

卒業論文

クラウドワークの模倣モデルと投票作業の一致率を 考慮した結果集約手法

太田 奈那

2023年2月8日

岐阜大学 工学部 電気電子・情報工学科 情報コース
鈴木研究室

本論文は岐阜大学工学部に
学士（工学）授与の要件として提出した卒業論文である。

太田 奈那

指導教員：

鈴木 優 准教授

クラウドワーカーの模倣モデルと投票作業の一致率を 考慮した結果集約手法*

太田 奈那

内容梗概

本研究では、クラウドワーカーの模倣モデルによる予測結果と複数人の作業に対する多数決にて付与された評価ラベルとの一致率を用いて、クラウドソーシングの質を向上させるような結果集約手法を提案する。クラウドソーシングは、インターネットを介してワーカーを募集するためスパムワーカーが一定数存在する。スパムワーカーが行った作業結果も含めて集約を行うと、作業依頼者の求める結果が得られないことがある。そこで、全データを集約した全体の結果と個人の作業結果との一致率が高いほど品質の高いワーカーであるとし、品質の低いワーカーの影響力を抑える結果集約手法を提案する。本手法では、ワーカーの模倣モデルによって得られた結果をもとに算出した一致率の値を用いて、ワーカーごとに票の重みを設定して多数決を行う。ワーカーが実際に行った作業データを用いて訓練済み BERT モデルをファインチューニングすることにより、ワーカーの模倣モデルを構築する。模倣モデルを構築することにより、一つのデータに対するワーカーを擬似的に増やすことができる。そのため、時間や費用を抑えた上で多くの意見を集めることが可能になると考えた。また、票の重みをワーカーごとに変えることによって品質の低いワーカーの影響力を抑えることができ、品質の高いワーカーの作業結果が反映されやすくなると考えた。それにより、複数人の作業に対する多数決にて付与された評価ラベルと比較して、手法を適用させることによって付与された評価ラベルの方が作業依頼者の求めるものに近い作業結果を得ることができると考えた。そこで、本手法の有効性を確かめるために実験を行った。その結果、品質が低いワーカーによる作業の重みが小さくなり、評価ラベルを付与する際に及ぼす影響力が抑えられていた。実験を通して、本手法を適用することによって品質の低いワーカーの影響力を抑え、作業依頼者の求め

*岐阜大学 工学部 電気電子・情報工学科 情報コース 卒業論文, 学籍番号: 1193033029, 2023 年 2 月 8 日.

る結果に近い作業結果が得られることを確認した.

キーワード

機械学習, クラウドソーシング, 一致率

目次

図目次	v	
表目次	vi	
第 1 章	はじめに	1
第 2 章	基本的事項	4
2.1	クラウドソーシング	4
2.2	ファインチューニング	4
2.3	BERT	4
2.4	結果集約手法	5
2.5	評価指標	5
第 3 章	関連研究	7
第 4 章	提案手法	9
4.1	分類器の構築	9
4.2	結果集約手法	10
4.2.1	一致率によるワーカーの選定	11
4.2.2	一致率による票の重み設定	13
4.2.3	一致率による無効票の決定	15
第 5 章	評価実験	17
5.1	データセット	17
5.2	実験 1：一致率によるワーカーの選定	19
5.2.1	実験手順	19
5.2.2	結果・考察	20
5.3	実験 2：一致率による票の重み設定	22
5.3.1	実験手順	22
5.3.2	結果・考察	23

5.4	実験 3：一致率による無効票の決定	25
5.4.1	実験手順	25
5.4.2	結果・考察	26
第 6 章	おわりに	29
	謝辞	31
	参考文献	32
	発表リスト	34

図目次

4.1	一致率 A_i によるワーカの選定の概要	11
4.2	8 クラス分類で最大の閾値を 0.35 としたときの閾値の推移	13
4.3	一致率 A_i による票の重み設定の概要	14
4.4	一致率 A_i による無効票の決定の概要	15

表目次

2.1	評価値導出の混同行列	6
5.1	評価ラベルの内容と正解ラベルの数	17
5.2	ワーカを選定して多数決を取り直した時の <i>Accuracy</i>	20
5.3	複数人の作業に対する多数決にて付与された評価ラベルと実験 1 で付与された評価ラベルに変化があったツイート例	21
5.4	票の重みを設定して多数決を取り直した時の <i>Accuracy</i>	23
5.5	複数人の作業に対する多数決にて付与された評価ラベルと実験 2 で付与された評価ラベルに変化があったツイート例	24
5.6	一致率を確率として多数決を取り直した時の <i>Accuracy</i>	26
5.7	作業数が多い上位 5 人の各模倣モデルによる評価予測の一致率	27
5.8	複数人の作業に対する多数決にて付与された評価ラベルと実験 3 によって付与された評価ラベルに変化があったツイート例	28

