

Web 情報を利用した駅クラスタリング手法の提案

Clustering Stations based on Web Information

清 雄一 小池 亜弥 白井 康之

Yuichi Sei Aya Koike Yasuyuki Shirai

(株)三菱総合研究所

Mitsubishi Research Institute, inc.

ユーザの位置情報に基づいてサービスを提供するロケーションベースサービスが普及しつつある。ユーザの正確な位置情報が悪意のある第三者に渡るとユーザの居住地や職場が特定されてしまう脅威が存在するため、位置情報を含めたユーザの個人情報を保護できる能力を持ったサービス提供者でなければ、ロケーションベースサービスの提供はできない。だが、このような能力を持った事業者がユーザの位置情報を適切に管理し、ユーザを一意に特定できないよう匿名化された位置情報を作成することができれば、より多くの事業者がロケーションサービスの分野に参入できると考えられる。このように、ユーザを一意に特定できないように匿名化を行う指標の一つとして、 k -匿名性が提案されている。これは、ユーザの位置情報を曖昧化し、曖昧化された領域に少なくとも k 人のユーザが存在するように匿名化することで、個人を一意に特定させない手法である。位置情報を曖昧化することにより、位置情報を利用したサービスの品質低下が予想されるが、既存研究では、 k 人以上のユーザが存在する領域の面積をできるだけ最小化することによって、品質低下を防いでいる。だが、提供するサービスの種類によっては、ユーザが存在する可能性がある位置の領域が最小化されることよりも、ユーザがどのような属性を持つ場所にいるか（たとえば、繁華街にいるのかビジネス街にいるのか）についての情報が失われないことのほうが重要である場合もある。本論文では、場所の属性に基づいた匿名化を行うことにより、場所の属性を利用するサービスの品質低下を防ぐ手法を提案する。また、場所の属性を設定する手法として、Web 情報を用いた手法を提案する。

1. はじめに

ユーザの位置情報を GPS 機能がついた携帯電話等から取得し、その情報に基づいてサービスを提供するロケーションベースサービスが普及しつつある。ユーザの位置情報を取得し、位置情報に応じて広告などの情報を通知するサービス [5] が考えられる。また、香港の企業が、ユーザの日々の位置情報を取得し、他のサービス提供者に向けて公開している例もある [7]。ユーザの正確な位置情報が悪意のある第三者に渡ると、ユーザの居住地や職場が特定されてしまう脅威が存在する。したがって、位置情報を適切に管理できるサービス提供者でなければ、ロケーションベースサービスを提供することはできない。だが、位置情報を適切に管理できるサービス提供者が、位置情報からユーザを一意に特定できないように匿名化を施すことによって、匿名化された位置情報を第三者に提供することが可能となる。匿名化された位置情報は、他の事業者が位置情報に基づいたサービス提供を行ったり、データマイニングを行う研究のために利用することができる。

ユーザを一意に特定できないように匿名化を行う指標の一つとして、 k -匿名性が提案されている [9, 8]。たとえば、ユーザの位置が (x, y) であるとする。 $x_1 < x < x_2, y_1 < y < y_2$ を満たす (x_1, x_2, y_1, y_2) を用意し、ユーザが (x_1, x_2, y_1, y_2) の頂点によって表される矩形領域に存在するという情報のみを第三者に提供することにより、ユーザの正確な位置を隠蔽する。このとき、 x_1, x_2, y_1, y_2 の値は、 (x_1, x_2, y_1, y_2) の領域に k 人以上のユーザが存在するように設定する。この領域内に k 人以上のユーザが存在するので、この情報を通知する元になったユーザを第三者は一意に特定することができない。 k -匿名性を対象とする既存研究は、 k 人以上のユーザが存在する領域の面積をできるだけ最小化することをめざしている。 k 匿名性を満たす領域の面積を最小化する問題は、NP 困難であることが示

