

ABCPRec : ユーザの消費者としての役割と創作者としての役割の適応的対応付けによるユーザ生成コンテンツ推薦

佃 洸撰 深山 覚 後藤 真孝

産業技術総合研究所

{k.tsukuda, s.fukayama, m.goto}@aist.go.jp

概要 ユーザ生成コンテンツを扱う Web サービス上では、一人のユーザが消費者としての役割と創作者としての役割の両方を持つことがある。コンテンツ推薦を行う既存モデルの大半は、ユーザの消費者としての役割しか考慮しておらず、推薦精度を改善するためにこれら二つの役割をどのように活用できるかは十分に検討されてこなかった。本稿では、ユーザ生成コンテンツ推薦における state-of-the-art 手法である CPRec (consumer and producer based recommendation) に基づき、ABCPRec (adaptively bridging CPRec) を提案する。二つの役割が全ユーザで統一的对対応付けられている CPRec とは異なり、ABCPRec では各ユーザの消費者・創作者それぞれの性質の類似度に基づいて適応的に二つの役割を対応付ける。これにより、ユーザの消費者・創作者としての特徴をより柔軟にデータから学習し、コンテンツのより高い精度での推薦が可能になる。二種類の実世界のデータセットを用いて推薦精度の評価実験を行い、提案手法の性能が CPRec を含む比較手法よりも統計的に有意に上回ることを示した。

キーワード 情報推薦, ユーザ生成コンテンツ, Bayesian Personalized Ranking

1 はじめに

ユーザ生成コンテンツ (user-generated content, UGC) とは、プロのクリエイターではない一般の人々が創作し、主に Web 上で公開されるコンテンツのことである。動画や写真の他にも、掲示板に投稿される文章など、多様なタイプの UGC が存在しており、それらは YouTube や Flickr, Reddit といった様々な Web サービス上で公開されている。UGC の特徴のひとつとして、非 UGC に比べてコンテンツの創作されるペースが速い点あげられる [1]。それゆえ、ユーザが望むコンテンツを膨大な UGC の中から発見することを支援するために、コンテンツの推薦は極めて重要である。UGC のもうひとつの特徴は、非 UGC を扱う Web サービス上ではユーザは消費者としての役割しか持たないのに対して、UGC を扱う Web サービス上では一人のユーザが「消費者」と「創作者」の二つの役割を持つことがあるという点である。ユーザはコンテンツを消費するだけでなく、Web サービス上でコンテンツを公開することで他のユーザにコンテンツを見てほしいといったことを動機として、積極的にコンテンツの創作も行う [2]。

こうした特徴をふまえて、Kang ら [3] は行列分解を用いた CPRec (consumer and producer based recommendation) と呼ばれる UGC 推薦手法を提案した。CPRec では、各ユーザは消費者と創作者のそれぞれに対応する二つのベクトルを持つ。これらのベクトルは、そのユーザの総合的な性質を表す「コアベクトル」を線形変換することで得られる。ユーザのコンテンツに対する好みの

度合いはユーザとコンテンツの相性および、ユーザとコンテンツの創作者の相性に基づいて計算される (詳細は 2.2 節で述べる)。UGC に関するデータセットを用いて、Kang らは CPRec が state-of-the-art 手法を上回る推薦精度であることを示した [3]。

CPRec の有用性は示されたものの、CPRec ではユーザの消費者ベクトルと創作者ベクトルが常にそのユーザのコアベクトルと対応付けられるため、モデルの柔軟性に欠けるという問題点がある。たとえば CPRec において同じ消費者ベクトルを持つユーザは、同じ創作者ベクトルを持つことになる。しかし現実世界では、消費者として類似している二人の消費者が、創作者としても類似しているとは限らない。CPRec が持つ問題点の詳細は 2.3 節で詳しく述べる。

本稿では、CPRec の持つ問題点を解決するため、ABCPRec (adaptively bridging CPRec) という手法を提案する。ABCPRec でも、各ユーザは二つのベクトル、すなわち消費者としての役割に対応するベクトルと、創作者としての役割に対応するベクトルを持つ。ただし CPRec とは異なり、モデルの柔軟性を高めるため、二つのベクトルはコアベクトルからは生成されず、ユーザはコアベクトルを持つこともない。その代わりに、「ユーザの消費者としての性質と創作者としての性質が類似していれば、そのユーザの消費者ベクトルと創作者ベクトルは類似すべきである」という制約を加える。二つの性質の類似度が高くなるほど、二つのベクトルはより類似する。逆に、二つの性質が全く違っていれば、二つのベクトルは互いに独立に値を持つことができる。つま