第1章 はじめに

クラウドソーシング [1] とは、インターネット上で募集した不特定多数の人に作業を代わりに行ってもらうことである。クラウドソーシングを行うことによって、複雑な作業や一人で行うには困難な作業を複数の人に分担して行ってもらうことができる。そのため、一人ひとりの作業は単純になり、一人にかかる負担を減らすことができる。また、インターネット上で作業が行われるため、誰でも簡単に好きな時間に好きな場所で作業を行うことができるという利点がある。

しかし、不特定多数の人に作業を依頼しているため意見のばらつきが見られることもある。インターネットを介してワーカを募集しているため、スパムワーカと呼ばれる品質の低いワーカが一定数存在する。スパムワーカとは、数をこなすためだけに故意に不適切な作業をするワーカである。スパムワーカが存在する状態で結果集約を行うとスパムワーカの影響を受けてしまい、作業依頼者の求める結果が得られないことがある。スパムワーカは他のワーカと違う結果を出すため、スパムワーカを取り除くことができるのではないかと考えられる。そのためには、作業に対する正しい評価が必要になる。そして、その評価は作業依頼者が求める結果に近い作業結果となっていなければならない。

作業依頼者が求める結果を得るためには、より多くのワーカによる意見が必要であると考えた。しかし、より多くのワーカの意見を集めるためには時間や費用がかかってしまう。そこで我々は、ワーカの作業結果をもとにワーカの作業を模倣するような分類器を構築することにより、擬似的にワーカを増やすことができるのではないかと考えた。擬似的にワーカを増やすことができれば、時間や費用を抑えた上で多くのワーカの意見を集めることが可能になる。そして、その意見を集約することによって作業依頼者の求める結果を得ることができると考えられる。

模倣モデルを構築するワーカの中には、先ほど述べたスパムワーカが存在している可能性がある。スパムワーカの作業結果を含めて結果集約を行った場合、作業依頼者の求める結果を得ることは難しい。そこで、ワーカの品質を考慮した結果集約を行うことによって、スパムワーカの影響を抑えることができるのではないかと考えた。ワーカ個人が付与した評価ラベルと複数人の作業に対する多数決にて付与された評価ラベルとの一致率をもとにワーカの品質を求める。その品質を用いて結

果集約を行い、作業依頼者の求める結果を得ることを目指す。

本研究では、クラウドソーシングの質を向上させるための異なる三つの結果集約手法を提案する。一つ目は、一致率をもとにワーカの選定を行って多数決をとる手法である。品質の低いワーカを取り除いて品質の高いワーカの作業のみを使用することによって、作業依頼者の求める結果が得やすくなると考えた。二つ目は、一致率をもとにワーカの品質を求め、その品質を用いて各ワーカの票に重み付けして多数決をとる手法である。票の重みをワーカごとに変えることによって、品質の低いワーカの影響力を抑えることができると考えた。三つ目は、一致率の値をワーカが一つのデータに対して投票できる確率として多数決をとる手法である。品質の低いワーカが行った作業の中にも適切な評価ラベルが付与されたデータも存在するため、全ての作業結果を取り除くと有用なデータが失われてしまう。そのため、集約時にデータを使用するかどうかを確率的に決定することによって、有用なデータが使用される可能性があると考えた。

これら三つの手法で共通している点は、ワーカの一一致率をもとにした多数決を行うことである。そのため、一致率が高いワーカの評価は多数決の結果に反映されやすく、低いワーカの評価は多数決の結果に反映されにくくなる。これにより、模倣モデルを構築したワーカの中にスパムワーカが含まれていたとしても、その影響は少なくなる。そして、作業依頼者の求める結果に近い作業結果を得られる可能性が高くなるのではないかと考えられる。

提案手法がクラウドソーシングの質を向上させるために有効であるのかどうかを確かめるために評価実験を行った。先ほど述べた三つの手法を用いて付与される評価ラベルに加えて、複数人の作業に対する多数決にて付与された評価ラベルが、作業依頼者の求める評価ラベルとどの程度一致するかを比較することにより手法の有効性を確かめた。

三つの手法を用いてそれぞれ実験を行ったところ、二つ目の手法を用いた場合に作業依頼者が求める結果と一致したデータが最も多いという結果が得られた。一つ目の手法や三つ目の手法を用いると、品質の低いワーカの作業データを取り除くことは可能であるが、そのワーカの有用な作業データまで取り除かれてしまうことがある。また、三つ目の手法では品質の高いワーカの作業結果であっても取り除かれてしまうことがあり、正しい評価ラベルを得ることが難しいと考えられる。一方、

二つ目の手法では品質の低いワーカーであっても作業結果が残るため、有用な作業データが取り除かれることもなく、評価ラベルを決める際に少なからず影響を与えることが可能である。そのため、品質が低いワーカーの影響を抑えつつ、正しい評価ラベルを得ることが可能であったのではないかと考えられる。

