

ユーザの興味オントロジ生成に基づくブログコミュニティ作成支援

中辻 真[†] 三好 優[†] 大塚 祥広[†]

[†]日本電信電話株式会社 NTT ネットワークサービスシステム研究所
〒180-8585 東京都武蔵野市緑町 3-9-11
E-mail: †{nakatsuji.makoto,miyoshi.yu,otsuka.yoshihiro}@lab.ntt.co.jp

あらまし 近年、ユーザ興味を発信する手段としてのブログ利用が目覚しい。しかし、ユーザの興味情報を簡易に詳細化する手段がないため、大量のブログエントリが日々発信されているにも関わらず興味に即したエントリやユーザを発見できない事が多い。そこで、ユーザ興味をクラス階層表現し、各クラスの興味度をも備える興味オントロジを導入する。そして興味オントロジを簡易に生成するため、まず、サービスドメイン毎のオントロジへユーザの蓄積ブログエントリを分類する事で興味オントロジを自動生成し、次に、各ユーザが自身のオントロジを調整する、ファーストトップダウン/セカンドボトムアップなアプローチに基づくオントロジ生成手法を提案する。さらに、興味度やクラストポロジを考慮したオントロジ近似度計測手法を提案し、コミュニティ形成へと適用する。また、ブログポータル Doblog における大規模な実ブログデータと、音楽サービスドメインオントロジを用いた検証より、本手法が適切な興味オントロジ自動生成と、コミュニティ形成支援を実現できることを示した。
キーワード オントロジ生成、オントロジマッピング、情報推薦、コミュニティ形成、セマンティック Web、ブログ

Formalization Support of Weblog Community based on Automatic Generation of User Interest Ontology

Makoto NAKATSUJI[†], Yu MIYOSHI[†], and Yoshihiro OOTSUKA[†]

[†]NTT Network Service Systems Laboratories, NTT Corporation
9-11 Midori-Cho 3-Chome, Musashino-Shi, Tokyo, 180-8585 Japan
E-mail: †{nakatsuji.makoto,miyoshi.yu,otsuka.yoshihiro}@lab.ntt.co.jp

Abstract Recently, the use of Weblog is remarkable as the means to publish the user interest. However, there is no means to informationize the user interest in detail, it is difficult to find suitable information resources in spite of a large amount of weblog entries are published every day. In order to resolve above problems, we firstly classify user entries to service domain ontologies and create interest ontology which expresses user's interest as a hierarchy of classes with interest strength by top-down approach. Nextly, with bottom-up approach, users modify their interest ontologies to reflect their interest more in detail. Furthermore, we propose similarity measurement technique between two ontologies based on the interest strength and class topology, then try to adopt it to formalization of weblog communities. We evaluate the performance of our proposed automatic interest ontology generation and community formalization support based on large-scale weblog entries on Weblog portal Doblog and music domain ontologies.

Key words Ontology Creation, Ontology Mapping, Interest Score, Community Formalization, Semantic Web, Weblog

1. はじめに

近年、インターネット上でユーザの興味対象を発信しユーザ間での議論を促進するブログサービスや互いに友人として承認し合ったユーザ間で興味対象を議論する Social Networking (ソーシャル・ネットワーキング) サービスが注目されており、今後ますますユーザ数やこれらを利用したサービスは拡大していくと考えられる [13]。そして、この種の情報流通サービスは、

ユーザが自身の興味に近いユーザの発信記事やコミュニティでの議論内容を閲覧する事を通じ、各自の興味対象を拡大する基盤となる可能性を持つため、興味深い。

しかし、現状のブログサービスにおける情報検索は、goo^(注1)などの Web ページ検索エンジンや、RDF Site Summary (RSS) という簡単なメタデータ記述を利用したキーワード検索でしか

(注1) : <http://www.goo.ne.jp>

い、更に、個人の興味情報を簡易かつ詳細に生成する機能を備えていないため、自身の興味に即した検索目的語を適切に構成する必要があり、検索キーワードの選択に手間がかかる。また、事前に検索対象をある程度把握していないとキーワード自体を構成できないため、興味を持つ可能性があるがキーワードを特定できない場合は、情報検索自体ができないことも多い。

このようなユーザの興味情報生成や情報推薦に関する研究は、従来の Web 検索においても様々な研究が行われている [4,9]。例えば、個人の興味を検索ログや個人が分類した Web ページに対するブックマーク情報などから推定し検索に活用する研究がある [4]。しかし、個人のブックマークは様々な分野の興味を混在しているため、情報検索に適用する際の精度を阻害する。一方、Web 上の情報に対しその背景となる意味情報を機械処理可能なオントロジとして記述する事で、様々なソフトウェアが自動でオントロジによる処理を実行する事を旨とするセマンティック Web 技術 [2] をブログ検索に適用する研究として、Semblog [9] ではユーザの興味情報をオントロジとして構築し、オントロジマッピングに基づく検索フレームワークの実現を試みている。このように、セマンティック WEB 技術をブログ検索に適用する事の有効性が唱え始められているが、その核となるオントロジ構築が難しい。