り、ユーザの消費者ベクトルと創作者ベクトルがコアベクトルを介して常に対応付けられる CPRec とは異なり、ABCPRec では消費者としての役割と創作者としての役割の類似度に応じて、二つのベクトルが適応的に対応付けられる。

一般に公開されている二種類の UGC データセットを用いた評価実験を行い、ABCPRec の推薦精度が CPRec を含む比較手法の推薦精度よりも高いことを示した。より具体的には以下の二点を示した：(1) CPRec からコアベクトルを除いてモデルの柔軟性を高めることの有用性、(2) ユーザが消費・創作したコンテンツに基づいて計算される、ユーザの消費者としての役割と創作者としての役割の類似度に応じて適応的に制約を加えることの有用性。

2 モデル

本章ではまず、UGC 推薦における state-of-the-art 手法である CPRec [3] の説明をする。次いで CPRec の問題点を述べ、最後にそれを解決するための提案手法について述べる。

2.1 記号の定義

U をユーザ集合、 I をコンテンツ集合とする。ユーザ $u \in U$ によって消費されたコンテンツ集合を I_u^+ によって表す。また、一つ以上のコンテンツを消費したユーザ（消費者）の集合を $C \subseteq U$ とする（つまり、 $C = \{u \mid u \in U \wedge |I_u^+| > 0\}$ ）。UGC を扱う Web サービスでは、全てのコンテンツはユーザによって創作されるため、一つ以上のコンテンツを創作したユーザ（創作者）の集合を $P \subseteq U$ と定義する。以上の定義に基づく我々の目的は、各ユーザ u に対して、 u がまだ消費していないコンテンツ集合 $I \setminus I_u^+$ から、パーソナライズされたコンテンツランキングを生成することである。

2.2 CPRec

上述の目的を達成するため、Kang らは CPRec と呼ばれる、行列分解に基づく手法を提案した [3]。CPRec では、任意のユーザ u は K 次元のコアベクトル γ_u を持つ。さらにユーザ u は、 u の消費者としての役割を表す K 次元の潜在ベクトル γ_u^c および、 u の創作者としての役割を表す K 次元の潜在ベクトル γ_u^p も持つ。これらのベクトルは、 γ_u に基づいて次のように生成される。

$$\gamma_u^c = W^c \gamma_u, \quad \gamma_u^p = W^p \gamma_u. \quad (1)$$

ここで、 W^c および W^p は $K \times K$ 行列であり、それぞれ u のコアベクトルを γ_u^c と γ_u^p に変換するために使用される。ただし、 W^c と W^p はいずれも全ユーザ間で同一である。ユーザ u のコンテンツ i に対する好みの度合

いは次式により計算される。

$$\hat{x}_{ui} = \alpha + \beta_u + \beta_i + \langle \gamma_u^c, \gamma_i \rangle + \langle \gamma_u^c, \gamma_{p_i}^p \rangle. \quad (2)$$

α は全ユーザに共通のオフセットであり、 β_u はユーザ依存のバイアスを、 β_i はコンテンツ依存のバイアスを表す。 γ_i はコンテンツ i の K 次元の潜在ベクトル、 p_i は i の創作者である。 $\langle \gamma_u^c, \gamma_i \rangle$ はベクトルの内積を表す。

非 UGC の推薦を目的とした大半の手法がユーザの消費者としての役割のみ考慮している [4] のに対して、CPRec は消費者とコンテンツの相性 ($\langle \gamma_u^c, \gamma_i \rangle$) に加えて、消費者と創作者の相性 ($\langle \gamma_u^c, \gamma_{p_i}^p \rangle$) も考慮してユーザのコンテンツに対する好みを推定する。CPRec ではパラメータの値を Bayesian Personalized Ranking (BPR) [5] により求める。BPR はペアワイズなアプローチによるランキング最適化のためのフレームワークであり、ユーザがコンテンツに対する好みを 5 段階でレーティングするような明示的なフィードバックではなく、ユーザがコンテンツに対して「like」をしたり、コメントを投稿したりするような暗黙的なフィードバックを扱う点に特徴がある。

