部
Department of Science and Engineering Doshisha University

In this paper, we propose the article recommendation method using the user's Kansei model that is derived by interactive Genetic Algorithm. In this research, we assume that each user has own Kansei model and user's preference is related to the function of Kansei model. In the proposed method, a novel neighborhood definition and degree of similarity between words by the concept network in the Kansei model is proposed. Using the real data of web report, an example of neighborhood definition and degree of similarity are described and some results are examined.

1. はじめに

情報化技術の発展により、Web上に膨大な量の情報がテキスト（文書）として公開されている。これらの大規模なデータから情報を閲覧するユーザが、求めるコンテンツを的確に得ることは難しい。そのため、一部のニュースサイトなどではユーザ個別に推薦内容を変化させるパーソナライズされた推薦システムが組み込まれている [Das 07]。個人の感性をモデル化することで利用者個別の嗜好に合わせた推薦が可能であると考えられており、様々な研究が行われている [Balabanović 98]。

本研究では、専門性の高い用語が頻出するレポート（テキストデータ）群で構成される Web 記事において、ユーザの感性に近い推薦を行うことを目標とする。そのため、レポートの特徴ベクトルを設計変数とした対話型遺伝的アルゴリズムの要素を取り入れる。

従来の遺伝的アルゴリズムを用いた記事推薦手法では、文書中に出現する全単語を設計変数とし、重みをユーザの行動履歴から遺伝的操作により最適化することで、ユーザの求めるトピックを予測する手法を用いている [Sheth 93, Atsumi 97]。しかしこの手法では、設計変数の数が膨大になるため、解探索の精度が悪化することが考えられている。本研究では、設計変数として用いる単語数に制限を設けることで、解探索精度の悪化を防ぐ。通常、異なる設計変数同士には互いに関連の定義がされていないために、遺伝的操作を行うことができない。そのため、概念語ネットワークによる単語間の類似度を近傍の定義とすることで、異なる次元間での遺伝的操作を可能とする手法を提案する。

2. 推薦システム

2.1 協調フィルタリング

Amazon[Linden 03] や Google[Das 07] などを用いられており、多数のユーザの中から行動履歴の類似したユーザを抽出することで、そのユーザの参照したコンテンツを推薦し合う方式

である。他のユーザとの好みの類似性を基本としたアルゴリズムであるため、ユーザが未知の商品が推薦される場合がある。また、推薦・予測にコンテンツ自体の先験情報が必要ないため手軽に導入が可能であり、多分野のコンテンツが入り交じるシステムでの利用が可能である。しかし、システムの条件としてすべてのコンテンツをユーザが評価する必要があるため、ユーザが多数いることが必要となる。そのため、誰も評価していないコンテンツは推薦される可能性が低くなり、推薦されるコンテンツが集中する問題点もある [Herlocker 00]。

2.2 内容ベースフィルタリング

ユーザの行動履歴とコンテンツに含まれるメタ情報（著者・出版社・内容など）を特徴ベクトルとしてマッチングさせる手法である。推薦システムがユーザ評価を必要としないため、全コンテンツを公平に推薦対象とすることができ、小規模なシステムへの導入が可能である。

2.3 推薦システムのパーソナライズ

ユーザ個人の感性をモデリングする手法として前述の内容ベースフィルタリングが有用であると考えられる [Pazzani 07]。ユーザの嗜好に応じてパーソナライズされた推薦を行うためには、ユーザの嗜好を表す感性モデルを何らかの方法で学習・予測する必要がある。代表的な手法としては、確率推論に基づく手法であるベイジアンネットワークや隠れマルコフモデル、ユーザの求める要素（感性）のパラメータを推定・最適化する手法である対話型遺伝的アルゴリズムなどが挙げられる [Tanaka 09, Tanaka 10]。

3. 対話型遺伝的アルゴリズム

3.1 概要

対話型遺伝的アルゴリズム (iGA: interactive Genetic Algorithm) は、多点探索の最適化アルゴリズムである GA をベースとした対話型最適化手法である。人間の感性のモデルを設計変数空間のランドスケープとして捉え、その空間における最良点、もしくは最良域を探索する。iGA を実装したシステムは、ユーザに対して、多数の候補解を提示し、ユーザは感性や好みに基づいてそれらを評価し、その評価値を用いてシステムは遺

