

ウェブからのランドマーク抽出に基づくクエリフリーな地域情報閲覧

Query-Free Browsing of Local Information based on Landmark Extraction from the Web

手塚 太郎[△] 田中 克己[△]

Taro TEZUKA Katsumi TANAKA

地域的に限定されたウェブ検索は今後ウェブの重要な利用形態のひとつとして普及していくことが予想される。一方、地域情報の場合、特定の目的を定めずに情報を閲覧するという利用も一般的である。ガイドブックを閲覧する行為、あるいはドライブや車窓からの風景を眺める行為などは特定の情報の取得を目的としているわけではない。このような利用を念頭に、ウェブ上の文書情報と地理情報システム(GIS)を結合することによって、地域情報のクエリフリーな閲覧を可能にさせるシステム「車窓」を開発した。このシステムではウェブマイニングによって地域のランドマーク情報を事前に抽出し、ユーザが閲覧する際の暗黙的なクエリとして使用する。ランドマークはオーソリディ型とハブ型の二種類が取得される。ユーザは地図インタフェース上で描画した経路に沿ってウェブ空間を仮想的に移動し、地域に関連するページをクエリフリーな形で閲覧できる。

The local web search is one of the promising applications of the World Wide Web (WWW). However, in regional information there also exists a demand for browsing overall information of the area, without specific interest. Reading a guidebook or viewing scenery from a car or train window are examples of such regional browsing. Based on this idea, we created a query-free regional web browsing system SHASOU. The system is based on a combination of web meta-search and a geographic information system (GIS). SHASOU extracts landmark information from the Web and uses them as implicit queries for regional information while the user performs browsing. Two types of landmarks, the authorities and hubs, are obtained from web mining. The user can move through a path drawn on a map and obtain regional web pages without sending specific queries.

1. はじめに

地域的に限定されたウェブ検索、いわゆるlocal web searchは今後、ウェブの重要な利用形態のひとつとして普及していくことが予想される[1]。しかし、地域的な情報の利用においては特定の情報を見つけ出す「検索」だけでなく、地域全体の概略を「閲覧」したいという要求も存在する。例として、ガイドブックは特定の情報を調べるためだけでなく、地域の概略を知るためにも用いられる。この場合、クエリは

対象地域そのものであり、利用の目的はその地域内で推薦される情報を閲覧することである。さらに別の視点から考えれば、ドライブや電車の車窓から風景を眺める行為もまた、地域的な情報を閲覧することに相当する。

本研究ではこのような地域情報のクエリフリーな閲覧をウェブ上で実現するシステムについて論じる。このシステムにおいては、ユーザが地域空間内を仮想的あるいは現実的に移動するに伴い、近接するランドマークの情報が順次提供されていく。このような機能を実現するにあたっては、対象地域におけるランドマークの情報が必要になる。ランドマークは多くのユーザにとって認知的に顕著な地理オブジェクトであり、移動時に目標となる建物である[2]。すなわち、ランドマークはユーザが地域の概略を知りたい場合に優先的に提供されるべき情報である。

ランドマークのような認知的な地理情報は既存のGIS(地理情報システム)のデータとして含まれているとは限らず、また、含まれていたとしてもそれは複数のレイヤに分けられている程度で、「ランドマーク性の強さ」といった定量的な値が与えられていることは少ない。

そこで本研究ではウェブ上に存在する大量の文書情報からランドマークに関する情報を抽出し、それに基づき地域情報のクエリフリーな閲覧を可能にするシステムを実現する。

2. 関連研究

空間情報をもとにランドマークを抽出する手法に関して、すでにいくつかの研究が行われている。小磯らは3D地図のデータから、視覚的な大きさ、および、属性値の周囲のオブジェクトとの相違を指標としてランドマークを抽出した[3]。また、Burnettらは被験者に経路の案内文を記述させ、そこから手作業によってランドマークを抽出した[4]。Raubalらは2D地図・交差点の写真・建築データベースなど多様なソースに基づき、建物のファサード・色調・長方形からのずれなど多様な属性をランドマーク性の要因として仮定し、検証を行った[5]。Brennerらは航空写真と空中からのレーザースキャンによって建造物の位置と高さを取得し、ID3およびクラスタリングのアルゴリズムを用いてランドマークを取得する手法を開発した[6]。Eliasはランドマーク性の一般化された要件について論じた[7]。さらに、Sorrrowsらは地理空間ならびに電子的空間におけるランドマーク的要素を満たすべき条件を挙げ、広範な定義を行った[8]。

