

論文の内容の要旨

論文題目 オンラインプラットフォームにおける
 バイアスを考慮したコンテンツの質の定量化

Quantifying the Unbiased Quality of the Contents
in Online Platforms

氏 名 福馬 智生

インターネット技術の爆発的な発展に伴い、様々な種類のウェブサイトやフォーラムが誕生している。それらには個人の意見、コメント、レビューが大量に集積されており、その結果情報の収集や他者との情報交換はかつてより遥かに容易になった。その結果個人が探したい情報が見つからないといった弊害が発生しており、それらは「情報過多 (information overload)」といった名前で広く知られている。そのようなユーザー体験の低下はオンラインプラットフォームの存続に致命的であり、昨今ではプラットフォーム運営側で様々な試みを用いて膨大な情報源の中からユーザに適した情報を見つけ出し提示するための「コンテンツの有用度の自動評価技術」に注目が集まっている。

コンテンツの質の自動評価に関する既存の試みとして「集合知」に基づく試みと、「機械学習」に基づく試みが存在する。前者はコンテンツの有用性をユーザー自身が投票するシステムを用い、他ユーザーによる評価の合計の多い順に並べ替えることで、容易に有用性の高いコンテンツに到達できるように誘導している。後者はPageRankのように開発者がヒューリスティックに設計した目的関数を最適化させる自己教師あり学習に基づく手法と、コンテンツの質を定量化した教師データに基づいて機械学習モデルを構築する教師あり学習に基づく手法に更に大別される。

本研究ではまずこれらの手法に潜む「バイアス」の観点から議論を行う。例として、プレゼンテーションバイアスを用いる。ユーザーはWEB上の全ての情報に目を通すことはなく、システムから提示される上位の一部しか目を通さないことが知られている。その結果見られた情報はクリックされる可能性があるが、見られないものは、本当は有用または興味がある内容だとしてもクリックされないというプレゼンテーションバイアス

が発生する。その結果、集合知に基づく試みでは、上位のコンテンツが票を独占し、真に価値があるコンテンツでも見られないため評価されないといったケースが起きる。加えて教師あり学習に基づく試みでは、それらユーザーのクリック情報を用いて最適化されることがほとんどであり(推薦システムなど)、同様のバイアスをシステム自体がアルゴリズム内に抱えることになる。

これらバイアスはユーザーに繰り返し提示することで自己フィードバックを起し、更にそのバイアスを強め意見の極化が生じることが知られている。意見の極化による弊害としては、人気度と質の相関が低くなることやフィルターバブルの発生、フェイク情報の拡散などが知られている。つまり既存の自動価値推定技術はバイアスを含んでおり、真にモノの価値を表しているとは言い難い。

そこで本研究では、オンラインプラットフォーム上における様々な意見の集合から、「真のモノの価値」の推定、またどれだけ既存のプラットフォームの評価がそれらバイアスによって歪められているかの定量化を目的とする。本研究は大きく二部構成になっており、第一部は投票行動におけるバイアスの除去、第二部は推薦システムにおけるアルゴリズムバイアスの緩和、第三部はプラットフォームにおける過大評価・過小評価の可視化がテーマになっている。

本論文の貢献は上記課題に対し、細分化された以下のResearch Questionに答える形で行われる。

Research Question 1: 複数のバイアスの影響を受けているユーザーの行動データから、「バイアスがなければどのように行動したのか」という反実仮想のもと、バイアスの影響を取り除いた真のコンテンツの質を測定できるか?

本研究では、「いいね」や「役に立った」といった投票形式の評価行動に焦点を当て、事前のバイアスが及ぼす影響について明らかにする。またそれら情報がなければ人はどう評価していたのかといった反実仮想な現象に対する推定技術の提案について述べる。

提案手法では、コンテンツの評判スコアのような認知バイアス情報の有無に基づく、ランダム化比較試験を用いて収集したデータセットについて、人々の投票行動の違いに着目する。続いて新たに人間の投票行動のモデリングを提案し、投票行動の違いについて機械学習モデルを用いて説明する。具体的には、ニューラルネットワークを用いて、周囲の認知バイアス情報によって「そのコンテンツが何倍有益に見えているか」を定量化する。さらに、認知的バイアスの影響を割り引くことで、認知バイアスの影響を取り除