されているため [6, 1]、より簡潔でありかつ良い解を与えるアルゴリズムが提案されている [3, 4]。

だが、提供するサービスの種類によっては、ユーザが存在する可能性がある位置の領域が最小化されることよりも、ユーザがどのような属性を持つ場所にいるか（たとえば、繁華街にいるのかビジネス街にいるのか）についての情報が失われないことのほうが重要である場合もある。たとえばある地域において、駅の東側と西側で様相が異なる街を考える。ユーザがどのような属性を持つ場所にいるかの情報を考慮せずに匿名化を行うと、ユーザが駅の近くにいる場合、ユーザが駅の東側にいるのか西側にいるのかの情報が失われる恐れがある。一方、ユーザが駅の東側にいるのか西側にいるかの情報を追加情報として第三者に提供すると、 k 匿名性が失われる恐れがある。したがって、ユーザがどのような属性を持つ場所にいるかの情報が第三者にとって重要である場合は、この属性情報を考慮した上で k 匿名性を満たすアルゴリズムを提案する必要がある。

本論文では、場所の属性に基づいた匿名化を行うことにより、場所の属性を利用するサービスの品質低下を防ぐ手法を提案する。また、場所の属性を設定する手法として、Web 情報を用いた手法を提案する。後者の提案については、マピオン (図 1) の Web サイトを利用した。マピオンの Web サイトには、各駅の周辺に、飲食店や公園等がそれぞれいくつ存在しているかの値が記載されている。このデータを用いて駅のクラスタリングを行うことにより、駅周辺の状況 (ビジネス街、住宅街、繁華街等) の把握を行った。

2. 場所属性情報を考慮した匿名化

2.1 場所属性の記述方法の定義

場所の属性の集合を C 、各属性を C_i とおく。たとえば、場所の属性としてビジネス街、住宅街、繁華街の 3 つが想定されている場合は、 $C = \{ \text{ビジネス街, 住宅街, 繁華街} \}$ であり、 $C_1 = \text{ビジネス街}$ となる。地点 (x, y) が属性 C_i を保有している度合を $V_i(x, y)$ と表現する。地点 (x, y) がビジネス街



図 1: マピオン

に完全に属している場合は、 $V_1(x, y) = 1$ であり、全く属していない場合は $V_1(x, y) = 0$ である。このように場所の属性情報として、ある地点 (x, y) がある属性を保持しているかどうか 2 値で表されている場合以外に、連続値で表現されている場合も考えられる。この場合は、 $V_1(x, y)$ の値は 0 から 1 までの実数値を取る。

2.2 匿名化指標の提案

多くの既存研究はデータを曖昧化することによって k 匿名性を実現しており、曖昧化されたデータの有用性を失わないようにするため、曖昧化の幅をできるだけ小さくすることを目標としている [4, 2]。たとえば、あるユーザの位置が (x, y) であるとすると、k 人以上が存在し、かつ面積が最小となるような矩形の領域を算出し、その頂点 (x_1, x_2, y_1, y_2) を第三者に提供する。一方、本論文においては、面積が最小となる領域を算出することが目標ではなく、場所属性も考慮した上で算出する必要があるため、新たな指標が必要である。

矩形の面積をできるだけ小さくするという目標に加え、矩形の頂点を (x_1, x_2, y_1, y_2) としたとき、その中に含まれる各点 (s, t) において、各場所属性の値 $V_i(s, t)$ の最大値と最小値の幅を小さくすることも目標となる。前者と後者のどちらを重要視するかは、曖昧化されたデータを利用する第三者が設定する。

したがって、 (x_1, x_2, y_1, y_2) で表現される領域を S とおき、最小化すべき指標を T とすると、

$$T = |x_2 - x_1| + |y_2 - y_1| + w \cdot \sum_{i=0}^{|C|-1} \left\{ \max_{(s,t) \in S} V_i(s, t) - \min_{(s,t) \in S} V_i(s, t) \right\} \quad (1)$$

となる。 S 内に k 人以上のユーザが存在し、かつ、式 1 を最小化するように、曖昧化された領域 S を設定する。

3. Web 情報を用いた場所属性の抽出

マピオンの Web サイトには、各駅の周辺に、飲食店や公園等がそれぞれいくつ存在しているかの値が記載されている。このデータを用いて駅のクラスタリングを行う。クラスタリングを行った結果に対して、それぞれの属性（ビジネス街、住宅街等）の付与については、手動で行う。