こうした問題を解決するため本研究では、ユーザ興味をクラス階層表現し、各クラスの興味の強さを示す値である興味度をも備えた興味オントロジを自動生成する手法や、興味度やクラストポロジを考慮した興味オントロジの近似度計測手法を提案し、ブログコミュニティ生成への適用を試みる。そして、個人が興味を持つ可能性が高い未知な情報を意外な情報と定義し、一般的な Web 利用者が意外な情報を自然と発見し、新たな興味に組み込むことを支援する基盤を提供することを目的とする。

具体的には、個人の興味情報を簡易かつ詳細に抽出するため、まず、近年急速に利用が進んできたブログを利用し興味オントロジを自動生成する。つまりブログは一般的な Web ページや BBS と異なり、個人としての興味を記述していることが多いため、個人の興味特定に利用する。なお、ブログでは例えば音楽と映画など複数サービスドメインに跨り興味が混在した形態で自由記述されている事が多いため、各サービスドメイン毎の情報をクラス階層として記述したドメインオントロジをユーザの興味オントロジを作成する雛型となるオントロジ (雛型オントロジ) として予め用意し、それに対し個人のブログエントリを分類する事で、興味オントロジをサービスドメイン毎に分離しトップダウン的に生成する。その上で、ユーザは興味オントロジを、自身の興味により適合したものへとボトムアップ的に修正を行う。そして、複数ユーザの興味オントロジ間でクラス毎の興味度やクラスの接続形態であるトポロジを考慮し、オントロジ間の近似度を計測する。さらに、近似度が高いオントロジ間で一部トポロジが異なるクラスに属するエントリを、意外な興味エントリとしてユーザ推薦することで、ユーザの興味幅の拡大と、他ユーザ間とのコミュニケーション促進を狙う。

また、ブログポータル Doblog^(注2)における大規模データ (約 5 万 5 千ユーザ, 160 万エントリ) を用い提案手法の検証を行い、本提案が高精度な興味オントロジ生成やブログコミュニティ解析、およびユーザ毎の興味に即したエントリ推薦によるコミュニティ形成に対し有効性を持つことを確認した。

以下、2. 章では、本論文の背景となるブログの概要説明とその問題点について述べ、関連研究の紹介も行う。3. 章では、雛型オントロジを用い、ブログユーザ毎の興味オントロジを自動生成する手法を提案し、4. 章において、興味オントロジ間の近似度計測手法の提案と、ユーザ興味に即したコミュニティ作成支援への適用について述べる。5. 章では、実ブログデータを用いた興味オントロジ生成とユーザ分布解析、および推薦エントリの検証を行い、6. 章の結論と将来の課題で結ぶ。

2. ブログの概要と関連研究

ブログの定義は明確ではないが、主に MovableType^(注3)などのブログ作成ツールや、Doblog などの Weblog ホスティングサービスを用い個人が発信するニュースサイト、または日記サイトという位置づけで理解されることが多い。その特徴を説明すると [11,15]、(1) 個人が Web 上で発信する記事 (エントリ) の集合体であり、(2) エントリが時系列に表示されている。そして、(3) ブログ作成ツールを用いることで作成したエントリを簡単に公開できる。また、ブログ作成ツールや Weblog ホスティングサービス毎にいくつか異なる点はあるが、(4) サイト内外のエントリを元情報にしたエントリを発信する *trackback* といわれる機構を持つ。これにより、ブログサイトを跨りエントリー (スレッド) を構成する事ができ、ユーザ間のコミュニケーションやエントリ毎にテーマを絞った議論が行いやすい。それ以外に、(5) RDF Site Summary (RSS) と呼ばれるメタデータをエントリ記述の際に生成し、ブログの更新情報を集め提供しているサーバである *ping* サーバや Weblog ホスティングサービスの運営するサーバに登録することで、エントリの存在や簡単な内容、更新情報などを他のユーザへ公開できる。RSS は、Web サイトの各ページのタイトル、アドレス、要約などをメタデータ記述できるものであり、公開 RSS に対し RSS フィードというサービスを用いることで、多数 Web サイトの更新情報を効率的に把握できる。

このように RSS はメタデータを構成し、流通させる仕組みを持つため、Semantic Web におけるオントロジ普及に期待されている。しかし、RSS はユーザがブログを公開するときに最低限必要なメタデータのみを提供するものであり、Semantic Web で期待される詳細なクラス関係を持つオントロジを簡単には生成できない。RSS フィードを用いた検索でも、メタデータが上記単純なものであるため、ブログエントリの発見には、ユーザが検索キーワードを予め構成する必要があることに変わりない。そのため、キーワードが明確でない限り興味に即したエントリであっても発見できない事が多い。

(注2) : <http://www.doblog.com/weblog/PortalServlet>

(注3) : <http://www.movabletype.org/>

これに対し Semblog [9] では、パーソナルオントロジというツリー構造を持つカテゴリ体系をユーザが生成し、自身が記述もしくは収集したエントリを分類する。つまり、ユーザの興味体系であるパーソナルオントロジをボトムアップアプローチで生成することで、他ユーザの持つパーソナルオントロジや各種トピックディレクトリとのマッピングができる。一方、本研究はサービス毎に用意する雛型オントロジを用い興味オントロジをトップダウンアプローチで生成するため、ユーザによるオントロジ設計・構築の手間が少ない。さらに、興味対象クラス・インスタンスに対する興味度を与え、各ユーザの興味の強さに即した情報推薦の実現を試みる。