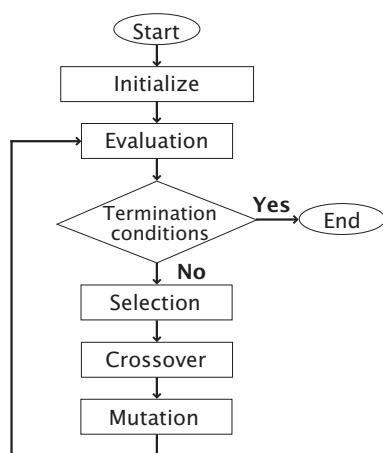


図 1: 遺伝的操作の流れ

伝的操作を適用する。これらの操作を繰り返すことで、集団全体をユーザの好むものへと変化させる。iGA は感性による評価を必要とするアプリケーションに利用されている。

3.2 アルゴリズムの設計

iGA では、他の最適化問題と同様に、最適化の対象とする候補解を設計変数として表現する。例えば、服飾デザイン支援システムであれば、デザインする服の形状や色、装飾などが設計変数として定義され、各解は、その設計変数のベクトルによって構成される。最適化を行う遺伝的操作のフェーズでは、さらにこの設計変数を 01 のビット列や遺伝子の型と実数値などの染色体に修正して用いる。

図 1 に遺伝的操作の流れを示す。まず、試行の最初に、染色体を多数含む母集団を初期化する。そして、この染色体一つ一つに対してユーザによる評価を行う。評価値の高い染色体を、親個体として選択し、これらに交叉として情報を組み替える操作を与えることで、より高い評価が期待される子個体を生成する。また、探索途中に局所解に陥ることを防ぐために確率的に突然変異を行う。これらの選択、交叉、突然変異を一連の流れを 1 世代の操作とする。何世代か繰り返すことで、徐々に評価の高い集団へと進化させていく。

4. 提案推薦システム

4.1 提案システム概要

本研究では 3 章で述べた iGA に基づく記事推薦システムを提案する。iGA により、ユーザの感性パラメータを予測することで、ユーザの感性に近い推薦を行うことを目標としている。提案手法の流れを以下に示す。

1. 推薦対象のコンテンツを特徴ベクトルとなる単語列に分解
2. ユーザの利用履歴から感性パラメータ候補を生成
3. 感性パラメータ候補に類似するコンテンツの提示
4. 手順 2,3 を繰り返す

記事推薦における iGA を用いたシステムのコンテンツ評価フェーズでは、通常の iGA とは違い、記事タイトルのみがユーザに推薦コンテンツとして提示されるため、提示されたすべての記事の内容を閲覧せずにコンテンツをユーザに評価してもら

表 1: 机と椅子の一次属性

| 概念 | 一次属性 |
|----|-------------------------------|
| 机 | (学校, 0.6) (勉強, 0.3) (本棚, 0.1) |
| 椅子 | (勉強, 0.5) (教室, 0.3) (木, 0.2) |

うことは難しい。そのため、提案手法では提示された記事タイトルの中からユーザが次に遷移・閲覧したコンテンツを、遺伝的操作における評価とする。コンテンツの特徴ベクトルの定義には先行研究で多く用いられている手法と同様に、文書内に出現する単語を特徴ベクトルとし、重み付けは TF・IDF 法を用いる [Salton 75, Pazzani 07]。本研究においては、それらの組み合わせ及びパラメータをユーザの感性パラメータの候補とする。しかし出現単語数が膨大な数になることから、解探索に悪影響を及ぼすことが予想される。そのため、膨大な量の設計変数を機械的に扱いやすい形に変形することで、解探索の性能を向上させる研究が行われている。その手法として、設計変数を主成分分析により、別の主成分へと写像することで次元数を削減する手法 [Tanaka 09, Tanaka 10] や、初期個体の生成時に予め SVM によるユーザ嗜好の学習を行わせることで個体の収束を早める手法 [Amamiya 09] などが考案されている。

本研究では、設計変数間に重みとは異なる関連度を定義することで、別次元同士の設計変数での遺伝的操作を可能とする手法を用いる。それにより、各コンテンツが全種類の設計変数を保持する必要がなくなる。単語間の関係性を数値で表す概念ベースを用いることで単語ネットワークを作成し、それらの組み合わせ及びパラメータを最適化対象とする。