ウェブのクエリフリーな閲覧に関しては、一般的なウェブページやテレビ映像ストリームを対象にいくつかの研究が行われている。現在のウェブ検索では検索キーワードの入力やハイパーリンクのクリックといったユーザ側からの動作が常に要求されており、場合によってはユーザにとって負担となっている。一方、テレビやラジオといったメディアにおいてはユーザはほとんど操作を要求されず、少ない動作で利用することができる。ウェブにおいても今後このような利用形態が普及していくことが考えられる。Henzingerら[9]や、馬ら[10][11]は、字幕付きのテレビ映像ストリームそのものをクエリとして、関連のあるウェブページを連続的に検索するアルゴリズムやシステムを提案した。また、灘本らは閲覧中のニュースサイトのウェブページ(ニュース記事)全体をクエリとして異なるニュースサイトの類似記事を自動検索して提示する同時比較型ブラウザを開発した[12]。

これらのシステムは、閲覧中のウェブページやテレビ映像ストリームそのものをクエリと見なして自動的に検索処理

[△] 正会員 京都大学大学院情報学研究科
{tezuka,tanaka}@dl.kuis.kyoto-u.ac.jp

を行うものであるが、本研究では、GISにおける経路情報そのものをクエリとして自動処理して地域情報を連続的に提示する手法を提案する。

3. ランドマーク抽出

3.1 地理空間認識とランドマーク

人間の地理空間認知におけるランドマークの重要性を見るために、GISのデータ構造と比較する。GISにおいては各地理オブジェクトは座標によって保存されている。一方、人間の地理空間認知は多くの場合、顕著なランドマーク同士が多様な関係性によって結び付いたグラフ構造として記憶されていると考えられる(図1)。

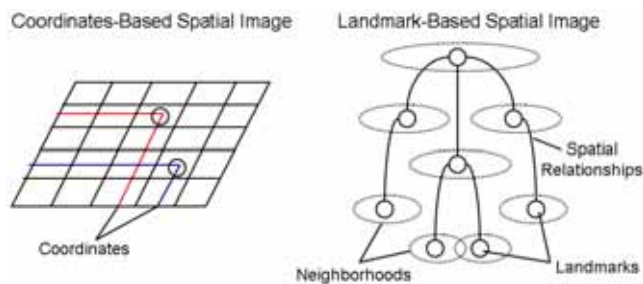


図1 空間情報のモデル

Fig.1 Models for Spatial Information

ここで、ランドマーク間の関係は非対称である。たとえば関係 $R(a,b)$ が「移動者が地理オブジェクトaを目印にして地理オブジェクトbを見つける」を表すとした時、対称性が成り立たないことは明らかである。また、関係に重要性などの重みを与えると、ランドマーク集合をノードの集合 V 、関係の集合をエッジ集合 E とした重み付き有向グラフ $G(V,E)$ と考えることができる。

3.2 ランドマークの分類

人間の地理空間認知をグラフによってモデル化するという考え方に基づき、ランドマークの集合に二つの重要なタイプを定義する。

ウェブのリンク構造の解析として、ハブ-オーソリティモデルが知られている[13]。Kleinbergは多くのエッジの起点となるノードをハブ、多くのエッジの到達点となるノードをオーソリティと捉え、ウェブページ集合から二つの特徴的なタイプのページ群を抽出した。

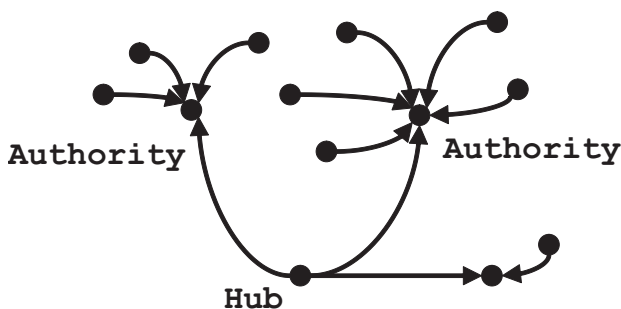


図2 有向グラフにおけるハブとオーソリティ

Fig.2 Hub and Authority in a Directed Graph

本研究ではランドマーク集合にオーソリティ型とハブ型を定義する。

オーソリティ型のランドマークは人々が頻繁に目的地とする建物、あるいは人々によって広く知られている建物に相当する。一方、ハブ型のランドマークは空間的な位置関係を把握する際に重要となる地理オブジェクトであり、経路を説明する場合に頻繁用いられるものである。ハブ型のランドマークは必ずしもユーザの移動の目的地であるとは限らない。