その他、個人の登録ブックマークや保持フォルダなどの階層構造と、ブックマークやフォルダの格納ファイルに基づき個人の興味情報を階層的に構築し、協調フィルタリングに基づきユーザにとって興味の近い別ユーザの興味階層に属する情報を推薦する研究がある [4]。しかし、ブックマークなどの階層構造には様々な分野の興味が混在する事も多いため、ブログコミュニティ形成に適用するには、分野毎の詳細な情報を適切に切り出す技術が必要となる。本研究では、様々な対象を記述しているブログエントリに対し、サービスドメイン毎に適切な粒度のオントロジを用意する事で、ドメイン毎に分離した興味オントロジを生成できる。

また、トピックディレクトリに対し Web ページを分類する研究 [3,5,10] が存在する。本研究の興味オントロジ生成の際の雛型オントロジに対するエントリ分類法との違いは、本研究では (1) オントロジのみを利用し、従来手法に必要な大量な Web ページより構成されるディレクトリやページ間のリンクを必要としない点、(2) オントロジの持つクラス特性を利用し分類誤りを除去する事で高精度な分類を実現できる点、および (3) ブログに適用しユーザ毎の興味オントロジを自動生成できる点である。

一方、エントリ間のリンク構造を分析し関係のあるブログサイト集合をコミュニティと捉え、抽出を試みる研究がある [1,6,15]。これらは、従来の Web コミュニティ抽出手法 [5] をブログへ適用するものが多い [1,6]。上記研究のコミュニティ形成支援への適用課題は、抽出されるコミュニティの興味対象が明示できない点と、既にリンク関係が形成されているページ集合を抽出するためユーザに対し意外な情報を推薦するとは言えない点である。

また、オントロジ間の近似度計測やマッピングに関する研究がある [7,8]。例えば、文献 [7] では、クラスの接続形態であるトポロジを考慮した近似度計測手法を提案している。それに対し本稿では、トポロジに加え、ユーザ毎の興味度も考慮した手法を提案している。また、近似度の高い興味オントロジ間で共起するクラスやトポロジを分析し、あるユーザの興味オントロジには出現しないが、そのユーザと近似度の高いオントロジに頻出するクラスを意外情報として抽出しユーザ推薦する事でブログユーザ間のコミュニティ形成を狙う。

3. 興味オントロジ自動生成手法

本章では、まず、音楽や映画といったサービスドメイン毎の雛型オントロジ設計法について説明し、ユーザの蓄積ブログエントリを雛型へ分類することによる興味オントロジ生成手法を提案する。

3.1 雛型オントロジの設計

本節では、OWL(Web Ontology Language) [12] を基にオントロジの説明をした上で、雛型オントロジの設計法を説明する。

OWLにおけるクラスは、同様の性質を持つ個体をグループ化しその性質を論理的に表現するための機能を提供する。クラスは、クラスの持つ個体であるインスタンスの列挙などのクラス表現を用い定義される。また個体同士の関係や個体とデータ値の関係を定義するプロパティを、特定のクラスとともに使用することでクラスの特徴を詳細に記述することができる。さらに、クラスのインスタンスに関する公理を用い、例えば `owl:sameAs` により 2 つのインスタンスが同値である事を記述できる [12]。

このように OWL を利用すれば、ドメイン毎のオントロジを詳細に設計できる。さらに、OWL を人間が手書きで記述するのは困難なため、*protege*^[14] などのオントロジ記述サポートツールの研究も進んでいる。とはいえ、やはり詳細なオントロジ設計や記述を一般ユーザが行うのは負担が大きく、オントロジ生成・流通を阻害する。そのため、本研究ではまずは雛型オントロジを、OWL 記述法則の中でもクラスの階層関係 (`subClassOf` 記述) とクラスに所属するメンバーであるインスタンスの列挙 (`oneOf` 記述)、階層構造の基準となるメタデータを指定するプロパティ記述のみを用いるライトウェイトなオントロジ [14] として設計する。さらに、ユーザの興味オントロジは、雛型オントロジへのユーザエントリ分類を通じ自動生成し、ユーザによるオントロジ記述は行わない。

なお、雛型オントロジの設計には、クラス間の階層関係やユーザ興味を細やかに反映するための末端クラスの粒度調整が必要である。幸い、*goo* 音楽^[15]等のポータルサイトにおけるトピックディレクトリは詳細化が進んでおり、例えば、音楽サービスドメインのジャンルを例に挙げると Web で公開されるジャンルの階層情報はユーザ興味に従う検索を考慮し、粒度を細かく設定している。そのためまずは、これらのトピックディレクトリを基に雛型オントロジを構築し、本研究における分析を通じ、適切な粒度考察を進める。

以下、雛型オントロジの設計手順を図 1 に示す例を基に説明する。まず、(1) 設計者は興味オントロジとしてどのサービスドメインのオントロジを生成するかを選択する。その上で、(2) そのドメインにおいてユーザ興味を反映するメタデータを選択する。選択材料としては、掲示板などの既存コミュニティの傾向を分析すればよい。例えば、音楽ドメインは、ジャンル・アーティストなどでコミュニティが生成されていることを考慮し、上記メタデータがユーザ嗜好を反映すると想定し、選択する。次

