

アカウントティング情報を用いた利用者の専門性分析

神奈川大学理学部情報科学科

田代 淳一[†] 堀 幸雄[‡] 後藤 英一[§]

現在アカウントティング情報は主に利用者への課金目的に利用されている。しかしアカウントティング情報からはその使用傾向や行動パターンが分析可能であり、個人の知識管理に利用できると考えられる。本研究ではアカウントティング情報を用いて利用者の動向分析を行い利用者の類似性と専門性の評価を行なった。クラスター分析を用いて類似性の判定を行なった結果、同じ専門性を持つ利用者を正しく分析し、その要素を限定できることを確認した。

Expertise analysis of the user using accounting information

Department of Information Science,
Faculty of Science, Kanagawa University

Junichi Tashiro[†] Yukio Hori[‡] Eiichi Goto[§]

Accounting information is mainly used for a user's fee collection purpose. However, it is thought that the use tendency and action pattern can be analyzed and it can use for knowledge management of an individual from accounting information. In this study, trend analysis of a user was performed using accounting information, and evaluation of a user's similarity and speciality nature was performed. As a result of judging similarity using cluster analysis, the user with the same speciality nature was analyzed correctly, and it checked that the element could be limited.

1 はじめに

本研究着手にあたる背景は、次の2つの点に着眼したことである。1. アカウントティング情報が利用者を与える有用性が極めて少ない。2. IT技術者の計算機スキルを体系化した指標がこれまで統一されていない。

アカウントティング情報は容易に取得することができ、個人の計算機に関する情報が網羅されているにも関わらず利用者向けの情報ではなく管理者向けの情報として認識、利用されているため、利用者に対して特に有用されていないのである。

また計算機スキルの指標の統一されていないことも大きな問題点であり、統一化は非常に重要なことである。近年IT産業の売上げ及び利益がハードからソフト・サービスにシフトするに伴い、戦略的・体系的なスキル管理・育成の重要性が高まってきている。メインフレームの時代には、スキル伝承手段としてOJT(On the Job Training)が機能していた。しかし各個人のスキルの明確化が情報サービスを提供している企業にとって課題となっている今日においては、従来のスキル管

理を見直す必要がある。すなわち、市場において求められるスキルが多様化・深化するに伴い誰が何をできるのか客観的に指標化し各個人もまた把握することがスキルアップのうえで重要となってくるのである。国の対応としても2002年12月に経済産業省がITサービスに必要とされる実務能力を体系化した指標である「ITスキル標準」^[1]が策定段階であり、国内における能力の指標化の整備が必須課題であることが明らかである。

また商業においても購買履歴といった顧客情報をマーケティングや販売戦略に活かすようにアカウントティング情報も利用者情報としての価値を有用化し組織の活性化につなげたいと著者は考える。

本研究の目的は、上記のような側面を統合的に考慮しながらアカウントティング情報から利用者の専門性が分析できることを示すことである。

2 アカウントティング情報と専門性

2.1 アカウントティング情報とは

アカウントティングシステムは、主にユーザに課金するため、システムリソースの使用量を記録するよう

[†]tashijun@goto.info.kanagawa-u.ac.jp

[‡]horiyuki@goto.info.kanagawa-u.ac.jp

[§]goto@goto.info.kanagawa-u.ac.jp

に設計されている。一般的にUNIXでは、ユーザベ-
 ースのプロセスの情報を記録・管理している。すなわち、
 オペレーティングシステムでは、実行される各プロセ
 スについて、ユーザ含む統計値を含むことにより、シス
 テムの使用量が記録される。アカウントシステムによ
 って収集されたデータは、システムパフォーマンスの
 監視にも使用することができる。アカウントイン
 グレコードには、システムで実行する各プロセスに
 関する次のようなデータが収められている [2]。

