

放棄セッションのユーザ操作に着目した モバイル検索カードの順位付け

川崎 真未^{1,a)} Inho Kang^{2,b)} 酒井 哲也^{1,c)}

受付日 2018年3月8日, 採録日 2018年7月21日

概要: 本研究は、与えられたモバイルクエリに対し、より適切なカード（従来研究におけるパーティカルと呼ばれるものに近い）の順位付けを目的とする。ウェブ検索における従来研究において、URL への適合性ラベル自動付与を目的とし、クリックデータを用いて URL を順位付けするものがある。彼らの手法を本研究の目的に適用できると考えられるが、モバイル検索では、クリックせずにそれを閲覧するだけで所望の情報が得られる「良い放棄」が PC での検索よりも多く起こることが知られている。したがって、モバイル検索においては、クリックデータだけでなく放棄セッションでのユーザの操作もカードの順位付けに役立つ可能性が高い。そこで本稿では、クリックセッションに加えて放棄セッションも用いた、カードの順位付け方法を提案する。評価には、韓国で最も普及している検索エンジン NAVER の実際のモバイルクエリログを使用した。評価データは、3 人の評価者により作成された、992 のユニーククエリを含む 2,472 組のカードの種類どうしのプリファレンスデータである。放棄セッションを用いることにより、クリックセッションのみを用いた場合に比べ、5.0 ポイント高いプリファレンス精度を達成できた。また、放棄セッションで、よりクエリに適したカードを選ぶ際は、カードが一部分でも映った時間の合計（time）よりも、カード全体のうちどれくらいの部分が何度閲覧されたか（completeness）の方が有効な手がかりとなることを確認した。

キーワード: クリック, 放棄, モバイル検索, パーティカル, ランキング

Mobile Search Card Ranking based on User Actions in Abandoned Sessions

MAMI KAWASAKI^{1,a)} INHO KANG^{2,b)} TETSUYA SAKAI^{1,c)}

Received: March 8, 2018, Accepted: July 21, 2018

Abstract: We consider the problem of ranking rich verticals, which we call “cards,” for a given mobile search query. Prior art in web search suggests that it is possible to effectively rank URLs based on click data for the purpose of automatically assigning relevance grades to the URLs. While the click-based approach is applicable to our card ranking problem, it is known that *good abandonment* (i.e., the user successfully obtains the information needed without clicking on the search engine result page) happens more often in mobile search than in desktop search. Hence, to infer an appropriate card ranking for a given query, it is desirable to leverage not only click data but also user actions on abandoned sessions. In this paper, we propose a method for card ranking that utilises abandoned sessions in addition to clicked sessions. Using a real mobile query log of a Korean search engine (NAVER), we constructed a data set containing 2,472 pairwise card type preferences covering 992 distinct queries, by hiring three independent assessors. By using abandoned sessions, our proposed method outperformed a click-only baseline by 5.0 points in preference accuracy. Moreover, among the simple features that we used for selecting the most informative card from an abandoned session, we found that *completeness* (i.e., how much of the card was shown on the mobile phone screen) is more useful than *time* (i.e., for how long the card was shown to the user during the session).

Keywords: click, abandonment, mobile search, vertical, ranking

1. はじめに

スマートフォンによるモバイル検索は、ユーザが分からないことがあるときにいつでも簡単に情報を得る手段として、現代では欠かせないものとなっている。たとえば2015年にGoogleは以下のように述べている。“more Google searches were completed on mobile devices than desktop computers.”*1 このため、モバイル検索の質はつねに向上が求められている。たとえばKadotamiら[8]はユーザが検索結果に満足するよう、ヤフーのモバイル検索におけるパーティカル[2], [12], [16] (情報の種類)の順位付けを行っている。適切なパーティカルの順位付けはモバイル検索の質の向上のための1手段である。

本研究では、韓国で最も普及している検索エンジンNAVER*2のモバイル検索結果における「カード」の順位付けを題材として、モバイル検索の質の向上に取り組む。NAVERの「カード」とは、検索結果ページ上で罫線により囲まれた、クエリに関するある情報を表示する領域のことであり、従来研究におけるパーティカル[2], [12], [16]やアンサー[3], [15]と呼ばれるものに近い。図1にカードの例を示す。左のカードは地域の1週間の天気をアイコンを使って示し、選択した日の天気を地図上に示している。中央のカードは連絡先や地図などの店舗情報を示し、右のカードは映画の情報を示している。図1に示したとおり、カードは画像や文章で構成され、検索結果ページ上を占める面積もそれぞれ異なる。

ウェブ検索の研究において、Agrawalら[1]は、与えられたクエリに対し、各URLに適合性ラベル(すなわち正解ラベル)を自動的に付与する取り組みについて報告している。これは、与えられたクエリに対するクエリログをもとに、「クリックされたURLは他のURLよりもユーザが好むものであった」という仮定のもと、URLをノードとする有向グラフを構築するものである。ここで、ノード間を結ぶ有効エッジはURL間のプリファレンス(選好)を表す。Agrawalらの手法は、上記有向グラフをもとにURL間に半順序関係を与えることができる。したがって、URLに対するクリックの代わりにNAVERのカードに対するクリックを用いることにより、与えられたモバイルクエリに対し適切なカードの種類順位付けを行う我々のタスクに、彼らの手法を応用できる可能性がある。

上記のアプローチにとどまらず、一般に検索エンジンの最適化はユーザのクリックに強く依存している(たとえば

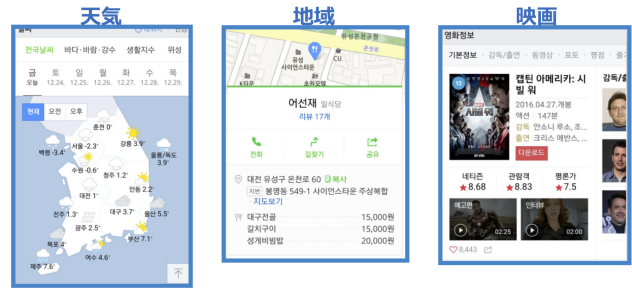


図1 カードの例
Fig. 1 Examples of cards.