2.3 ABCPRec

CPRec では、各ユーザがコアベクトルを持つと仮定しているため、 γ_u^c は W^c を介して、 γ_u^p は W^p を介して、コアベクトル γ_u に常に対応付けられる。これによりモデルの柔軟性が下がり、以下のような問題を引き起こす。 W^c および W^p の各要素の値は学習を通じて求められるため、両者とも非常に高い確率で逆行列を持つ。これは、 $\gamma_u^p = W^p (W^c)^{-1} \gamma_u^c$ のように、ほぼ常に γ_u^c が γ_u^p に線形変換できることを意味する。加えて、 W^c および W^p は全ユーザで共通であるため、消費者と創作者の二つの役割が全ユーザで統一的に対応付けられる。つまり、同じ消費者ベクトルを持つユーザは、同じ創作者ベクトルを持つことになる。しかし現実世界では、消費者として類似している二人の消費者が、創作者としても類似しているとは限らないため、妥当性に欠ける。同様に、同じ創作者ベクトルを持つユーザは同じ消費者ベクトルを持つことになる。こうした制約のために、学習を通じてベクトルが適切な値を持っていないということも起こりうる。この問題点を解決するために、我々はユーザのコンテンツに対する好みを推定する新たな手法 ABCPRec (adaptively bridging CPRec) を提案する。

2.3.1 スコア予測

ABCPRec においても、ユーザ u は消費者としての役割に対応する K 次元ベクトル ν_u^c と、創作者としての役割に対応する K 次元ベクトル ν_u^p を持つ。ただし、CPRec とは異なり、 u はコアベクトルを持たず、 ν_u^c と ν_u^p の間には線形な関係はない。これにより、モデルの

柔軟性を高めている。ユーザのコンテンツに対する好みの度合い \hat{x}_{ui} は次式により求める。

$$\hat{x}_{ui} = \alpha + \beta_u + \beta_i + \langle \nu_u^c, \gamma_i \rangle + \langle \nu_u^c, \nu_{p_i}^p \rangle. \quad (3)$$

これはベクトルの表記方法を除いて、式2と同一である。

2.3.2 パラメータ学習

パラメータの学習にはBPR [5]を用いる。BPRでは、パラメータの最適化をする際に用いる学習データ \mathcal{D} は次式のように定義される。

$$\mathcal{D} = \{(u, i, j) \mid u \in \mathcal{U} \wedge i \in \mathcal{I}_u^+ \wedge j \in \mathcal{I} \setminus \mathcal{I}_u^+\}. \quad (4)$$

三つ組 (u, i, j) はユーザ u がコンテンツ j よりもコンテンツ i の方を好んでいることを意味する。 \mathcal{D} に基づく目的関数は次式で与えられる。

$$\sum_{(u, i, j) \in \mathcal{D}} \ln \sigma(\hat{x}_{uij}) - \lambda_{\Theta} \|\Theta\|^2. \quad (5)$$

σ はシグモイド関数であり、 $\Theta = \{\beta_i, \gamma_i, \nu_u^c, \nu_u^p\}$ はモデルの全パラメータを、 λ_{Θ} はハイパーパラメータを表す。 \hat{x}_{uij} は u の i と j それぞれに対する好みの度合いの差であり、 $\hat{x}_{uij} = \hat{x}_{ui} - \hat{x}_{uj}$ により定義される¹。

式5では ν_u^c と ν_u^p の間に制約はない。一般に、モデルの複雑さ（パラメータ数）が増加すれば、モデルが過学習してしまう可能性も高くなる。そのため、単純にCPRecからコアベクトルを除くことでモデルの複雑さを高めるだけではモデルとして適切ではないかもしれない。パラメータの値をより適切に学習するためには、より現実に即した制約を加える必要がある。そこで、次のような仮定を設ける： u の消費者としての性質が創作者としての性質と類似していれば、 ν_u^c と ν_u^p の値も近くなる。さらに、二つの性質の類似度が高くなるほど、 ν_u^c と ν_u^p はより近づく。つまり、二つの性質の類似度に応じて、 ν_u^c と ν_u^p の間に適応的に制約が加えられる。この考えに基づき、式5を拡張した以下のような新たな目的関数を提案する。

$$\sum_{(u, i, j) \in \mathcal{D}} \ln \sigma(\hat{x}_{uij}) - \lambda_{\Theta} \|\Theta\|^2 - \lambda_s \sum_{u \in \mathcal{U}} \text{sim}(u) \|\nu_u^c - \nu_u^p\|^2, \quad (6)$$