- コマンドイメージ名
- 使用したCPU タイム
- プロセスが完了するまでの経過時間
- プロセスが開始された時刻
- 対応するユーザID とグループID
- 存続時間でのメモリ使用量
- 読み書きした文字数
- 読み書きしたディスク I/O ブロック数
- プロセスが起動されたTTY
- プロセスに関連するアカウントingフラグ
- プロセスの終了ステータス

上記のデータを用いて、専門性分析を行う際にいかな
 るデータ処理が必要かについて次節で論ずる。

2.2 専門性分析手法

前節で述べたデータをアカウントing情報から取
 得し、とくに使用したプロセスを重み付けを行う。こ
 のとき成分を v 、サンプルサイズを n 、重み付け係数を
 w とし、

$$v_j = \sum_{i=1}^n w^{(i-1)}$$

で表した。 w は、0.99 とした。これはプロセスの単純
 な線形総和の場合、特徴空間で出現頻度の莫大な差が
 あるためこの範囲を 0 ~ 100 までの間に正規化してい
 る。すなわち、このデータがアカウントing情報から
 クラスタ分析を行う際の対象である。クラスタ
 分析では、クラスタ間の距離 (非類似度) 関数に基づ
 き、最も距離の近い二つのクラスタを逐次的に併合す
 る。そして、この併合を、すべての対象が一つのクラ
 スタに併合されるまで繰り返すことで階層構造を獲得
 する。この階層構造はデンドグラムによって表示する。

デンドグラムは、各終端ノードが各対象を表し、併合
 されてきたクラスタを非終端ノードで表した二分木
 である。非終端ノードの横軸は、併合されたときのク
 ラスタ間の距離を表す。クラスタを形成する際の対象
 間の距離を求めるにはいくつかの手法があるがウォー
 ド法 (Ward's method) [6] を適用した。ウォード法は以
 下のような一般式で表せる。

$$D(C_1, C_2) = E(C_1 \cup C_2) - E(C_1) - E(C_2)$$

$$\text{ただし、} E(C_1) = \sum_{x \in C_1} (D(x, c_1))^2$$

ウォード法は、各対象から、その対象を含むクラ
 スタのセントロイドまでの距離の二乗の総和を最小化
 しクラスタを形成していく手法である。これはクラ
 スタ分析の中でも階層的方法の一つであるが、実用性
 が高い手法であり、鎖効果が起きにくい。鎖効果とは、
 ある一つのクラスタに対象が順に一つずつ吸収されて
 クラスタが形成されていく現象である。このようなデ
 ンドグラムが得られた場合には、どの距離で切っても
 あるクラスタとその他の対象一つずつで構成されたク
 ラスタに分かれることになり、グループに分けたこと
 にならない。

ウォード法により形成されたクラスタは、クラスタ
 間の結合距離の近いもの同士で成り立っているが本研
 究では要素の次元数が 1000 以上と膨大なためクラ
 スタ間のプロセス成分の相関関係を求めるのは難しい。
 そこで、クラスタ間に主成分分析を行い軸の回転を利
 用して対象を記述するのに最小次元を獲得、すなわち
 クラスタ間で特化しているのか成分を抽出する。

専門性の定義は、クラスタ分析によって分類され
 たクラスタで代表となる成分やクラスタ間で相関の強
 いものとする。代表となる成分とは、クラスタ内のユー
 ザが共通して使用したプロセスの中でも特に頻度の高
 いものことであり、相関が強いものとはユーザの第一、
 第二主成分得点である傾向を示している成分のこと
 である。

3 アカウントing分析実験

3.1 方法と目的

前節で述べた手法を用いて、アカウントing情報
 から利用者の類似性を見つけ、専門性となりうる要素
 を限定することを目的とする。

今回は、神奈川大学の教育用共同計算機システムで、
 平成 14 年 3 月から平成 14 年 12 月までの本校の利用
 者 (理学部情報科、化学科、生物科、経営学部国際経
 営学科、学生と教職員) のアカウントing情報を対