文献[5], [14]). しかしながら、モバイル検索、特にNAVERのようにカードという形態で視覚的に情報を検索結果ページに提示するモバイル検索においては、ユーザが検索結果を放棄(ユーザがクリックをせずに検索セッションを終えること)するケースが少なくない。放棄には、ユーザが検索結果を閲覧するだけで所望の情報を得て満足する「良い放棄」と、ユーザが検索結果に失望しセッションを終了する「悪い放棄」があり、モバイル検索では、良い放棄がPCでの検索よりも多く起こることが知られている[11]。このことから、クリックセッション(クリックがあったセッション)をもとにクリックされたカードに関するプリファレンスを推定できるのと同様に、放棄セッション(クリックがなかったセッション)においても、モバイル検索結果画面中でユーザが情報を得たと思われるカードを推定することにより、プリファレンスを推定できるのではないかと考えた。したがって、本研究では、カードの順位付けのために、クリックセッションに加えて放棄セッションを用いることを提案する*3。なお、本研究では良い放棄と悪い放棄とを区別せずに、放棄セッションをカード間のプリファレンス推定に利用する。Liら[11]が手作業により良い放棄を抽出していることから分かるように両者の自動分類は極めて難しく、また、「悪い放棄におけるユーザの操作はカード間のプリファレンス推定に役立たない」とは必ずしもいえないからである。

2. 関連研究

以下、2.1節で、検索エンジンの正解データ作成を目的とした、クリックセッションに基づくURLの順位付けに関する従来研究について説明する。本研究は、この手法を放棄セッションを使用できるよう拡張したものである。2.2節では、放棄セッションがよい放棄であったか悪い放棄であったかを推定する従来研究について述べる。2.3節では、モバイル検索におけるビューポート(検索結果ペー

¹ 早稲田大学基幹理工学研究所
Graduate School of Fundamental Science and Engineering,
Waseda University, Shinjyuku, Tokyo 169-0072, Japan

² Naver Corporation
Naver Corporation, Seongnam-si, Gyeonggi-do, Korea

a) marerhg@ruri.waseda.jp
b) once.ihkang@navercorp.com
c) tetsuyasakai@acm.org

*1 <https://techcrunch.com/2015/10/08/mobile-searches-surpass-desktop-searches-at-google-for-the-first-time/>

*2 <https://www.naver.com>

*3 本研究は、ACM CIKM 2017で発表したshort paper[9]に追加実験と考察を加えたものである。

ジ全体のうちモバイルアプリケーション端末の画面に表示される領域)とユーザの注視の相関を調査した従来研究について述べる. 本研究では, ここでの知見を参考に, 放棄セッションにおいてユーザが注視したカードを推定する.

2.1 クリックデータに基づく URL の適合性ラベルの付与

Joachims ら [7] はユーザのクリックに基づいて URL のプリファレンス対を予測した. また彼らはアイトラッキングを用いた実験で, クリックされた URL の1つ下のランクの URL は, およそ 50% と高い確率で見られていると報告している.

Agrawal ら [1] は Joachims らのアプローチを拡張し, クリックデータに基づいて, 与えられたクエリに対し, 自動的に各 URL の適合性ラベルを付与する手法を提案した. Agrawal らの手法は, 与えられた特定のクエリに対するクリックデータをもとに, 以下の処理を行う.

ステップ 1 ノードを URL とし, エッジとその重みが URL 対のプリファレンスを表す有向グラフ (プリファレンスグラフ) を作成する. 重みはクリックデータの集計により決定する.

ステップ 2 エッジの重みに基づいて, グラフのノード (URL) の順位付けを行う.

ステップ 3 順位付け結果をいくつかの領域に区切り, 各 URL に適合性ラベルを付与する.

本研究の目的は, URL への適合性ラベル付与ではなく, モバイル検索におけるカードの種類を適切に順位付けすることである. そこで, 提案手法では Agrawal らの手法を URL ではなくカードを有向グラフのノードとしたうえで拡張し, また上記ステップ 3 は行わない. 以下, 提案手法において拡張する上記ステップ 1 および 2 について詳述する.

ステップ 1 を説明する. 基本的なアイデアは, 「クリックされた URL はそれまでに見ていた URL よりもユーザが好むものであった」という仮定 [7] に基づき, 各クリックセッションから URL 対のプリファレンス (どちらの URL が好まれたか) を得て, それらを集計することである. ある URL 対に対しプリファレンスが得られた場合, グラフ上では, それに対応するエッジを生成し, 重みを 1 とする. 以後, 同じ URL 対に対し同じプリファレンスが得られた場合には, エッジの重みをインクリメントする. エッジはより好まれた URL を示すノードから, そうでない URL を示すノードの向きに引く.

たとえば, 図 2 のように, $URL_A, URL_B, URL_C, URL_D$ がこの順番で提示され, ユーザが URL_B のみをクリックしたセッションがあるとすると. このとき, Agrawal らの手法では, URL_B をクリックしたユーザが, クリックをする前に他の各 URL を見たと思われる確率 (推定閲覧確率) を利用する. たとえば, URL_B すぐ上およびすぐ下に位置する URL_A, URL_C を見た確率はそれぞれ 1, それより下の

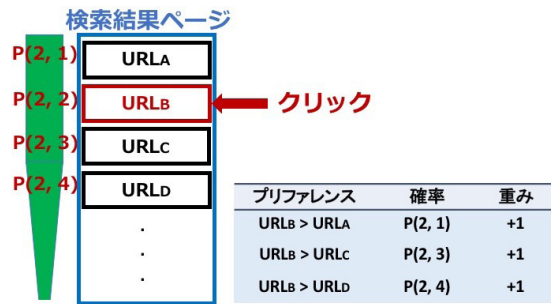


図 2 Agrawal らの確率的なプリファレンスルール
Fig. 2 Probabilistic preference rules of Agrawal et al.