λ_s はハイパーパラメータ、 $\text{sim}(u)$ は類似度関数である。 $\text{sim}(u)$ の値が大きいほど ν_u^c と ν_u^p の距離が小さくなるべきであることを表しており、これがモデルの適応性につながっている。

2.3.3 類似度関数

$\text{sim}(u)$ を計算するために、以下の二つの仮説を立てる。

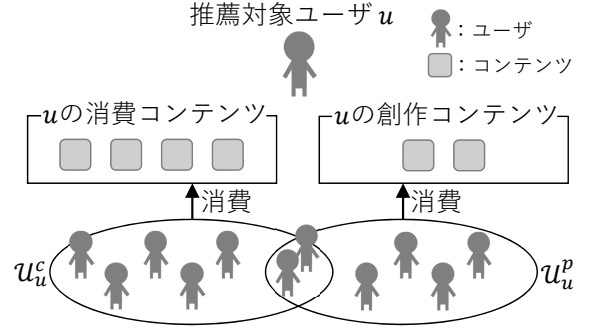


図1 推薦対象ユーザ u に対するユーザ集合 U_u^c および U_u^p .

仮説1: u によって消費されたコンテンツを好むユーザが、 u によって創作されたコンテンツも消費していれば、 ν_u^c と ν_u^p は類似するべきである。

仮説2: u によって創作されたコンテンツを好むユーザが、 u によって消費されたコンテンツも消費していれば、 ν_u^c と ν_u^p は類似するべきである。

仮説1では u が消費したコンテンツに焦点をあてているのに対して、仮説2では u が創作したコンテンツに焦点をあてている。 U_u^c を u によって消費されたコンテンツを一つ以上消費したユーザ集合、 U_u^p を u によって創作されたコンテンツを一つ以上消費したユーザ集合とする(図1)。すなわち、 U_u^c は u によって消費されたコンテンツを好むユーザの集合、 U_u^p は u によって創作されたコンテンツを好むユーザの集合であると仮定する。そのうえで、仮説1を反映した類似度関数を次式により与える。

$$\text{sim}(u) = \frac{|U_u^c \cap U_u^p|}{|U_u^c|}. \quad (7)$$

また、仮説2を反映した類似度関数を次式により与える。

$$\text{sim}(u) = \frac{|U_u^c \cap U_u^p|}{|U_u^p|}. \quad (8)$$

3 評価実験

本章では、実験を通して以下の三つの疑問に答える。

RQ1 CPRecにおいてモデルの柔軟性の低下を引き起こす原因となっている、ユーザのコアベクトルを除くことは、モデルを改善するために有用であるか。

RQ2 ユーザの二つの役割の類似度に応じて ν_u^c と ν_u^p の間に適応的に制約を加えることは有用であるか。

RQ3 RQ2が真である場合、2.3.3項で述べた二つの仮説のうち、どちらの仮説に基づく類似度関数の方がより有用であるか。

¹ \hat{x}_{uij} を計算する際に α および β_u はキャンセルされるため、両者は Θ には含まれない。

表1 データセットの統計情報.

	Flickr	Reddit
ユーザ数 ($ \mathcal{U} $)	99,060	52,654
アイテム数 ($ \mathcal{I} $)	557,286	336,743
総アイテム消費数 ($\sum_u \mathcal{I}_u^+ $)	11,256,457	1,786,032
消費者の割合 ($ \mathcal{C} / \mathcal{U} $)	99.83%	99.60%
創作者の割合 ($ \mathcal{P} / \mathcal{U} $)	40.82%	87.24%
消費者兼創作者の割合 ($ \mathcal{C} \cap \mathcal{P} / \mathcal{U} $)	40.65%	86.85%

3.1 データセット

本実験では以下にあげる二件の UGC の公開データセットを用いた.