象とした。クラスタ分析はワード法の自作実験プログラムを用いて、ユーザ ID 同士のクラスタリングを行う。最終的に一つのクラスタに全てのクラスタが統合されたら、終了する。その結果を用いてデンドグラムを形成する。また、その後適当な部分クラスタを選択し、そのクラスタ内での成分値を実験プログラムで主成分分析を行い、クラスタにおける成分の主成分得点係数と対象ユーザの主成分得点を出し、変量とユーザについて散布図を作る。

3.2 実験結果および考察

上記手続きを実行した結果、全てのケースにおいて、デンドグラム、散布図を得ることができた。その中で特に注目すべき結果を図1～図5に載せた。これは情報科学科3回生のクラスター分析と主成分分析によるデンドグラムと散布図である。特に3回生を利用した理由として、授業カリキュラムで主に情報系専攻科目が大半を占めており、プログラミングを用いた演習や実験等が多く、他学年の情報科学科あるいは他学科に比べ計算機システムの利用頻度も高いからである。

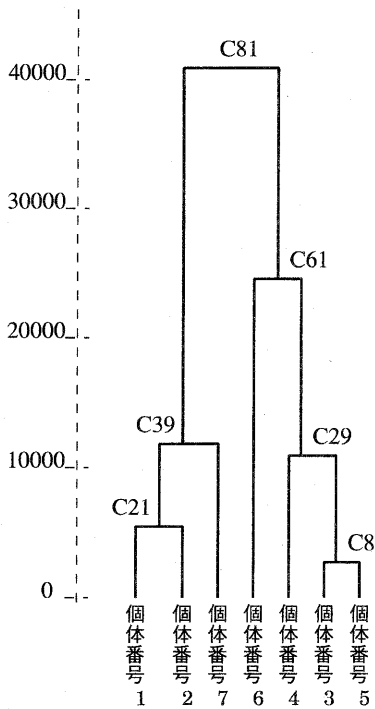


図1: クラスタ (C81) のデンドグラム

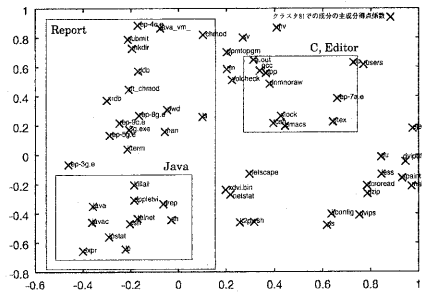


図2: クラスタ (C81) の変量プロット

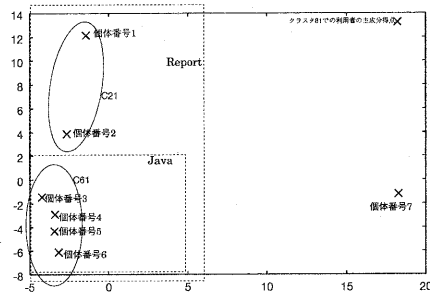


図3: クラスタ (C81) のサンプルプロット

図1のデンドグラムにある個体番号1～7は、次のサンプルプロットと変量プロットの個体番号のことである。ここで変量プロットについては、主成分分析対象である成分からあらかじめ説明変量となりうる成分を著者の勘案で選別し、どの利用者でもログインしているときに用いるであろう成分については除いてあることを述べておく。変量プロットは第1主成分得点係数を横軸、第2主成分得点係数を縦軸とした成分の散布図である。またサンプルプロットは、各変量に割り当てられた主成分得点係数を用いて利用者の主成分得点を出したものであり、横軸を第1主成分得点、縦軸を第2主成分得点とした散布図である。

この変量プロットとサンプルプロットから興味深い発見が幾つかある。まず変量プロットを見ると、このクラスタにおいてある程度成分がその専門性毎に分かれている。横軸の第1主成分得点係数からの側面で見ると正の領域にはC言語関連のcpp,a.out, それにgccといった成分がある。負の領域には逆にJava言語を用いる際のプロセスであるjava.javac それにjavaアプレットを実行する際のappletviがある。他の成分も含まれているが、明らかにC言語系の成分とJava言語関連の成分とで正負に分かれているところから、第1主成分