URL_D 以下を見た確率は 1 未満で, 同図の左側に視覚的に示したように, 徐々に小さくなるよう定める. 推定閲覧確率は, Joachims らのアイトラッキングを用いたユーザ実験結果を参考に定められたものである. 図中では, たとえば検索結果中第 2 位の URL をクリックした場合の第 1 位の URL に対する推定閲覧確率を $P(2,1)$ のように表している. 「クリックされた URL はそれまでに見ていた URL よりユーザが好むものであった」という仮定より, この例からは, 以下のプリファレンスが考えられる. $URL_B > URL_A, URL_B > URL_C, URL_B > URL_D$. そして, たとえばプリファレンス $URL_B > URL_A$ は, 確率 $P(2,1)$ で有向グラフに反映させ, $URL_B > URL_D$ は確率 $P(2,4)$ (比較的小さい値) で反映させる. ここで, 反映させるとは, ノード間に新たなエッジを設けて重みを 1 とするか, 既存のエッジの重みをインクリメントすることを意味する. 上記の例では, 1つのセッションにおける1つのクリックから得たプリファレンスを有向グラフに反映させる手順について示したが, 実際には与えられたクエリに対応する複数セッションの各クリックをもとに同様の処理を行い, そのクエリに対する有向グラフを完成させる.

図 2 の例では, プリファレンス $URL_B > URL_A$ の確からしさを推定閲覧確率 $P(2,1)$ で, $URL_B > URL_D$ の確からしさを $P(2,4)$ で表しており, この確からしさを考慮したうえで有向グラフを構築していることになる. 前述のとおり, Agrawal らはこの確からしさを, すなわち推定閲覧確率をユーザ実験を参考に定めている.

Kadotami ら [8] は, ヤフーのモバイル検索におけるバーティカル (本研究におけるカードに相当) の順位付けのために Agrawal らのアルゴリズムを適用しており, Agrawal らと同様に推定閲覧確率に依存した処理を行っているが, 推定閲覧確率の利用の効果が見られなかったことを報告している (“it is not clear if a decaying probability curve is necessary for this particular task.”). そこで, 本研究では, 推定閲覧確率の代わりに, モバイルクエリログから得られるユーザのビューポートデータ, すなわちユーザが検索結果のどの部分を実際に画面に表示していたかを活用する (3.1 節). なお, Kadotami らの研究では放棄セッショ

ンを扱っていない。

ステップ2を説明する。本ステップではエッジの重みに基づいてグラフのノードの順位付けを行う。Agrawalらが試した手法のうち、最も単純かつ効果的であった Δ -order [1]を説明する。ステップ1により、たとえばプリファレンス $URL_B > URL_A$ は、有向グラフ上では、 URL_B を示すノードから URL_A を示すノードへの有向エッジおよびその重みにより表現されている。そこで、各ノードのスコアを、そのノードから出ていくエッジの重みの和から、そのノードに入ってくるエッジの重みの和を除くことにより算出する。そして、グラフ中の全ノードを上記スコアにより順位付ける。本研究では、ステップ2については Δ -orderをそのまま採用しているため、具体例について3.2節で改めて説明する。

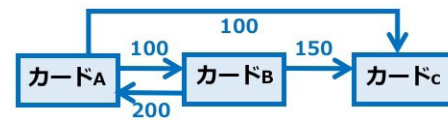


図3 プリファレンスグラフ

Fig. 3 A preference graph.

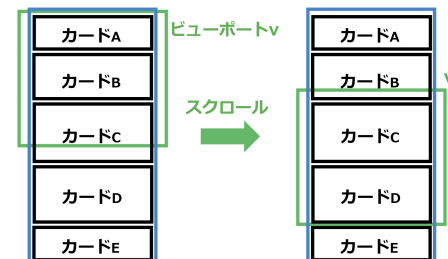


図4 セッション中のスクロール

Fig. 4 Scrolling within a session.

2.2 良い放棄と悪い放棄

良い放棄と悪い放棄を自動的に判別する問題に取り組んだ研究 [4], [13], [15]がある。この中で、Williamsら [15]の研究では、検索結果中最初の視覚的なアンサー（本研究のカードにほぼ相当）の推定閲覧時間、最初のアンサーのうちビューポートに表示されているピクセル数の割合などがユーザの検索に対する満足度と相関があると報告している。

また、Lagunら [10]や、Guoら [6]は、モバイルユーザの満足度推定に、ビューポート時間（検索されたアイテムが表示されていた時間）が有用であることを示した。

提案手法では、良い放棄と悪い放棄の判別は行わないが、上記の研究はカードがユーザに与えた情報量を予測するのに役に立つと考えられる。

2.3 ビューポート情報とユーザの注視の関係

Lagunら [10]は、セッション時間に対して、ユーザが検索結果ページ上のどの結果をどのくらい見たかについて、ビューポートデータから予測した場合と、アイトラッキングデータで決めた場合との相関を報告した。ビューポートデータにはビューポート時間、検索されたアイテムのうちどれくらいの面積がユーザに見えていたか、検索されたアイテムがビューポートの面積中どれくらいを占めていたかの3つがあり、それぞれを用いた場合より、3つすべてを用いた場合の方が相関が高いと述べている。

提案手法は、カードがユーザに与えた情報量を予測するために上記3つの要素を使用する。

3. 提案手法

本研究の提案手法は2.1節で説明した手法を拡張したものであり、与えられたクエリに対するモバイルクエリログを用いて、以下の処理を行う。

ステップ1 セッションデータ中のユーザの操作ログをもとに、与えられたクエリについて、どのカードがどの

カードより好まれるか（プリファレンス）を推定し、これを集計した有効グラフ（プリファレンスグラフ）を作成する。

ステップ2 プリファレンスグラフのエッジの重みに基づいて、グラフのノード（カード）の順位付けを行う。

3.1 ステップ1：プリファレンスグラフの作成

図3に、ある与えられたクエリに対し作成されるプリファレンスグラフの簡単な例を示す。プリファレンスグラフのノードはカードであり、それぞれのエッジの向きはカードどうしのプリファレンスを表し、エッジの重みはそのプリファレンスがセッションデータより何回推定されたかを表す。この例の200と書かれたエッジは、ある与えられたクエリに対して「カードBはカードAよりもユーザが好むものであった」とセッションデータ中のユーザの操作ログより200回推定したことを示す。

提案手法は、「ユーザに多くの情報を与えたカードは、その他のユーザが閲覧したカードよりも好まれた」という仮定のもとプリファレンスを推定する。多くの情報を与えたカードとユーザが閲覧したカードの定義を以下で述べる。