いた形でも有用性を推定する手法を提案する。最後に、大手Q&AサイトであるStack Exchangeの行動ログを用いて、我々の手法の有効性を実験的に示し、モデルが学習した知識について解釈を行った。我々の研究は、オンラインプラットフォームにおける中心的な問題である、バイアスの影響を取り除いたコンテンツの有用性をどのように測定するかに取り組んでいる。

Research Question 2: ユーザーのレビュー投票行動から不良ユーザー・レビューを発見し、信頼度に基づいた意見の集約によって従来よりも良い集合知を獲得できるか？

例えばある商品のレビューのように、既存の星について単に平均といった集約方法は、必ずしも良い集合知に繋がるとは限らない。近年では、やらせレビューやアンチレビューのような評価を操作するためのユーザーが存在するためである。

そこで、リンク構造に基づくレビューの信頼性評価アルゴリズムによって得られた「信頼度」によって意見を集約することで得た集合知を、既存の平均といった集約方法と、有識者の意見との類似度という観点から比較を行った。レビューの信頼性評価にはKumar et al.によって提案されたREV2を用いる。レビューのデータには東京都内のラーメン屋のレビューを独自に収集した。結果として信頼性に基づいて集約することが有識者に近いよりよい集合知が獲得されることを確認した。また同章では、信頼性に基づいた評価によって平均と比べて評価が大きく変化した店舗についての特徴を分析した。さらにREV2によって判断される信頼性とは何かといったモデルの解釈性についても検証を行った。

Research Question 3: 推薦システムは履歴の少ないユーザーやアイテムの組み合わせについても正しくそれぞれの間の嗜好度をモデリングすることは可能か？

Research Question 4: 推薦システムは予測の不確実性をモデリングすることは可能か？

推薦システムは、ユーザーとアイテムのインタラクションの履歴データを元に、未知のユーザーとアイテムのインタラクションを予測する。それらは個人ごとに嗜好度または質を推論している点において、多元的なコンテンツの質の評価を行っているとは本研究では捉える。

従来の代表的な協調フィルタリング手法は、十分にユーザー・アイテムの嗜好度に関する履歴が手に入る場合には上手くモデリングができることが知られている。一方で、疎なデータ、つまりユーザー・アイテムの嗜好度に関する履歴がほとんど得られていないケースでは十分にそれらの関連性をモデリングできないことが知られている。しかし一般的に学習に用いられる履歴データは、プレゼンテーションバイアスの影響により、ごく一部のユーザーやアイテムのみ履歴データが豊富にあり、残りの殆どには履歴が僅かといった頻度での観点でバイアスのかかったデータであることが知られている。

また推薦システムの欠点として紹介したフィルターバブルの解消のためには、今後探索と活用を能動的に繰り返しながら、ユーザーの嗜好度を能動的に学習しながら推薦を行うといったアプローチが求められると考えられる。それらはつまり既存のオフラインデータに対してのみ、性能を向上させるようなモデルの開発と評価を行うことでは、それらの解消が根本的に困難であることを意味する。ユーザーのアイテムに対する嗜好度の予測とその予測の信頼度は、今後の能動的学習における探索と活用のバランスをとる上で非常に重要と考えられるが、既存の深層学習に基づく協調フィルタリング手法はそれらが行えないという欠点がある。

そこで本研究では、既存の潜在因子に基づく協調フィルタリング手法について、疎なデータからの学習と不確実性のモデリングを同時に可能にする汎用的な学習フレームワークMetaCFを提案した。提案手法は推薦システムをメタラーニングの一種であるNeural Processesの観点から定義し直すことで、従来のモデルにほとんど手を加えることなく、それらの予測を可能にした。

以上を総括すると、我々の提案手法・分析によって、以下のことが解決・明らかになった。1) 事前バイアスがなければ人はどのように評価したかに基づいて、コンテンツの質を認知バイアスの影響を取り除いた形で推定が可能になった。2) コンテンツの質を測るにあたり、個人の意見の集約は、評価者の過去の評価履歴などから判断した信頼性に基づいた形で行うことでより有識者に近い意見が得られる。3) 推薦システムにおいて疎なデータからの学習が可能になった。4) 推薦システムの予測に対して確信度をモデリングすることが可能になった。本研究の貢献により、オンラインプラットフォーム上におけるコンテンツの価値を全員にとって、または個人に対してより正確に測ることが可能になった。