

観測値の類似度を考慮した協調フィルタリング

On a Collaborative Filtering Method using Improved Similarity of Observed Value

†大島 敬志
Keishi Ohshima

‡鈴木 誠
Makoto Suzuki

†平澤茂一
Shigeichi Hirasawa

Abstract—Recommender systems apply knowledge discovery techniques to the problem of making personalized recommendations for information, products or services during a live interaction. These systems especially collaborative filtering ones, are achieving widespread success. Conventional collaborative filtering systems use methods that are based solely on observed scores for items.

In this paper we will propose a new collaborative filtering method using improved similarity of value that are based on observed binary scores for items.

Keywords—Collaborative Filtering, Recommender Systems, Memory Base Reasoning

1 はじめに

近年、膨大なデータから利用者や消費者等のユーザの要求を自動的に推定し、その要求を満たす情報を積極的に推薦するシステムが研究されている。例えば流通業では書籍やCDについて、過去のユーザの評価データからユーザの要求を満たす商品を推薦している。本論文ではこのような推薦システムの基本技術である協調フィルタリング [4] について考察する。

協調フィルタリングとは過去の購買履歴のデータから購買やアクセスパターンの類似するユーザを同定し、それらのユーザの嗜好から特定ユーザの嗜好を推定する手法である。協調フィルタリングの手法としてよく知られたアルゴリズムに相関係数法 [1] や二項関係学習法 [2] がある。これらの手法ではそのアイテム (商品) が好きか嫌いかという意味での多段階の「嗜好についての類似度」を用いて特定ユーザの嗜好を予測している。しかしアイテムが商品や映画の場合の推薦システムでは「嗜好についての類似度」とは別に、どのアイテムについて回答しているかという意味での 1, 0 で表される「回答パターンの類似度」も存在している。また従来手法では「嗜好についての類似度」を求める際に全てのアイテムを同等に扱って

る。しかしアイテムによっては嗜好に影響を与えないものや一部のアイテムには嗜好に強い影響を与えるものが存在している。

そこで本論文では相関係数法と二項関係学習法に対し、「回答パターンの類似度」を用いて嗜好に影響を与えているアイテムに対してウエイトづけを行い推薦システムの精度が向上する事を示す。

2 協調フィルタリング

2.1 協調フィルタリングとは [4]

協調フィルタリングとは、あるアイテムに対するユーザの評価を、そのユーザの別のアイテムに対する評価値と、他ユーザの評価データに基づいて予測することである。具体的な適用例として、映画 (アイテム) に対する多数のユーザの評価を入力とし、あるユーザがまだ見ていない映画に対する評価を、そのユーザの他の映画に対する評価と、他のユーザの評価に基づいて予測する。この時入力された評価の数はデータベース全体から見ると非常に少ない。協調フィルタリングの考え方は、評価対象ユーザとの類似度が高いユーザが高評価を与えたアイテムは、評価対象ユーザも高い評価を付ける可能性が高いというものである。

図 1 に協調フィルタリングの概念図を示す。図 1 では、評価対象ユーザ u_i のアイテム v_x に対する 5 段階評価値を予測している様子が示されている。ここで未観測の要素は空白で示す。図 1 に示された u_i 以外のユーザ u_1, u_2, u_3 のうち u_i と最も評価の傾向が類似しているユーザは u_2 であり、 u_2 の v_x に対する評価が高いことから、 u_i の v_x に対する評価も高いものと予測される。また u_1, u_3 はいずれも u_i とは評価値の傾向が異なっており、 u_i との類似度が低いことから、 v_x の評価値予測の際にはそれほど考慮されない。

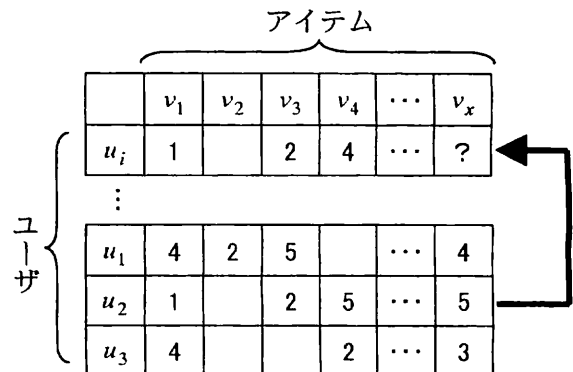


図 1 協調フィルタリング概念図

† 〒169-8555 東京都新宿区大久保 3-4-1 早稲田大学大学院理工学
研究科 経営システム工学専攻. Dept. of Industrial and
Management System Engineering School of Science and
Engineering, Waseda University, Okubo 3-4-1, Shinjuku-ku,
Tokyo, 169-8555 JAPAN.