3.1 マピオン Web サイトからのデータ収集

Mapion 電話帳には、「和食店」や「スーパーマーケット」などの各カテゴリについて、駅ごとにその施設についての情報が

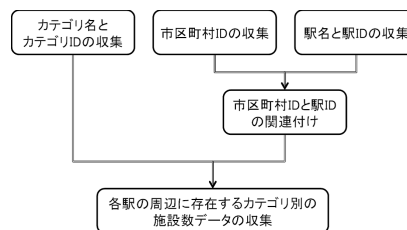


図 2: 駅単位のカテゴリ別施設データについての収集手順

記載されている。Mapion 電話帳から、この情報を収集し、駅ごとにおけるカテゴリ別施設数のデータを作成した。図 2 にデータ作成の手順を示す。(1) カテゴリ名とカテゴリ ID の収集、(2) 市区町村 ID の収集、(3) 駅名と駅 ID の収集では、マピオン Web サイトにおけるページ構造の把握を行っている。

1. カテゴリ名とカテゴリ ID の収集 Mapion 電話帳では、「和食店」などのカテゴリに分類して情報を記載している。このカテゴリ名の収集を行った。また、各カテゴリには ID が付加されており、この ID に基づいて、カテゴリ用の URL が決定されている。そのためカテゴリ ID の収集も同時に行った。たとえば、「和食店」のカテゴリ ID は「M01001」であり、「和食店」に関する情報は、<http://www.mapion.co.jp/phonebook/M01001/> に記載されている。
2. 市区町村 ID の収集 Mapion 電話帳では、「千代田区」など市区町村ごとに分類して情報を記載している。また、各市区町村には ID が付加されており、この ID に基づいて、市区町村用の URL が決定されている。そのため市区町村 ID の収集も同時に行った。たとえば、「千代田区」の市区町村 ID は「13101」であり、「千代田区」に関する「和食店」の情報は、<http://www.mapion.co.jp/phonebook/M01001/13101/> に記載されている。
3. 駅名と駅 ID の収集 Mapion 電話帳では、「大手町駅」などの駅ごとに分類して情報を記載している。それらの駅名の収集を行った。また、各駅には ID が付加されており、この ID に基づいて、駅用の URL が決定されている。そのため駅 ID の収集も同時に行った。たとえば、「大手町駅」の駅 ID は「ST22564」であり、「大手町駅」に関する「和食店」の情報は、<http://www.mapion.co.jp/phonebook/M01001/13101/ST22564/> に記載されている。ここで、「大手町駅」は「千代田区」に存在するため、URL には「千代田区」の市区町村 ID である「13101」も含まれている。
4. 市区町村 ID と駅 ID の関連付け上記 (3) で述べたように、各駅がどの市区町村に属するかの情報が必要となる。従って、各駅がどの市区町村に存在するかの関連付けを行った。
5. 各駅の周辺に存在する各カテゴリの施設数データ収集上記の (1) (2) (3) をまとめると、市区町村 A_i に属する駅 S_j の周辺に存在するカテゴリ C_k の施設に関する情報は、<http://www.mapion.co.jp/phonebook/Ck/Ai/Sj/> から得られる。ただし、この URL からは、駅 S_j の周辺に存在するカテゴリ C_k の施設についての情報を、最大 20 件までしか

表 4: 街カテゴリー一覧とクラスごとの街カテゴリー

街カテゴリー	クラス	街カテゴリー
ビジネス街	1	ビジネス繁華街
ビジネス繁華街	2	住宅街
大規模繁華街	3	住宅街
中規模繁華街	4	小規模繁華街
小規模繁華街	5	大規模住宅街
大規模住宅街	6	住宅街
中規模住宅街	7	住宅街
住宅街	8	ビジネス街
歓楽街	9	住宅街
公園街	10	中規模繁華街
官公庁・大使館街	11	住宅街
	12	大規模繁華街
	13	官公庁・大使館街
	14	官公庁・大使館街
	15	住宅街
	16	ビジネス街
	17	公園街
	18	学生街
	19	中規模繁華街
	20	住宅街
	21	歓楽街
	22	小規模繁華街
	23	公園街
	24	中規模住宅街
	25	中規模住宅街

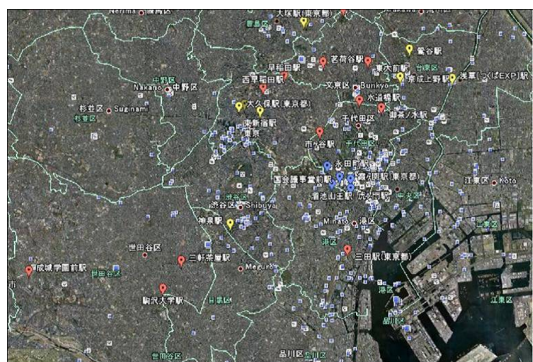


図 3: 街クラスターのプロット (官公庁街・学生街・歓楽街)