本論文における貢献は以下のとおりである。

- ワーカーの模倣モデルを構築してそのモデルの精度を確かめた。
- ワーカーの模倣モデルによる評価ラベルと投票作業から得た評価ラベルの一致率を使用したクラウドソーシングの質を向上させることができた。

本論文の構成は以下の通りである。2章では本論文にて用いた技術や手法の基本的事項について述べる。3章では関連研究について述べる。4章では本論文の提案手法について述べる。5章では提案手法を用いた評価実験の目的や手順、結果・考察などについて述べる。最後に6章では本論文のまとめと今後の課題について述べる。

第 2 章 基本的事項

2.1 クラウドソーシング

クラウドソーシングとは、一人で行うには困難なタスクやその一部をインターネットを介して募集した不特定多数の人に委託することである。クラウドソーシングを利用して複数の人に作業を分割して割り振ることにより、一人ひとりの作業は単純なものになり、一人にかかる負担を減らすことができる。また、クラウドソーシングを利用することによって自分が苦手な作業を得意な人に委託することができ、質の高い作業結果を得ることができる。さらに、インターネット上で作業を行うことができるため誰でも簡単に利用でき、好きな時間に好きな場所で作業を行うことができるという利点がある。実際に活用されている例として、ホームページやバナーのデザイン作成、アンケートなどの発注が挙げられる。

2.2 ファインチューニング

ファインチューニングとは、訓練済みモデルをベースに出力層などを変更したモデルを構築し、自身で用意したデータを使用してモデル全体のパラメータを学習させる手法である。ファインチューニングでは、訓練済みモデルのパラメータを初期値として利用し、自身で用意したデータや目的に合うようにパラメータを再学習させる。ファインチューニングは訓練済みモデルをもとに学習を行うため、自身が用意したデータが少量の場合でも精度の高いモデルを構築することができる。また、一から学習するよりも短時間でモデルを構築することができる。

2.3 BERT

BERT[2]とは、自然言語処理のための Transformer[3] モデルをベースとしたニューラルネットワークモデルのことである。事前学習として MLM(Masked Language Model) と NSP(Next Sentence Prediction) を行っている。MLM は入力の一部の単語を隠して元の単語を予測するタスクである。NSP は長い文章の中

から選んだ2文が連続しているかどうかを予測するタスクである。事前学習した後のモデルをファインチューニングすることにより、様々な自然言語処理タスクに応用することができる。

2.4 結果集約手法

結果集約とはアンケートなどの回答を集計することであり、結果集約手法は集計をするための方法を指す。集計方法には、単純な加算で計算できるため全体像の把握に適している「単純集計」や、単純集計よりも詳細な分析ができるため属性ごとの傾向を知るのに向いている「クロス集計」が存在する。クロス集計は、詳細なアンケート分析によって細かい傾向を見つけ出すことができるため、アンケートの回答を集計する主な方法として使用される。また、回答が自由記述で行われたデータを集計する「自由記述集計」もあり、類似単語などを見つけてカテゴライズしたり関連性を見つけて集計したりする方法である。さらに、多数決をとることによって一つのタスクに対してただ一つの評価を与えるといった集約方法もある。本研究で使用するデータセットを作成する際は、多数決をとる方法を用いた。

2.5 評価指標

本研究では、手法の評価を行うために *Accuracy* を用いる。表 2.1 は評価用データの予測結果をまとめた混同行列である。正例とは、ある問題に対して正クラスである事例のことであり、負例とは、ある問題に対して負クラスである事例のことである。例えば、患者が癌であるかどうかを判定する問題においては、患者が癌であったら正クラス、癌ではなかったら負クラスとなる。*Accuracy* は、表 2.1 の混同行列をもとに式 (2.5.1) を用いて導出する。この数値は、単純な正解率を示しており、予測結果がどれだけ真の値と同じ結果であるかどうかを測ることができる。