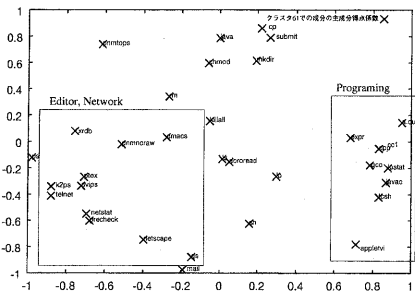


図4: クラスタ (C61) の変量プロット

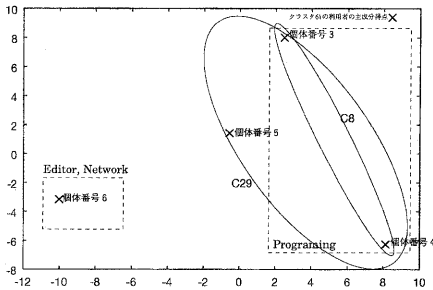


図5: クラスタ (C61) のサンプルプロット

得点係数は Java 言語と C 言語を大別する値であると解釈できる, すなわち利用者が Java 言語と C 言語のどちらの利用傾向があるかを識別しているのである. また, 第2主成分得点係数では, 正の領域に先頭が rep で始まる成分が顕著に分布している. これは, 計算機システムを用いた演習課題を課する講義の提出用実行プロセスである. 特にこの講義では通常, プログラミングは C 言語を用いる学生が既知であり, C 言語関連との正の相関があることは多いと考えられる. 従って横軸を基本として見ると第2主成分得点係数は, 提出用実行プロセスに用いた言語が Java 言語か C 言語かを判別する成分と考えることができる. また Java 言語関連の成分, C 言語関連の成分は第2主成分得点係数からみても第1主成分得点係数からみたと同じ関係にあり, 利用者の Java 言語か C 言語の専門性を大別するのは第2主成分得点係数からも全く言えることである. すなわちこれらを統合的に考慮すると, 第1主成分得点係数と第2主成分得点係数が正の場合, 普段比較的 C 言語を用いて演習などを行っていると考えることができ, 第1主成分得点係数, 第2主成分得点係数ともに負の場合, 日常的に C 言語よりは Java 言語の

ほうが利用頻度が高いことが言える.

上記の変量プロットの結果をサンプルプロットに反映させた結果, 図2のようになる. この結果を見てわかるようにクラスタ同士で散布がわかれており他のクラスタに比べ比較的クラスタ間のプロセスの専門性がでている傾向があると考えられる. クラスタ分析によりクラスタ化されたものが類似性が近く, 得る専門性がクラスタによってちがうという結果のよい例といえる. 個体番号1と2はクラスタ21で, 個体番号3~6はクラスタ61で結合されていて, クラスタ61は, どの利用者も第1, 第2主成分得点ともに負の領域にある. これは先の変量プロットをみてもわかるように, Javaを普段から利用していることが高いと思われる. したがってこのクラスタ81内においてはクラスタ61はJavaの専門性があると断定できる. また, それより第2主成分得点が正の領域では, クラスタ21の利用者が分布されている. これはクラスタ61と比べ, 横軸第1主成分得点に関してはほぼ同じ値を取るものの, 第2主成分得点に関しては, 正負で全く反対の値をとっている. すなわち, クラスタ21はJavaを普段から他の成分と比較して多く使っているという専門性はいえないが, レポート用の実行プロセスに関してはクラスタ61に比べ多く行っていることが明らかであり, 授業の課題に対する積極性が高い傾向がわかる. 今回は, レポート課題の実行プロセスだったため, 専門性とは言いがたいが課題に対する積極性が高いということも一つのアプリティとして考えれば, 一専門性として捉えることもできる.