まず、多くの情報を与えたカードを定義する。クリックセッションの場合は、クリックされたカードを多くの情報を与えたカードとする。放棄セッションの場合は、ユーザが最も注視したカードを多くの情報を与えたカードとする。具体的には、ある放棄セッションにおいて、ビューポートにその一部もしくは全体が表示されたカードの集合 C があるとき、各カード $c \in C$ について以下の方法により $cardscore(c)$ を計算し、最も値の高いカードを多くの情報を与えたカードとする。ここでビューポートとは、検索結果ページに対するモバイルアプリケーション端末の表示領域のことである。各セッションデータは、図4に示すように、ユーザのスクロールにより定義される一般に

複数のビューポートにより構成される．そこで，あるセッションを構成するビューポートの集合を V とするとき， $cardscore(c)$ を以下のように算出する．

$$cardscore(c) = \sum_{v \in V} score(c, v), \quad (1)$$

$$score(c, v) = time(c, v) \times dominance(c, v)$$

$$\times completeness(c, v),$$

$$time(c, v) = \frac{c \text{ を } v \text{ に表示した時間}}{\text{セッションの時間}},$$

$$dominance(c, v)$$

$$= \frac{c \text{ の } v \text{ に表示されている部分の高さ}}{v \text{ の高さ}},$$

$$completeness(c, v)$$

$$= \frac{c \text{ の } v \text{ に表示されている部分の高さ}}{c \text{ の高さ}}.$$

次に，ユーザが閲覧したカードであるが，クリックセッションにおいては，各クリック以前にビューポートに一部でも表示されたカードと定義する．一方，放棄セッションにおいては，セッション中でビューポート内に一部でも表示されたカードすなわち前述の $c \in C$ と定義する．提案手法では，モバイルアプリケーション端末の小さいビューポートに表示されたカードのみをプリファレンスの対象とするため，Agrawalら，Kadotamiらのような推定閲覧分布を考慮せずとも信頼性の高いプリファレンスが収集できると考えられる．

図 5 の上半分は，クリックセッションからプリファレンスを推定する様子を表している．このように，クリックセッションにおいては，クリックされた各カード（一般にはセッション中に複数存在する）を多くの情報を与えたカードと見なし，各クリック以前に閲覧されたカードに対するプリファレンスを推定する．一方，同図の下半分は，放棄セッションからプリファレンスが生成される様子を表している．この場合は，このセッションに対し $cardscore(c)$ が最大のカードを選定し，これをユーザに多くの情報を与えたカードと見なす．そして，同セッション中にビューポートに表示された全カードを対象にプリファレンスを推定する．なお， $cardscore(c)$ が最大のカードが複数枚ある場合は，それらのカードそれぞれについて，同様にプリファレンスを推定する．

3.2 ステップ 2: プリファレンスグラフのノードの順位付け

各クエリに対して構築したプリファレンスグラフからカード間の半順序関係を得るために，2.1 節で説明した Δ -order を使用する．図 6 は，図 3 のプリファレンスグラフが与えられた場合に，ノードを順位付ける様子を表している．図の上半分には，カードごとにスコアを記載して

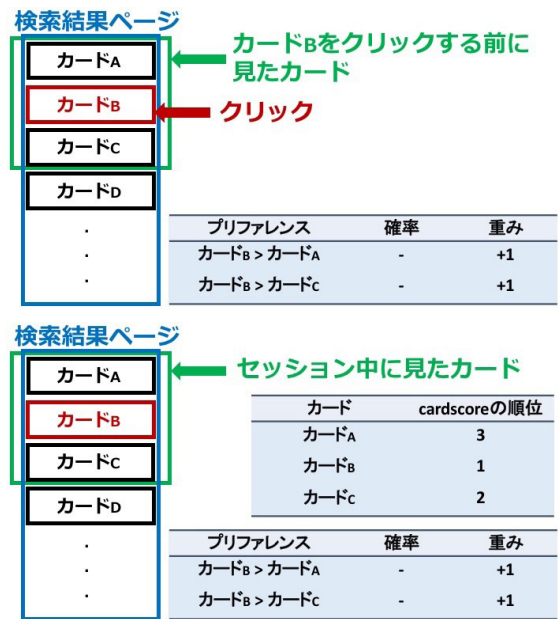


図 5 クリックセッション (上) と放棄セッション (下) でのプリファレンスの推定方法

Fig. 5 Preference estimation from a clicked session (top) and from an abandoned session (bottom).

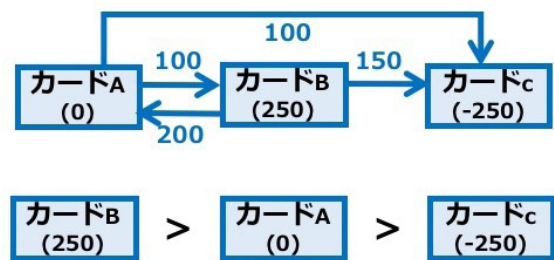


図 6 ノードの順位付け

Fig. 6 Ordering the nodes.

いる．たとえばカード_A のについては，ノードから出ていくエッジの重みの和が 100 + 100 で 200，ノードに入ってくるエッジの重みの和が 200 であるので，スコアは 0 である．同様にカード_B，カード_C について計算を行うと，同図の下半分のように，スコアに基づいたカードの順位付け結果が得られる．

4. プリファレンスによる評価

提案手法の有効性を確認するため，NAVER のモバイルクエリログを用いて正解つきデータセットを作成し，評価実験を行った．

4.1 使用するモバイルクエリログ

評価実験に使用したデータは，2016 年 12 月 12 日から 18 日までの韓国の検索エンジン NAVER のモバイルクエリログの一部である．NAVER のモバイル検索結果ページにはカードが並んでいる．モバイルクエ

リログの形式は, USER_ID, SERP_ID, QUERY, UNIXTIME, INTERACTION_TYPE, VISIBLE_ITEMS, CLICKED_CARD のようになっており, VISIBLE_ITEMS には, そのときにビューポートに表示されているカードの種類や, 各カードのビューポートに表示されている部分の高さが記録されている.

NAVER のカードの種類は数百にもものぼるが, 正解つき評価用データ作成のためには, カードの種類とクエリ数がある程度限定する必要がある. 今回は, クリック率 > 0.8 のカード, $0.8 \geq$ クリック率 ≥ 0.2 のカード, クリック率 < 0.2 のカードから, それぞれ 10 種類を選定し, 合計 30 種類のカードを対象とした. なお, クリック率はあるカードの, 最初のページに表示されたセッション数に対するクリックされた回数の割合であり, その期間は評価実験に使用したデータと同じである. 対象のカードの例を図 7 に示し, クリック率を図 8 示す.