表1 ランドマーク種別

Table 1 Landmark Types

ランドマーク種別	特徴	指標
オーソリティ的	目標以外の用途	文書頻度
ハブ的	空間的目標物	周辺共起度

ハブ型のランドマークからは多様なオーソリティに到達することができる。オーソリティ型のランドマークには多様なハブから到達できる。

本研究では二つの型のランドマークを別個の指標によって抽出することを試みる。

3.3 オーソリティ型ランドマーク

オーソリティ型のランドマークは一般的な意味でよく知られた地理オブジェクトであり、ウェブ上でも頻繁に言及されていると予測される。そのため、テキストマイニングにおいて広く用いられている指標である文書頻度を用いて評価を行う[14]。文書頻度 df の計算式は以下の通りである。

$$df(p_i) = |\{d \in D \mid p_i \in s \wedge s \in d\}|$$

但し、 p_i は対象となる地名、 D は文書集合、 s は文を表す。

3.4 ハブ型ランドマーク

ランドマークを基準として地域の概略や建物間の位置関係を把握させるような場合には、ハブ型のランドマークを使用するのが望ましい。その指標として、周辺共起総数という値を定義する。周辺共起総数は、各地名がウェブ上で自らの周辺地名と共起する頻度を用いて計算された値である。これは各地名が空間的文脈の中でどれだけ重要であるかを表している。高い周辺共起総数を持つ地名は、都市空間の理解において重要な地名、すなわちハブ型のランドマークであると推測できる[15]。

周辺共起総数は、対象となる地名が周囲の地名と共起する頻度によって計算される。具体的な計算は、以下の手順に従って行われる。まず、対象地名 p から十分近い周辺地名の集合、 $P(p)$ を定義する。この定義には二つの手法が考えられる。ひとつは、閾値となる距離を設け、それより近い地名はすべて周辺地名とみなす手法である。周辺地名の定義の一例は、以下の式によって表される。

$$P'(p) = \{p_i \mid p_i \in P_{all} \wedge \delta(p, p_i) \leq R\}$$

但し、 p は対象となる地名、 $P'(p)$ は固定的な距離によって定義される周辺地名の集合、 P_{all} は対象地域に含まれるすべての地名である。関数 δ は二つの地名間の距離を与える。

R は距離の閾値である。だが、この定義を用いる場合、ある地名の周囲に多数の地名が存在する場合もあれば、ほとんど存在しない場合もあり、固定された距離ですべての地域に対応できないという問題が生じる。そこで本手法では固定的な距離の代わりに、「もっとも近い n 個の地名」を周辺地名として定義する。

$$P(p) = \{p_i \mid p_i \in P_{all} \wedge \delta(p, p_i) \leq \delta(p, p_{i+1}) \wedge 1 \leq i \leq n\}$$

これに基づき、周辺共起総数 rc は以下のように定義される。

$$rc(p_i) = \sum_{p_j \in P(p_i)} \kappa(p_i, p_j)$$

但し、 $\kappa(p_i, p_j)$ は p_i と p_j を共に含む文書数、すなわち p_i と p_j の共起数である。周辺共起数は地名の共起関係をエッジとみなしてグラフを描き、接続するエッジ数が多いノードを抽出することに相当する。

4. 地域情報のクエリフリーな閲覧

4.1 ランドマーク情報の利用

前節で述べられた評価基準をもとにランドマークを抽出し、地域情報のクエリフリーな閲覧を可能にするシステム「車窓」の実装を行った。車窓はウェブ上での利用を可能にするクライアント・サーバ型のアーキテクチャを採用している。クライアントはJAVA appletならびにJSPを用いて実装し、通常のウェブブラウザで動かすことができる(図3)。バックエンドのデータベースシステムとしてはPostgreSQLを使用している。対象地域として京都市を選び、地図データならびにランドマーク抽出にしようされる地名集合として、株式会社ゼンリンによって提供される住宅地図のデータを使用した[16]。

車窓システムのインタフェースにおいては、表示中の地図の範囲の中心からもっとも近いランドマーク名がクエリとしてGoogle Web APIに送られ、結果として取得されたウェブページが表示される。ユーザがクエリを入力する必要なしに、地域に関連するページを取得することができる。