(注4) : <http://protege.stanford.edu/>

(注5) : <http://music.goo.ne.jp/>

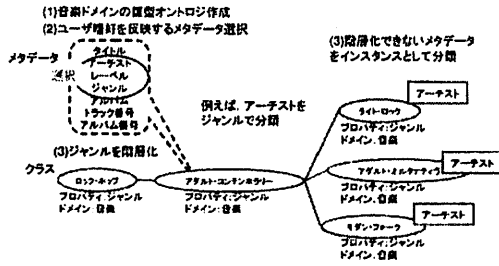


図1 雛型オントロジ構築手順

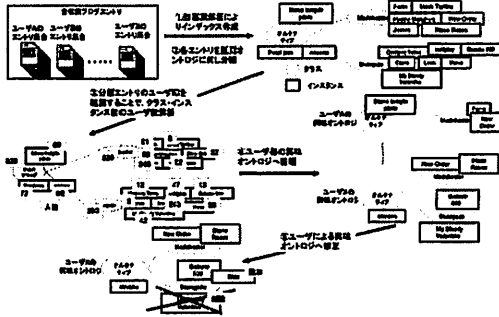


図2 ユーザ分布解析と興味オントロジ自動生成手順

に、(3) クラスツリーを作成可能なメタデータを選択し、クラス階層を形成する。この際、選択されたメタデータをクラスの性質を制約するプロパティとしてクラス階層間で継承する。例えば、ジャンルをプロパティとして継承するクラス階層を構築しアーティストなどをインスタンスとして各クラスに分類する。

本研究ではまずは雛型オントロジをトップダウンアプローチにより設計し、それに基づき次節で述べるユーザの興味オントロジを自動生成する。そして、興味オントロジによる情報推薦サービスなどを普及する事で、興味オントロジの流通やユーザ主導のボトムアップアプローチによる興味オントロジ修正の実現を狙う。なお、雛型オントロジは、サービス設計者の意図するサービス領域から構築していくことができるため、サービスプロバイダの提供サービスとの関係を構築しやすいという利点があり、さらにサービスの拡充に伴い徐々に増やして行けばよい。

3.2 ユーザの興味分布解析と興味オントロジ自動生成

図2に示す雛型オントロジ例を用い、ユーザの興味分布解析と興味オントロジ自動生成手順としてベーシックアルゴリズム(BA)を説明する。

まず、(1) ping サーバなどを通じ収集した全ブログエントリーに対し形態素解析を行いインデックスを作成する。ここで、収集されたブログエントリーは、一意なユーザIDを持つとする。

その上で、(2) 全ブログエントリーを雛型オントロジに対し分類する。分類方法としては、あるエントリー内の記述に雛型オントロジのあるクラス C_i の名前属性があれば、そのエントリーを C_i に分類し、また、 C_i に所属するインスタンス $I_i \in C_i$ の名

前属性があれば、エントリーをクラス C_i のインスタンス I_i に分類する。なお、エントリーが複数クラスに分類されても良い。例えば、図2において、エントリー内の記述に“Charlatans”という文字列がある場合、そのエントリーはクラス“Madchester”のインスタンス“Charlatans”に分類される。

次に、(3) 雛型オントロジを形成する最下層クラス C_r の持つ各インスタンスに対し興味を持つユーザ数を計測する。なお、クラス C_r のインスタンスに興味を持つユーザ数を算出する際、同一ユーザが複数エントリーにおいて同一インスタンスを記述していたとしても、ユーザ数は1と計測する。次に上記計測を最下層クラスに対しても実施し、最下層クラスに興味を持つユーザ数を、最下層クラス配下の全インスタンスに興味を持つユーザ数と最下層クラス C_i 自身に興味を持つユーザ数の総和で計測する。この場合も、同一ユーザが複数インスタンスに興味を持っていたり、最下層クラスとそのクラスに所属するインスタンスに同時に興味を持つとしても、ユーザ数は1と計測する。このようにしてユーザ数をルートクラスまで再帰的に計測する事で、そのドメインに興味を持つユーザ分布を計測できる。

そして、(4) 分類結果からユーザIDの一致するエントリーの分類体系のみを抽出すれば、そのユーザに対する興味オントロジを生成できる。例として、図2にユーザAのエントリー集合がインスタンス“stone temple pilots”や“New Order”, “Farm”を記述している場合に生成される興味オントロジを示す。

最後に、(5) 自動抽出された興味オントロジをユーザが閲覧し、興味クラス・インスタンスの追加や、興味の無いクラス・インスタンスを削除することで、各個人に適切な興味オントロジを生成する。また、こうして得られたユーザの興味オントロジ修正情報を収集し、雛型オントロジへフィードバックすることも可能である。