次に, 評価データに含めるクエリは, 上記 30 種類のカードを対象とすることを前提として以下のように選定した.

- 信頼性の高いプリファレンスグラフの作成には, 与えられたクエリに対するある程度のセッション数が必要であるため, セッション数の合計で上位 70% に入るヘッドクエリとする.
- 提案手法は検索結果ページに現れるカードを順位付けるため, 1 つの検索結果ページに対象のカードが 2 つ以上現れるセッションを 1 つ以上持ち, 対象のカードが 1 つ以上現れるセッションを複数持つクエリとする.
- 実験の目的が放棄セッションを使用する有用性を示すことであるため, クリック率が 50% 以下のクエリとする.

上記のとおりクエリを絞り, 最終的に 992 のユニーククエリで合計 720,764 セッションが得られ, そのうち 315,759 (約 44%) が放棄セッションであった. このデータセットを使って, プリファレンスグラフを作成した. なお, 上記のような絞り込みを行わない場合でも NAVER モバイルクエリログにおける放棄セッション数の割合は一般に 25% 程度あり, 我々のデータセットは実データと極端に乖離したものではない.

4.2 正解データ

上述の 992 件のクエリに対する正解データを以下のように作成した*4. 本研究は日本で行われたが, 実験で用いた NAVER のカードデータは韓国語で書かれているため, クラウドソーシングサイト Lancers*5を通して 2 名のネイ

*4 実際に作成したプリファレンスデータは 1,000 のユニーククエリに対する合計 7,476 件のプリファレンスデータであったが, エクセルでデータを整理した際に, 8 のユニーククエリを変換 (たとえば「34-1」というクエリを「Jan-34」と変換した) してしまい, そのまま判定者に表示していたことが分かったため, 20 件のデータ, すなわち合計 60 件のプリファレンスデータを取り除いた.

*5 www.lancers.jp



0.8 < クリック率



0.2 ≤ クリック率 ≤ 0.8



クリック率 < 0.2

図 7 本実験で対象とする 30 種類カードの例

Fig. 7 Examples from the 30 card types used in our experiments.

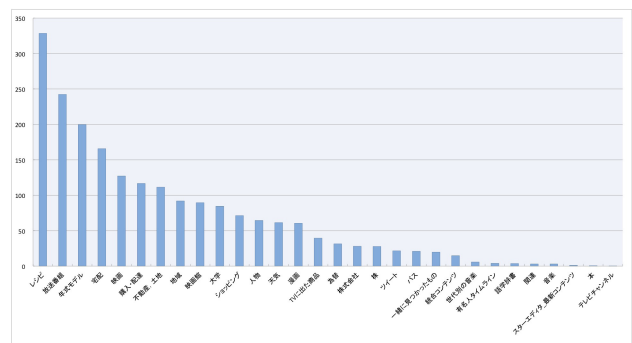


図 8 本実験で対象とする 30 種類カードのクリック率

Fig. 8 Clickthrough rates for the 30 card types used in our experiments.

ティブの韓国人と 1 名の韓国語が話せる日本人を雇い, 各カード対に対し判定者 3 名によるプリファレンスを独立に収集した. ここで, 個々のプリファレンスは, 与えられたクエリに対し 2 種類のカードのうちいずれのカードがより適しているかの判断結果を表す. 以上により合計 7,616 件のプリファレンスを収集し, 各カード対について 3 名のプリファレンスの多数決をとることにより最終的に 2,472 件の正解プリファレンスを得た.

図 9 に判定者に提供した正解データ作成用画面のスク



図 9 正解データを作成するための画面. 表示されている 2 種類のカードの内容は上のクエリと関係がない

Fig. 9 An interface for obtaining gold preferences. The contents of the displayed cards are not related to the query shown at the top.

表 1 正解データラベルの内訳

Table 1 Statistics of the gold preference labels.

N1-N2-NN	判定数	正解カード	判定数
L-L-L	993 (40.17%)	L	1,742
L-L-R	506 (20.47%)		(70.47%)
L-R-L	105 (4.25%)		
R-L-L	138 (5.58%)		
R-R-R	347 (14.04%)	R	730
R-R-L	177 (7.16%)		(29.53%)
R-L-R	92 (3.72%)		
L-R-R	114 (4.61%)		
合計	2,472 (100%)	total	2,472 (100%)

リールショットを示す. クエリ (“香港為替”) が一番上に表示され, 2 種類の異なるカードが左右に表示されている. 本研究の目的は, 個々のカード間ではなく, カードの種類間のプリファレンスを推定することであるため, 図に示したように, クエリとは内容的に直接関係がないカードを左右に並べて提示している. また, 種類が同じであっても, 実際の表示形式が若干異なるカードがあるため, 表示形式が複数ある場合にはスクロールによりそれらを閲覧できるようにしている.

判定者への指示は日本語で書いたドキュメントにより与えた. この中で, “提示された「検索ワード」に対して, 左と右のどちらのカードの形式が検索結果として適切かを判定してください.” と記載した. また, 判定者はブラウザ上で上記の画面 (図 9) を使用してリモートで判定作業を行った.

表 1 に 3 名の判定者による判定ラベルの統計を示す. 1 列目についてだが, **N1**, **N2** はどちらも韓国語がネイティブである判定者 2 名を示し, **NN** は韓国語がネイティブでない判定者 1 名を示す. また, たとえば L-L-R は, **N1** が左のカード, **N2** が左のカード, **NN** が右のカードを正解

表 2 判定者間の一致度. 3 名の場合は Fleiss' κ を, 2 名の場合は Cohen's κ を使用した

Table 2 Inter-assessor agreements. Fleiss' κ is used for three assessors; Cohen's κ is used for two assessors.

	κ	95% CI
3 名全員	0.562	[0.539, 0.584]
N1 と N2	0.570	[0.551, 0.588]
N1 と NN	0.199	[0.178, 0.219]
N2 と NN	0.245	[0.224, 0.265]

データ作成用の画面で選んだことを示す. 3 列目の “正解カード” は上述のとおり多数決で選んだ. 2 列目の判定数の太字は, **NN** がネイティブの 2 名の判定と大きく異なることを示しており, このことから, **NN** の判定内容が信頼できない可能性があるといえる. 表 2 に判定者間の一致度を Fleiss' κ と Cohen's κ [17] により示す. Fleiss' κ で 3 名の判定者間の一致度を測り, Cohen's κ で 2 名ずつ判定者間の一致度を測る. 95%信頼区間も算出する.

