

評価傾向の差異を考慮した分散表現による協調フィルタリング

1X15C039-1 後藤 亮介
指導教員 後藤 正幸

1 研究背景と目的

近年、情報技術の発展に伴い、EC サイトが取り扱うアイテム数が増加し、利用するユーザの嗜好も多様化している。そのため、EC サイトを運営する企業にとって、過去の購買履歴データから各ユーザの嗜好を把握し、嗜好に合致したアイテムを提示するための推薦システムは重要な Web マーケティング技術の 1 つになっている。

推薦手法の 1 つとして、評価履歴データを用いて学習し、未評価アイテムの中で予測評価値の高いアイテムを推薦するアイテムベース協調フィルタリングがある。評価履歴データを用いることで、選択したアイテムが好きか嫌いかを考慮することができるため、ユーザの嗜好をより反映した推薦システムを実現することができる。評価履歴データを用いたアイテムベース協調フィルタリング手法の 1 つに、Kuzmin の手法 [1] がある。この手法は、Item2vec[2] をベースとした未評価アイテムの評価値予測モデルであり、Item2vec に入力する学習データの前処理として、評価値をもとにデータを分割している。過去に評価を付けた全アイテムに対する評価値の平均値をユーザ毎に算出し、平均値以上の評価値が付与されたアイテムを“高評価アイテム群”，平均値未満の評価値が付与されたアイテムを“低評価アイテム群”と定義している。そして、Item2vec の学習を行う際、同じユーザに評価され、かつ同じアイテム群に分類されたアイテムの分散表現が類似するように学習を行う。その際、Kuzmin の手法では、“高評価アイテム群”と“低評価アイテム群”を同時に Item2vec で学習し、単一の意味空間を構成している。しかし、ユーザが対象アイテムに対して高評価する理由と、低評価する理由は異なると考えられる。そのため、それらを単一のモデルで表現するのではなく、別々のモデルで表現することで、ユーザの嗜好をより捉えたアイテムの分散表現が可能になると考えられる。

そこで本研究では、各ユーザに対して与えられる“高評価アイテム群”と“低評価アイテム群”を、それぞれ独立に学習し、高評価ベクトル空間と低評価ベクトル空間を得る。そして、これらを統合した評価値予測を行う手法を提案する。また、提案手法を映画評価履歴データに適用し、従来手法と比較することで、提案手法の有効性を検証する。

2 準備

2.1 Item2vec

自然言語処理の分野において、単語を低次元の意味空間上のベクトルで表現する手法として Word2vec[3] が知られている。ここで Word2vec をアイテムの分散表現モデルに援用した手法が Item2vec である。ユーザが過去に購買した全アイテムを購買系列として、あるアイテムは、同じ購買系列中の他のアイテムから予測できると仮定する。そして、注目アイテムベクトルと周辺アイテムベクトルとの内積を大きくするように各アイテムベクトルの学習を行う。また、学習の進行を高速化し、分散表現の精度を高めるために、全アイテムから確率的にアイテムのサンプリングを行うノイズ分布をあらかじめ定義し、注目アイテムに対してそのアイテム系列とは関係なく獲得されるアイテム（ネガティブアイテム）をこのノイズ分布に従ってサンプリングする。そして、ネガティブアイテムと注目アイテムとの内積を小さくするように各アイテムベクトルを更新する。

いま、アイテム数を N 、ユーザ数を M とし、全アイテム集合を $\mathcal{I} = \{i_n : 1 \leq n \leq N\}$ 、全ユーザ集合を $\mathcal{U} = \{u_m : 1 \leq m \leq M\}$ とする。また、ネガティブサンプリング数を K とする。ユーザ u_m が j 番目に購買したアイテム $x_{m,j} \in \mathcal{I}$ を表す意味空間上のベクトルを $\mathbf{v}_{m,j}$ とし、ある

1 つの周辺アイテム $x_{m,l} (l \neq j) \in \mathcal{I}$ を表すベクトルを $\mathbf{v}_{m,l}$ とする。また、ネガティブアイテム $y_{m,j}^k \in \mathcal{I} (1 \leq k \leq K)$ のベクトルを $\mathbf{s}_{m,j}^k$ で表す。このとき、注目アイテム $x_{m,j}$ に対する周辺アイテムへの損失関数は式 (1) で定義される。

$$\text{Loss}(x_{m,j}, x_{m,l}) = \log(\sigma(\mathbf{v}_{m,j}^T \cdot \mathbf{v}_{m,l})) + \sum_{k=1}^K \log(\sigma(-\mathbf{v}_{m,j}^T \cdot \mathbf{s}_{m,j}^k)) \quad (1)$$

ただし、 $\sigma(\cdot)$ はシグモイド関数を表す。ユーザ u_m の全購買アイテム数を e_m とすると、学習データ全体における損失項は以下の式 (2) で定義される。

$$\text{Loss}_{All} = \sum_{m=1}^M \sum_{j=1}^{e_m} \sum_{l \neq j}^{e_m} \text{Loss}(x_{m,j}, x_{m,l}) \quad (2)$$

