

SNS からの重要意見抽出のためのデータセット構築及び LLM による分類検証

矢口 一晟¹櫻井 恵理子²櫻井 義尚³¹明治大学大学院 先端数理科学研究科
cs233017@meiji.ac.jp²産業能率大学 経営学部
SAKURAI_Eriko@hj.sanno.ac.jp³明治大学 総合数理学部
sakuraiy@meiji.ac.jp

概要

本研究では、重要意見抽出のためのデータセット構築および大規模言語モデル（以下：LLM）を用いた重要意見抽出の手法を提案し、その有効性を検証した。そこで、人間がラベリングを実施せずに SNS から企業がマーケティング活動に必要な重要意見抽出モデルが構築可能であることを示す。これにより、重要意見データセットの作成や抽出にかかる時間的かつ金銭的コストの削減が期待できる。データセットの構築検証ではワークショップを開催することで重要意見に関連する議論を実施し学生による重要意見データの作成が可能であるかを検証した。また重要意見の抽出検証では、LLM による Zero-shot 分類および擬似データを作成した後、ファインチューニングによる重要意見分類の精度を比較した。

1 はじめに

Twitter（現 X）などの SNS にはユーザーから発信される投稿データが蓄積されており、その中には企業のサービス向上や改善などに有益となる情報が存在し、マーケティング調査に活用されている。一方で、投稿データの中でも有益な情報が含まれている割合は少なく、人手でそれらを抽出することは大変な人的コストを要する。そこで本研究では、LLM を用いることで重要な意見を含む投稿の抽出を効率化する手法を提案し、その有効性を検証する。

具体的には、まず重要意見の抽出検証の先駆けとして X と Google レビューデータを用いた重要意見ワークショップを実施することで検証用データとなる重要意見データの作成を行う。ワークショップでは計 14 名の学生とともに 4 つの観点に対して、個人およびグループワークを行い、重要意見データセットの構築を図った。次に、ワークショップで作成した重要意見データから、得られる特徴を LLM によって抽出する。最後に抽出した特徴を用いて、2 つ

の手法による重要意見の抽出効果を検証する。

2 関連研究

SNS にはユーザーから発信される膨大な投稿データが蓄積されており、その中には企業のマーケティング調査に有益な情報が存在する。そのため、SNS を用いた意見抽出の研究が進められている [1]。

また、これまでの自然言語処理分野において、意見の中でも「要望」「不満」「現状認識」などの意見自体の役割を担う部分の抽出は扱われてこなかった。そこで、山本ら [2]の研究では人々が何を求めているかを調査するため機械学習モデルである

Support Vector Machine を用いて、アンケート回答に記載されている意見の中から要望が述べられている「要望文」を抽出する研究を行った。山本らの研究では「要望」に着目し、横浜環状線北西線の建設に関する業務における要望文を抽出するシステムを提案しているが、本研究では人々の共感性に着目し、意見データの中から企業がマーケティング調査を行う上で必要となる重要意見の抽出を目指す。

さらに、近年 LLM に関する研究が盛んに行われており Brown ら [3]は GPT-3 の開発を行い、LLM に特定のタスク説明文と少量の教師データを入力することで高い精度を達成し、日本語タスクにおいても高い性能を発揮できる場合があることが報告されている。そこで、藤井ら [4]は LLM によって生成した擬似データを用いて学習した小規模モデルが 6 個の日本語タスクにおいて LLM との比較検証を実施した結果、フォーマルなテキストを入力としたときの分類タスクにおいて Zero-shot および Few-shot Prompting を上回る結果となった。このアプローチは、LLM による重要意見抽出といった特定のタスクにおいても有効性を持つことが考えられる。

本研究では、大規模言語モデルである GPT-4o を用いた分類および品質の良い擬似データセットの作成による分類手法を提案し、重要意見の抽出タスクに

おける有効性を検証する。

3 データセット

本研究では、重要意見データセットの構築におけるワークショップおよび重要意見抽出効果の検証において分類用の検証データでは、Twitter（現 X）と Google レビューから合計 196 件のデータを用いた。また、本研究は重要意見抽出の対象をディズニー施設としているため、ディズニーに関するサービス向上・改善に繋がる重要意見抽出を行う。データの内訳は表 1 の通りである。