- **Flickr** は写真共有サービスであり, ユーザは他のユーザが創作 (撮影および投稿) した写真の中で気に入った写真に「星マーク」を付けることができる. 我々は Cha ら [6] によって公開されている, 2006 年と 2007 年に収集されたデータセットを使用した. 本実験では「星マーク」を付ける行動をコンテンツの消費とみなす. Flickr では各写真は一人のユーザにより創作されている.
- **Reddit** は掲示板スタイルのインタフェースを介して, ユーザ同士で様々なトピックについて議論ができるオンラインコミュニティである. ユーザはニュース記事へのリンクや, テキストなどを元にスレッドを作成できる. また, 作成されたスレッドに対してコメントを投稿することもできる. 我々は Reddit により公開されているデータセット²を用いた. このデータセットには 2017 年 3 月に作成された全てのスレッドと, それらに投稿された全てのコメントが含まれている. 本実験ではスレッドの作成をコンテンツの創作, スレッドに対するコメントの投稿をコンテンツの消費とみなす. 各スレッドは一人のユーザにより創作されている.

学習により求められるモデルパラメータの信憑性を高めるために, コンテンツの消費数と創作数の少なくとも一方が 10 件未満のユーザと, 消費ユーザ数が 10 人未満のコンテンツを除いた. 表 1 にデータセットの統計情報を示す. 表にあるように, 消費者の割合は二つのデータセットでほぼ同じ数値である. 一方で, 創作者の割合および消費者兼創作者の割合は Flickr データセットの方が Reddit データセットよりも低いという特徴がある.

3.2 比較手法

Kang ら [3] と同様に, 本実験では以下のベースライン手法を用いた³.

²<https://www.reddit.com/comments/6607j2>

³UGC を扱う Web サービス上では, ユーザ同士のソーシャル・ネットワーク情報が常に存在するとは限らないため, SBPR [7] 等の

- **PopRec**: この手法はパーソナライズされていない手法であり, コンテンツを人気度の高い順にランキングする. コンテンツの人気度はそのコンテンツを消費したユーザ数とする.
- **BPR** [5]: この手法では消費者とコンテンツの相性だけに基づいて, ユーザのコンテンツに対する好みを次式により求める.

$$\hat{x}_{ui} = \alpha + \beta_u + \beta_i + \langle \nu_u^c, \gamma_i \rangle. \quad (9)$$

- **Vista** [8]: この手法では, 各ユーザは ν_u^c と ϕ_u の二つの潜在ベクトルを持つ. \hat{x}_{ui} は次式により計算される.

$$\hat{x}_{ui} = \langle \nu_u^c, \gamma_i \rangle + \langle \phi_u, \phi_{p_i} \rangle. \quad (10)$$

つまり, この手法では消費者 u とコンテンツ i の相性を求める際と, 消費者 u と創作者 p_i の相性を計算する際で, 使用する u のベクトルが異なる.

- **Factorization Machines (FMs)** [9]: この手法では, 消費者・創作者・コンテンツの交互作用を考慮する. 二次の交互作用に基づいた, ユーザのコンテンツに対する好みは次式で与えられる.

$$\hat{x}_{ui} = \alpha + \beta_u + \beta_i + \beta_{p_i} + \langle \nu_u^c, \gamma_i \rangle + \langle \nu_u^c, \nu_{p_i}^p \rangle + \langle \gamma_i, \nu_{p_i}^p \rangle. \quad (11)$$

以上のベースライン手法では, ベクトル間に制約は存在しない. 本実験では以下の比較手法も用いる.

- **CPRC** [3]: この手法は UGC 推薦における state-of-the-art 手法である. 手法の詳細は 2.2 節で述べた.
- **NBCPRC**: この手法では ν_u^c と ν_u^p の間に制約を設けない. つまり, ユーザのコンテンツに対する好みは式 3 により求められ, 目的関数は式 5 により与えられる. NBCPRC は no-bridging CPRC を意味する.

我々の提案手法に関しては, 類似度を式 7 により求める手法を **ABCPRC^{H1}**, 式 8 により求める手法を **ABCPRC^{H2}** と表す.

3.3 評価指標

各ユーザについて, \mathcal{I}_u^+ をトレーニングデータ, バリデーションデータ, テストデータの三つに分ける. そのために, ユーザ u が消費したコンテンツの中から, バリデーションデータ (\mathcal{V}_u) とテストデータ (\mathcal{T}_u) を 1 件ずつランダムに選択した. 選択されなかった消費コンテンツユーザ同士のソーシャルな関係に基づく手法は比較手法として用いない. 本実験では, UGC の推薦精度を改善するために, ユーザの二つの役割をどのように活用するかという点により焦点をあてる.

