

# 協調フィルタリング推薦によって 誤推薦されたコンテンツに関する一考察

山内 一騎† 當間 愛晃†

†琉球大学工学部情報工学科

## 1 はじめに

現在、情報推薦システムのほとんどは協調フィルタリング (CF) を用いている。この CF による推薦手法では、推薦の精度を最適化することに重点を置いているため、コンテンツのバリエーションという点においてはコンテンツベースより良いものの未だに悪いという問題があった [1]。つまり、CF 推薦は再現率という点では悪い推薦である。

そこで、本研究では、CF 推薦において誤推薦されたコンテンツの特徴を発見することで、再現率の向上を図り、ユーザの満足度の向上を目指す。本研究では、誤推薦されたコンテンツの特徴を発見するために、正/誤推薦されたコンテンツの違いを調査した。誤推薦された特徴を発見することにより、その特徴に重みを与え、再現率の向上を図ることができるからである。

調査の際、情報推薦のデータセットとして MovieLens を用いた。今回はその中のユーザ数 943、映画数 1682、評価数 100,000 を用いた。

MovieLens を用いて、正/誤推薦されたコンテンツの違いを調査したところ、データセットに含まれる情報源からの直接的な要因は観察できなかった。これは MovieLens が実際の映画の情報と比較すると映画のタグ情報に欠落があることが原因だと考えられる。MovieLens では、ある映画に対するタグが付加されているどうかの 1 か 0 である。一方、実際の映画はタグ情報には関連度があると考えられる。例えば、ある映画はタグ A と B に属され、それぞれが付加されているが、この映画はタグ B よりはタグ A の方が関連度が高いなどである。そこで、本稿では映画のタグの関連度を求めた。また、MovieLens ではユーザのタグも欠落しているため、ユーザの視聴履歴からユーザのタグを推測し、映画と同じように関連度を求めた。そして、これらの新しく求めた情報を用いて、正/誤推薦されたコンテンツの違いを調査した。

A statistical study on negative recommendation sets using collaborative filtering algorithm

†Kazuki YAMAUCHI †Naruaki TOMA

†Department of Information Engineering, Faculty of Engineering, University of the Ryukyus

## 2 提案手法

本稿では、コンテンツの情報を詳細に検証するため、タグの関連度を求めた。そして、この関連度をランク (式 (2)) で表し、その情報の違いを正/誤推薦されたコンテンツで比較し調査した。下記にその調査手順を示す。

Step1 : ユーザへのタグとそのランクを求める

Step2 : Step1 を用いてコンテンツへのタグのランクを求める

Step3 : Step1,2 を用いて、CF によって正/誤推薦されたコンテンツへのランク検証

下記にこれら 3 つの step の詳細を示す。

Step1 : ユーザへのタグとそのランクを求める

ユーザ集合を  $U = \{U_1, U_2, \dots, U_n\}$ 、コンテンツ集合を  $C = \{C_1, C_2, \dots, C_m\}$ 、タグ集合を  $T = \{T_1, T_2, \dots, T_l\}$  とする。まず、ユーザが評価したコンテンツのタグをカウントすることで、ユーザのタグを求める (式 (1))。

$$T_{(k,U_i)} = \sum_{j=1}^m T_{(k,U_i,C_j)} \quad (1)$$

$T_{(k,U_i)}$  はユーザ  $i$  に対するタグ  $k$  のカウントである。 $T_{(k,U_i,C_j)}$  はユーザ  $i$  が評価した中のコンテンツ  $j$  にタグ  $k$  が付加されていれば 1 の値、そうでなければ 0 の値を取り、これらの合計をカウントすることで  $T_{(k,U_i)}$  を求める。なお、 $T_{(k,U_i)}$  は 1 ユーザに対して、 $l$  (タグの総数) 個求める。次に  $T_{(k,U_i)}$  のそれぞれのタグを降順に並び替えランクを付ける (式 (2))。

$$R(T_{(k,U_i)}) = \text{rank}(\text{sort}(T_{(k,U_i)})) \quad (2)$$

$R(T_{(k,U_i)})$  は式 (1) でカウントした値を基に降順ソートした結果を用いてランク (順位そのもの) を返す関数である。これによって、ユーザの嗜好のタグにランクをつけることができ、ユーザの嗜好の中で最も好みのタグを求めることができた。