$$Accuracy = \frac{TP + TN}{TP + FP + FN + TN} \quad (2.5.1)$$

表 2.1 評価値導出の混同行列

		予測値	
		正例ラベル	負例ラベル
真の値	正例ラベル	TP	FN
	負例ラベル	FP	TN

TP 正例と予測されたデータのうち、実際のラベルも正例であるデータの数

FP 正例と予測されたデータのうち、実際のラベルは負例であるデータの数

FN 負例と予測されたデータのうち、実際のラベルは正例であるデータの数

TN 負例と予測されたデータのうち、実際のラベルも負例であるデータの数

第3章 関連研究

クラウドソーシングの質の向上に関する研究はいくつか存在する。西ら [4] の研究では、ソーシャルネットワークを用いたワーカーの品質向上を目指している。作業を行うワーカーは他一人に作業を委託することができ、報酬は作業に正解したワーカーとそのワーカーに作業を委託したワーカーに支払われる。こうすることによって、能力の低いワーカーは能力の高いワーカーに作業を委託するようになり、品質の高いワーカーの作業結果を得られる。芦川ら [5][6] の研究では、ワーカーに作業の適性があるかどうかのフィルタリングを行うことによって、クラウドソーシングの質の向上を目指している。ワーカーの作業前、作業中、作業後に加えて、得られた結果を用いて推測された未知データの結果精度を用いてフィルタリングを行う。フィルタリングを行ってワーカーを絞ることにより、クラウドソーシングの精度を向上させている。Halpin ら [7] の研究では、スパムワーカーの検出手法を提案している。ワーカーごとに作業数や一つの作業を行うのにかかった平均時間などの特徴を用い、機械学習を行うことによってスパムワーカーの検出を行っている。松原ら [8] の研究では、ワーカーに適した作業の割り当てを行っている。提案されている手法では、複数のタスクをワーカーに割り当てる場合において、まず各ワーカーに希望するタスクの優先順位をつけさせる。各ワーカーがつけたタスクの優先順位をもとに各ワーカーに対して作業を割り当て、虚偽の順位をつけたワーカーに対して不利益が生じるように設定する。そのため、ワーカーは真実の優先順位をつけることとなり、ワーカーごとに適した作業を割り当てることによって精度の高い結果を得ている。上記の研究ではワーカーの品質に着目することによってクラウドソーシングの質の向上を目指している。

本研究では、ワーカーごとの模倣モデルを構築し、模倣モデルを用いて作業データを増やすことによってクラウドソーシングの質の向上を目指す。一つのデータに対する作業結果を増やすことによって多くの意見を得ることができ、作業依頼者の求める結果に近い作業結果を得ることができるのではないかと考えた。そのため、上記の研究とはクラウドソーシングの質の向上を目指す点では同じである。一方、質を向上させるためのアプローチとして作業データを増やすという点で異なっている。

また、結果集約手法に関する研究もいくつか存在する。Dawid ら [9] の研究では、

EM アルゴリズムを用いたラベル付与の手法について提案している。ワーカが各ラベルを回答したときの正解率を EM アルゴリズムを用いて推定する。推定して得られた正解率が最も高いラベルを真のラベルとしてデータに付与するという手法である。小山ら [10] の研究では、高精度なラベル統合方法について提案している。行った作業の処理結果をどの程度確信しているのかをワーカに申告してもらい、その確信度をもとにワーカの作業結果がラベルを付与する際に必要な情報であるかどうかを確率的に判断している。また、自己申告した確信度はワーカの正解率と相関があると考えており、確信度を用いることによって高い精度で適切なラベルを付与している。

本研究では、クラウドワーカの模倣モデルを構築することによって求めた一致率を用いてワーカごとに票の重みを設定したり、一致率の値をワーカが一つのデータに対して投票できる確率としたりして多数決をとり、ラベルを付与する。そのため、結果集約の質を高めるために品質の低いワーカの作業結果が集約結果に現れにくくなるという点で上記の研究と異なっている。また、我々が一致率を求めることによってワーカが自己申告する手間がなくなるため、ワーカへの負担が少ないという利点もある。

第4章 提案手法

本研究は、クラウドワーカの模倣モデルによる予測結果と複数人の作業に対する多数決にて付与された評価ラベルとの一致率 A_i を用いた結果集約手法によって、クラウドソーシングの質を向上させることを目的とする。模倣モデルによる予測結果と複数人の作業に対する多数決にて付与された評価ラベルとの一致率 A_i が高いほど品質の高いワーカであると考えられる。そのため、一致率 A_i を用いた結果集約手法を用いることによってより良いワーカの結果を反映させることができ、作業依頼者の求める結果に近い作業結果が得られるのではないかと考えた。

そこで本研究では、以下の異なる三つの結果集約手法によってクラウドソーシングの質を向上させることを目指す。

- (1) 一致率 A_i をもとにワーカの選定をして多数決をとる方法
- (2) 一致率 A_i をもとにワーカごとに票の重みを設定して多数決をとる方法
- (3) 一致率 A_i の値を一つのデータに対して投票できる確率として多数決をとる方法

これら三つの結果集約手法については 4.2.1 項から 4.2.3 項で詳しい説明を行う。

4.1 分類器の構築

本研究では、各クラウドワーカの作業を模倣するような分類器を構築する。クラウドワーカの模倣モデルを構築する際に、東北大学の乾・鈴木研究室で構築された訓練済み日本語 BERT モデル*を使用して、テキストデータに対して評価ラベルを予測するような分類器を構築する。このとき、BERT モデルの最終層のパラメータのみを更新するように設定し、ファインチューニングを行うことによって模倣モデルを構築する。ここで、以下の異なる二種類の方法を用いて BERT モデルのファインチューニングを行う。

