

位置履歴からのユーザ属性の推定

松尾 豊^{†1} 岡崎 直観^{†2,†3} 中村 嘉志^{†1}
西村 拓一^{†1} 橋田 浩一^{†1} 中島 秀之^{†1,†4}

近年、ユビキタス情報環境の中でユーザの位置情報を取得する技術が進展している。本論文は、ユーザの位置履歴からユーザに関する属性を推定する汎用的な方法を提案し、実証実験によって得られたデータに基づいて検証する。センサが各ユーザを何回検出したかという情報は、ユーザ-センサ行列として表すことができる。これは、情報検索における文書-語行列と同様の形式になり、ユーザ属性を推定する問題は文書分類の問題に置き換えることができる。本論文では、センサの検出回数から属性を構築する複数の方法を比較しながら、SVMを用いてユーザ属性を推定する。また、センサの重要度を定義する方法を提案し検証する。

User Modeling Based on Location History

YUTAKA MATSUO,^{†1} NAOAKI OKAZAKI,^{†2,†3} YOSHIYUKI NAKAMURA,^{†1}
TAKUICHI NISHIMURA,^{†1} KÔITI HASIDA^{†1}
and HIDEYUKI NAKASHIMA^{†1,†4}

Recent development for location detection techniques enables us to obtain location histories for users in a ubiquitous environment. This paper describes a new method to infer user attributes from a user's location history. Using the number of counts each sensor detects each user, we can obtain a user-sensor matrix, which is similar to document-term matrix in the context of information retrieval. The problem to detect a user's attributes can be reduced into text classification problem, to which support vector machine can be effectively applied. We also propose a method to measure the importance of sensors.

1. はじめに

ユビキタス情報環境において、ユーザを取り巻く状況を把握し、個々のユーザに応じたサービスを提供することは、中心的な課題の1つである^{1)~3)}。ユーザの位置情報を取得し情報支援に活用するシステムは、位置情報システム⁴⁾と呼ばれる。位置情報を用いたシステムとして以前から有名なものに、Active Badgeがある⁵⁾。Active Badgeは、赤外線デバイスを埋め込んだ名刺を使ってユーザの位置情報を交換するシステムで、「あるユーザがどこにいるか」「あるユーザは誰といるか」などの問合せに答えたり、ユーザに電話を転送したりするなどの支援が可能である。近年では、さ

まざまなデバイス技術の進展にともなって、GPSやRFIDタグを用いてユーザの位置を推定する研究、家の中に各種のセンサを配置しユーザの行動を把握する研究⁶⁾などが行われている。たとえば、お年寄りや要介護者などの行動記録や行動意図の把握⁷⁾、また日常的な行動の文脈把握とそれに応じた行動の支援^{6),8)}などがあげられる。多くの人が訪れる万博会場やテーマパークで混雑管理や迷子防止のために位置情報が用いられている例もある^{9),10)}。

個々のユーザに応じた情報支援を行うためには、個々のユーザに関する知識を用い、ユーザに適した情報提示やインタラクションを行う必要がある。個々のユーザをモデル化する技術はユーザモデリングと呼ばれ、システムの挙動に関連したユーザに関する明示的な知識をユーザモデルと呼ぶ^{11)~13)}。ユビキタス情報技

†1 産業技術総合研究所

National Institute of Advanced Industrial Science and Technology

†2 東京大学

The University of Tokyo

†3 マンチェスター大学

University of Manchester

†4 公立はこだて未来大学

Future University, Hakodate

文献 13) には、ユーザモデルは次のように定義されている。A user model is a knowledge source in a system which contains explicit assumptions on all aspects of the user that may be relevant to the behavior of the system. These assumptions must be separable by the system from the rest of the system's knowledge.

術におけるユーザモデリング技術の重要性は近年認識が高まっている。たとえば Heckmann は、できるだけドメインに依存しない汎用的なユーザモデリングのために、統一的に用いることのできるオントロジと記述法を提案している¹³⁾。ユーザモデルにはさまざまな表現の形式があるが、本研究では、ユーザの特徴をユーザ属性とその属性値によって表す方法をとる。この方法は、Heckmann の attribute-value のペアによるユーザの表現と親和性が高い。

本論文では、ユビキタス情報環境におけるユーザの位置履歴から、さまざまなユーザ属性の属性値を推定するユーザモデリングの方法を提案する。ユビキタス情報環境におけるユーザモデリング技術の重要性に鑑み、ドメインに依存せず汎用的に適用可能であることを目指している。センサのユーザ検出情報を、行列の形で表し、情報検索の分野で行われている文書分類の手法を適用する。その際、センサの検出回数をどう処理するかによって、分類の精度が変わることを示す。また、ユーザ属性の推定のためのセンサの重要度を測る指標を提案する。本論文では、2004年に産業技術総合研究所において行われた実験で得られたデータを基に、分析と評価を行っている。

なお、本手法では、特定のアプリケーションのためのユーザモデリングではなく、広く適用可能なユーザモデリング手法の構築が目標である。そのため、データの処理やアルゴリズムの設計を、できるだけドメインに依存しない形に簡単化している。こういった手法をベースラインとして用いることで、さらに精度の高い方法やドメインの性質を利用した方法の構築を行うことができ、今後のユビキタス情報環境におけるユーザモデリング技術の発展につながっていくと考えている。本論文で念頭においているアプリケーションは、たとえば次のようなものである。

- オフィス環境における情報支援：たとえば来客対応システム、ミーティングの予約システムなど、ユーザの所属するグループや身分、コーヒーを飲むかどうかといったユーザ属性を考慮してユーザの操作を簡便化する。

- 駅や博物館での情報支援：年齢や喫煙するかどうかで案内の経路を変える。
 - ショッピングセンターにおける位置情報の取得：年齢や性別などの情報とあわせて立ち寄る店や広告効果などの分析を行う。
- などである。

本論文は次のように構成されている。次章で関連研究について述べたあと、本論文で用いる位置情報について3章で説明する。4章でユーザ属性の推定問題および提案手法について説明し、データの分析と考察を5章に示す。さらに6章でセンサの重要度の計算についての提案および評価を行う。7章で議論を行い、結論を述べる。

2. 関連研究

情報支援に利用可能なユーザ属性にはさまざまなものがある。ユーザ属性には、現在の位置、混雑状況やリソースの使用状況など外部環境に関するもの、ユーザの意図や目的、同行者や携帯しているものなどユーザの状態に関するもの、興味や嗜好などユーザの性質に関するもの、ユーザの日常的な行動などがある。また、消費者行動に関する研究では、店舗内における行動の分析が行われているが^{14),15)}、こういった分析では、年齢や性別、職業などのデモグラフィック属性という属性でユーザを特徴づけることがある。