表2 ベクトルの次元数 K の値を 20, 50, 80 としたときの AUC の比較結果. PopRec, BPR, Vista, FMs, CPreC, NBCPreC, ABCPreC^{H1} との有差 ($\alpha = 0.01$) をそれぞれ †, ‡, *, †, ‡, ♣, ◇, ♠ により表す.

Dataset	K	PopRec	BPR	Vista	FMs	CPreC	NBCPreC	ABCPreC ^{H1}	ABCPreC ^{H2}
Flickr	20	0.6737	0.8698	0.8436	0.8764	0.8563	0.8839 †**♣	0.8861 †**♣	0.8900 †**♣◇
	50	0.6737	0.8772	0.8435	0.8822	0.8664	0.8937 †**♣	0.8949 †**♣	0.8992 †**♣◇♠
	80	0.6737	0.8777	0.8394	0.8810	0.8712	0.8955 †**♣	0.8988 †**♣	0.9028 †**♣◇
Reddit	20	0.6392	0.8713	0.8829	0.8960	0.9138	0.9209 †**♣	0.9296 †**♣◇	0.9340 †**♣◇♠
	50	0.6392	0.8721	0.8918	0.8999	0.9201	0.9302 †**♣	0.9346 †**♣◇	0.9391 †**♣◇♠
	80	0.6392	0.8709	0.8946	0.9001	0.9211	0.9322 †**♣	0.9376 †**♣◇	0.9408 †**♣◇♠

ツは全てトレーニングデータ (\mathcal{R}_u) として使用した. この操作を各ユーザについて 5 回繰り返し, 5 件のトレーニング・バリデーション・テストのデータの組を作成した. 3.4 節では, これら 5 件の平均精度を報告する. 手法間で精度を比較する際の公平性を保つため, 5 件のトレーニング・バリデーション・テストのデータの組は全ての手法で同じものを使用した. コンテンツの推薦精度は次式で表される AUC (Area Under the ROC Curve) により評価する.

$$AUC = \frac{1}{|\mathcal{U}|} \sum_{u \in \mathcal{U}} \frac{1}{|\mathcal{D}_u|} \sum_{(i,j) \in \mathcal{D}_u} \delta(\hat{x}_{ui} > \hat{x}_{uj}). \quad (12)$$

$\mathcal{D}_u = \{(i, j) \mid i \in \mathcal{T}_u \wedge j \in \mathcal{I} \setminus \mathcal{I}_u^+\}$ であり, $\delta(z)$ は z が真のとき 1, 偽のとき 0 の値をとる. いずれの手法でも, Adam optimizer [10] を用いた Tensorflow [11] によりモデルのパラメータを学習した. 正則化項のハイパーパラメータの値は $\{0.0001, 0.001, 0.01, 0.1, 1\}$ の中からバリデーションデータを用いて AUC の観点から最適なものを選択した. 学習率は 0.01 とした.

3.4 結果

表 2 に, ベクトルの次元数 K の値を 20, 50, 80 としたときの各手法の精度をテューキーの検定の結果と共に示す. ν_u^c と ν_u^p の間に制約を持たない NBCPreC であっても, いずれのデータセットでも CPreC を有意に上回る精度であった. この結果から, CPreC からコアベクトルを除くことは UGC の推薦精度を改善するうえで有用であることが言え, これが RQ1 に対する回答となる. また, ABCPreC^{H1} と ABCPreC^{H2} は共に, CPreC と NBCPreC よりも推薦精度が優れている. したがって RQ2 に対する回答は次のとおりである: CPreC のように二つの役割を常に対応づけたり, NBCPreC のように二つの役割を全く対応づけなかったりするよりも, 適応的に ν_u^c と ν_u^p を対応づける方が有用である.

いずれのデータセットでも, 全ての K で ABCPreC^{H2} が ABCPreC^{H1} を上回っていることから, RQ3 に対する回答は次のようになる: ユーザの消費者としての性質と創作者としての性質の類似度を計算する際は, ユーザが創作したコンテンツに焦点をあてた仮説 2 の方が

表 3 $sim(u)$ の分布.