Model 1 ワーカ個人の作業結果を集めたデータセットを使用して、BERT モデル

*<https://github.com/cl-tohoku/bert-japanese>

をファインチューニングする。

- Model 2 (i) 複数人の作業に対する多数決にて評価ラベルを付与したデータセットを使用して，BERT モデルをファインチューニングする。
- (ii) ワーカー個人の作業結果を集めたデータセットを使用して，(i) で構築したモデルをファインチューニングする。

Model1 では，ワーカーのデータ数が少なかった場合に精度の高い分類器を構築することは難しい。そこで，Model2 を用いることによってデータ数の少なさを補うことができ，どのようなワーカーも精度の高い分類器を構築することができるのではないかと考えた。そのため，上記に示した二種類の方法を用いて各クラウドワーカーの模倣モデルを構築する。

4.2 結果集約手法

本節では，異なる三つの結果集約手法について説明する。三つの結果集約手法に共通している，模倣モデルによる予測結果と投票作業による作業結果の一致率 A_i を算出する手順について説明する。

クラウドワーカーの集合を W として各クラウドワーカーを $w_i \in W$, $i = 1, 2, \dots, n$ とし，各クラウドワーカーが実際に作業を行ったデータを集めたデータセットを D_i とする。ここで n はクラウドワーカーの人数とする。そして，複数人の作業に対する多数決にて評価ラベルを付与したデータセットを D_{all} とする。また，評価ラベルを付与するテキストデータを d_j , $j = 1, 2, \dots, N$ とする。ここで N は評価ラベルを付与するデータの総数とする。複数人の作業に対する多数決にて評価ラベルを付与したデータセット D_{all} の中から，ワーカー w_i が実際に作業していないデータ d_j を抽出してデータセット $D'_i = \{d_j | d_j \in D_{all}, d_j \notin D_i\}$ を作成する。構築した模倣モデルにてデータセット D'_i に含まれるデータ d_j の評価ラベルを予測することによって，データ d_j に対して評価ラベルを付与する。このとき，ワーカー w_i の模倣モデルによって付与されたデータ d_j の評価ラベルを $l_i(d_j)$ と表す。評価ラベル $l_i(d_j)$ と複数人の作業に対する多数決にて付与された評価ラベル $l_{all}(d_j)$ を比較して，一致

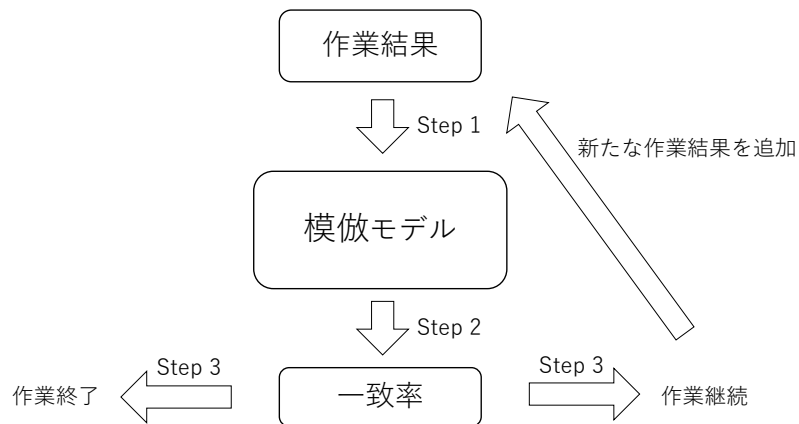


図 4.1 一致率 A_i によるワーカーの選定の概要

率 A_i を求める．ワーカー w_i の一致率 A_i は以下のように算出する．

$$A_i = \frac{\sum_{j=1}^N M(l_i(d_j), l_{all}(d_j))}{N} \quad (4.2.1)$$

$$M(l_i(d_j), l_{all}(d_j)) = \begin{cases} 1 & (l_i(d_j) = l_{all}(d_j)) \\ 0 & (l_i(d_j) \neq l_{all}(d_j)) \end{cases} \quad (4.2.2)$$

式 (4.2.2) は，ワーカー w_i の模倣モデルによって付与された評価ラベルと複数人の作業に対する多数決にて付与された評価ラベルを比較したとき，一致していたら 1，異なっていたら 0 を出力する関数である．式 (4.2.1) を用いて算出された一致率 A_i を使用して三つの結果集約手法で多数決をとる作業を行う．

4.2.1 一致率によるワーカーの選定

構築した模倣モデルによって得られた一致率 A_i をもとにワーカーを選定し，多数決をとる手法について説明する．この手法では逐次的にワーカーの作業の質を確認して，品質の低いワーカーを除去する．手法の概要は図 4.1 に示す通りである．

Step 1 ワーカー w_i が行った作業結果を訓練データとして模倣モデルを作成する．

Step 2 作成した模倣モデルを用いてワーカー w_i が作業を行っていないデータの評価ラベルを予測し、複数人の作業に対する多数決にて付与された評価ラベルとの一致率 A_i を算出する。