さらに、ユーザ属性の変化を考えると、変化しにくいもの（長期的なもの）とそうでないものがある。たとえば、外部環境やユーザの意図や目的などは変化しやすい。ユーザの興味や嗜好などは比較的变化しにくい。また、年齢や性別、職業などは（数年にわたって使うシステムでなければ）基本的には変化しないと考えることができる。

本研究では、ユーザの長期的な属性を対象とする。年齢、職業といったデモグラフィック属性と、ユーザの性質に関する属性のうち長期的なもの、そしてユーザの日常的な行動を表す属性の一部であり、具体的には、年齢、職業（チーム、身分）、嗜好（コーヒーや喫煙）、日常的な行動（出勤頻度、居室、通勤方法）である。ユーザの長期的な属性はさまざまなドメインにわたって有効であるのに対し、ユーザの意図や目的、状態など、短期的に変わりうるユーザ属性はドメインへの依存性が高いためである。

ユーザモデリングの研究で有名な Kobsa は、ユーザモデリングを用いたサービスについてまとめており、その中で、ユーザの長期にわたる属性として、知識、嗜好、能力などをあげている¹¹⁾。しかし、長期的な

ユビキタス情報環境におけるユーザモデリング技術のワークショップがいくつか開催されている。Workshop on User Modeling for Ubiquitous Computing (2003) や Workshop on Personalized Context Modeling and Management for UbiComp Applications (2005)、最近では UbiqUM2006 (Workshop on Ubiquitous User Modeling) などである。また、2005年には Journal of User Modeling and User-Adapted Interaction において User Modeling in Ubiquitous Computing という特集が組まれている。

ユーザ属性にはどのようなものがあり、どう分類できるかに関して一般的な知見はなく、本研究では、他のドメインでもできるだけ利用できるように、かつ実験の範囲内で取得可能な長期的なユーザ属性を選択している。

ユーザの長期的属性は、アンケートやユーザ登録などによって得られる場合もあるが、位置情報から推定することが効果的な状況も考えられる。たとえば、1日から数日にわたって開催される展示会や見本市では、参加者のデモグラフィック属性や興味の記入を促すことも多い。プレゼントなどの賞品を用意して参加者のアンケートをとるキャンペーンも、ユーザ属性とあわせて意見を集めようということである。したがって、ユーザが自発的に登録しなくても、ユーザ属性がある程度の精度で推定できれば、その後のデータ分析や個人に合わせた情報の配信などに便利である。また、大規模なショッピングモールで客の位置情報を取得し、その一部のユーザ属性が分かっているのであれば、それを学習データとして本手法を適用し、他の客のユーザ属性を推定することが可能である。

これまで、位置情報をはじめとするセンサ情報からユーザの行動を推定するさまざまな研究が行われている。Wilsonらは、家の中でのセンサデータ（接触センサと圧力センサ）を用い、ユーザの行動を表す一連のデータの系列（エピソード）を分離するシステムを構築している^{7),16)}。Narrotorというシステムを用いることで、行動のラベル付けを容易にし、行動の要約を得ることができる。文献17)では、GPSのデータをクラスタリングし、予測可能なモデルを作る手法を提案している。クラスタリングすることで、ユーザにとって重要な「場所（location）」を特定したあと、その関係をマルコフモデルで表現し、予測可能にしている。また、Hightowerらは、GPS、WiFi（無線LAN）、GSM（デジタル携帯電話の通信方式）を使って、1カ月の間の行動を記録し、どこに行くかを学習し予測する方法を提案している¹⁸⁾。

以上の研究では、ユーザが何をしているかというユーザの行動を分析の対象としている。それに対して、本論文は、ユーザの年齢や職業などのユーザ属性の推定が目的である。いずれのアプローチも、目的はユーザの文脈を得るということであるが、ユーザの行動の推定はしばしば暗黙に仮定されたユーザ属性に依存しており、お互いに補完するものであると考えられる。

3. 本論文で取り扱う位置情報

3.1 カード型 CoBIT

本論文では、ユーザの位置と同時にユーザのIDを取得できるセンサを対象とする。具体的に使用したのは、カード型 CoBIT^{19),20)}である。カード型 CoBITは、赤外線 ID タグと液晶シャッタを用いて、利用者の位置と方向に応じてIDを発信する。ただし、今回用いたカード型 CoBITは、簡便化のため液晶シャッタを省略したバージョンである（図1）。3～5mの距離までユーザのIDを発信することができ、タグ検出器（図2）を環境中に設置することで、ユーザの位置と同時にIDを取得することができる。CoBITの検出間隔は約3秒（ID発信の連続的な衝突を避けるためにランダムに間隔が変化する）である。カード型 CoBITは、名札代わりに気軽に用いることができるので、学会や展示会といったイベント空間を対象とした情報支援で実際に用いられている。2003年度および2004年度人工知能学会全国大会では、学会支援システムとし



図1 カード型 CoBIT
Fig. 1 Card-type CoBIT.

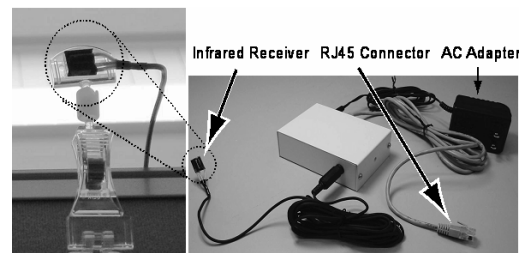


図2 2種類のタグ検出器（センサ）
Fig. 2 Two types of tag detectors (sensors).

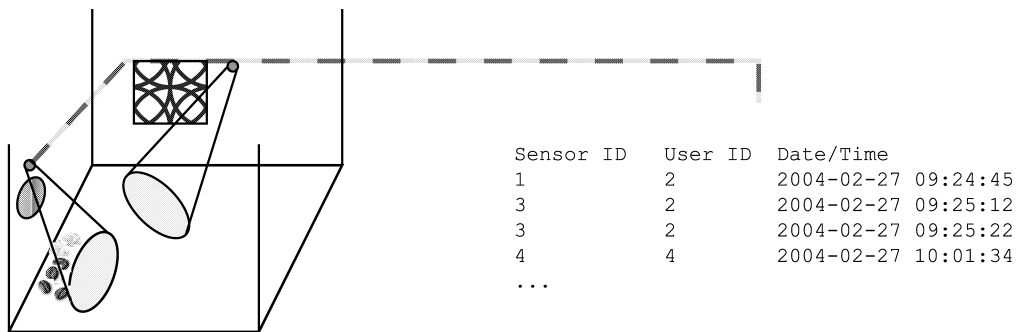


図 3 空間に配置されたセンサと蓄積される検出履歴

Fig. 3 Allocated sensors and their detection history.