‡ 〒251-8511 神奈川県藤沢市辻堂西海岸 1-1-25 1410-3 湘南工科大
学 情報工学科. Shonan Institute of Technology, Tujido
Nisikaigan 1-1-25 Fujisawa-shi Kanagawa-ken, 251-8511
JAPAN.

2.2 従来手法

2.2.1 相関係数法 [1]

行列 $M = (M_{iz})$ は行がユーザを列がアイテムを表す行列として、その (i, z) 要素が i 番目のユーザの z 番目のアイテムの評価値とする。未観測の要素 (空白) については欠損値として扱い、 $M_{iz} = *$ とする。通常、行列 M は欠損値を多く含み、ほとんどの M_{iz} には値が与えられない。相関係数法では評価対象ユーザ i の z 以外のアイテムに対する評価値と、他のユーザの評価値に基づいて未知の評価値 M_{iz} を以下のように算出する。

$$M_{iz} = M_i + \frac{\sum_j C_{ij} (M_{jz} - M_j)}{\sum_j |C_{ij}|} \quad (\text{式 1})$$

ここで M_i と M_j はユーザ i と j の評価値の平均値

$$M_i = \frac{\sum_x M_{ix}}{\sum_x 1} \quad (\text{式 2}), \quad M_j = \frac{\sum_x M_{jx}}{\sum_x 1} \quad (\text{式 3})$$

である。ここでの \sum_x は評価値が既知のアイテムだけについてとる。また C_{ij} は評価対象ユーザ i とデータベース中のユーザ j の相関係数

$$C_{ij} = \frac{\sum_x (M_{ix} - M_i)(M_{jx} - M_j)}{\sqrt{\sum_x (M_{ix} - M_i)^2 \sum_x (M_{jx} - M_j)^2}} \quad (\text{式 4})$$

である。ここでの \sum_x は評価対象ユーザ i とデータベース中のユーザ j の両方で評価値が既知のアイテムについてとる。相関係数法にはその他、ユーザの類似度ではなくアイテムの類似度を用いる手法 [3] 等も研究されている。

2.2.2 二項関係学習法 [2]

今、 w_{ij} を行 i から j への重みとする。評価値が入力されている要素を参照し、重み w_{ij} を以下のように更新する。

weight update: $M = (M_{jz})$ が欠損値でない全ての (j, z) について

$$\begin{cases} w_{ij} = (2 - \gamma)w_{ij} & \text{if } M_{jz} = M_{iz} \\ w_{ij} = \gamma w_{ij} & \text{if } M_{jz} \neq M_{iz} \end{cases} \quad (\text{式 5})$$

ここで γ はユーザが任意に与えた値 (ただし $0 < \gamma \leq 1$) である。

観測データを正しく予測するのに寄与している行の重みを $2 - \gamma$ 倍し、誤って予測するのに寄与している行の重みを γ 倍する。つまり w_{ij} は以下の式と同義である。

$$w_{ij} = (2 - \gamma)^a \gamma^b \quad (\text{式 6})$$

ただし a はユーザ i に対して j が同じ評価を示した回数、 b はユーザ i に対して j が異なった評価を示した回数である。実際に予測を行う場合は w_{ij} の和が最大になる値を予測値とする。

$$M_{iz} = \arg \max_{a \in A} \sum_{j: M_{jz} = a} w_{ij} \quad (\text{式 7})$$

ただし A は評価値の集合である。

2.3 従来手法の問題点

従来手法では類似度を計算する際に全てのアイテムを同等に扱って類似度を計算している。しかしアイテムによっては類似度に影響を与えるアイテムとそうでないアイテムが存在している。例えば、ホラー映画の評価値の予測を行いたいのに恋愛映画の評価値を用いて類似度の予測を行っても精度は向上しない。その問題点を改善するにはアイテムのコンテンツが分かればいいが、コンテンツの分類にはさまざまな概念からさまざまな分類ができてしまい、一意に正解・不正解が定まらない。またコンテンツによるフィルタリングを導入すると、人の嗜好を考慮するのはなく同じコンテンツのアイテムに同じ評価値が予測される傾向が強くなり問題である。