Step2: Step1 を用いて映画へのタグのランクを求める

まず、ユーザは評価したコンテンツを自身の嗜好の中で最も好きなタグ目線でコンテンツを見ていると仮定する。すると、最も多くのタグ目線で見られているコンテンツはそのタグによる関連度が大きいと考えられる (式 (3))。

$$T_{(k,C_j)} = \sum_{i=1}^n R(T_{(k,U_i,1)}) \quad (3)$$

$T_{(k,C_j)}$  はコンテンツ  $j$  に対するタグ  $k$  のカウントである。 $R(T_{(k,U_i,1)})$  は step1 によって求めたユーザ  $i$  のタグのランクの 1 位を示している。これを、step1 と同様にランク付けをする (式 (4))。

$$R(T_{(k,C_j)}) = rank(sort(T_{(k,C_j)})) \quad (4)$$

$R(T_{(k,C_j)})$  は式 (3) を式 (2) と同様に求めた関数である。

Step3: Step1,2 を用いて、CF によって正/誤推薦された映画へのランク検証

step1,2 で求めたランクを用いて、正/誤推薦されたコンテンツのランクはユーザのランクからすると、どうなっているかを検証するために、 $x$  軸を式 (2)、 $y$  軸を式 (4)、 $z$  軸を式 (5) とし、3次元で作図した。

$$Z = \sum_{i=1}^n \sum_{j=1}^m \sum_{k=1}^l T_{(x,y)} \quad (5)$$

$$x=R(T_{(k,U_i)}) \cdots \text{式 (2)}, \quad y = R(T_{(k,C_j)}) \cdots \text{式 (4)}$$

$T_{(x,y)}$  は  $x$  と  $y$  の条件でタグが付加されていれば 1 の値、そうでなければ 0 の値を取る。

### 3 調査結果

step3 によって作図した正/誤推薦されたコンテンツの調査結果の考察、及びユーザ毎に誤推薦されたコンテンツの調査結果 (図 1) と考察を下記に示す。

#### 3.1 正推薦されたコンテンツのタグ

この調査結果では、正推薦された全てのコンテンツを step3 の通りに合計して求めた。図は省略する。この省略した図は図 1 の  $x$  軸、 $y$  軸が共に 1 位に近づくにすれ、 $z$  軸が高くなっていき、 $x$  軸、 $y$  軸が共に 1 位のところで  $z$  軸が圧倒的に高くなっていた。このことから、CF 推薦によって正推薦されたコンテンツは、タグの関連度が 1 番大きいタグとユーザが最も好むタグが同じであることが多いということである。これは、人間が推薦を行う際も、最も一般的で推測がしやすい推薦手法と同じである。

#### 3.2 誤推薦されたコンテンツのタグ

ここでは、誤推薦された全てのコンテンツを step3 の通りに合計して求めた。ここでも図を省略する。省略した図では、正推薦されたコンテンツのような規則はなく、 $z$  軸はまばらになっていた。このことから、逆に、正推薦されたコンテンツのように、人間が推測しやすい推薦手法で推測されなかったコンテンツが CF でも上手く推薦されなかったことが考えられる。今回用いた一般的な CF 推薦であるユーザ間 CF 推薦は、人間をフィルターに通しているの、人間が行う推薦と関係性があるのは、当然である。今回は、それを新たな面で再認識できたともいえる。しかし、本来の目的である誤推薦されたコンテンツの特徴が発見できなかったため、ユーザ毎に誤推薦されたコンテンツのタグを調査した。

##### 3.2.1 ユーザ毎に誤推薦されたコンテンツのタグ

図 1 は、ユーザ数 943 人の内、良い結果が得られたユーザを 1 人選んだ。図 1 からわかるように、ピークが複数観測され、このユーザの誤推薦されたコンテンツには特徴があることがわかる。そのうちの最も高いピークは  $x$  軸が 3、 $y$  軸が 10 であり、ここから、このユーザが 10 番目に好きなタグがコンテンツの 3 番目に関連度をもつタグと同じ場合、評価を高くすることがわかる。つまり、この特徴に該当するコンテンツに重みを付加することで誤推薦されたコンテンツが正推薦できる。

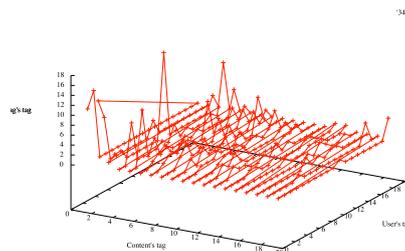


図 1: ユーザ毎の調査結果

### 4 今後の予定

今後の予定として、ユーザ、映画へのタグの重みの求め方の改善、及び統計解析の改善が第一にあげられる。また、評価値の低いデータの調査、コンテンツ間の共起関係、クラスタの検討などもあげられる。

### 参考文献

[1] MICHAEL J.PAZZANI, "A Framework for Collaborative, Content-Based and Demographic Filtering", Artificial Intelligence Review 13: 393-408, 1999.