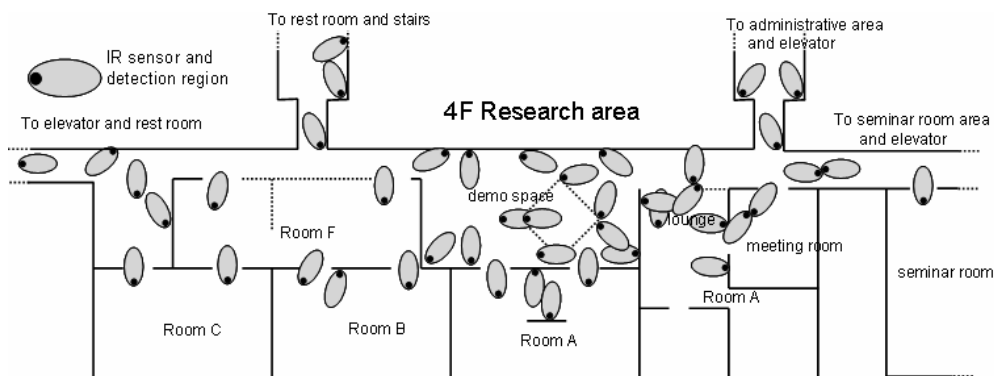


図 4 センサの配置図の一部

Fig. 4 A part of the sensor allocation map.

て利用された^{21),22)}。

タグ検出器（以下ではセンサという）を環境中に配置し、得られるデータを模式的に示したものが図3である。センサのIDと、センサが検出したユーザID、検出時間が記録されたデータである。提案手法では、検出回数だけを扱い、検出時間は用いない。本論文で仮定するこのデータ形式は一般性のあるものであり、たとえば、ICタグ（RDIFタグ）で取得されるユーザの位置情報はこの形式に加工することができる。また、GPSのように広範囲に連続的な位置情報が取得できるデータは、位置のクラスタリング¹⁷⁾の前処理を行うことで、このタイプのデータ形式に変換できる。また、位置データを集約することができれば、環境側にタグを配置し、ユーザ端末で位置を読み取る、たとえばRFIDリーダ搭載型の端末でも本手法は適用可能である。なお、本研究ではユーザ属性を推定することが目的であるので、完全に匿名のセンサを用いることは想定していない。

3.2 位置情報取得の実験

著者らは、産業技術総合研究所臨海副都心センターにおいて、2004年2月16日（月）～20日（金）の1週間、位置情報取得のための実験を行った。研究センター内の1階および4階にセンサを計94個設置し、4Fに勤務する内部スタッフおよび当期間中の来館者すべてにカード型CoBITを配布した。4階のオフィス環境に配置されたセンサの配置図を図4に示す。この領域は約1,000m²であり、センサを設置した領域全体の約3分の1である。内部スタッフは研究員、事務職員、学生、補助スタッフなどであり、対象となる66名全員にカード型CoBITを配布し、そのうち47名のデータを取得することができた。来館者については、打合せや見学などの来客、物品の納品業者など、計170名の位置情報を取得した。

センサの総検出回数は24,317回であり、そのうち20,273回が内部スタッフであった。1人あたりの平均は、内部スタッフで431.3回、来館者で23.8回である。来館者には、来館目的や訪問先など簡単なアンケートに記入してもらった。また内部スタッフに対しては、年齢、来所頻度、居室、喫煙するかなどのユーザ属性

文献 23) のように、ID の取得できるセンサと組み合わせて用いることはできる。

を手で調べ、ユーザ属性のデータとした。このユーザ属性は、他のドメイン（たとえば駅や博物館）でも使えるような汎用性と、オフィス環境における実験という制約の両方を考慮し選んだ。なお、内部スタッフ、来館者とも、実験の目的を簡単に説明するとともに、実験に関する協力の了承を得ている。本論文では、このうちスタッフに関するデータを用いて分析を行う。

4. ユーザ属性の推定方法の提案

ユーザの位置履歴とユーザ属性があらかじめ分かっていたら、位置履歴からユーザ属性を推定するモデルを作ることができる。本章では、位置履歴からユーザ属性を推定する手法について述べる。

4.1 ユーザ属性の推定と文書分類への帰着

各センサはユーザの ID を検出できるので、それぞれのユーザに対する検出回数を求めることができる。これは、センサとユーザをそれぞれ行、列とする行列で表すことができ、ユーザ-センサ行列と呼ぶことにする。ユーザ-センサ行列 W はユーザの数を n 、センサの数を m とすると、 $n \times m$ 行列であり、その要素 W_{ij} はセンサ s_j においてユーザ u_i が検出された回数となる。たとえば、センサが 4 個、ユーザが 3 人 とすると、次のような行列が一例である。

$$W = \begin{pmatrix} 1 & 2 & 2 & 4 \\ 1 & 0 & 2 & 0 \\ 3 & 2 & 0 & 0 \end{pmatrix}. \quad (1)$$

次に、ユーザ属性を考えよう。たとえば、コーヒーを飲むかどうかというユーザ属性を考え、この属性値が { 飲む, 飲まない } という 2 つ (それぞれ 1 と 0 とする) をとりうるとする。3 人のユーザ属性がそれぞれ 1, 0, 1 であるとすると、次のように表すことができる。

	s_1	s_2	s_3	s_4	coffee
u_1	1	2	2	4	1
u_2	1	0	2	0	0
u_3	3	2	0	0	1

ここで、 $u_1 \sim u_3$ は各ユーザ、 $s_1 \sim s_4$ は各センサ、coffee はコーヒーを飲むかどうかのユーザ属性である。つまり、 $s_1 \sim s_4$ というセンサの検出回数をもとにした属性から、coffee という属性が 1 であるか 0 であるかを分類するモデルを作ればよい。これは機械学習の問題となる。

仮に、新たなユーザ u_4 がいて、各センサの検出回数が以下のとき、このとき、このユーザの coffee 属性

は 1 であろうか、0 であろうか。

	s_1	s_2	s_3	s_4	coffee
u_4	1	2	0	0	?