3.3 分類誤りのフィルタリング手法

しかしベーシックアルゴリズムでは、例えば図3において、クラス“Madchester”配下のインスタンス“Farm”などの多義語に対しては、“Madchester”というジャンルのアーティストである“Farm”でなく、農場という意味の“Farm”を記述するエントリーをも、クラス“Madchester”のインスタンス“Farm”に分類してしまい誤りが多い。そこで、本研究では、オントロジの持つ(1) 同一クラスに所属するインスタンスは同一の性質を持つという特性と、(2) クラス階層の近いクラス間の性質は近く、両者のインスタンス間の性質も近いという特性を利用し、分類誤りを除去するフィルタリングアルゴリズムを2種類提案する。

以下、フィルタリングアルゴリズムを説明する。

ベーシックアルゴリズムの手順(2)を細分化し、(2-1) あるユーザのあるエントリー E_i 内に雛型オントロジのあるクラス C_i に所属するインスタンス $I_i \in C_i$ の名前が記述されている場合、そのユーザの蓄積する全エントリーに対し、 C_i に所属する I_i 以外のインスタンス $I_k \{I_k \in C_i\}$ や C_i の記述があるかどうかをチェックする。そして、(2-2) 記述がある場合にエントリー E_i はクラス C_i に所属するインスタンス I_i を話題にするエントリーとして分類し、ない場合は誤りとする。図3を用い説明すると、“Farm”に対する記述が、あるユーザのエントリー E_i に存在し、

共通インスタンスを $I_i \in I$ と定義する。特に、トポロジ T_1 に所属するクラス集合を $C(T_1)$ 、トポロジ T_2 に所属するクラス集合を $C(T_2)$ とする。また、共通クラスに対する興味一致度を $I(C_i)$ とし、共通インスタンスに対する興味一致度を $J(I_i)$ とし、共通クラス配下の興味一致度を $I_i(C_i)$ とする。

(1) まず、 O_A と O_B 間で共通クラスを分析し、トポロジ T_1 を形成する共通クラスと、トポロジ T_2 を形成する共通クラスを抽出する。ここで、両トポロジを形成する共通クラスが存在することに注意する。例えば、図 5 では、共通クラス $a1$, $b1$, $b2$ はトポロジ T_1 を形成し、 $b2$, $b3$, $c4$ はトポロジ T_2 を形成する。

(2) 共通クラス C_i が両オントロジ間で共通インスタンス I_i を持つとすると、共通インスタンス I_i に対する興味一致度 $I(I_i)$ は、両オントロジのインスタンス I_i に対する興味度のうち小さい方の値とする。例えば、図 5 では、共通クラス $c3$ 配下の共通インスタンス a に対する興味一致度は 2 となる。

(3) 同様に、両オントロジ O_A と O_B における共通クラスに関する興味一致度 $I(C_i)$ は、両オントロジのクラス C_i に対する興味度のうち小さい方の値とする。例えば、共通クラス $C3$ に対する興味一致度は 0 となる。

(4) 次に、クラス $C_i \in C(T_1)$ 配下の興味一致度 $I_i(C_i)$ は、 O_A と O_B における C_i 配下の subclasses の積集合を $N(C_i)$ とし、和集合を $U(C_i)$ とすると、 $I_i(C_i) = \frac{\sum_{C_j \in N(C_i)} I(C_j)}{|U(C_i)|}$ で与える。例えば、共通クラス $b2$ 配下に対する興味一致度は $(9+3)/2 = 6$ となる。そして、興味一致度 $I_i(C_i)$ をトポロジ T_1 を形成する全共通クラスで足し込んだ値 $\sum_{C_i \in C(T_1)} I_i(C_i)$ をトポロジ T_1 を形成する共通クラス集合に対する興味一致度 $S(T_1)$ とする。

(5) 一方、共通クラス $C_i \in C(T_2)$ に対し、 O_A における C_i 配下のインスタンス集合を $I_A(C_i)$ とし、 O_B における C_i 配下のインスタンス集合を $I_B(C_i)$ とし、 C_i の興味一致度 $I_i(C_i)$ を、 $\frac{\sum_{I \in C_i} I(I_i)}{|I_A(C_i) \cup I_B(C_i)|}$ で与える。例えば、共通クラス $c3$ 配下に対する興味一致度は $((2+0+3+0)/4) = 5/4$ となる。そして、興味一致度 $I_i(C_i)$ をトポロジ T_2 を形成する全共通クラスで足し込んだ値 $\sum_{C_i \in C(T_2)} I_i(C_i)$ をトポロジ T_2 を形成する共通クラス集合に対する興味一致度 $S(T_2)$ とする。

(6) $S(T_1)$ および $S(T_2)$ 、両トポロジに対する重要度に応じた評価関数 $f(X)$ を用いオントロジ間の近似度 S_0 を $S_0(AB) = S(T_1) + f(S(T_2))$ で与える。

4.2 コミュニティ形成支援への適用

次に、オントロジ間の近似度計測手法を意外情報の推薦やコミュニティ形成支援へと適用する。

まず、近似度計測アルゴリズムをユーザ A とその他のログユーザ集合 $u \in U$ との間で総当りで行う。そして、ヒューリスティックな閾値 δ を用い、 $S_0(Au) > \delta$ を満たすユーザグループ G_U を導出し、グループ G_U に属する各ユーザの興味オントロジとユーザ A の興味オントロジの差分クラスとインスタンスを分析する。そして、オントロジ間の近似度が近いにも関わらずユーザ A に存在しないクラス、インスタンスに関するエン