	仮説 1		仮説 2	
	Flickr	Reddit	Flickr	Reddit
$sim(u) = 0$	62.70%	13.77%	62.87%	14.17%
$0 < sim(u) \leq 0.1$	36.46%	47.34%	2.93%	0.82%
$0.1 < sim(u) \leq 0.2$	0.54%	10.68%	4.92%	1.79%
$0.2 < sim(u) \leq 0.3$	0.08%	6.54%	5.25%	2.67%
$0.3 < sim(u) \leq 0.4$	0.02%	4.39%	4.78%	2.97%
$0.4 < sim(u) \leq 0.5$	0.01%	2.93%	4.53%	2.76%
$0.5 < sim(u) \leq 0.6$	0.0%	2.93%	4.71%	5.56%
$0.6 < sim(u) \leq 0.7$	0.0%	2.2%	4.07%	6.21%
$0.7 < sim(u) \leq 0.8$	0.0%	1.63%	3.21%	6.79%
$0.8 < sim(u) \leq 0.9$	0.0%	1.33%	2.02%	9.68%
$0.9 < sim(u) \leq 1.0$	0.17%	6.25%	0.72%	46.59%

有用である. なぜ仮説 1 よりも仮説 2 の方が優れているのかを分析するため, 各仮説で計算された類似度の値の分布を表 3 に示す. いずれのデータセットでも, ν_u^c と ν_u^p の間に制約を持たないユーザ, つまり $sim(u) = 0$ のユーザの割合は両仮説でほぼ同じであった. 残りのユーザに関しては, 仮説 1 では残りのユーザの大半が $0 < sim(u) \leq 0.1$ に分布している. それに対して仮説 2 では, Flickr では $0 < sim(u) \leq 1.0$ の範囲にある程度万遍なく分布しており, Reddit では大部分のユーザが $0.9 < sim(u) \leq 1.0$ に分布している. そのため, 仮説 2 では ν_u^c と ν_u^p の間により強い制約が適切に加えられ, これが推薦精度の改善につながったと考えられる. この分析に加えて, U_u^c と U_u^p それぞれに含まれるユーザのコンテンツに対する好みの類似性に関する分析も以下のように行った. U_u^c 中のあらゆる二ユーザのペアについて, 両方のユーザが共通に消費したコンテンツ数を数えた. この値が大きいほど, 両ユーザの消費コンテンツに対する好みの類似度が高いことを意味する. この処理を全てのユーザの U_u^c に対して行った. さらに, U_u^p についても同様の処理をした. 図 2 に Reddit におけるユーザペアの割合の分布を示す. U_u^c の場合, 全ユーザペアのうち, およそ半数のユーザペアは共通のコンテンツを全く消費していなかった. U_u^p の場合, 共通のコンテンツを 2 個以上消費しているユーザペア (両ユーザのコン

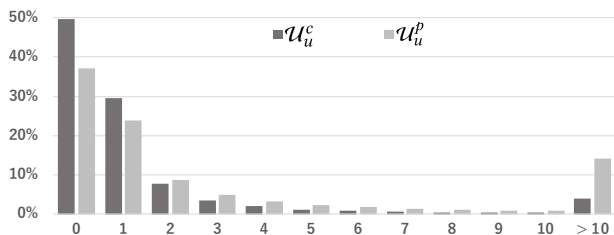


図2 Reddit データセットにおけるユーザペアの割合の分布 (x 軸: 二人のユーザによって共通に消費されたコンテンツの数)。

コンテンツの好みの類似度が比較的高いことを意味する)の割合は U_u^c よりも常に高く、14.06%ものユーザペアが共通のコンテンツを 11 個以上消費していた。同様の傾向は Flickr データセットでも見られた。以上の結果から、 U_u^c に含まれるユーザは消費コンテンツの好みに関して多様性が高いのに対して、 U_u^p に含まれるユーザは消費コンテンツの好みが似ていることがわかる。そのため、 U_u^p に焦点をあてた仮説 2 の方がより信頼性の高い制約を加えることができ、仮説 1 を上回る結果となった。