もし s_2 がコーヒーマシンのある場所に設置されたセンサであるとする、 s_2 が大きければ coffee 属性が 1 であるという分類ルールが学習でき、したがって u_4 は coffee 属性が 1 であろうと推定することができる。

このように、本論文では、センサがユーザを検出した回数を属性として、ユーザ属性の値を学習する。それによって、新たなユーザに対しても、センサの検出回数からユーザ属性の値を推定することができる。

情報検索の分野では、各文書に各語が何回出現したかという、文書-語行列がよく用いられる。前述の行列 W は、 n 個の文書と m 個の単語と考えたときの、文書-語行列と同じ形である²⁴⁾。文書分類のタスクでは、文書に出現する語の出現頻度をもとに文書の属するカテゴリを推定する。位置履歴からのユーザ属性の推定問題を文書分類と同様にとらえ、各ユーザが検出されるセンサの検出回数をもとにユーザ属性の推定を行う。

本論文では、各種のタスクで良好な結果を出すことが報告されている Support Vector Machine (SVM) を適用する。SVM ではデータを 2 つの種類に分離するために、各データ点との距離 (マージン) が最大となる分離平面を求め²⁵⁾。マージンを最大化する分離平面を構築することで、高い汎化性能を持つ。線形に分類することが難しい場合は、カーネルトリックによって入力空間をより高次の特徴空間に写像し、そのうえで線形分離を行う。そのためのカーネルとして、多項式カーネル、RBF (radius basis function) カーネルなどがある。SVM は情報検索や自然言語処理の分野でよく使われる手法であり、詳しくはたとえば文献 26), 27)などを参照していただきたい。なお、位置履歴は時系列の情報であるが、ユーザ-センサ行列を作る時点で時間に関する情報は失われている。これについては 7 章で議論する。

4.2 センサの検出回数からの属性構築

センサの検出回数をそのまま属性として用いる方法のほかに、検出回数に何らかの処理を施したうえで属性として用いる方法がある。ユーザによってはセンサに多く検出されるユーザと、そうでないユーザがあり、検出回数ではなく、すべてのセンサによる検出回数の総数と比較しての割合を用いた方がよいかもれない。また、検出回数を使うのではなく、ある閾値を設けて、センサに検出された/されていないの 2 値で表すこと

も考えられる。

情報検索の分野では、 $tf \cdot idf$ という指標がよく用いられる。文書の特徴づける索引語としては、文書内によく出てくる語も重要であるが、他の文書にあまり出てこない語も重要である。 $tf \cdot idf$ は、語が他の文書にどのくらい出てくるかを考慮して語の重みを定義するものである。この場合にあてはめると、センサが他のユーザをどのくらい検出しているかを加味して属性を構成することになる。どのユーザも検出するセンサ（たとえば建物の入り口に置かれたセンサ）は、ユーザ属性の推定という目的では大きな情報を持たない。逆に、コーヒーサーバや喫煙所など、特定のユーザ属性の人だけを検出する場所はユーザ属性の推定という観点からは情報量が多い。ユーザ u_i におけるセンサ s_j の属性値を次のように計算する。

$$tfidf(s_j, u_i) = freq(s_j, u_i) \times idf(s_j) \quad (2)$$

ただし、 $freq(s_j, u_i)$ は s_j による u_i の検出回数であり、 $idf(s_j)$ は次式で定義される。

$$idf(s_j) = \log(n/uf(s_j)) \quad (3)$$

n はユーザ数、 $uf(s_j)$ は s_j が検出したユーザの数であり、どのユーザも検出するセンサは $uf(s_j)$ が大きくなるので、 $idf(s_j)$ は小さい値になる。

まとめると、本論文では、以下の 8 種類の方法でユーザ u_i に対するセンサ s_j の属性 a_{ij} を構築する。

- 検出回数： $a_{ij} = freq(s_j, u_i)$
- バイナリ：検出回数を 2 値化したもの

$$a_{ij} = \begin{cases} 1 & \text{if } freq(s_j, u_i) \geq thre \\ 0 & \text{otherwise} \end{cases}$$

ただし、 $thre$ は閾値であり本論文では予備実験により 1 としている。

- IDF：

$$a_{ij} = \begin{cases} idf(s_j) & \text{if } freq(s_j, u_i) \geq thre \\ 0 & \text{otherwise} \end{cases}$$

- TFIDF：

$$a_{ij} = tfidf(s_j, u_i)$$

- 上記 4 つの方法を、各ユーザごとに 1 になるようにそれぞれ正規化したもの。ただし m はセンサ数である。

$$a_{ij}^{normalized} = \frac{a_{ij}}{\sum_{i=1}^m a_{ij}} \quad (4)$$

次章では、これらの方法の比較を行う。

5. 分析と考察

本章では、3.2 節で述べた位置情報取得の実験によ

るデータを用いて、ユーザ属性を推定する。その結果を示し、考察を行う。

5.1 ユーザ属性推定の精度

本研究で用いたユーザ属性を表 1 に示す。各ユーザ属性ごとに、属性値を判別する分類問題を作る。SVM は 2 値の分類が基本であるので、ユーザ属性の各属性値ごとに分類問題を生成している。たとえば年齢であれば、とりうる属性値が 5 つあるので、計 5 個の分類問題が生成される。各分類問題ごとに SVM で学習を行い、Leave-one-out 法により評価を行う。なお、予備実験により最もパフォーマンスの良かった RBF カーネルを用いた。

分類精度を表 2 に示す。センサの検出回数からの属性の構築法 (8 種類) による Recall と Precision、および F 値を表している。F 値とは、Recall と Precision の相乗平均であり、次式で表される。

$$F = \frac{2\text{Recall} \cdot \text{Precision}}{\text{Recall} + \text{Precision}}$$

表中の値は全分類問題に対する平均である。たとえば、検出回数を用いたときには、Recall は 73.56 であり、あるユーザ属性が特定の値をとる人の 73.56% を見つけることができることを示している。また、Precision が 37.75 なので、分類器が特定の属性値であると出力した人のうち 37.75% が正解であったということになる。全般に、Recall が 70% 程度、Precision が 50% 弱であり、それほど高いわけではない。これは、ユーザ属性の中には位置履歴からの確に推測できるものもあれば、そうでないものも含まれているためである。

検出回数からの属性の構築法ごとの違いを見ると、検出回数を正規化したものが最も F 値が高く、ついで TFIDF を正規化したものが高い値となっている。正規化を行う方が、行わない場合よりもつねに良い結果となっており、ユーザによるセンサの検出回数の違いを吸収するために正規化が重要であることを示している。最も良いものと悪いものでは、10 ポイント以上の差があり、センサの検出回数をどう処理するかが、ユーザ属性の推定のうえで重要であることが分かる。

5.2 ユーザ属性ごとの考察

次に、各ユーザ属性ごとに詳しく F 値を調べたものが表 3 である。センサの検出回数からの属性の構築法として、正規化した検出回数を用いている。各ユーザ属性ごとに F 値に大きな差があることが分かる。

ここで注意しなければならないのは、ユーザ属性の

多項式カーネル、線形カーネルでも全般的な傾向は変わらず、F 値が数ポイント程度悪化する。

表 1 ユーザ属性
Table 1 User properties.