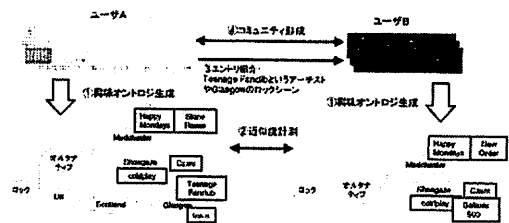


図 6 興味オントロジの近似度計測による実現サービスイメージ

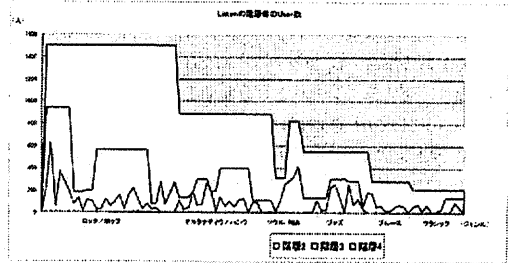


図 7 雛型オントロジに対するログユーザの興味分布

リをユーザ A の興味に即している意外な情報として推薦する。

オントロジ間の近似度計測手法に基づくコミュニティ形成支援サービスのイメージを図 6 を用いて説明する。ユーザ間の興味オントロジの近似度を計測し、近似度の高いオントロジ間で共起するクラスやインスタンスを分析することで、例えば、“Madchester”等のクラスや“Happy Mondays”等のインスタンスに興味を持つユーザは“Glasgow”クラスやその“Teenage Fanclub”というインスタンスにも興味を持つ可能性が高い事が分かる。こうしたトポロジが異なるにも関わらず興味を持つ可能性が高い情報を他ユーザのエントリを介して意外な情報としてユーザに推薦できる。

5. 提案手法の実装と評価

本章では、ブログポータルサイト Doblog における大規模データ (約 5 万 5 千ユーザ, 160 万エントリ) に対し、興味オントロジの自動生成とユーザ全体の興味分布の検証を行う。また、コミュニティ形成支援プロトタイプの実装を通じ、ユーザ毎の興味オントロジを用いたエントリ推薦の有効性を検証する。

実験にあたり音楽ドメインの雛型オントロジを作成した。作成した雛型オントロジは、goo 音楽などの Web ポータルの公開情報を参考とし、114 ジャンルをクラスとし、末端クラスに約 4300 のアーティストをインスタンスとして分類し、図 2 に一部示すような雛型オントロジを作成した。図 2 ではクラス階層のみを表示しているが、実際には末端クラスにインスタンスを配置している。なお、各クラス、インスタンスには名前属性を複数与えている。例えば、“verb”というインスタンスには、“ヴァーヴ”と“verv”という 2 つの名前属性を与える。このようにし、4300 のアーティストに対し約 7600 の名前属性を与えた。

そして、ベーシックアルゴリズム (BA)、フィルタリングアル

表1 提案手法による興味オントロジの精度 (FA2, ホップ数2)

	Rock	Jazz・Classic・その他	Total
正解数	911	1440	2351
適合率	911/1001=91.0%	1440/1520=94.7%	2351/2521=93.2%
クラス数	36	78	114
インスタンス数	2133	2158	4291

表2 1単語の名前属性を持つインスタンスの分類精度

	Rock	Jazz・Classic・その他	Total
正解	138	97	235
適合率	138/204=67.6%	97/125=77.6%	235/329=71.4%
1単語よりなるインスタンスの割合	455/2133=21.3%	458/2158=21.2%	913/4291=21.2%

リズム (FA1, FA2) により生成される興味オントロジの精度を測定した。なお、検証方法としては、興味オントロジを構成するクラス・インスタンスに分類されるユーザのエントリを人手で確認した。正解の導出根拠としては、実際にそのクラス・インスタンスの名前属性の記述があるエントリを正解とした。精度の尺度としては、正解数と分類結果中の正解数の割合 (適合率) を用いた。正解数が多いほど、ユーザが記述した興味がカバーされるが、適合率が低いと興味オントロジに誤りが含まれ、ユーザへの推薦情報の信頼性が落ちるため本稿では適合率向上がまず必須と考える。更に、1単語から形成される名前属性を持つインスタンスやクラスが特に多義語となる可能性が高い事を考慮し、そうしたインスタンスやクラスにフィルタリングアルゴリズムを適用した。なお、ログエントリのインデックス作成には全文検索エンジン Namazu^(注6)を用いた。

FA2の精度を、分類結果の1/10のデータをランダムに抽出し検証した(表1)。これによると、適合率は90%以上にまで達しており、フィルタリングが、エントリ分類や興味オントロジ生成に効果を持つことが確認できた。また、表2に1単語の名前属性を持つクラス・インスタンスへの分類精度を示す。こうした語は多義語である可能性が高く適合率が落ちている。