一方、従来手法ではそのアイテムが好きか嫌いかという意味での多段階の「嗜好についての類似度」を用いて特定ユーザの嗜好を予測している。しかしアイテムが商品や映画の場合の推薦システムでは「嗜好についての類似度」とは別に、どのアイテムについて回答しているかという意味での 1, 0 で表される「回答パターンの類似度」も存在している。

次節ではこの問題点を踏まえ、「回答パターンの類似度」を用いてアイテムごとにウエイト付けを行い、予測精度の向上を目指す。

3 提案手法

3.1 提案手法の概要

提案手法ではアイテム z の嗜好に大きな影響を及ぼすアイテム x とは、アイテム z に評価値を入力しているユーザがアイテム x に評価値を入力しており、逆にアイテム z に評価値を入力していないユーザが評価値を入力していないアイテムだと考える。なぜなら同じような傾向を持っているアイテムであればユーザがそのアイテムに興味を持っていれば両方に回答し、逆に興味を持っていなければ両方に回答しないからである。このアイテム間の影響を考えたアイテムの「回答パターンの類似度」を用いて 2 つの手法を提案する。

①相関係数法に用いることでアイテム z の嗜好に強く影響を与えるアイテムを考慮した「嗜好についての類似度」を求める。

②二項関係学習法に用いることでアイテム z に対してユーザのウエイトを更新する際にアイテム z の嗜好に強く影響を与えているアイテムほどウエイトを大きく更新する。

提案手法のアルゴリズムの概要は以下の通りである。

ステップ 1: アイテムについて回答しているかという意味での 1, 0 で求められる評価値を用いてアイテムの「回答パターンの類似度」を求める。

ステップ 2: アイテムの「回答パターンの類似度」を考慮し、①では「嗜好についての類似度」を、②で

はウエイトを求める。
 ステップ3: ①では「嗜好についての類似度」を、②ではウエイトを用いて未知の評価値 M_{iz} を予測する。

3.2 提案手法① (相関係数法)

従来手法ではユーザの相関係数 C_{ij} を求める際に全てのアイテムを同等に扱っている。そこで前節で述べたアイテムの「回答パターンの類似度」を用いることでより正確な「嗜好についての類似度」を求める。

(ステップ1) 行列 $M' = (M'_{iz})$ は行がユーザを表し、列がアイテムを表す行列として、その (i, z) 要素が i 番目のユーザが z 番目のアイテムに評価値が入力されていれば1, そうでなければ0を入力する。

$$M'_{iz} = \begin{cases} 0 & \text{if } M_{iz} = * \\ 1 & \text{if } M_{iz} \geq 1 \end{cases} \quad (\text{式 8})$$

行列 M' の値を用いてアイテム間の相関を求める。ここでアイテム間の類似度は内積を用いる。

$$d_{xz} = \frac{\sum_j M'_{jx} \cdot M'_{jz}}{\sum_j M'_{jx} \sum_j M'_{jz}} \quad (\text{式 9})$$

d_{xz} はアイテム z に評価値を入力しているユーザが評価値を入力しているアイテムほど値が大きくなる。 d_{xz} はアイテムの「回答パターンの類似度」でありこの値が大きいほどアイテム z の嗜好に影響を与えていると考えられる。

(ステップ2) d_{xz} の値を用いてアイテム z の予測を行う際の「嗜好についての類似度」 C'_{ijz} を以下のように修正する。

$$C'_{ijz} = \frac{\sum_x d_{xz} (M_{ix} - M_i) (M_{jx} - M_j)}{\sqrt{\sum_x d_{xz} (M_{ix} - M_i)^2 \sum_x d_{xz} (M_{jx} - M_j)^2}} \quad (\text{式 10})$$