属性	とりうる属性値	備考
年齢	25 歳未満, 24~29 歳, 30~34 歳, 35~39 歳, 40 歳以上	SC は幹部である
身分	SC, 常勤職員, 非常勤職員, 補助スタッフ, 一時的なスタッフ	
チーム	研究グループ A, B, C, D, 研究系事務, 事務部門	コーヒーを飲む頻度
出勤頻度	高, 中, 低	
コーヒー	高, 中, 低	
喫煙	する, しない	
居室	A, B, C, D, E, F	常駐する部屋
通勤方法	駅 A, 駅 B	利用できる駅が 2 つある

表 2 属性の計算方法による分類精度の違い

Table 2 Classification performance depending on various feature weighting.

	F 値 (%)	Recall (%)	Precision (%)
検出回数	44.45	73.56	37.75
バイナリ	43.92	65.83	41.62
TFIDF	44.28	71.62	37.38
IDF	44.37	68.45	45.33
検出回数 (正規化)	54.46	68.83	49.23
バイナリ (正規化)	40.73	63.80	40.97
IDF (正規化)	41.23	61.02	41.46
TFIDF (正規化)	53.00	65.50	47.88

値ごとに正例の数 (その属性値をとるユーザの数) が異なることである。正例の数が少ない場合には高い F 値が出にくい難しい分類問題となる。逆に、ほとんどが正例であれば、極端にはすべてを正と出力する分類器でも高い F 値となる。たとえば、喫煙しないという smoking0 の属性がこれに該当する。したがって、すべてを正と出力する分類器における F 値をベースライン (BL) とし、F 値がベースラインと比べて 10 ポイント以上上回っているものに*印をつけている

各ユーザ属性ごとの考察を以下に述べる。

年齢 「25 歳未満」「40 歳以上」という両端の属性値は F 値で 60 以上とある程度的確に推測できている。しかし、中間的な属性値に対しては、ベースラインとほとんど変わらない。

身分 「常勤職員」「非常勤職員」といった区別は名目的なものであり、非常勤職員であっても常勤している人もある。このため、この分類は難しい。一時的に来ているスタッフは、比較的いろいろな場所で検知される特徴があり特定しやすい。

チーム 「研究系事務」「事務部門」は高い F 値となっている。これは、いつも同じ場所で座っており、定型的な業務が多いからであると考えられる。それに対して、研究スタッフは、出入りが多く、打合せなどで移動することも多い。研究グループと居室は必ずしも一致していないが、比較的固まっているグループ (グループ A や B) は推定がし

表 3 各ユーザ属性による F 値の違い

Table 3 F-value for each user property.

属性	F 値	(BL)	属性値	
年齢	(平均)	41.79	(32.95)	
	age0*	63.15	(25.53)	25 歳未満
	age1	17.39	(25.53)	24~29 歳
	age2	48.49	(45.28)	30~34 歳
	age3	19.05	(29.17)	35~39 歳
age4*	60.87	(39.22)	40 歳以上	
身分	(平均)	41.33	(32.66)	
	position0	13.33	(20.00)	SC
	position1	41.03	(44.83)	常勤職員
	position2*	60.61	(42.11)	非常勤職員
	position3	25.00	(36.36)	補助スタッフ
position4*	66.67	(20.00)	一時的なスタッフ	
チーム	(平均)	59.18	(28.50)	
	team0*	58.83	(34.04)	研究グループ A
	team1*	60.00	(30.43)	研究グループ B
	team2	21.05	(22.73)	研究グループ C
	team3*	46.15	(26.67)	研究グループ D
	team4*	83.33	(26.67)	研究系事務
team5*	85.71	(30.43)	事務部門	
出勤頻度	(平均)	48.75	(49.42)	
	atd0	50.00	(60.00)	高
	atd1	48.65	(51.06)	中
	atd2*	47.62	(37.21)	低
コーヒー	(平均)	57.22	(48.90)	
	coffee0	58.83	(56.41)	低
	coffee1	66.67	(60.00)	中
	coffee2*	46.15	(30.30)	高
喫煙	(平均)	64.13	(57.58)	
	smoking0	91.89	(93.33)	しない
	smoking1*	36.36	(22.22)	する
居室	(平均)	67.00	(28.29)	
	room0*	72.73	(23.08)	B
	room1*	77.78	(26.42)	C
	room2*	66.67	(19.61)	D
	room3*	62.50	(29.63)	A
	room4*	46.15	(29.63)	E
room5*	76.19	(41.38)	F	
通勤方法	(平均)	61.86	(63.74)	
	station0	83.72	(83.72)	駅 A
	station1	40.00	(43.75)	駅 B
総計	54.46	(38.38)		

やすい。

出勤頻度 ベースラインとほとんど変わらない。事前に予想していた結果と反対であった（正規化した検出回数以外の属性構築でも同じであった）。このユーザ属性は、普段の出勤頻度を表すものであるが、センサから得られるのは当該の1週間における検出回数であり、普段は多く来ている人でもたまたま出張と重なったり、センサに検出されにくかったりという理由で推測できないためであると考えられる。逆に、普段からあまり来ない「低」の属性値については、ベースラインよりも若干良い結果となっている。

コーヒー 全般にあまりF値は高くないが、ベースラインと比較すると、コーヒーをよく飲む人は検出しやすいことが分かる。

喫煙 喫煙する人は、ベースラインと比較して高いF値になっている。これは、喫煙所付近のセンサで判断できるためである。しかし、喫煙所はラウンジに隣接しており、飲み物を飲んでいるだけでも喫煙所で検出される可能性もある。また、喫煙するときにはカード型 CoBIT を外していたという人もおり、こういった要因が精度を下げる要因であると考えられる。

居室 居室は、最も位置情報と関連したユーザ属性であり、いずれも70%近い高い値となっている。居室Eは低いですが、これはきちんと区切られた部屋ではないオープンスペースであり、検出するのが難しい。

通勤方法 ベースラインと変わらない。1日の最初あるいは最後にどちらの入り口で検出されるかで推定できるのではないかと予想していたが、良い結果にはならなかった。時間情報を用いれば改善される可能性がある。