ここで, 下位クラスタの比較をしてみる. 図4, 図5はクラスタ81の下位クラスタのそれぞれ変量プロットとサンプルプロットである. クラスタ61はクラスタ81で, 特にJavaの専門性が高かった. クラスタ61での変量プロットでは第1主成分は, 正の領域に主に gcc, a.out, cpp や javac, appletvi といった C 言語や Java 言語を用いる際のプロセスが分布しており, 負の領域には dvips, emacs や ptex といった文書作成の Latex に用いるプロセスやエディタのプロセスが分布している. また, netscape や netstat といったプロセスも負の領域にあり, 文書作成などのエディタとの相関があると考えられる. この変量の主成分得点係数を用いたサンプルプロット図5を見るとクラスタ61もクラスタ

	第一主成分 (%)	第二主成分 (%)
クラスタ 81	35.17	20.56
クラスタ 61	44.11	29.40

表1: 主成分の寄与率

81のサンプルプロットと同様、構成しているクラスタ毎に散布が分かれている。個体番号は図1のデンドグラムと図3のデンドグラムとも対応している。すなわちクラスタ61の結合でみると個体番号3と個体番号4は第1主成分得点が特に正の値を示しており、図4の変量プロットから考えると主にプログラミング言語関連のプロセスが高い頻度で利用されていることが推測でき、特に他のプロセスに比べプログラミングの専門性が高いと思われる。また個体番号6は明らかに負の領域に位置しており、他のプロセスに比べエディタやネットワークのプロセスが高く、特に頻度が高いことが分かる。また、全ての主成分分析対象成分の中での寄与している割合をみる寄与率というのがある。これをクラスタ81とクラスタ61の主成分得点係数に対応させて第2主成分までを表1に示す。これを見ると、クラスタ81は、第1主成分が35%クラスタ61では第1主成分が44%となっており、クラスタ61のほうが第1主成分が全体から見た影響度が高い。しかし寄与率を累積した場合、両方のクラスタはともに第2主成分までが50%を越えているものの、元データの全体を集約しているといえる80%を越えていないため、全ての情報を集約はしていないことがわかる。この点から累積寄与率が80%を越えている主成分までを含めた測定が次回以降、必要となる。

アカウント情報を知識として捉え、利用者の専門性かつ類似性の判定ができるかを検証した。クラスタ分析と主成分分析を融合させることでトレンド発見できることは一般的なことであり本研究でもこの手続きを適用したわけだが、必ずしもどのクラスタにおいても共通の傾向が発見できたとは言いがたい。しかし、特定の部分クラスタによっては第二主成分までの散布を見ると似た傾向が出ているクラスタもあることから、距離とコマンドの使用傾向は相関性が高い。

4 今後の課題

本研究は、アカウント情報から計算機の利用履歴を分析し個人の計算機スキルを導き出し、計算機スキルの体系化を行った。すなわちアカウント情報はアカウントを持った利用者の知識の潜在的に持っていたのである。この知識の応用としてナレッジコミュニティへのアプローチを課題として考えてみたい。アカウント情報から専門性が導き出せることは、ナレッジコミュニティにおいて専門家の特定する際の要素として用いることができないかということである。以下にそのナレッジコミュニティについて概略を説明する。

4.1 ナレッジコミュニティ

デフレの進行やB2C市場を中心としたITバブルの崩壊、中国メーカーなどが低価格攻勢を強めるなか、IT技術を業務効率化、合理化に用いることはもはや必然的なことであり他社との差別化にはならない。今や企業競争力の底上げには社内外の知識財産を一元管理し重要な経営資源として経営に活かすことが今企業がすべき重要課題の一つである。すなわち財務諸表の資本と並んで知識資本の良し悪しが企業の競争力の決め手となる時代が始まったのである^[7]。モノ作り大国であったわが国企業のホワイトカラーは、米国企業に比べて生産性が三分の一程度であるといわれていた。しかし知識社会が到来するに及んで、生産性の低いわが国の多くのホワイトカラーもいかに効率的に知識資産を経営資源に活用するかが問われているのである。