2.2 従来手法: Kuzmin の手法 [1]

Kuzmin の手法では、各ユーザの平均評価値を閾値として、全アイテムを高評価アイテム群と低評価アイテム群に分割する。そして、同一アイテム群を 1 つの購買系列と捉えて Item2vec の学習を行う。これにより、同じユーザに選択されていて、かつ似た評価がされているアイテム同士が、類似したベクトルとなるように分散表現が学習される。そして、得られたアイテムベクトル同士の類似度を用い、式 (3) によって予測評価値を算出する。

$$\hat{r}_{m,x} = \eta_x + \frac{\sum_{\tilde{x} \in \tilde{\mathcal{I}}_m} \text{Sim}(\tilde{x}, x) \times (r_{m,\tilde{x}} - \eta_{\tilde{x}})}{\sum_{\tilde{x} \in \tilde{\mathcal{I}}_m} |\text{Sim}(\tilde{x}, x)|} \quad (3)$$

ここで、 $\hat{r}_{m,x}$ はユーザ u_m のアイテム $x \in \mathcal{I}$ に対する予測評価値、 η_x は全ユーザのアイテム x に対する平均評価値、 $\text{Sim}(\tilde{x}, x)$ はアイテム \tilde{x} と x を表すベクトルのコサイン類似度、 $\tilde{\mathcal{I}}_m \in \mathcal{I}$ はユーザ u_m が評価した中でアイテム x と類似度が高いアイテム集合とする。そして、最終的に各ユーザに対して予測評価値が上位のアイテムを推薦することで、評価値を考慮した推薦システムを実現している。

3 提案手法

3.1 概要

従来 Kuzmin の手法では、評価値の高いアイテムと低いアイテムを同一のモデルで学習するため、得られる意味空間は単一である。しかし、ユーザの評価行動の背後には多様な嗜好が存在している。例えば、1 人のユーザに関して考えた場合も、あるアイテムに対して好きと認識するに至るまでのメカニズムと、嫌いと認識するに至るまでのメカニズムは異なると考えられ、これらを単一のモデルで表現することは困難である可能性がある。

そこで本研究では、高評価アイテム群と低評価アイテム群に関して、それぞれ独立の意味空間を学習する（以下、高評価ベクトル空間、低評価ベクトル空間と呼ぶ）。すなわち、高評価アイテム群のみを学習させ、高評価アイテム群に特化した分散表現を獲得できるモデルと、低評価アイテム群のみを学習させ、低評価アイテム群に特化した分散表現を獲得できるモデルを獲得する。各アイテム群に対して独立に意味空間を構成することで、評価による意味の違いを捉えたベクトルが得られると考えられる。そして、それぞれのモデルから得られた分散表現を統合する形で評価値の予測を行うことで、より高精度な推薦システムの実現を図る。

3.2 提案手法のアルゴリズム

提案手法のアルゴリズムを以下に示す。

STEP1: アイテム評価系列の分割

各ユーザに対して、ユーザが付与した全評価値の平均を算出する。そして、算出した平均評価値よりも高く評価されているアイテムと低く評価されているアイテムに分割する。

STEP2: 意味空間の構成

STEP1 で得られた各アイテム群を学習データとし、Item2vec を学習する。この際、高評価アイテム群と低評価アイテム群はそれぞれ独立のモデルで学習を行う。

STEP3: 2つの意味空間での予測評価値の算出

各意味空間上のアイテムベクトル同士の類似度を用いて、式(3)により別々に予測評価値を算出する。

STEP4: 予測評価値の算出

式(3)によって求めた各予測評価値の加重平均により、最終的な予測評価値を式(4)で算出する。

$$\hat{r}_{m,x} = \alpha \times \hat{r}_{m,x}^{(low)} + (1 - \alpha) \times \hat{r}_{m,x}^{(high)} \quad (4)$$

ここで、 $\hat{r}_{m,x}$ はユーザ u_m のアイテム x に対する予測評価値、 $\hat{r}_{m,x}^{(low)}$ は低評価ベクトル空間上で得られる予測評価値、 $\hat{r}_{m,x}^{(high)}$ は高評価ベクトル空間上で計算される予測評価値、 α は評価値の混合比とする。

4 実験

4.1 実験条件

評価値の予測実験には、公開データセット MovieLens の映画評価データを用いる。評価値データ数は 100,000 件、ユーザ数 942 人、アイテム数 1,683 本、評価値は 1~5 の 5 段階を実験対象データである。本実験では評価指標として式(5)の MAE(平均絶対誤差)を用い、5 分割交差検証法によりテストデータに対して得られる各予測精度の平均を求める。

$$MAE = \frac{1}{D} \sum_{m=1}^M \sum_{x \in \mathcal{I}} |\hat{r}_{m,x} - r_{m,x}| \delta_{m,x} \quad (5)$$

ここで、 $r_{m,x}$ はユーザ u_m のアイテム x に対する真の評価値、 D はテストデータ数である。また、 $\delta_{m,x}$ はインジケータ関数であり、ユーザ u_m がアイテム x を評価している場合は 1、それ以外は 0 を示す。