4.2 特徴単語の重みの決定

レポートの特徴ベクトル項目として抽出した単語の重みを決定する必要がある。その手法として、TF・IDF 法を用いる。TF・IDF 法とは、単語の出現頻度 (Term Frequency: TF) 及び逆文書頻度 (Inverse Document Frequency: IDF) を用いた文書内での単語の重要度を表す指標であり、それぞれ式 1~式 3 として表される。

$$tf(i, j) = \frac{n_{i,j}}{\sum_k n_{k,j}} \quad (1)$$

$$idf(i) = 1 + \log_2 \frac{|D|}{|\{d : d \ni t_i\}|} \quad (2)$$

$$tfidf = tf \times df \quad (3)$$

$n_{i,j}$ は単語 i の文書 j における出現回数、 $|D|$ は総ドキュメント数、 $|\{d : d \ni t_i\}|$ は単語 i を含むドキュメント数である。同一ドキュメント中で頻出な単語は TF 値が高くなり、多くのドキュメントで用いられている単語は IDF 値が低くなる。

4.3 概念ベースによる単語ネットワークの作成

概念ベースとは、ある単語の意味 (概念) をその単語に関連の深い単語群 (属性) で定義した単語空間のことである。例えば、人間は [学校] から「教師」や「生徒」などの単語を連想できる。この場合、概念ベースでは [学校] の概念を {教師: 0.4, 児童: 0.3, 生徒: 0.1, 学生: 0.1, ...} のように概念の重みあり属性群として格納している。また、その属性を一次属性と呼び、さらにその属性語の概念を展開したものを二次属性と呼ぶ。例として、表 1、表 2 に机と椅子の一次属性及び二次属性を展開した表を示す。上記の属性を再帰的に展開することで、すべての単語における関連度を表したネットワークが表現され

表 2: 机と椅子の二次属性

| 一次属性 | 二次属性 |
|------|-------------------------------|
| 学校 | (大学, 0.4) (校舎, 0.4) (木造, 0.2) |
| 勉強 | (予習, 0.5) (試験, 0.3) (本, 0.2) |
| 本棚 | (図書, 0.6) (書物, 0.3) (本, 0.1) |
| 教室 | (教師, 0.4) (校舎, 0.4) (生徒, 0.2) |
| 木 | (森林, 0.5) (木造, 0.4) (葉, 0.1) |

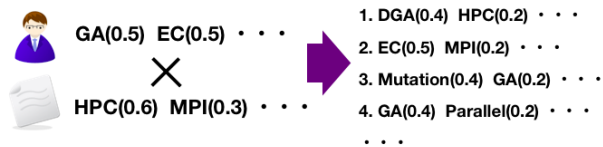


図 2: 感性パラメータの生成

る。本稿では、この概念ベースによる単語ネットワークを、概念語ネットワークと呼ぶ。概念ベースの自動構築手法なども考案されており [Watabe 01]、本研究ではレポートデータから、出現単語の共起確率を用いて作成する。解析対象とするレポートには専門性の高い一般的ではない名詞が含まれているため、単語辞書に専門用語を追加する必要がある [Hiroaki 03]。本研究では予め重要な専門用語を登録している。

4.4 感性パラメータからの推薦レポートの決定

本研究における感性パラメータは任意の数の単語及び重みで構成される。ユーザがレポートへのアクセスを行った際に、ユーザの感性パラメータと閲覧レポートのもつ特徴パラメータから新たに感性パラメータの候補となる複数の単語の組み合わせを生成する。図 2 に候補パラメータ生成の例を示す。

図 2 では現在のユーザの感性パラメータが [GA(0.5), EC(0.5)] であるとき、[HPC(0.6), MPI(0.3)] の特徴ベクトルを持つレポートを閲覧した状況を表している。GA と HPC, EC と MPI という単語の組み合わせから新たに別の単語を生成し、それらを感性パラメータの候補としている。

レポートへの初回アクセス時には、レポートが持つ特徴パラメータを重みの合計が 1 になるよう、正規化した状態で入力される。図 3 に交叉処理、図 4 に突然変異処理の例を示す。図 3 の親単語 A・B の交叉では単語の関係ネットワーク上での最短経路上からルーレット選択により選択を行う。この際の重みは A・B のノード間で線形となるように与える。また、パラメータ候補から推薦レポートへは、すべての出現単語で構成されるベクトル空間上のユークリッド距離で最も類似しているレポートを提示する。