ユーザ属性ごとに、よく推定できるものもあればそうでないものもあり、特にどういった種類のユーザ属性か（デモグラフィック属性やユーザの性質に関する属性など）による大きな違いは見られなかった。よく推定できるユーザ属性かそうでないかは事前の予想とは異なっている場合も多かった。SVMによって学習されたモデル（センサの検出回数からの属性の重み）を見ると、この様子をよりの確に理解することができる。たとえば、これまでに説明の例でも用いてきたように、コーヒーというユーザ属性の予測に最も寄与するのは、コーヒーサーバの前のセンサであると考えていた。しかし、実際にはコーヒーサーバ付近の2つのセンサの寄与は5位、9位であり、意外にもコーヒーサーバの

前のセンサはわずかではあるが負の影響であった（つまり、コーヒーサーバの前で検出されれば、コーヒーを飲む頻度が低いと推定される確率が高まる）。最も寄与の高かったセンサは、お茶に集まる共用の机の前のセンサであった。たしかに、コーヒーサーバの前は、コピー機やドアがあり、そこで検出されたからといって必ずしもコーヒーを取りにいったとは限らないこと、また実際にコーヒーを取った人は共用の机で飲むか、そこにお茶菓子を取りに行くことが多いことから、この結果は後から考えれば妥当であった。

このように、本手法では、単純に事前に予想されるような場所とユーザ属性の結び付きではなく、学習データに基づいたユーザ属性の推定が行われている。

6. センサの重要度の提案

前章では、ユーザ属性とユーザの位置履歴があらかじめ分かっているとき、推定しやすいユーザ属性があることを述べた。しかし現実には、ユーザの位置履歴とユーザ属性が分かっており、その関連が学習できる状況は少ないであろう。

では、取得できるユーザ属性が分からないときに、どのセンサがユーザ属性の推定のために潜在的に有効であるか調べる方法はないだろうか？本章では、前章までのように学習データがある場合ではなく、学習データがない場合に、ユーザ属性の推定に寄与するセンサの重要度を測る方法について述べる。

6.1 センサの重要度の提案

最も簡単なセンサの重要度の定義として、センサの総検出回数がある。1回もユーザを検出しないセンサは、設置しておいてもユーザ属性の推定に有効でないことは明らかである。ほかに、総検出回数が多いセンサではなく、多くのユーザを検出するセンサ、また他のセンサで検出しないようなユーザを検出するセンサを重要と考えるべきかもしれない。

本章では、センサの重要度として次のものを定義して比較する。

- 総検出回数 (freq): $w(s_j) = \sum_{j=1}^n \text{freq}(s_j, u_i)$
- 検出ユーザ数 (user-freq): $w(s_j) = uf(s_j)$
- TFIDF 合計 (tfidf): $w(s_j) = \sum_{j=1}^n \text{tfidf}(s_j, u_i)$
- 上記3つの方法について、ユーザごとに正規化したうえで全ユーザについて足し合わせたもの (normalized と表記する): 式(4)を用いて

$$w(s_j) = \sum_{j=1}^n a_{ij}^{\text{normalized}}$$

さらに、次のような重要度を考える。他のセンサで

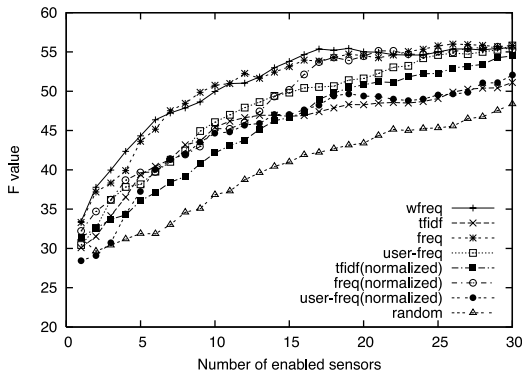


図 5 各重要度ごとのセンサ数と F 値

Fig. 5 Number of enabled sensors by each weighting versus F-value.

検出していないようなユーザを検出しているセンサは重要であろう。したがって、各ユーザがどのくらい他のセンサで検出されているかを考慮した検出回数を定義する。

- 重みづけ検出回数 (wfreq) :

$$w(s_j) = \sum_{i=1}^n freq(s_j, u_i) \times \log(m/sf(u_i))$$

ここで、 $sf(u_i)$ はユーザ u_i がいくつのセンサで検出されたかを表す。これは、TFIDF の行と列を逆にしたものととらえることもできる。

以上の 7 つの重要度について、比較を行う。

6.2 比較

センサの重要度の高いものから順に用いたときに、ユーザ属性の推定精度はどのように上がっていくであろうか。図 5 は、前章で行った SVM による推定精度が、用いるセンサ数の増加とともにどのように上がっていくを示している。最終的にすべてのセンサを用いると、結果は前章と一致する。検出回数からの属性の構築には正規化した検出回数を用い、ユーザ属性には、位置履歴からの推定が有効なもの（表 3 中で*印のついている 18 個の属性）だけを用いている。図中の数字は*印のついたユーザ属性に対しての F 値の平均である。センサの重要度はどのユーザ属性を用いるかに密接に関わるが、ここでは位置履歴からの推定の効果のある 18 個のユーザ属性すべてに対して適用している。このような方法は、文書分類における語の重みの計算法の比較^{28),29)}でよく行われる方法である。

センサ数が少数のときは F 値は低いが、センサ数が増えるとともに F 値が徐々に上がっていく。一番下の折れ線 (random) は、ランダムにセンサを並べたものであり、F 値の上昇が一番遅い。最も良いもの

が、重みづけ検出回数 (wfreq) と検出回数 (freq) であり、ユーザ属性の推定に寄与するセンサを最も効果的に特定していることが分かる。他の方法 (検出ユーザ数 (user-freq), TFIDF 合計 (tfidf) やその正規化したもの) はランダムよりは良いが、上位 2 つには劣る。特に、5 章で良い結果を示した、正規化した検出回数は、センサの重要度としては有効でないという結果であった。

センサの重要度を計算することで、ユーザ属性の推定に潜在的に寄与するセンサを特定することができる。したがって、環境中のセンサの数を絞りたいときに重要度の低いセンサを取り除いたり、またセンサを追加するときに重要度の高いセンサの近辺に配置したりするといった形で利用することができる。

7. 議論

センサで得られた情報は、いくつかの前処理を行う必要がある。たとえば、センサの検出間隔によっては、同じ場所で連続して何回も検出されてしまうことがある。この場合はある時間内に同じ場所で検出された場合は 1 回と数えるなどの工夫が必要である。また、地点 A-B-C がつながっている場合など、A-C と検出されれば、実は A-B-C と検出されるべきところが検出漏れがあったと推測することができる。こういった前処理は本論文では行っていないが、より詳細な分析を行う際には必要となると考えられる。

本論文におけるユーザ-センサ行列では、時間の情報を用いていない。時間情報はユーザ属性を知るうえで重要な情報となりうるが、次の 2 つの理由のために本論文では捨象している。その理由は次の 2 つである。