(ステップ3) この C'_{ijz} を用いて未知の評価値 M_{iz} を以下のように算出する。

$$M_{iz} = M_i + \frac{\sum_j C'_{ijz} (M_{jz} - M_j)}{\sum_j |C'_{ijz}|} \quad (\text{式 11})$$

3.3 提案手法② (二項関係学習法)

従来手法では重みを学習する際に γ というユーザが任意に与える値を用いていた。しかし本論文では γ のかわりに (式9) で求めたアイテムの類似度を用いる。

(ステップ1) 相関係数法と同じく (式9) から d_{xz} を求める。

(ステップ2) 重み w'_{ijz} を以下のように更新する。
 weight update: $M = (M_{jz})$ が欠損値でない全ての $(i \neq j)$ についてアイテム z の予測を行う際のウエイト w'_{ijz} について

$$\begin{cases} w'_{ijz} = \frac{1}{2} - (1 - d_{xz}) w_{ij} & \text{if } M_{jz} = M_{iz} \\ w'_{ijz} = (1 - d_{xz}) w_{ij} & \text{if } M_{jz} \neq M_{iz} \end{cases} \quad (\text{式 12})$$

とする。

(ステップ3) 実際に予測を行う場合は w'_{ijz} の和が最大になる値を予測値とする。

$$M_{iz} = \arg \max_{a \in A} \sum_j M_{jz=a} w'_{ijz} \quad (\text{式 13})$$

4 提案手法の評価と考察

4.1 評価データ

評価データには協調フィルタリングの検証用データとしてよく用いられる MovieLens [5] データと EachMovie [6] データを用いた。

(1) MovieLens データは 943 人の 1682 の映画に対する 5 段階評価データである。評価されている項目は全部で 100000 個あり欠損率は 93.7% である。ここではそのうち 80000 個をランダムに取り出し学習データとみなし、残りの 20000 個のデータで推定誤差を評価した。

(2) EachMovie データは 72916 人の 1649 の映画に対する 6 段階評価のデータである。本実験ではそのうち 1800 の 1649 の映画に対してシミュレーションを行った。評価されている項目は全部で 107803 個あり欠損率は 96.4% である。ここではそのうち 87803 個をランダムに取り出し学習データとみなし、残りの 20000 個のデータで推定誤差を評価した。

両データともに同じ条件でシミュレーションを 5 回行い、その平均について評価した。

4.2 評価基準

提案手法を MAE と F 値で評価した。

MAE とは平均絶対誤差のことであり以下の式で表される。

$$\text{MAE} = \frac{\sum_{i=1}^N |p_i - q_i|}{N} \quad (\text{式 14})$$

ただし N はテストデータ数、 p_i は予測値、 q_i は正解値とする。

F 値とは推薦されるべきアイテムが正しく推定される割合を表しており以下の式で表される。ここでは評価値が 4 を超えるアイテムを推薦されるべきアイテムとした。

$$F = \frac{2PR}{P+R} \quad (\text{式 15})$$

ただし P (再現率) は

$$\frac{\text{推薦された適合アイテム数}}{\text{全テストデータ中の全適合アイテム数}} \quad (\text{式 16})$$

R (正解率) は

$$\frac{\text{推薦された適合アイテム数}}{\text{推薦されたアイテム数}} \quad (\text{式 17})$$

である。再現率は適合アイテムを漏れなく推薦できる度合い、正解率は適合アイテムだけを推薦できる度合いを表している。なお再現率と正解率はトレードオフ関係になっており、このふたつを同時に評価する評価基準として F 値は代表的な評価基準である。

4.3 評価結果と考察

4.3.1 相関係数法の評価と考察

表1, 2に評価結果を示す. 表1, 2より MovieLens データでは MAE, F 値が従来手法に比べそれぞれ 0.021, 0.012 向上した. EachMovie データでも MAE, F 値が従来手法に比べそれぞれ 0.032, 0.016 向上した. これは嗜好に影響を与えているアイテムを考慮した「嗜好についての類似度」が改善されたためだと考えられる. 実際にテストデータを予測する際に求めた類似度の式全てについて

$$|M_i + (M_{jz} - M_j) - M_{iz}| \leq 0.5 \cap C'_{ijz} - C_{ij} \geq 0 \quad (\text{式 18})$$

または

$$|M_i + (M_{jz} - M_j) - M_{iz}| > 0.5 \cap C'_{ijz} - C_{ij} < 0 \quad (\text{式 19})$$

になる類似度の総数を調べた(ただしここでの M_{iz} は既知の値). ここで(式18)は正解に寄与しているユーザの「嗜好についての類似度」を従来手法より増加している場合である. 逆に(式19)は不正解に寄与しているユーザの「嗜好についての類似度」を従来手法より減少させている場合である.(式18)もしくは(式19)を満たしている類似度の式が全体の59%を占めており提案手法①が全体的に類似度の式を向上させていると考えられる.