たとえば、グルメ施設やカラオケ施設等の施設それぞれについて、それらが存在する密度情報を属性として付与することができる。つまり、ある地域におけるグルメ施設の存在密度を属性 C_1 などとする。これにより、あるユーザがグルメ施設が多い地域にいるのか、カラオケ施設が多い地域にいるのかの情報をできるだけ失わないように匿名化を行うことが可能となる。

各施設が存在する密度情報の算出方法として、各施設が存在する地点を中心に二次元正規分布を描き、その和を取ることが考えられる。ある施設が地点 (μ_x, μ_y) に存在する場合、地点 (x, y) へ与える属性値の大きさは

$$N(x, y, \mu_x, \mu_y) = \frac{1}{2\pi} \exp\left(-\frac{(x - \mu_x)^2 + (y - \mu_y)^2}{2}\right) \quad (3)$$

となる。属性 C_i を与える施設が n 個存在し、それぞれの位置が (μ_{x_j}, μ_{y_j}) である場合、地点 (x, y) における属性値 $V_i(x, y)$ は、

$$V_i(x, y) = \sum_{j=0}^{n-1} N(x, y, \mu_{x_i}, \mu_{y_i}) \quad (4)$$

となる。

4. おわりに

本論文では、場所の属性情報を用いた k 匿名化手法の提案を行った。また、場所の属性情報の決定手法として、Web 情報を用いた手法を提案し、実際にマピオン Web サイト上の情報を用いて、駅のクラスタリングを行った。今後、マピオン Web サイト等から得られた情報を用いて場所の属性情報を細

かく設定し、提案した k 匿名化手法についての評価を行う予定である。

参考文献

- [1] G. Aggarwal, T. Feder, K. Kenthapadi, R. Motwani, R. Panigrahy, D. Thomas, and A. Zhu. Anonymizing tables. *Database Theory-ICDT 2005*, pp. 246–258.
- [2] R. J. Bayardo and R. Agrawal. Data privacy through optimal k -anonymization. In *ICDE '05: Proceedings of the 21st International Conference on Data Engineering*, pp. 217–228, Washington, DC, USA, 2005. IEEE Computer Society.
- [3] K. LeFevre, D. DeWitt, and R. Ramakrishnan. Incognito: Efficient full-domain k -anonymity. In *Proceedings of the 2005 ACM SIGMOD international conference on Management of data*, p. 60. ACM, 2005.
- [4] K. LeFevre, D. DeWitt, and R. Ramakrishnan. Mondrian multidimensional k -anonymity. In *Data Engineering, 2006. ICDE'06. Proceedings of the 22nd International Conference on*, pp. 25–25, 2006.
- [5] J. Meyerowitz and R. Roy Choudhury. Hiding stars with fireworks: location privacy through camouflage. In *MobiCom '09: Proceedings of the 15th annual international conference on Mobile computing and networking*, pp. 345–356, New York, NY, USA, 2009. ACM.
- [6] A. Meyerson and R. Williams. On the complexity of optimal k -anonymity. In *Proceedings of the twenty-third ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*, p. 228. ACM, 2004.
- [7] M. Terrovitis and N. Mamoulis. Privacy preservation in the publication of trajectories. In *IEEE International Conference on Mobile Data Management*, pp. 65–72, Los Alamitos, CA, USA, 2008. IEEE Computer Society.
- [8] R. Yarovsky, F. Bonchi, L. Lakshmanan, and W. Wang. Anonymizing moving objects: how to hide a MOB in a crowd? In *Proceedings of the 12th International Conference on Extending Database Technology: Advances in Database Technology*, pp. 72–83. ACM, 2009.
- [9] P. Zacharouli, A. Gkoulalas-Divanis, and V. S. Verykios. A k -anonymity model for spatio-temporal data. In *ICDEW '07: Proceedings of the 2007 IEEE 23rd International Conference on Data Engineering Workshop*, pp. 555–564, Washington, DC, USA, 2007. IEEE Computer Society.