Step 3 算出した一致率 A_i が閾値を上回っていたら作業を継続し、下回っていたらその時点で作業を終了する。

Step3 で算出した一致率 A_i が閾値を上回っていた場合は作業を継続し、新たな作業結果を追加して Step1 から Step3 の手順を繰り返し行う。ワーカーの作業の質を逐次的に確認する作業を繰り返し行うのは、最初は真面目に作業を行っていても途中から適当にこなしていたり、判定することが簡単なデータが作業序盤に割り振られていることにより作業の質が一時的に良くなっていたりする可能性があるためである。最終的に一致率 A_i が閾値を下回らなかったワーカーの作業結果を使用して多数決をとる。ここで、Step1~3 の手順を繰り返し行っている回数を m 、 k クラス分類のチャンスレートを C 、閾値の最大値を θ_{max} として閾値 θ_m を以下の式を用いて定める。このとき、最大の閾値 θ_{max} は $C < \theta_{max} \leq 1$ の範囲で値をとる。

$$\theta_m = 2 \times (\theta_{max} - C) \times \text{sigmoid}\left(\frac{m}{5}\right) - \theta_{max} + 2 \times C \quad (4.2.3)$$

$$\text{sigmoid}(x) = \frac{1}{1 + \exp(-x)} \quad (4.2.4)$$

$$C = \frac{1}{k} \quad (4.2.5)$$

式 (4.2.4) はどのような入力値も 0.0~1.0 の範囲の数値に変換して出力することができるシグモイド関数である。式 (4.2.3) を用いて閾値を定めることによって、入力値が正の数のおきにチャンスレートから閾値の最大値までの範囲の中で数値を出力することができる。閾値の最大値とは、手順を繰り返し行うごとに増加する θ_m が最終的に収束する値のことである。この値を自身で定めることにより、使用するデータセットの難易度に合わせた閾値を定めることができるようになる。また、下限としてチャンスレートをを用いることにより、無作為に選んだ場合よりも一致率 A_i の低い分類器を取り除き、ワーカーの選定をすることができる。図 4.2 では閾値の推移の例を表している。この図は、8 クラス分類の分類器を構築する際に最大の閾値を 0.35 としたときの閾値の推移を表したものである。

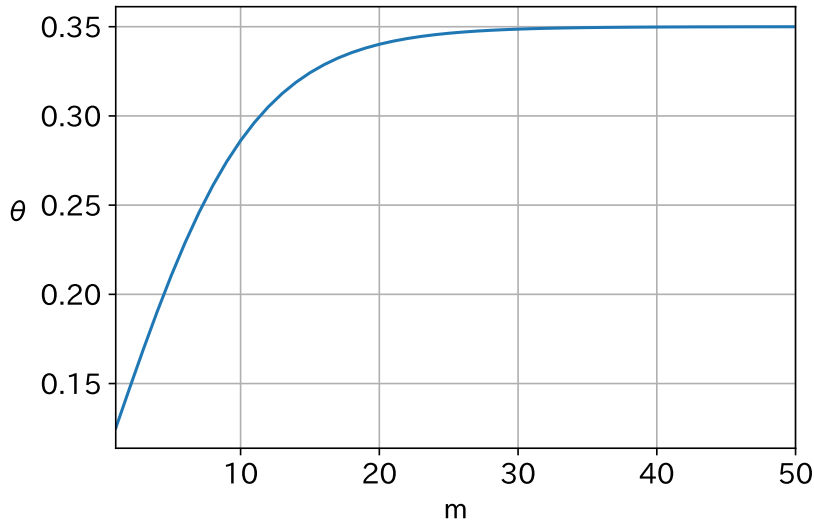


図 4.2 8 クラス分類で最大の閾値を 0.35 としたときの閾値の推移

4.2.2 一致率による票の重み設定

クラウドワーカーの模倣モデルによって得られた一致率 A_i をもとに、ワーカーごとに票の重みを設定することによって多数決をとる手法について説明する。この手法では、品質の低いワーカーの作業結果を集約結果に反映されにくくするために、多数決をとる際にワーカーごとに票の重みを設定して集約をする。手法の概要は図 4.3 に示す通りである。ワーカー w_i の一致率 A_i を算出する流れについては StepA, 算出した一致率 A_i をもとに票の重みを設定して多数決をとる流れは StepB の手順で行う。

Step A-1 データセット D_i から無作為に抽出したデータを訓練データとして学習を行い、ワーカー w_i の模倣モデルを構築する。

A-2 模倣モデルを用いてデータセット D'_i に含まれるテキストデータ d_j の評価ラベルを予測し、複数人の作業に対する多数決にて付与された評価ラベルとの一致率 A_i を算出する。