表1 MovieLens データの相関係数法での結果

	MAE	再現率	正解率	F 値
従来手法	0.712	0.724	0.705	0.714
提案手法①	0.691	0.737	0.715	0.726

表2 EachMovie データの相関係数法での結果

	MAE	再現率	正解率	F 値
従来手法	0.859	0.654	0.750	0.699
提案手法①	0.827	0.671	0.764	0.715

4.3.2 二項関係学習法の結果と考察

表3, 4に評価結果を示す. 表3, 4より従来手法の場合, MovieLens データでは $\gamma=0.8$ の時に MAE の値が, $\gamma=0.9$ の時に F 値の値が最も良くなった. EachMovie データでは $\gamma=0.6$ の時に MAE の値が, $\gamma=0.7$ の時に F 値の値が最も良くなった. 提案手法②では従来手法が最も良い γ を用いたものよりも MovieLens データでは MAE, F 値が従来手法に比べそれぞれ 0.016, 0.009 向上した. EachMovie データでも MAE, F 値が従来手法に比べそれぞれ 0.021, 0.017 向上した. これは従来手法ではどのアイテムに対しても一定のウエイトを与えていたために, 嗜好に影響を与えないアイテムが原因で嗜好の似ているユーザの類似度を下げていたのが改善されたためだと考えられる.

表3 MovieLens データの二項関係学習法での結果

	MAE	再現率	正解率	F 値
従来手法 $\gamma=0.7$	0.766	0.756	0.698	0.726
従来手法 $\gamma=0.8$	0.762	0.767	0.694	0.729
従来手法 $\gamma=0.9$	0.763	0.782	0.688	0.732
提案手法②	0.744	0.775	0.710	0.741

表4: EachMovie データの二項関係学習法での結果

	MAE	再現率	正解率	F 値
従来手法 $\gamma=0.5$	0.898	0.744	0.700	0.721
従来手法 $\gamma=0.6$	0.894	0.751	0.699	0.724
従来手法 $\gamma=0.7$	0.894	0.754	0.696	0.724
提案手法②	0.873	0.772	0.712	0.741

5 むすび

本論文では「回答パターンの類似度」を用いて嗜好に影響を与えているアイテムに対してウエイトづけを行う手法を提案した. またシミュレーションを行い, 従来手法に比べ, 精度が向上する事を示した.

今後の課題として, どのようなデータの傾向があれば, 「回答パターンの類似度」という概念を用いるのが有効なのかを考えたい. また人の嗜好というあいまいなものに対してどのような傾向があるのかを考え手法に反映させたい.

謝辞

著者の一人大島は, 研究を進めるにあたり御討論, ご助言を頂いた湘南工科大学小林学先生に感謝します.

<参考文献>

- [1] P.Resnick N.Iacovou M.Suchak P.Berstom "GroupLens:an open architecture for collaborative filtering of netnews", 1994, in Proc of CSCW94, pp.175-186
- [2] A.Nakamura N.Abe, "Collaborative filtering using weighted majority prediction algorithms", 1997, in Proc.of ICML98, pp.395-403
- [3] B.Sarwar G.Karypis J.Konstan Riedl, J, "Item-based Collaborative Filtering Recommender Algorithms.", 2001, in proc of the 10th international world wide web conference
- [4] 福原知宏「協調フィルタリングに関する研究動向」
http://ai-www.aist-nara.ac.jp/doc/people/tomohi-f/Docs/cf_review.html, 1998
- [5] <http://www.cs.umn.edu/Research/GroupLens>
- [6] <http://www.research.digital.com/SRC/eachmovie>