- たとえば、滞在している、通過している、複数回訪れているなどの情報をどのように用いるか、さまざまな選択肢が存在する。ヒューリスティックなルールを多く作って精度を上げるよりも、本論文では、まずは一般的に適用可能な方法を構築することを目的とした。
- 情報検索では bag of words と呼ばれる単語の出現順序を考慮しないモデルがよく使われ、有効であることが知られている。時間情報を捨象することで、情報検索におけるさまざまな手法が適用可能になる。

ユーザの位置が検出された時間を考慮し、位置の時系列的な情報を用いることも可能である。この場合に比べて、提案手法は次のようなメリットがある。

- 処理が簡単であること。ベースラインとして用いることができる。

- センサの検出漏れに対してロバストである。
- センサ間の同期をとるのが難しい場合でも利用可能である。

しかし、時間的な情報を用いることでさらに詳細なユーザ属性の推定が可能になる。たとえば、ユーザがある場所を定期的に訪れている場合には、短時間に同じ回数だけそこで検出されるよりも、ユーザにとって重要な場所である可能性が高い。これを扱うには、長時間の間隔があった場合にはセンサ検出の重みを上げる、時間帯や曜日などと合わせて検出情報を処理するなどが必要である。このように、ドメインにできるだけ依存しない形で、時間情報を考慮して精度を上げていくことは今後の課題である。

本論文では、ユーザの位置履歴がすべて得られている状況でユーザ属性を推定する方法を示している。ユーザ属性とあわせて行動分析や事後の情報支援には本論文の手法は有効であるが、ユーザ属性の推定をただちに情報支援に生かそうとすると、リアルタイムでのユーザ属性の推定が必要である。こういった手法についても、時間情報の利用と同時に検討していく必要があると考えている。

本論文で用いている位置履歴データやユーザ属性は、1回の実験によるものである。本来は、複数回、さまざまな状況を想定した実験を行うことが望ましいが、環境にセンサを配置するための物品の準備や配置・配線計画、配置許可や実際の配置作業、また、デバイスを常時管理したり参加者に実験の説明を行ったりする人員の確保など、作業負荷が非常に大きいことがネックになる。本実験では、3カ月以上の計画期間を経て、10人以上の協力者によってセンサの設置作業を行った。将来的に、センサを設置したユビキタス情報環境が増えてくればよいが、こういった研究を進めていくには、当面はいかに環境を簡単に設置するかという要素も重要であろう。

最後に、2005年から個人情報の保護に関する法律が施行され、個人情報に対する意識が高まっている。ユビキタス情報環境においてユーザ個人に合わせた支援を目指すとして、研究から実証段階において個人情報の取扱いに注意する必要がある。本研究では、参加者に研究の目的を簡単に説明し、データを分析する旨の許可を得ている。一方で、現時点では技術的に、ユーザの行動のデータや各種センサのデータから、ユーザに関する情報がどの程度推測されるのか明らかではない。今後、本研究を進めていくことで、利用者にとっても取られたくないセンサ情報、出してもよいセンサ情報を判断する参考になるのではないかと考えている。

8. 結 論

本論文では、ユーザの位置情報をもとにユーザ属性を推定する手法について述べた。本論文で提案した手法は、ユーザ-センサ行列に基づく手法であり、情報検索におけるさまざまな手法が適用できる。SVMによるユーザ属性の推定法を示し、特にセンサの検出回数をユーザごとに正規化する方法、もしくは正規化した $tf \cdot idf$ を用いる方法が良いことが分かった。さらに、センサの重要度を計算する方法を提案し、総検出回数、もしくは重みづけ検出回数がユーザ属性を推定するうえで良い指標であることを示した。

本論文で提案した手法は非常にシンプルな方法であり、位置情報からユーザ属性を推定するという研究において、利用しやすいベースラインを提供すると考えている。今後、時間情報の加工、空間の意味的な情報の付加によって、ユーザ属性の精度をどのように向上させることができるか、検討していく予定である。

謝辞 本研究は、総務省戦略的情報通信研究開発推進制度研究主体育成型研究開発「建物内の位置履歴からのユーザモデリングに関する研究」の助成を得て行ったものである。

参 考 文 献

- 1) 総務省：何でもどこでもネットワークの実現に向けて—「ユビキタスネットワーク技術の将来展望に関する調査研究会」報告書(2002)。
- 2) 総務省：平成16年版情報通信白書世界に広がるユビキタスネットワーク社会の構築(2004)。
- 3) 中島秀之, 橋田浩一, 森 彰, 伊藤日出男, 本村陽一, 車谷浩一, 山本吉伸, 和泉 潔, 野田五十樹: 情報インフラに基づくグラウンディングとその応用—サイバーアシストプロジェクトの概要—, コンピュータソフトウェア, Vol.18, No.4, pp.48-56 (2001)。
- 4) Hightower, J. and Borriello, G.: Location Systems for Ubiquitous Computing, *IEEE Computer*, Vol.34, No.8, pp.57-66 (2001)。
- 5) Want, R., Hopper, A., Falcao, V. and Gibbons, J.: The Active Badge Location System, *ACM Trans. Information Systems*, Vol.10, No.1, pp.91-102 (1992)。
- 6) 美濃導彦: ユビキタスホームにおける生活支援, 人工知能学会論文誌, Vol.20, No.5, pp.579-586 (2005)。
- 7) Wilson, D.H.: The Narrator: A Daily Activity Summarizer Using Simple Sensors in an Instrumented Environment, *Proc. UbiComp 2003* (2003)。
- 8) 西田佳史, 本村陽一, 山中龍宏: 乳幼児事故予防