上記のように, ネイティブの判定者とネイティブでない判定者のラベルの一致度が低いことが分かったため, ネイティブ 2 名の判定結果のみを用いた第 2 の正解データを作成した. こちらでは多数決が適用できず, ネイティブ 2 名の判定結果が一致したデータのみを使用したため, 正解プリファレンス数は 2,023 件となった.

4.3 評価尺度

提案手法の評価指標として, 以下のように正解率と精度を定義する.

$$\text{プリファレンス正解率} = \frac{n}{N} \quad \text{プリファレンス精度} = \frac{n}{N'}$$

N : 用意した正解プリファレンスの総数.

N' ($\leq N$): 上記のうち, 各手法が作ったランキングにより推定できたプリファレンス総数. 提案手法でビューポートに表示されたカードどうしのみについてプリファレンスを推定するため, 正解プリファレンスが提案手法によるカードのランキングに含まれない場合がある. また, プリファレンスグラフに Δ -order を適用した結果, 2 つのカードの順位が同じになってしまう場合もある. 以上により, 一般に N' は N より小さくなる.

n ($\leq N'$): 推定プリファレンスのうち, 正解プリファレンスと一致するものの個数.

4.4 比較する手法

実験で比較する手法は以下のとおりである. なお, C (Click) はクリックセッション, A (Abandonment) は放棄セッションを示し, score は多く情報を与えたカードを提案手法の cardscore により選択することを, random は多く情報を与えたカードをビューポートに表示したカードの中からランダムに選択することを示す.

表 5 各手法のプリファレンス精度の差について Tukey HSD 検定を行った結果 (判定者 3 名).
有意水準 $\alpha = 0.05$ で有意な p 値を太字で示す

Table 5 Tukey HSD test results for the differences in preference precision (with three assessors) for all methods. p -values smaller than $\alpha = 0.05$ are shown in bold.

手法	精度の差	同時信頼区間	p 値
C+A(score) A(random)	0.2067	[0.1601, 0.2534]	0.0000
A(score) A(random)	0.1907	[0.1428, 0.2386]	0.0000
C(score)+A(score) A(random)	0.1813	[0.1357, 0.2269]	0.0000
C A(random)	0.1566	[0.1056, 0.2076]	0.0000
C+A(random) A(random)	0.1180	[0.0713, 0.1646]	0.0000
C+A(score) C+A(random)	0.0888	[0.0438, 0.1337]	0.0000
A(score) C+A(random)	0.0728	[0.0265, 0.1191]	0.0001
C(score)+A(score) C+A(random)	0.0633	[0.0195, 0.1072]	0.0006
C+A(score) C	0.0501	[0.0006, 0.0996]	0.0454
C C+A(random)	0.0387	[-0.0108, 0.0882]	0.2256
A(score) C	0.0341	[-0.0166, 0.0848]	0.3911
C+A(score) C(score)+A(score)	0.0254	[-0.0185, 0.0693]	0.5643
C(score)+A(score) C	0.0247	[-0.0238, 0.0732]	0.6956
C+A(score) A(score)	0.0160	[-0.0303, 0.0623]	0.9231
A(score) C(score)+A(score)	0.0094	[-0.0358, 0.0546]	0.9914

C+A(score) [提案手法] クリックセッションと放棄セッションを使用し, 放棄セッションで多く情報を与えたカードは提案手法の cardscore により選択する.

C+A(random) クリックセッションと放棄セッションを使用し, 放棄セッションで多く情報を与えたカードはランダムに選択する.

C クリックセッションのみを使用する.

A(score) 放棄セッションのみを使用し, 放棄セッションで多く情報を与えたカードは提案手法の cardscore により選択する.

A(random) 放棄セッションのみを使用し, 放棄セッションで多く情報を与えたカードはランダムに選択する.

C(score)+A(score) クリックセッションと放棄セッションを使用し, いずれの場合においても, 多く情報を与えたカードは提案手法の cardscore により選択する. すなわち, クリックセッションにおいて, クリックされたカードではなく, cardscore の最も高いカードを多く情報を与えたカードとし, セッション中で一部もしくは全体が表示されたカードを閲覧したカードとする.

4.5 実験結果と考察

表 3 と表 4 に, 各手法のプリファレンス精度と正解率を示す. 表 3 は判定者 3 名の判定結果から正解データを作った場合, 表 4 はネイティブ 2 名の判定結果から正解データを作った場合の結果である. また, 表 5 に, 判定者 3 名による正解データを使用した場合の, 各手法のプリファレンス精度の差について Tukey HSD 検定*6を行った結果を示す. 以降, 各手法により推定されたプリファレンスの性質

表 3 各手法のプリファレンス精度と正解率 (判定者 3 名)

Table 3 Preference precision and accuracy for each method (with three assessors).

手法	精度	正解率	n	N'	N
C+A(score)	0.6932	0.5138	1270	1832	2472
C+A(random)	0.6045	0.4482	1108	1833	2472
C	0.6431	0.3354	829	1289	2472
A(score)	0.6772	0.4490	1110	1639	2472
A(random)	0.4865	0.3135	775	1593	2472
C(score)+A(score)	0.6678	0.5489	1357	2032	2472

表 4 各手法のプリファレンス精度と正解率 (判定者が 2 名)

Table 4 Preference precision and accuracy for each method (with two assessors).

手法	精度	正解率	n	N'	N
C+A(score)	0.7175	0.5348	1082	1508	2023
C+A (random)	0.6287	0.4696	950	1511	2023
C	0.6820	0.3584	725	1063	2023
A(score)	0.6884	0.4652	941	1367	2023
A(random)	0.4789	0.3134	634	1324	2023
C(score)+A(score)	0.6842	0.5645	1142	1669	2023

の違いについて考察を行うため, 主としてプリファレンス精度に着目する.