高評価ベクトル空間と低評価ベクトル空間での予測評価値の混合比 α は、学習データをさらに 2 分割し、交差検証法により最適化する。具体的には、混合比 α を 0 から 1 まで 0.1 ずつ変化させた全パターンに対して、設定用学習データから得られる 2 つの意味空間を用いて設定用テストデータに対する MAE を求める。そして、交差検証法により分割された 5 つの設定用テストデータに対して求めた MAE の平均が最も低い α を最適な混合比とする。また、高評価アイテム群と低評価アイテム群の分散表現モデルの違いを確認するために、ベクトル空間が異なることによるアイテム間の類似度の差異についても検証する。

4.2 実験結果と考察

図 1 は、混合比 α を 0 から 1 まで変化させた場合の設定用テストデータに対して求めた MAE の推移を表す。

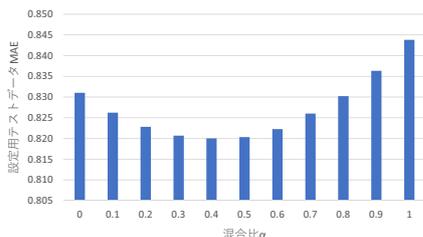


図 1: 混合比を変化させた際の MAE の推移

図 1 より、 $\alpha=0.4$ のときに最も MAE が小さくなった。そこで、最適な混合比 α を 0.4 と設定した時の評価値予測実験の結果を表 1 に示す。また、比較手法として従来手法と高評価ベクトル空間、低評価ベクトル空間の片方のみを用いて評価値予測をした結果を示す。

表 1: 各手法における MAE

手法	MAE
従来手法 (Kuzmin の手法)	0.8107
高評価ベクトル空間のみを用いた場合 ($\alpha = 0.0$)	0.8110
低評価ベクトル空間のみを用いた場合 ($\alpha = 1.0$)	0.8204
提案手法 ($\alpha = 0.4$)	0.8026

表 1 より、提案手法は従来手法と比べて予測精度が向上している事が分かる。高評価ベクトル空間のみを用いた場合と低評価ベクトル空間のみを用いた場合において予測精度が低下したのは、学習データを高評価と低評価に分割する必要があり、モデルの学習に用いるデータが相対的に少なくなってしまうためと考えられる。しかし、これらのモデルを統合して評価値を予測することで、アンサンブル効果により従来法よりも予測精度を向上させることが可能である。

表 2 に、高評価ベクトル空間と低評価ベクトル空間での類似度の差が小さいアイテムの組み合わせ、表 3 に類似度の差が大きいアイテムの組み合わせの一例を示す。

表 2: 類似度の差が小さいアイテムの組み合わせ

映画タイトル	高評価	低評価
スタートレック 4・スタートレック 6	0.8540	0.9068
Sweet Hereafter・地球が静止する日	0.2239	0.2698
フリッパー・パーフェクトワールド	0.7697	0.7990

表 3: 類似度の差が大きいアイテムの組み合わせ

映画タイトル	高評価	低評価
ニューシネマパラダイス・StarWars5	0.7017	0.3456
ジュマンジ・エイリアン 4	0.5874	0.9267
トップガン・バットマンリターンズ	0.4928	0.8012

表より、例えばスタートレック 4 と 6 は、高評価ベクトル空間と低評価ベクトル空間で、共に似たような類似度を示した。すなわち、スタートレック 4 を高く評価しているユーザは、スタートレック 6 を同時に高く評価している。また、4 を低く評価しているユーザも、同じく 6 を低く評価していることがわかる。一方、ニューシネマパラダイスを高く評価しているユーザは、StarWars5 も高く評価している傾向があるが、どちらかの作品を低く評価しているユーザは、もう一方の作品は MovieLens 上で評価を行っていない傾向がわかる。

このように、高評価アイテム群の背後に存在するユーザの嗜好と、低評価アイテム群の背後に存在するユーザの嗜好を独立に学習することで、ユーザの嗜好を詳細に把握することを可能とし、予測精度の向上に寄与していると考えられる。

5 まとめと今後の課題

本研究では高評価ベクトル空間と低評価ベクトル空間を独立に学習し、各空間上から推定される評価値の加重平均で評価値予測を行うモデルを提案した。また、ベンチマークデータを用いた実験により提案手法の有効性を示した。

今後の課題として、提案手法を他のデータにも適用し有効性を検証する。また、評価値予測だけでなく TopN 精度でアイテムに対する推薦精度の測定をすることが挙げられる。

参考文献

- [1] Vitali Kuzmin, “Item2Vec-based Approach to a Recommender System”, 2017.
- [2] Oren Barkan, Noam Koenigstein, “Item2vec: Neural item embedding for collaborative filtering” 2016 IEEE 26th International Workshop on Machine Learning for Signal Processing, 2016.
- [3] T. Mikolov, I. Sutskever, K. Chen, G. Corrado, and J. Dean, “Distributed Representations of Words and Phrases and their Compositionality”, *Advances in NIPS26*, pp. 3111-3119, 2013.