最後に、PopRec を除いた全手法の AUC の推移を図 3 に示す。いずれのデータセットでも、ABCPRC^{H2} は全ての K において最も高い精度であった。その精度は Flickr データセットで K が 80 のとき、Reddit データセットで K が 60 のときにほぼ上限に達した。FMs は ν_u^c と ν_u^p の間に制約を設けず、かつ消費者とコンテンツの相性および、消費者と创作者の相性を考慮しているという点で NBCPRec と類似している。しかし、NBCPRec の精度は常に FMs を上回った。このことから、UGC 推薦においては、コンテンツと创作者の相性 (式 11 中の $\langle \gamma_i, \nu_{p_i}^p \rangle$) は考慮しない方が効果的であると言える。3.1 節で述べたように、Flickr データセットは创作者の割合および消費者兼创作者の割合が低い。そのようなデータセットでは、CPRC の精度はベースライン手法 (BPR および FMs) よりも低かった。それに対して、我々の提案モデルはいずれのデータセットでも最も優れた精度であり、このことは提案モデルの頑健性を示している。

4 まとめ

本稿では、ユーザの消費者としての役割と创作者としての役割を考慮し、その類似度に応じて適応的に二つの役割を対応付ける UGC 推薦手法を提案した。実験の結果は提案手法の有用性を示しており、特に推薦対象のユーザが創作したコンテンツを消費したユーザ集合に着目して類似度を計算する手法が有用であった。今後の方針として、ユーザの二つの役割の類似度を計算する際に、そのユーザが消費および創作したコンテンツの外見 (色など) やテキストデータ (タグなど) を利用すると

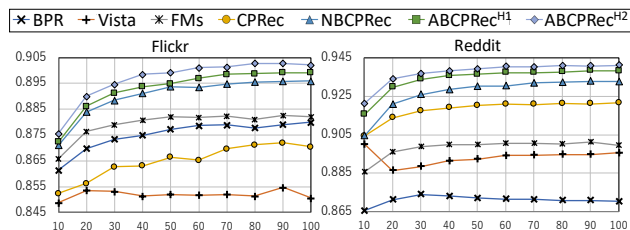


図3 ベクトルの次元数 K の値に応じた AUC の推移 (x 軸: K , y 軸: AUC)。

いった拡張が考えられる。他にも、二人のユーザの创作者としての類似度に基づいて、ユーザ間の创作者ベクトルに制約を加えるなどのように、新たな制約を加えた手法の拡張も予定している。

謝辞

本研究の一部は JSPS 科研費 (17K12688) および JST ACCEL (JPMJAC1602) の支援を受けた。

参考文献

- [1] M. Cha, H. Kwak, P. Rodriguez, Y.-Y. Ahn and S. Moon: “I tube, you tube, everybody tubes: analyzing the world’s largest user generated content video system”, IMC, pp. 1–14 (2007).
- [2] A. Wilson, H. Murphy and J. C. Fierro: “Hospitality and travel: The nature and implications of user-generated content”, Cornell Hospitality Quarterly, **53**, 3, pp. 220–228 (2012).
- [3] W. C. Kang and J. McAuley: “Learning consumer and producer embeddings for user-generated content recommendation”, RecSys, pp. 407–411 (2018).
- [4] J. Bobadilla, F. Ortega, A. Hernando and A. Gutiérrez: “Recommender systems survey”, Knowledge-Based Systems, **46**, pp. 109–132 (2013).
- [5] S. Rendle, C. Freudenthaler, Z. Gantner and L. Schmidt-Thieme: “BPR: Bayesian personalized ranking from implicit feedback”, UAI, pp. 452–461 (2009).
- [6] M. Cha, A. Mislove and K. P. Gummadi: “A measurement-driven analysis of information propagation in the flickr social network”, WWW, pp. 721–730 (2009).
- [7] T. Zhao, J. McAuley and I. King: “Leveraging social connections to improve personalized ranking for collaborative filtering”, CIKM, pp. 261–270 (2014).
- [8] R. He, C. Fang, Z. Wang and J. McAuley: “Vista: A visually, socially, and temporally-aware model for artistic recommendation”, RecSys, pp. 309–316 (2016).
- [9] S. Rendle: “Factorization machines”, ICDM, pp. 995–1000 (2010).
- [10] D. P. Kingma and J. Ba: “Adam: A method for stochastic optimization”, ICLR (2015).
- [11] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard, M. Kudlur, J. Levenberg, R. Monga, S. Moore, D. G. Murray, B. Steiner, P. Tucker, V. Vasudevan, P. Warden, M. Wicke, Y. Yu and X. Zheng: “Tensorflow: A system for large-scale machine learning”, OSDI, pp. 265–283 (2016).