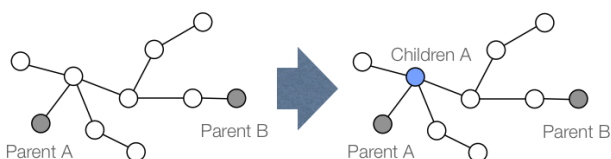


図 3: 交叉の例

5. 実験

5.1 実験目的および実験環境

本システムによってユーザの履歴から嗜好を学習し、類似するキーワードを主題とする記事が推薦結果に現れることを明らかにする。

実験は IS Report システムの公開レポートデータを対象に行った。IS Report システムは同志社大学医療情報研究室が管理する研究レポート公開システムである*1。

公開されているレポートとして研究に関する基礎知識や文献調査、研究報告などの内容であり専門性の高い用語が多数含まれている。

5.2 実験内容

本システムによるユーザのパラメータ推定を加えた推薦レポートと、パラメータ推定を加えないレポートの持つ主キーワードのみを特徴量として用いる手法を比較する。本実験では提案手法において、類似するキーワードを主題とする記事が推薦結果に現れることを明らかにするため、概念語ネットワーク導入以外の要素を極力排除した条件で行う。感性パラメータとなる子個体生成時の突然変異率は 0、学習に用いる感性パラメータ数は 1、推薦レポート数は 5 とした。なお、学習させるパラメータ数が 1 のため、感性パラメータの重みが 1 に正規化され、単語の選択確率は同一となる。

実験には「The Compact Genetic Algorithm (主キーワード: GA)」レポートを用い、「PC クラスタ管理システム Sun Grid Engine の概要 (主キーワード: ジョブ)」を学習済みのレポートとして用いた。

5.3 実験結果

表 3 に学習を行わない手法、表 4 に提案手法においての「The Compact Genetic Algorithm」レポート閲覧時の推薦レポート及び、レポートが主とするキーワードの例を示す。提案手法においては「PC クラスタ管理システム Sun Grid Engine の概要」レポートを学習済みデータとして与えており、提案手法によって推定された感性パラメータも記載する。表 3 に示される学習システムを組み込んでいない推薦結果では、本来の「The Compact Genetic Algorithm」レポートの主キーワードである [GA] のみをキーワードに全レポートから類似度の高い順に並べた結果である。そのため、遺伝的アルゴリズム (GA: Genetic Algorithm) に関連深いレポートが推薦結果として現れている。対して、表 4 に示される学習システムを組み込んだ推薦結果では、学習済みレポートの主キーワードである [ジョブ] と閲覧中のレポートの主キーワードである [GA] について提案手法によって感性モデルとして新たなキーワードを生成した上で、そのパラメータに類似するレポートを推薦結果として表示している。

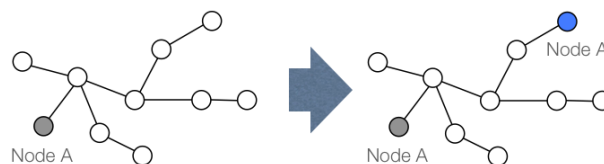


図 4: 突然変異の例

*1 <http://www.mis.doshisha.ac.jp/isreport>

表 3: 学習なし推薦レポート

| 推薦レポートタイトル | 主キーワード |
|---|------------------|
| PSA/AT(GA) の (Ala)10 への適用 自作 SGA と ga2k との解探索の性能比較 SGA の作成と動作確認 | AT ga2k GA |
| MGG と sGA (simple GA) の性能比較 環境分散遺伝的アルゴリズムの検証 | MGG 分散 GA |

表 4: 学習あり推薦レポート

| 推薦レポートタイトル | 主キーワード | 感性モデル |
|--|-----------------|--------------------|
| PSA/AN(GA) における遺伝的操作 PTH における力場パラメータの検討 | AN パラメータ | GA パラメータ |
| 大規模ジョブ問題におけるコーディング cron によるプログラムの自動実行 ランキングにおける角度パラメータ | ジョブ 実行 個体 | ジョブ 実行 パラメータ |