まず, 表 3 より C の精度は 0.6431, A(score) の精度は 0.6772 であり, 表 5 よりこれらの間に統計的有意差はない. このことから, 提案した cardscore(c) に基づく放棄セッションからのプリファレンス推定により, クリックセッションからのプリファレンス推定結果と遜色のない結

*6 3 つ以上のシステム間比較において, 各システム対の差に関する適切な p 値を与える多重比較法の 1 つである [17].

果が得られていることが分かる。

次に、表 3 と表 4 において、C+A(score) は C より精度が高く、推定できたプリファレンス数 N' もより多くなっている。さらに、表 3 における両者の精度の差は表 5 が示すように統計的に有意である ($p = 0.0454$)。よって、プリファレンスグラフを作成してカードを順位付ける際、クリックセッションに加えて放棄セッションを用いることは有効といえる。一方、表 5 によれば、C+A(score) と A(score) の精度の差は統計的に有意ではない。以上を総合すると、C+A(score) の高い精度に大きく貢献しているのは放棄セッションでのプリファレンス推定のほうであると考えられる。

表 3 と表 4 において C+A(score) は C(score)+A(score) よりも精度が高い。このことは、クリックセッションでは、提案した *cardscore(c)* を利用するよりも、実際にクリックしたカードに基づきプリファレンスを作成したほうがよいことを示唆する。ただし、表 5 によれば両者の差は統計的に有意ではないので ($p = 0.5643$)、今回の結果からただちに上記を結論付けることはできない。

表 3 と表 4 において C+A(score) は C+A(random) より精度が高く、表 5 より $p \approx 0.0000$ で有意差がある。また、同様に A(score) は A(random) より精度が高く、表 5 より $p \approx 0.0000$ で有意差がある。このことから、放棄セッションにおけるユーザに多く情報を与えたカードの選ぶ際に、提案した *cardscore(c)* を用いることは有効といえる。

表 3 と表 4 を比較してみると、A(random) を除いて、全般的に表 4 におけるプリファレンス精度および正解率のほうが高いことが分かる。このことと、表 2 の判定者間一致度の結果を総合すると、今回実験した手法は韓国語ネイティブの判断に比較的近い結果になっていると思われる。

4.6 ユーザに多くの情報を与えたカードの選定に用いた各特徴量の効果に関する考察

前節では *cardscore* に基づきユーザに多くの情報を与えたカードを選定し、プリファレンスを推定する提案手法の有効性を示した。本節では、*cardscore* を構成する特徴量である time, dominance, completeness (3.1 節参照) の効果について考察する。表 6 に、*cardscore* で使用する特徴量のうち 1 つもしくは 2 つを用いない場合、および、3 つすべてを用いた場合 (通常の *cardscore*) の手法 A(score) のプリファレンス精度を示す。判定者 3 名による正解データを使用した場合である。なお、表の t, d, c はそれぞれ time, dominance, completeness を示しており、たとえば td は time と dominance を用いたことを示す。また、この結果に対応する Tukey HSD 検定の結果を表 7 に示す。第 1 列において、たとえば “c t” は completeness のみを用いた場合と time のみを用いた場合の差を意味する。

表 6 より、completeness のみを使った場合が一番精度が

表 6 *cardscore* で使用する特徴量変えた場合の A(score) のプリファレンス精度 (判定者 3 名)

Table 6 Preference precision (with three assessors) of A(score) for different combinations of the features used for *cardscore*.

手法	time	dominance	completeness	プリファレンス 精度
t	使用	-	-	0.6236
d	-	使用	-	0.6486
c	-	-	使用	0.6928
td	使用	使用	-	0.6571
dc	-	使用	使用	0.6687
tc	使用	-	使用	0.6777
tdc	使用	使用	使用	0.6772

表 7 A(score) の *cardscore* に使用する特徴量を変えた場合のプリファレンス精度の差について Tukey HSD 検定を行った結果 (判定者 3 名)。有意水準 $\alpha = 0.05$ で有意な p 値を太字で示す

Table 7 Tukey HSD test results for the differences in preference precision (with three assessors) for A(score) with different combinations of *cardscore* features. p -values smaller than $\alpha = 0.05$ are shown in bold.

手法	精度の差	同時信頼区間	p 値
c t	0.0692	[0.0202, 0.1183]	0.0006
tc t	0.0541	[0.0051, 0.1031]	0.0195
tdc t	0.0537	[0.0047, 0.1026]	0.0211
dc t	0.0451	[-0.0038, 0.0941]	0.0940
c d	0.0442	[-0.0044, 0.0929]	0.1034
c td	0.0358	[-0.0130, 0.0845]	0.3157
td t	0.0335	[-0.0155, 0.0825]	0.4054
tc d	0.0291	[-0.0196, 0.0777]	0.5733
tdc d	0.0286	[-0.0200, 0.0772]	0.5905
d t	0.0250	[-0.0239, 0.0740]	0.7409
c dc	0.0241	[-0.0245, 0.0728]	0.7674
tc td	0.0206	[-0.0281, 0.0693]	0.8754
tdc td	0.0202	[-0.0285, 0.0688]	0.8859
dc d	0.0201	[-0.0285, 0.0687]	0.8867
c tdc	0.0156	[-0.0331, 0.0643]	0.9654
c tc	0.0152	[-0.0336, 0.0639]	0.9700
dc td	0.0116	[-0.0370, 0.0603]	0.9924
tc dc	0.0090	[-0.0397, 0.0576]	0.9982
tdc dc	0.0085	[-0.0401, 0.0572]	0.9986
td d	0.0085	[-0.0402, 0.0571]	0.9987
tc tdc	0.0004	[-0.0482, 0.0491]	1.0000

高く、time のみを使った場合が一番精度が低く、表 7 よりこれらの間 (“c t”) には統計的有意差がある ($p = 0.0006$)。すなわち、カードが一部分でも映った時間の合計よりも、カード全体のうちどれくらいの部分が何度閲覧されたかの方が放棄セッションにおけるプリファレンス推定のために有効であるといえる。さらに、time と completeness を使った場合 (tc) とすべての特徴量を用いた場合 (tdc) は、ともに time のみを使った場合 (t) よりも精度が高く、有

意差がある (それぞれ $p = 0.0195$, $p = 0.0211$). このことも, completeness の有用性を示している. また, 表 7 が示すように, completeness のみを使った場合 (c) と提案手法 (tdc) の間には有意差がないことから, 各特徴量は必ずしも相補的な効果をもたらしていない可能性がある.