表 1: ワークショップ用データセットの内訳

	Twitter（現 X）	Google レビュー
データ数	160	36

4 提案手法

本章では、ワークショップによる重要意見データセットの構築および重要意見抽出の 2 つの手法を説明する。

4.1 重要意見データセットの構築

重要意見のラベリングは個人の主観や意見によって偏るという特性があるため、グループでの議論を実施し、役に立つ又は結果的に議論が巻き起こった意見をラベリングすることが重要意見のラベルになるという仮説のもと、ワークショップを実施した。具体的には、産業能率大学の学生 14 名と有識者をお招きし、4 つの観点（図 1）での 5 値（3 値）ラベリングおよびコメントワークを行い、196 件のデータセットをもとに重要意見データセットを作成した。最終的に、学生によるデータセット構築が、どの程度有識者が作成した重要意見データを網羅できているかを確認するために、図 2 の流れで 5 つの分析を実施した。

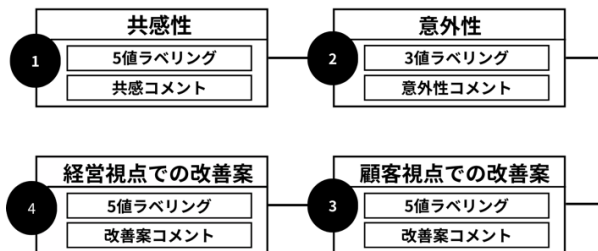


図 1: ワークショップの流れおよび 4 つの観点

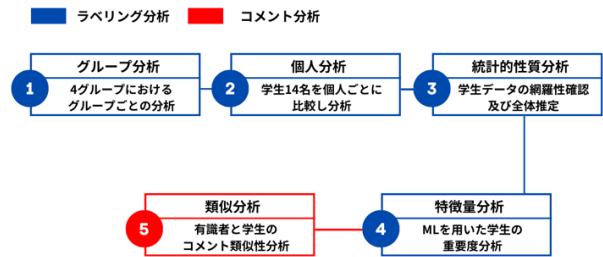


図 2: 5 つの分析フロー

4.2 LLM による重要意見抽出

LLM による重要意見の抽出検証では、GPT-4o を用いた Zero-shot 分類（図 3）と GPT-4o により作成した擬似データを、京都大学が公開している Deberta-v3-base-japanese に Fine-Tuning させることで重要意見抽出モデルによる分類（図 4）の 2 つの手法で検証した。後者の擬似データ作成においては同条件下で 1000 件、2000 件、3000 件の擬似データを作成し Fine-Tuning させ、比較する。検証時はワークショップで構築した重要意見データを検証用データとし、評価指標である Accuracy, Recall, Precision, F1-score を用いてモデル評価を行う。

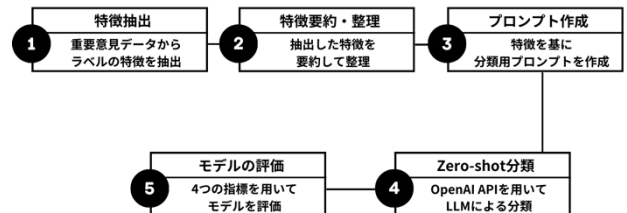


図 3: GPT-4o による Zero-Shot 分類手法

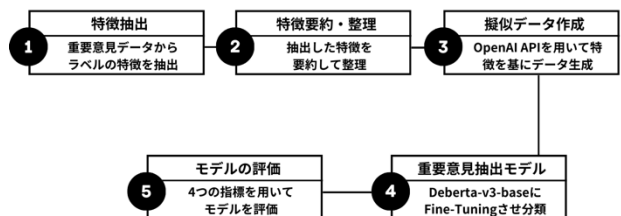


図 4: 擬似データ Fine-Tuning による分類手法

また、Deberta-v3-base-japanese は BERT [5] モデルを基盤としつつ Deverta-v3 [6] を日本語に特化させたモデルとなっている。

5 実験 1: 重要意見データセット構築

5.1 実験概要

本実験では、4つの観点からラベリングに関する個人ワークおよびグループワークを実施し、有識者意見と比較することで、学生によるデータセット構築が可能であるかを検証した。分析の先駆けとなるワークショップでは、図1の4つの観点から個人ワークとグループワークを行い、それぞれ5つのラベリングと、その観点から表2のラベル1と2に該当するデータに対してコメントを行うものである。また、分析に関しては図2のラベリング分析およびコメント分析に関する5つの分析を行うことで、学生による重要意見データの構築可否を検証した。