このような背景から最近、企業でデジタルネットワークを活用して社内外の専門家に技術相談できる仕組みが注目され始めており、導入企業も相次いでいる。このような仕組みを一般的にナレッジコミュニティという。具体的には、社内エキスパートや参加者の知識、得意分野をデータベースに登録し、質問を電子掲示板に書き込むと社内専門家や技術者がメールで回答するという流れである。このような仕組みを企業経営に活かす経営戦略手法をナレッジマネジメントといい企業の知識資産を活用におけるトレンドとなっている。ナレッジコミュニティは社内に埋もれている個人の潜在的な知識、ノウハウを共有可能となり会社の資産として活かすことができる。また同時にこれをイントラネットで構築することで社内で誰がどんな知識を持っているのかということも一元管理できる。特に率先してナレッジコミュニティを導入している先進企業においては、この知識管理の仕組みに人事考課の評価システムなどを併せて運用している。すなわち、どれだけナレッジを企業に還元したかということ人事考課制度として評価する仕組みである。このコラボレーションにより、社員がナレッジを表出するモチベーションにつなげているのである。

4.2 おわりに

知識社会における本研究の位置付けは、教育用計算機ルームでの計算機スキルの体系化と専門家検索である。例えば、大学の情報系の学科で学内教育用計算機システムを用いるときに利用者間の使用傾向の類似度を見ることができる。またある特定のコマンド群(言語関連、システム、ユーティリティ)などについて特に専門性が高い利用者を検索することも可能であり、ティー

チングアシスト (TA) 選出や、技術的な問題が発生したときの要員検索といった人員配置, 利用者が自分の能力を客観的に知ることに有効である. こういったアカウント情報情報が知識社会に与える有用性を本研究では, 特に取り上げたのである. これを具体化する提案として, 本研究をシステム化して大学の計算機への導入がある. 導入された環境下で利用者はコマンドラインから実行し, その時点での自分にはどのような専門性があるか, また他の利用者がどういった専門性があるかを確認でき, その分野に関する問題, 疑問がある場合に誰がその知識があるかを検索可能となり現在の Q&A 方式のナレッジコミュニティとは別の角度からの組織内の専門家に技術相談できる仕組みが構築できるであろう.

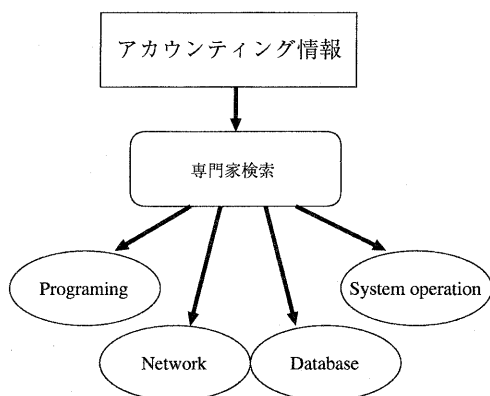


図 6: アカウント情報からの専門性発見

- [6] 神嶋敏弘: データマイニング分野のクラスタリング手法 (1), 人工知能学会誌, Vol.18 No.1 pp.59-65(2003)
- [7] 山崎秀夫 著: 未来型組織を支える企業ナレッジポータル. 野村総合研究所 (2002)

参考文献

- [1] 経済産業省 IT スキル標準 ver1.0(2002)
- [2] AEleen Frisch 著, 谷川哲司 監訳, 黒岩真吾 株式会社ユニテック 共訳: UNIX システム管理 改訂版オライリー・ジャパン (1998)
- [3] Everitt, B.S.: Cluster Analysis, Edward Arnold, third edition (1993)
- [4] Jain, A.K. and Dubes, R.C.: Algorithms for Clustering Data, Prentice Hall (1988)
- [5] Jain, A.K., Murty, M.N. and Flynn, P.J.: Data-Clustering: A Review, ACM Computing Surveys, Vol.31, No.3 (1999)