5.4 考察

感性パラメータの学習システムを組み込むことによって、推薦結果に違いが見られたが、その要因となっている感性パラメータの生成について考察する。本手法による学習とは、現在閲覧中のレポートのキーワードと、過去に推定された感性パラメータとなるキーワードを概念語ネットワーク上で最短経路となる単語の一つに、次世代の感性パラメータを遷移することを指す。本実験で感性パラメータ生成の際に用いた、[GA]と[ジョブ]の単語ネットワーク上での最短経路は [GA]-[パラメータ]-[実行]-[ジョブ] で構成されており、それぞれの選択確率はすべて同一である。そのため、今までの閲覧レポートの特徴を引き継いだ上で、その単語同士の概念上で間にあたる [パラメータ] や [実行] といったキーワードが用いられた推薦が行われている。しかし、本実験では学習単語数を主キーワード 1 つに制限していたため、極端な推薦結果が得られたと言える。そのため、複数の単語を学習させるとにより、複雑な特徴を受け継いだ推薦がされると期待される。

6. まとめ

本稿では、ユーザの感性を対話型遺伝的アルゴリズムによって学習し、その情報を用いることでパーソナライズされた記事推薦を実現するシステムを提案した。コンテンツに含まれるメタ情報を用いる内容ベースフィルタリングの推薦システムで問題となっていた出現単語を特徴ベクトルとした場合の一単語の寄与度の低さを、対話型遺伝的アルゴリズムを用いて解消するシステムを提案した。すべての単語を最適化対象とするのではなく、感性パラメータとなる特徴単語を、概念語ネットワークを用いることによって推定した。結果、一度別のレポートを閲覧することにより同一のレポートを閲覧した場合でも、閲覧済みレポートの特徴を残した推薦が行われることを確認した。

参考文献

[Amamiya 09] Amamiya, A., Miki, M., and Hiroyasu, T.: Interactive Genetic Algorithm using Initial Individuals Produced by Support Vector Machine, *The Science*

and Engineering Review of Doshisha University, Vol. 50, No. 1, pp. 34–45 (2009), [In Japanese]

[Atsumi 97] Atsumi, M.: Extraction of User's Interests from Web Pages based on Genetic Algorithm, *Intelligence*, Vol. 97, No. 108, pp. 13–18 (1997)

[Balabanović 98] Balabanović, M.: Exploring versus exploiting when learning user models for text recommendation, *User Modeling and User-Adapted Interaction*, Vol. 8, No. 1, pp. 71–102 (1998)

[Das 07] Das, A. S., Datar, M., Garg, A., and Rajaram, S.: Google news personalization: scalable online collaborative filtering, in *Proceedings of the 16th international conference on World Wide Web*, pp. 271–280 (2007)

[Herlocker 00] Herlocker, J. L., Konstan, J. A., and Riedl, J.: Explaining collaborative filtering recommendations, in *Proceedings of the 2000 ACM conference on Computer supported cooperative work*, pp. 241–250 (2000)

[Hiroaki 03] Hiroaki, Y., Tatsunori, M., and Hiroshi, N.: Term Extraction Based on Occurrence and Concatenation Frequency, *Journal of natural language processing*, Vol. 10, No. 1, pp. 27–46 (2003)

[Linden 03] Linden, G., Smith, B., and York, J.: Amazon.com recommendations: item-to-item collaborative filtering, *Internet Computing, IEEE*, Vol. 7, No. 1, pp. 76–80 (2003)

[Pazzani 07] Pazzani, M. and Billsus, D.: Content-based recommendation systems, *The adaptive web*, pp. 325–341 (2007)

[Salton 75] Salton, G., Wong, A., and Yang, C. S.: A vector space model for automatic indexing, *Commun. ACM*, Vol. 18, No. 11, pp. 613–620 (1975)

[Sheth 93] Sheth, B. and Maes, P.: Evolving agents for personalized information filtering, in *Artificial Intelligence for Applications, 1993. Proceedings., Ninth Conference on*, pp. 345–352 (1993)

[Tanaka 09] Tanaka, M., Hiroyasu, T., Miki, M., and Yokouchi, H.: Extraction of Design Variables using Collaborative Filtering for interactive Genetic Algorithms, *2009 IEEE International Conference on Fuzzy Systems Proceedings* (2009)

[Tanaka 10] Tanaka, M., Hiroyasu, T., Miki, M., Yasunari, S., and Yoshimi, M.: Automatic Generation Method to derive for the design variable spaces for interactive Genetic Algorithms, in *2010 IEEE World Congress on Computational Intelligence (WCCI 2010)* (2010)

[Watabe 01] Watabe, H. and Kawaoka, T.: The degree of association between concepts using the chain of concepts, in *Systems, Man, and Cybernetics, 2001 IEEE International Conference on*, Vol. 2, pp. 877–881, IEEE (2001)