表 2: 4つの観点における5値ラベルの説明

	共感性	意外性
ラベル1	とても共感	とても意外性あり
ラベル2	やや共感	やや意外性
ラベル3	わからない	意外性はなし
ラベル4	やや共感しない	
ラベル4	全く共感しない	
	顧客視点_改善案	経営視点_改善案
ラベル1	とても改善できる	とても改善できる
ラベル2	やや改善できる	やや改善できる
ラベル3	わからない	わからない
ラベル4	やや改善できない	やや改善できない
ラベル5	全く改善できない	全く改善できない

また、4つの観点を2値ラベルに変換する際には、表2の「ラベル1・2」に該当するデータを「ラベル1」、「ラベル3~5」に該当データを「ラベル0」とする。

さらに、実験を行う上で学生14名をそれぞれ Student1~14 に置き換え、有識者のラベリングデータを正解データと仮定し、検証を実施した。

5.2 実験結果

5.2.1 ワークショップの結果

ワークショップでは4つの観点に対し各1時間の時間を設けラベリングを実施した結果、「ラベル1・2」と判定された、つまり票数が多く議論が活発になったのは「共感性>意外性」「経営視点での改善案>顧客視点での改善」となった。また、有識者のラ

ベリングも同様の結果となり、学生および有識者ラベリングでは「共感性」のラベリング数が最多となった。

5.2.2 グループおよび個人分析

本実験では14名の学生を4つのグループに分け、個人ワークを実施した後、グループワークを実施した。結果として、ワークショップの結果に付随しグループでの議論も共感性に対する議論が最多であったが、グループ意見として出力された結果は個人ワークで得られた意見の要約であり新たな意見が創出された事例はごく僅かとなった。そのため、共感性に着目し学生14名と有識者を含めた各個人同士の分析を実施した結果、学生の共感票数(2値ラベルにおけるラベル1と判定されるもの)が多いものは有識者の共感を得ていることが分かった。

5.2.3 網羅性の確認および特徴重要度分析

学生の共感判定が有識者の正解データをどの程度網羅できているか、また Student1~14 の中で、どの学生が網羅性に寄与しているのかを確認した。そこで、14名分の学生票数の閾値を半数である7件に設定し比較検証した結果と、Random Forest を用いての予測および Feature importance での Student の重要度分析の結果を以下に示す。

表 3: 網羅性と Random Forest での予測

	網羅性	Random Forest
Accuracy	70.9%	84.8%
Precision	0.696	0.895
Recall	0.571	0.820
F1-score	0.627	0.831

表 4: Feature importance の値 (上位4名抜粋)

	値
Student 12	0.880
Student 13	0.400
Student 4	0.390
Student 9	0.270

5.2.4 類似分析

本実験では、学生と有識者が共感すると判定(2値ラベルにおけるラベル1)したコメントに着目し、両方のコメントを、Sentence BERT を用いて類似度分析を実施した。以下の表5は有識者が共感判定且つコメントをした82件/196件に対して、類似度の

閾値を0.5以上と0.6以上に分けて分析した結果である。

表 5 : Sentence BERT を用いたコメント類似度

閾値	該当件数/82 件
0.5 以上	43 件
0.6 以上	27 件

5.3 考察

本実験では、ワークショップの結果から共感性のラベリングおよびコメントが多く存在し、重要な指標であることが示された。そこで共感性に着目し、5つの分析を実施した。表3より網羅性の確認およびRandom forestを用いた正解データ予測ではAccuracyにて70~85%の精度となり高い割合で学生のラベリングが正解データを判定できているといえる。また、学生の中でも表4よりStudent12は正解データに重要となることが示された。ワークショップの事前アンケート結果からStudent12はディズニー施設内でのキャスト経験があるため、ドメイン知識の豊富さによる特定のキーワードへの理解度の深さが影響していると考えられる。最後に、コメントの類似度分析においては、閾値が0.5以上のとき1/2以上、0.6以上では1/3の精度となっており、学生のラベリングが有識者の正解データ（重要意見データセット）の構築に寄与できることが示された。