地図上に表示されるランドマークを選択するための指標としては文書頻度 df と周辺共起総数 $rc(p_i)$ のいずれかの値が高い n 件のランドマークが地図上に表示される。指標を切り替えによって、これによってオーソリティ型あるいはハブ型ランドマークを選択的に表示できる。地域の名所を知りたい旅行者はオーソリティ型、目的地と目標物の位置関係を把握したい地元住民にはハブ型、といった使い分けが可能である。

車窓システムでは地図の縮尺によって表示される n 件のランドマークの集合が異なるため、同じ地点を中心にしていても、縮尺が異なれば取得されるウェブページ集合が異なる。これは、地図の縮尺からユーザが求める情報の詳細度を推測しているとみなせる。

もし広域の地図を表示させているのであれば、ユーザは地域の概略を知りたいのだと推測し、広い範囲にわたって重要なランドマークが地図上に表示される。また、ブラウザの下部にはそれに関連するウェブページが表示される。

逆に、ユーザが狭い範囲の地図を表示させている場合は、ユーザは詳細な情報を求めていると考えられる。この場合は

局所的に重要なランドマークが地図上に現れ、それに関連するページが表示される。このように縮尺に応じて情報の詳細度が変化するのは、縮尺を暗黙的なクエリとして使用していることを意味する。これによってユーザが明示的に詳細度を指定することなしに、適切なレベルの地域情報が取得できる。

4.2 経路指定による閲覧

経路指定によって閲覧を行う場合、ユーザは地図インタフェース上に任意の移動経路を描く。その後、移動開始を選択すると、地図の表示範囲はその移動経路に沿って自動的に移動していく(図4)。この間も、地図の中心からもっとも近いランドマーク名がクエリとして送られるため、ウェブページ表示領域に現れるページは順次切り替わっていく。これによって車窓から風景を眺めるかのようにウェブ上の地域情報を閲覧できることが、「車窓」という名称の由来である。



図3 地域情報閲覧インタフェースの構成

Fig.3 UI for the Regional Information Browsing

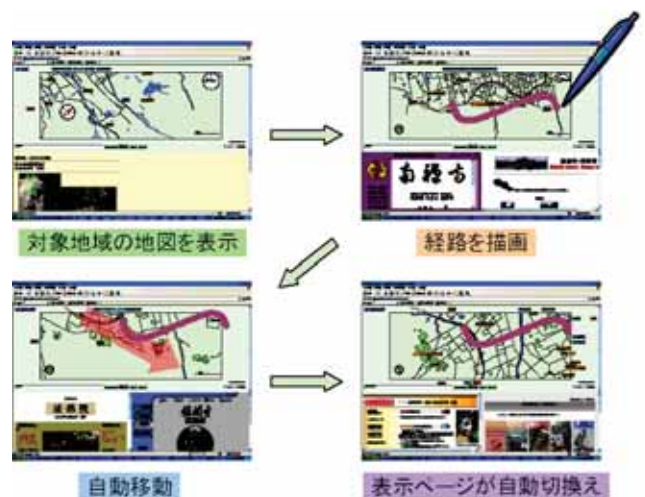


図4 経路指定による地域情報閲覧

Fig.4 Browsing of Regional Information by Path

移動経路はGPS などによって取得された軌跡、あるいは外部ソースからインポートされたデータなどを使用することも可能である。クエリフリーな問い合わせにおいてはコンテ

キストが重要な絞り込み情報として使われることが多いが、本システムの場合、地図の範囲と移動経路がそれに相当している。

5. まとめと今後の課題

本研究においては地域情報のクエリフリーな閲覧という利用形態に着目し、それを実現する手法としてのランドマーク抽出の手法、ならびに閲覧用ユーザインタフェースの実装を述べた。

今後の課題として、ランドマークは多様な側面を内包した概念であるため、ハブ型・オーソリティ型だけでは分類として不十分であると言える。そこでランドマークをさらに多様な種別に分類し、それぞれについて評価を行うことが求められる。

[謝辞]

本研究は、一部21世紀COEプログラム「知識社会基盤構築のための情報学拠点形成」による。ここに記し謝意を表します。また、本研究は一部文部科学省科学技術振興費プロジェクト「異メディア・アーカイブの横断的検索・統合ソフトウェア開発」(代表：田中克己)による。ここに記し謝意を表します。

[文献]