またグラフ7に、雛型オントロジ内の各クラスに対するユーザ分布をクラス階層毎に解析した結果を示す。なお今回作成した雛型オントロジはクラス階層が4階層からなるものが多いが、ソウルやブルースなどは末端クラスが3階層までしかないため、本グラフには4階層目は表示されていない。これによると末端クラスに所属するユーザ数であっても200程度存在する。末端クラスに分類されたエントリ集合を調査すると、そのクラスを特徴付ける語が多く頻出する事が確認できた。例えば、デス・メタル配下にはデスヴォイスという語などが頻出する事が分かった。雛型オントロジ作成においては、サービス毎に適切な粒度設計がポイントになると考えられるため、今後はこうしたクラスの特徴語や実験サービスを通じたユーザ数分析を通じた適切な粒度の定量化を試みる。

更に、フィルタリングアルゴリズムの性能を検証するため、ロックジャンルの1単語よりなる名前属性を持つ1/4のエントリをランダムに抽出し BA, FA1, FA2(ホップ数2)の正解数・

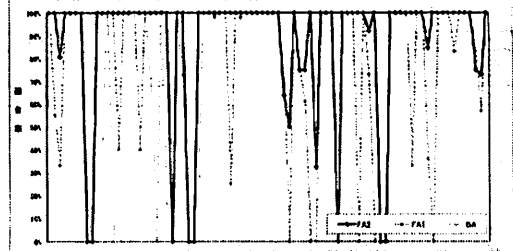


図8 1単語の名前属性を持つインスタンスの精度比較

表3 BA, FA1, FA2に対する正解数, 適合率の比較

	FA2	FA1	BA
正解	14	40	43
適合率	0.7	0.597	0.189

適合率を比較した。なお検証を公正にするため、FA2におけるインスタンスへの分類数が4以下のものは検証対象外とした。その結果83のインスタンスについて、正解数, 適合率を得る事ができた。なお、BAに対してはエントリ数が膨大であったためそのうち17インスタンスのみに検証を実施している。

図8は、83のインスタンスを横軸に各適合率を示す。また、表3にBAにおける17インスタンスに対する正解数, 適合率の比較を示す。本結果より、BA, FA1, FA2の順で適合率が向上することが分かる。また正解数は、FA1はBAよりそれほど落ちないが、FA2と比較するとFA2は大きく減少する。そのため今後は、ログの特徴である時系列の近いエントリやトラックバック等の関連エントリを利用し、FA2における分類決定要素を同一エントリのみでなく分類決定要素が出現する確率の高い上記エントリまで参照する機構を付加し、適合率を維持しつつ正解数を増やす事を試みる。なお、グラフよりFA2においても適合率をあげることができないインスタンスが8個存在し、全体の適合率を下げる事が分かる。そこでFA1からFA2にした際、急激に分類数が増えるインスタンスは、ユーザが興味を持ち多数エントリ記述しているにも関わらず、そのインスタンスの分類決定要素と全く共起しないため誤りである可能性が高いと考え、FA1とFA2で10倍以上結果が増えるインスタンスを自動抽出し分析した。結果、28個抽出でき、その内5個がFA2でも適合率が0%となる事が分かった。これより、FA1とFA2の比較を通じ急激に分類数が増えるインスタンスを分析し、雛型オントロジから削除する事が、適合率向上に有効と考える。

またFA2に対しホップ数を変化させ、精度を比較した(表4)。表4によると、ホップ0と2を比較するとホップ2が正解数, 適合率ともに良くなる。さらに、ホップ0と4を比較すると、ホップ4の方が若干正解数は増えるが、適合率は下がる。これは、今回用いた雛型オントロジが、例えばクラス“メタル”を例に挙げると“メタル”配下の末端クラスに“北欧メタル”や“ポップメタル”などがあり、“メタル”の親クラスが“ロック”であるなど、末端クラスとその親クラス間の概念間

(注6) : <http://www.namazu.org/>

表4 ホップ数変化による興味オントロジの精度比較

ホップ	正解数	適合率	検証数/検分類エントリ数	総分類エントリ数
ホップ0	475	475/533=89.1%	0.05	10882
ホップ2	495	495/544=91.0%	0.05	10888
ホップ4	477	477/557=85.8%	0.05	11133

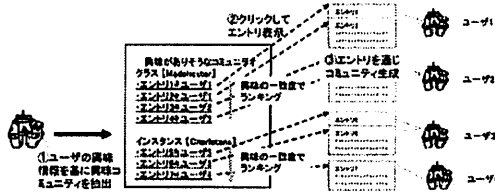


図9 コミュニティ形成支援プロトタイプ

表5 推薦エントリの精度

正解数	適合率
7164	7164/7552=94.8%

の距離と親クラスと祖父クラスの距離が遠くなるように設計されているため、ホップ0よりも分類決定要素が多いホップ2が最も適合率がよくなり、距離が遠くなるクラスのインスタンスまで分類決定要素に入れるホップ4が最も適合率が低くなると考えられる。また、ホップ数2と4の間で10倍以上分類エントリ数に変化がある1語からなる名前属性を持つインスタンスを自動抽出し分析した。結果、それらのインスタンスに分類された増加エントリはすべて誤りである事がわかった。例として“ヨーロッパ”というクラス“北欧メタル”配下のインスタンスを調べると、ホップ4では、クラス“メタル”の親クラスであるクラス“ロック”配下の“アダルトコンテンポラリー”配下のインスタンスも“ヨーロッパ”の分類決定要素に使われるようになるため、ヨーロッパ等の“ヨーロッパツアー”の話題などと共起してしまうような語を名前属性とするインスタンスの場合、誤りが急激に増える事がわかった。これより、ホップ数を変化させた際、急激に分類数が増えるインスタンスは、そのホップ数における分類は誤りである可能性が高いため、興味オントロジ生成に利用しない機構を付加する事で、正解数を増やしつつ適合率を維持できると考える。