5. 結論と今後の課題

本稿では, 与えられたモバイルクエリに対しより適切なカードの順位付けを目的として, 放棄セッションからもプリファレンスが得られる可能性に注目し, 放棄セッションでユーザにより多く情報を与えたカードを定義し, クエリセッションと放棄セッションの両方からプリファレンスを取得, プリファレンスグラフを作成した. プリファレンス精度による評価を行ったところ, 以下の知見が得られた.

- クリックセッションにおいてクリックに基づきプリファレンスを推定するのと同様に, 放棄セッションにおいてはユーザの注視したカードを推定することにより, プリファレンスを高精度に推定することが可能である.
- プリファレンスグラフを用いてカードを順位付けする際に, クリックセッションに加えて放棄セッションを用いることは有効である.
- 放棄セッションにおけるユーザに多く情報を与えたカードを選ぶ際に, 提案した time, dominance, completeness に基づく *cardscore* を用いることは有効である.
- 放棄セッションで多く情報を与えたカードを選ぶ際は, カードが一部分でも映った時間の合計 (time) よりも, カード全体のうちどれくらいの部分が何度閲覧されたか (completeness) の方が有効な手がかりとなる.

今後の課題を述べる. 今回は, 放棄セッション利用の有効性を検証するために意図的に放棄率の高いクエリを選んだので, 実際のモバイルクエリデータからの代表的サンプルに対する実験には必ずしもなっていない. したがって, 今回得られた効果の実用上の影響について検証する必要がある. また, 今回のモバイルクエリログ収集時点での元の検索結果と, 本研究における新たなカードの順位付けを適用した検索結果を, 実運用において直接比較評価することを検討したい.

謝辞 本稿に対し貴重なご意見をくださった査読者の先生方と編集委員の鷹野孝典先生に感謝申し上げます.

参考文献

- [1] Agrawal, R., Halverson, A., Kenthapadi, K., Mishra, N. and Tsaparas, P.: Generating labels from clicks, *Proc. ACM WSDM 2009*, pp.172-181 (2009).
- [2] Arguello, J., Diaz, F. and Callan, J.: Learning to aggregate vertical results into web search results, *Proc. ACM CIKM 2011*, pp.201-210 (2011).

- [3] Chilton, L.B. and Teevan, J.: Addressing people's information needs directly in a web search result page, *Proc. WWW 2011*, pp.27-36 (2011).
- [4] Chuklin, A. and Serdyukov, P.: Potential good abandonment prediction, *WWW 2012 Companion*, pp.485-486 (2012).
- [5] Craswell, N. and Szummer, M.: Random walks on the click graph, *Proc. ACM SIGIR 2007*, pp.239-246 (2007).
- [6] Guo, Q. and Song, Y.: Large-scale analysis of viewing behavior: Towards measuring satisfaction with mobile proactive systems, *CIKM 2016* (2016).
- [7] Joachims, T., Granka, L., Pan, B., Hembrooke, H., Radlinski, F. and Gay, G.: Evaluating the accuracy of implicit feedback from clicks and query reformulations in web search, *ACM TOIS*, Vol.25, No.2 (2007).
- [8] Kadotami, Y., Yoshida, Y., Fujita, S. and Sakai, T.: Mobile vertical ranking based on preference graphs, *Proc. ACM ICTIR 2017* (2017).
- [9] Kawasaki, M., Kang, I. and Sakai, T.: Ranking rich mobile verticals based on clicks and abandonment, *Proc. 2017 ACM on Conference on Information and Knowledge Management*, pp.2127-2130 (2017).
- [10] Lagun, D., Hsieh, C.-H. and Navalpakkam, D.W.V.: Towards better measurement of attention and satisfaction in mobile search, *Proc. SIGIR 2014*, pp.113-122 (2014).
- [11] Li, J., Huffman, S. and Tokuda, A.: Good abandonment in mobile and pc internet search, *Proc. ACM SIGIR 2009*, pp.43-50 (2009).
- [12] Ponnuswami, A.K., Pattabiraman, K., Wu, Q., Gilad-Bachrach, R. and Kanungo, T.: On composition of federated web search result page: Using online users to provide pairwise preference for heterogeneous verticals, *Proc. ACM WSDM 2011*, pp.715-724 (2011).
- [13] Song, Y., Shi, X., White, R. and Awadallah, A.H.: Context-aware web search abandonment prediction, *Proc. ACM SIGIR 2014*, pp.93-102 (2014).
- [14] Wang, K., Walker, T. and Zheng, Z.: PSkip: Estimating relevance ranking quality from web search clickthrough data, *Proc. ACM KDD 2009*, pp.1355-1364 (2009).
- [15] Williams, K., Kiseleva, J., Crook, A.C., Zitouni, I., Awadallah, A.H. and Khabza, M.: Detecting good abandonment in mobile search, *Proc. WWW 2016*, pp.495-505 (2016).
- [16] Zhou, K., Demeester, T., Nguyen, D., Hiemstra, D. and Trieschnigg, D.: Aligning vertical collection relevance with user intent, *Proc. ACM CIKM 2014*, pp.1915-1918 (2014).
- [17] 酒井哲也: 情報アクセス評価方法論: 検索エンジンの進歩のために, コロナ社 (2015).



川崎 真未

2018年早稲田大学基幹理工学研究科修士課程修了.



Inho Kang

Inho Kang obtained his Ph.D. in computer science from KAIST, Korea, and joined Samsung Advanced Institute of Technology in 2004. Then he did post-doctoral work at Language Technology Institute of Carnegie Mellon University from 2006 to 2007. In 2008, he joined Naver, Korea, and is a senior manager of the Clova NLP group.



酒井 哲也

1993年早稲田大学工業経営学専門分野修士課程修了。博士(工学)早稲田大学。Microsoft Research Asia等約20年間の企業経験の後、2013年に早稲田大学に着任。現在、情報理工学科教授・主任・国立情報学研究所客員教授。国際論文誌 Information Retrieval Journal (Springer) 共同編集長。本学会山下賞(2006)、論文賞(2006, 2007)、FIT 2005論文賞・2008船井ベストペーパー賞等受賞。

(担当編集委員 鷹野 孝典)