Step B-1 A-2 で算出した一致率 A_i をもとに票の重みを設定し、作業結果に対して重み付けを行う。

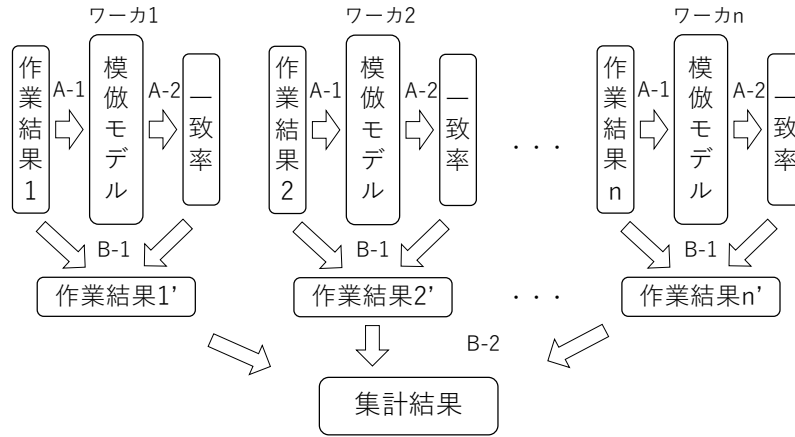


図 4.3 一致率 A_i による票の重み設定の概要

B-2 同じデータに対して評価を行ったワーカーのデータを集約し、結果を出す。

B-1 で行われるワーカー w_i の票の重み設定は以下の手順で行う。ワーカー w_i の品質を q_i として、最大値 1, 最小値 0 となるように一致率 A_i を正規化する。正規化をするときに k クラス分類のチャンスレート C を含めることによって、全ワーカーの一致率 A_i がチャンスレート C を上回っていたとき、全ワーカーが投票権を得ることができると考えた。ワーカーの品質 q_i は以下の式を用いて求める。

$$q_i = \frac{A_i - x_{min}}{x_{max} - x_{min}} \quad (4.2.6)$$

$$Q = \{x | x = C, A_i (i = 1, 2, \dots, n)\} \quad (4.2.7)$$

式 (4.2.7) は各ワーカーの一致率とチャンスレートの集合である。このようにして得たワーカーの品質 q_i をワーカー w_i の票の重みとする。

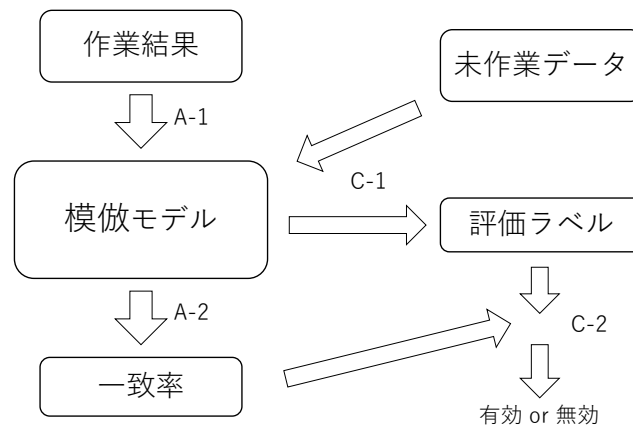


図 4.4 一致率 A_i による無効票の決定の概要

4.2.3 一致率による無効票の決定

クラウドワーカーの模倣モデルによって得られた一致率 A_i の値をもとに無効とする票を決めて多数決をとる手法について説明する。この手法では、品質の低いワーカーの作業結果を集約結果に反映されにくくするために、作業集約時に無効票を確率的に決定する。手法の概要は図 4.4 に示す通りである。ワーカー w_i の一致率 A_i を算出する流れについては StepA、算出した一致率 A_i をもとに無効票を決める流れは StepC の手順で行う。

Step A-1 データセット D_i から無作為に抽出したデータを訓練データとして学習を行い、ワーカー w_i の模倣モデルを構築する。

A-2 模倣モデルを用いてデータセット D'_i に含まれるテキストデータ d_j の評価ラベルを予測し、複数人の作業に対する多数決にて付与された評価ラベルとの一致率 A_i を算出する。

Step C-1 模倣モデルを用いてデータセット D'_i に含まれるデータ d_j の評価ラベルを予測する。

C-2 A-2 で算出した一致率 A_i の値をもとにワーカー w_i の票が有効か無効を決定する。

C-2では、A-2で算出した一致率 A_i の値をワーカ w_i のデータ d_j に対する票を有効とする確率としている。こうすることによって、ワーカごとにデータ d_j に対して票が有効であるか無効であるかを決定することができる。また、一致率 A_i が高いほど投票できる確率は高くなり、品質の高いワーカの作業が反映されやすいと考えられる。このようにして無効票を決定することによって多数決をとる。