更に、コミュニティ形成支援の検証として、図9に示すコミュニティ形成支援プロトタイプを作成し検証を行うことで、ユーザの興味オントロジに即し推薦されるエントリの精度を検証した。推薦精度の尺度としては、推薦エントリ中にユーザ毎の興味オントロジを形成するクラス・インスタンスに関する記述があるエントリを正解とし、正解数と推薦結果中の正解の割合(適合率)を用いた。結果を表5に示す。表によると、適合率は95%近くにまで達することが分かり、今回生成した興味オントロジに基づきユーザ毎に興味に即したエントリ推薦ができることが確認できた。

6. 結論と今後の課題

本稿では、ブログユーザのエントリからユーザ毎の興味オントロジ生成とクラス階層と興味度を考慮したオントロジ間の

近似度計測手法を提案し、意外情報の推薦やコミュニティ形成への適用を試みた。そして、ブログポータル Doblog の大規模データを基に高精度な興味オントロジ生成やコミュニティ解析およびコミュニティ形成の実現性を確認した。

今後、4.2章で述べたようなユーザ毎の興味エントリの統計処理に基づく意外なエントリ推薦の有効性確認のため、Doblog上で提案するコミュニティ形成支援実験サービス^(注7)を実施し、(1)意外なエントリ情報推薦によるコミュニティ活性化の有効性と(2)コミュニティのユーザ分布の時間変化の検証を進める。更に、ユーザのブログ閲覧履歴を用いた興味オントロジの拡張手法についても検討を深める。

謝辞

本研究の検証は、株式会社 NTT データのブログポータル Doblog のデータを利用させて頂いている。データ提供とコミュニティ形成サービスのプレインストーミングに快くご協力頂きました Doblog チーム及び株式会社ホットリンクには大変お世話になりましたことを感謝致します。

文献

- [1] Adar, E., Zhang, L., Adamic, L. and Lukose, R. M.: Implicit Structure and the Dynamics of Blogspace, *WWW 2004 Workshop on the Weblogging Ecosystem* (2004).
- [2] Berners-Lee, T.: An attempt to give a high-level plan of the architecture of the Semantic Web (1998).
- [3] Chakrabarti, S., van den Berg, M. and Dom, B.: Distributed Hypertext Resource Discovery through Examples, *Proceedings of the 25th VLDB*, pp. 375-386 (1999).
- [4] Jung, J. J., Yoon, J. S. and Jo, G.: Collaborative Information Filtering by Using Categorized Bookmarks on the Web, *INAP* (2001).
- [5] Kleinberg, J. M.: Authoritative sources in a hyperlinked environment, *Journal of the ACM*, Vol. 46, No. 5, pp. 604-632 (1999).
- [6] Kumar, R., Novak, J., Raghavan, P. and Tomkins, A.: On the bursty evolution of Blogspace, *Proceedings of the twelfth international conference on World Wide Web (WWW2003)*, pp. 568-576 (2003).
- [7] Maedche, A. and Staab, S.: Measuring Similarity between Ontologies, *In Technical Report, E0448, University of Karlsruhe* (2001).
- [8] Noy, N. F. and Musen, M. A.: Anchor-PROMPT: Using Non-Local Context for Semantic Matching, *In Proceedings of the Workshop on Ontologies and Information Sharing at the International Joint Conference on Artificial Intelligence (IJCAI)* (2001).
- [9] Ohmukai, I. and Takeda, H.: Metadata-driven Personal Knowledge Publishing, *Proceedings of the Third International Semantic Web Conference (ISWC2004)* (2004).
- [10] Page, L., Brin, S., Motwani, R. and Winograd, T.: The PageRank Citation Ranking: Bringing Order to the Web, Technical report, Stanford Digital Library Technologies Project (1998).
- [11] 奥村学, 南野朋之, 藤木裕明, 鈴木泰裕: blog ページの自動収集と監視に基づくテキストマイニング, 第6回人工知能学会セマンティック Web とオントロジー研究会 (2004.7).
- [12] 神崎正英: セマンティック・ウェブのための RDF/OWL 入門, 森北出版株式会社 (2005).
- [13] 総務省: ブログ・SNS の現状分析及び将来予測, http://www.soumu.go.jp/s-news/2005/050517_3.html (2005).
- [14] 橋本大也: ライトウエイト・メタデータの応用事例とその可能性, 第10回人工知能学会 SIGSWO 研究会 招待講演資料 <http://www.ringolab.com/note/daiya/archives/003614.html> (2005).
- [15] 谷口智哉, 松尾豊, 石塚潤: Blog コミュニティの抽出と分析, 第6回人工知能学会 SIGSWO 研究会 (2004).

(注7) : <http://music.doblog.com/>にて、Doblog ユーザ向けに6月中旬から開始する予定である。