- のための日常行動モデリング, 情報処理, Vol.46, No.12, pp.1373-1381 (2005).
- 9) 九州総合通信局: 無線技術を活用した顧客管理・情報提供システムに関する調査研究会—公開試験等の概要 (2005).
 - 10) 車谷浩一, 山下倫央, 和泉憲明, 幸島明男, 和泉潔: 愛・地球博グローバル・ハウス統合情報支援システム-CONSORTS アーキテクチャによる情報提供・会場運営支援システム, 情報処理学会誌, Vol.47, No.2, pp.105-108 (2006).
 - 11) Kobsa, A.: Generic User Modeling Systems, *User Modeling and User-Adapted Interaction*, Vol.11, pp.49-63 (2001).
 - 12) Brusilovsky, P.: Methods and techniques of adaptive hypermedia, *User Modeling and User-Adapted Interaction*, Vol.6, pp.87-129 (1996).
 - 13) Heckmann, D.: Ubiquitous Use Modeling, Ph.D. thesis, University of Saarland (2005).
 - 14) 小磯貴史, 服部可奈子, 吉田琢史, 今崎直樹: 歩行者動線分析システムを用いた大型家電量販店での行動分析, 情報処理学会ユビキタスコンピューティングシステム研究会研究報告, No.2003-UBI-002, pp.61-66 (2003).
 - 15) パコアンダーヒル: なぜこの店で買ってしまうのか ショッピングの科学, 早川書房 (2001).
 - 16) Wilson, D., Long, A. and Atkeson, C.: A ContextAware Recognition Survey for Data Collection Using Ubiquitous Sensors in the Home, *Proc. CHI 2005* (2005).
 - 17) Ashbrook, D. and Starner, T.: Using GPS to learn significant locations and predict movement across multiple users, *Personal and Ubiquitous Computing*, Vol.7, No.5, pp.275-286 (2003).
 - 18) Hightower, J., Consolvo, S., LaMarca, A., Smith, I. and Hughes, J.: Learning and Recognizing the Places We Go, *Proc. UbiComp 2005* (2005).
 - 19) 中村嘉志, 西村拓一, 伊藤日出男, 中島秀之: 無電源でユーザ属性と位置を発信する CHOBIT 端末の設計と実装, 情報処理学会論文誌, Vol.44, No.11, pp.2670-2680 (2003).
 - 20) Nakamura, Y., Nishimura, T., Itoh, H. and Nakashima, H.: ID-CoBIT: A Battery-less Information Terminal with Data Upload Capability, *Proc. IECON 2003* (2003).
 - 21) 西村拓一, 濱崎雅弘, 松尾 豊, 大向一輝, 友部博教, 武田英明: 2003 年度人工知能学会全国大会支援統合システム, 人工知能学会誌, Vol.19, No.1, pp.43-51 (2004).
 - 22) Nishimura, T., Nakamura, Y., Itoh, H. and Nakamura, H.: System Design of Event Space Information Support Utilizing CoBITs, *Proc. ICDCS 2004*, pp.384-387 (2004).
 - 23) Schulz, D., Fox, D. and Hightower, J.: People Tracking with Anonymous and ID-Sensors Using Rao-Blackwellised Particle Filters, *Proc. IJCAI-03*, pp.921-928 (2003).
 - 24) Manning, C. and Schütze, H.: *Foundations of statistical natural language processing*, The MIT Press, (2002).
 - 25) Vapnik, V.: *The Nature of Statistical Learning Theory*, Springer-Verlag (1995).
 - 26) ネロクリスティアニーニ, ジョンショー-テイラー (著), 大北 剛 (訳): サポートベクターマシン入門, 共立出版 (2005).
 - 27) 前田英作: 痛快! サポートベクトルマシン—古くて新しいパターン認識手法, 情報処理, Vol.42, No.7, pp.676-683 (2001).
 - 28) Joachims, T.: Text categorization with support vector machines, *Proc. ECML'98*, pp.137-142 (1998).
 - 29) Mladenic, D., Brank, J., Grobelnik, M. and Milic-Frayling, N.: Feature selection using linear classifier weights: interaction with classification models, *Proc. SIGIR 2004*, pp.234-241 (2004).

(平成 18 年 8 月 10 日受付)

(平成 19 年 3 月 1 日採録)



松尾 豊 (正会員)

1997 年東京大学工学部電子情報工学科卒業。2002 年同大学院博士課程修了。博士 (工学)。同年より、産業技術総合研究所情報技術研究部門勤務, 2005 年 10 月よりスタンフォード大学客員研究員。人工知能, 特に高次 Web マイニングに興味がある。人工知能学会, 言語処理学会, AAAI, INSNA の各会員。



岡崎 直観 (正会員)

2001年東京大学工学部電子情報工学科卒業。2003年同大学院情報理工学系研究科修士課程修了。2003年同研究科博士課程進学。2005年より英国国立テキストマイニングセンター (National Centre for Text Mining) のリサーチ・アシスタント。2006年帰国。2007年東京大学大学院情報理工学系研究科博士課程修了。現在、同大学院情報理工学系研究科コンピュータ科学専攻研究員 (科学技術振興特任教員)。文書自動要約、用語抽出を中心にテキストマイニングの研究を行っている。言語処理学会の会員。



中村 嘉志 (正会員)

1994年神奈川大学理学部情報科学科卒業。1996年電気通信大学大学院情報システム学研究科博士前期課程修了。1997年同大学院同研究科博士後期課程退学。同年同大学院同研究科助手を経て、現在、産業技術総合研究所情報技術研究部門研究員。2006年より芝浦工業大学客員助教授 (連携大学院) を併任。博士 (工学)。ロケーション・ウェアな情報支援システムの研究に従事。センサネットワークと実世界指向のヒューマン・コンピュータ・インタラクションに興味を持つ。電子情報通信学会、人工知能学会、IEEE 各会員。



西村 拓一 (正会員)

1992年東京大学工学系大学院修士 (計測工学) 課程修了。同年NKK (株)入社。X線、音響・振動制御関係の研究開発に従事。1995年RWCPに出向、1998年NKK (株)復帰。1999年RWCPつくば研究センタに所属。2001年産業技術総合研究所サイバーアシスト研究センターに所属、2005年同情報技術研究部門実世界指向インタラクショングループ長、筑波大学大学院知能機能システム専攻助教授 (連携大学院)、現在に至る。博士 (工学)。時系列データ検索・認識、実世界情報支援に興味を持つ。電子情報通信学会、人工知能学会、ヒューマンインタフェース学会、ACM 各会員。



橋田 浩一 (正会員)

1981年東京大学理学部情報科学科卒業。1986年同大学院理学系研究科情報科学専門課程修了。理学博士。同年電子技術総合研究所入所。1988年から1992年まで (財) 新世代コンピュータ技術開発機構に出向。現在電子技術総合研究所主任研究官。自然言語処理、人工知能等の研究に従事。編著書に『岩波講座 認知科学』(共編、岩波書店) 等。



中島 秀之 (フェロー)

1983年東京大学大学院情報工学専門課程修了 (工学博士)。同年電子技術総合研究所入所。人工知能を状況依存性の観点から研究。マルチエージェントならびに複雑系の情報処理とその応用に興味を持っている。2001年より産業技術総合研究所サイバーアシスト研究センター長。2004年より公立ほこだて未来大学学長。産業技術総合研究所情報技術研究部門研究顧問。認知科学会元会長、ソフトウェア科学会元理事、人工知能学会元理事、情報処理学会副会長、同フェロー。マルチエージェントシステム国際財団元理事。日本工学アカデミー会員、電子情報通信学会会員。