- [1] K. S. McCurley: "Geospatial Mapping and Navigation of the Web," The Tenth International World Wide Web Conference (WWW10), pp. 221-229, Hong Kong (2001).
- [2] Lynch, K.: The Image of the City, The MIT Press, Cambridge, Massachusetts (1960).
- [3] Koiso, K., Mori, T., Kawagishi, H., Tanaka, K. and Matsumoto, T.: "InfoLOD and LandMark: Spatial Presentation of Attribute Information and Computing Representative Objects for Spatial Data," International Journal of Cooperative Information Systems, Vol. 9, No. 1-2, pp. 53-76 (2000).
- [4] Burnett, G. E., Smith, D. and May, A. J.: "Supporting the navigation task: Characteristics of 'good' landmarks," Proceedings of the Annual Conference of the Ergonomics Society (eds. M.A. Hanson) Taylor & Francis (2001).
- [5] Raubal, M. and Winter, S.: "Enriching Wayfinding Instructions with Local Landmarks," Geographic Information Science (eds. M. Egenhofer and D. Mark), Lecture Notes in Computer Science 2478, Springer-Verlag, pp. 243-259 (2003).
- [6] Brenner, C. and Elias, B.: "Extracting Landmarks for Car Navigation Systems Using Existing GIS Databases and Laser Scanning," Proceedings "Photogrammetric Image Analysis," International Archives of the Photogrammetry, Remote Sensing and Spatial Information Sciences, Vol. XXXIV, Part 3/W8, Munchen (2003).
- [7] Elias, B.: "Determination of Landmarks and Reliability Criteria for Landmarks," Technical Paper, ICA Commission on Map Generalization, Fifth Workshop on Progress in Automated Map Generalization, IGN, Paris, pp. 28-30 (2003).
- [8] Sorrows, M. and Hirtle, S.: "The Nature of Landmarks for Real and Electronic Spaces," Spatial Information Theory: Cognitive and Computational Foundations of

- Geographic Information Science, (eds. C. Freska and D. Mark), Springer-Verlag, pp. 37-55 (1999).
- [9] Henzinger, M., Chang, B.-W., Milch, B. and Brin, S.: "Query-Free News Search", Proc. of the 12th International World Wide Web Conference, May 2003, Budapest, Hungary.
 - [10] Ma, Q. and Tanaka, K.: "WebTelop: Dynamic TV-content Augmentation by Using Web Pages", Proc. of IEEE International Conference on Multimedia & Expo (ICME2003), Vol.2, pp.173-176 (2003).
 - [11] 馬強, 田中克己: 話題構造に基づく放送と Web コンテンツの統合のための検索機構, 情報処理学会論文誌: データベース Vol.45 No.SIG 10 (TOD23), pp.18-36 (2004).
 - [12] Nadamoto, A. and Tanaka, K.: "A Comparative Web Browser (CWB) for Browsing and Comparing Web Pages", Proc. of the 12th International World Wide Web Conference, May 2003, Budapest, Hungary.
 - [13] Kleinberg, J.: "Authoritative Sources in a Hyperlinked Environment," Journal of ACM, 46 (1999).
 - [14] van Rijsbergen, C. J.: Information Retrieval - Second Edition, Butterworth & Co Publishers Ltd (1979).
 - [15] Tezuka, T., Yokota, Y., Iwaihara, M. and Tanaka, K.: "Extraction of Cognitively-Significant Place Names and Regions from Web-based Physical Proximity Co-occurrences," in Web Information Systems - WISE 2004 (eds. X. Zhou, S. Su, M. P. Papazoglou, M. E. Orłowska, and K. G. Jeffery), Lecture Notes in Computer Science 3306, pp. 113-124 (2004).
 - [16] Zenrin Co.,Ltd, <http://www.zenrin.co.jp/>

手塚 太郎 Taro TEZUKA

京都大学大学院情報学研究科博士研究員。2005年京都大学大学院情報学研究科博士後期課程修了、博士(情報学)。データベース、情報検索システムの研究に従事。日本学術振興会特別研究員。情報処理学会、日本データベース学会会員。

田中 克己 Katsumi TANAKA

京都大学大学院情報学研究科社会情報学専攻教授。1976年京都大学大学院前期博士課程修了、工学博士。主にデータベース、マルチメディアコンテンツ処理の研究に従事。IEEE Computer Society, ACM, 人工知能学会, 日本ソフトウェア科学会, 情報処理学会, 日本データベース学会会員。