6 実験 2 : LLM による重要意見抽出

6.1 実験概要

本実験ではGPT-4oによるZero-shot分類とGPT-4oにより作成した擬似データをFine-Tuningさせることで重要意見抽出モデルによる分類の2つの手法を検証した。2つの手法ではそれぞれOpenAIのAPIを用いており、GPT-4oによる分類では言語モデルに対して重要意見の特徴を含むプロンプトを与えることで、重要意見であるか否かを判定させた。この手法をBaseとする。また、Fine-Tuningによるモデル構築での分類では、言語モデルを使用することで1000件、2000件、3000件の擬似データを作成し、重要意見抽出タスクに特化させた形で重要意見判定を実施した。

3種類の件数でFine-tuningさせたモデルをそれぞれFT_1000, FT_2000, FT_3000として表記する。

6.2 実験結果

表 6 : 2 手法による重要意見抽出の検証結果

	Acc.	Pre.	Rec.	F1.
Base	69.90%	0.736	0.464	0.569
FT_1000	60.09%	0.603	0.255	0.331
FT_2000	59.69%	0.541	0.455	0.490
FT_3000	59.45%	0.556	0.358	0.431

6.3 考察

表6からFine-tuningさせたモデルは、学習件数を増やしてもBaseと比べて、4つの指標全てにおいて劣ることが分かった。そのため、擬似データのFine-Tuningによる重要意見抽出の有用性は検証できなかった。しかし、Baseによる分類では約7割の正解率となり、Precisionに関してはモデルがデータを重要意見であると予測したもののうち実際に重要意見である（正例）割合が高いことを示しており、予測の信頼性が高いといえる。また、Fine-Tuningモデルに着目するとFT_2000モデルでF1-scoreの値が高いことが分かる。しかし、データ数を増やしたFT_3000モデルでは全体的に値が下がっていることから、単にデータ数を増加させるだけでは精度が向上しないことが示された。これは、生成されたデータの中で類似データが多く含まれていることによる過学習への対応や適切なデータセットのサイズを選定する必要があると考える。

7 おわりに

本研究では、ワークショップの実施による重要意見データセットの構築およびワークショップから得られた特徴を用いて重要意見を抽出するシステムを提案した。実験1では、学生による重要意見ラベリングが正解データに対してAccuracyが70~85%の精度を達成することができ、コメントを含めた類似度では半数以上のデータを構築可能であることを検証できた。また、実験2では、重要意見指標を用いた2つの手法においてLLMによる分類にて正解率が70%の精度を達成できたが、擬似データ作成によるFine-Tuningでの分類では精度向上が見られなかった。これは、擬似データの質の低さが要因であると考えられる。そのため、今後は、LLMに入力するプロンプトの工夫やデータ増幅に対する過学習への対応などにより、モデルの分類精度向上を目指す。

謝辞

本研究は、ワークショップに参加してくださった産業能率大学経営学部櫻井ゼミの皆様のご協力により実現しました。心より感謝申し上げます。

本研究は JSPS 科研費 20K11960 の助成を受けたものです。

参考文献

- [1] 川島崇秀, 哲. 佐藤, 典. 神門, “Twitter からの消費者ニーズの抽出手法に関する提案,” 第 8 回データ工学と情報マネジメントに関するフォーラム, 2016.
- [2] 山本瑞樹, 孝. 乾, 大. 高村, 聡. 丸本, 裕. 大塚, 学. 奥村, “文章構造を考慮した自由解答意見からの要望抽出,” 言語処理学会 12 回年次大会併設ワークショップ「感情 評価 態度と言語」, 2006.
- [3] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child and A. Ramesh, "Language Models are Few-Shot Learners," 2020.
- [4] 藤井巧朗, 智. 勝又, “日本語タスクにおける LLM を用いた疑似学習データ生成の検討,” 言語処理学会 第 30 回年次大会 発表論文集, 2024.
- [5] J. Devlin, Ming-Wei Chang, Kenton Lee, Kristina Toutanova, “BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding,” 2019.
- [6] P. He, Jianfeng Gao, Weizhu Chen, “DEBERTAV3: Improving DeBERTa using ELECTRA-Style Pre-Training with Gradient-Disentangled Embedding Sharing,” 2023.