第5章 評価実験

本研究では，異なる三つの結果集約手法を提案した．そこで，それらの手法が有効であるかどうかを確かめるための実験を行った．

5.1 データセット

本研究では，クラウドソーシングによって構築された岐阜大学鈴木研究室のデータセットを使用する．このデータセットは605人のクラウドワーカーによって作業されたデータで構築され，全部で250,354件のデータが含まれている．このデータはワーカーID，ツイートID，ツイート内容，評価ラベルの四つのカラムで構成されている．またこのデータセットは，ツイート内に含まれる「笑」がポジティブな意味を持つのかネガティブな意味を持つのか，ネガティブの中でもどのような場面で使用されているのかを知ることを目的に作成されたものである．評価ラベルは表5.1に示した8種類のラベルを使用している．データセットを構築する際，評価ラベルは以下の手順で付与される．

- (1) 「笑」を含むテキストデータに対してポジティブ，ネガティブ，ニュートラル，ポジ+ネガ，その他の中から一つの評価ラベルを付与する．

表 5.1 評価ラベルの内容と正解ラベルの数

ラベル番号	ラベル内容	データ数 (件)
0	ポジティブ	8,955
1	ニュートラル	12,801
2	ポジ + ネガ	1,188
3	その他	706
4	攻撃性あり	1,382
5	攻撃性なし	3,116
6	自虐	1,524
7	判断できない	39

(2) (1) でネガティブが付与されたテキストデータに対して攻撃性あり，攻撃性なし，自虐，判断できないの中から一つの評価ラベルを付与する．

このように段階を踏むことによって，ネガティブの評価ラベルが付与された場合にさらに細かい評価ラベルを付与した．以上で述べた作業は，一つのデータに対して5人のワーカによって行われることが想定されている．

上記のデータセットには一つのデータに対して複数のワーカの作業結果が存在しているため，多数決をとることによって評価ラベルを一つに決定する．データセットの中には作業人数が5人に満たないデータが存在するため，5人の作業結果が残っているデータのみを使用して多数決をとることにする．複数人の作業に対する多数決にて評価ラベルを付与したデータセットの作成方法について説明する．クラウドソーシングによって構築されたデータセットの中から，作業人数が5人のテキストデータを抽出する．この抽出したテキストデータに評価ラベルを付与する際，上記で述べた手順(1)，(2)に倣い5人のワーカの評価を使用して多数決をとる．多数決によってただ一つの評価ラベルに決まったデータを集め，データセットを作成する．投票作業によってただ一つの評価ラベルが付与されたデータは全部で29,711件である．

5.2節から5.4節の実験では，投票作業によって得られたデータセット D_{all} と，実際に作業を行ったデータが500件以上あるクラウドワーカのデータセット D_i を使用する．ここで作業データが500件以上あるクラウドワーカを対象にした理由は，データ数が少ないとワーカの特徴を捉えたモデルを構築することが難しいと考えたためである．このクラウドワーカのデータセット D_i は，クラウドソーシングによって得られたデータの中から，ワーカIDを用いてデータを抽出することによって作成する．

作業データが500件以上あるクラウドワーカは127人であったため $n = 127$ とし，これらのデータセット $D_i (i = 1, 2, \dots, 127)$ を使用して三つの結果集約手法の有効性を確かめるための実験を行う．また，多数決をとるデータとして，ワーカ $w_i (i = 1, 2, \dots, 127)$ 以外のワーカの作業データを取り除いた後に残ったデータのうち，5人のワーカ w_i によって評価ラベルが付与されているデータを使用する．このデータの中から無作為に抽出した2,000件のデータを集めたデータセット D_{eva} を使用して多数決をとり，それぞれの手法の評価を行う．

手法の評価を行うにあたって、データセット D_{eva} に含まれるデータ d_j に対してあらかじめ著者が正解ラベルを付与した。この正解ラベルが付与されたデータセットを D_{true} とする。この付与された正解ラベルと、それぞれの実験で付与された評価ラベルや複数人の作業に対する多数決にて付与された評価ラベルをそれぞれ比較して *Accuracy* を算出する。ここで得られた *Accuracy* をもとに実験の結果を比較していく。

5.2 実験 1：一致率によるワーカの選定

本節では、4.2.1 項で説明した、一致率をもとにワーカを選定することによって多数決をとる手法の有効性を確かめるための実験について述べる。

5.2.1 実験手順

本実験では、5.1 節で説明したクラウドワーカ w_i のデータセット D_i を使用して模倣モデルを構築する。このとき、データに評価ラベルを付与した手順と同様の手順で評価ラベルの予測を行うため、二つの模倣モデルを構築する。一つ目は、ネガティブかポジティブかなどを判定するモデルである。二つ目は、ネガティブの中でもさらに細かい評価ラベルを判定するモデルである。この二つのモデルを組み合わせることによってデータ d_j に対する評価ラベルを予測する。また、模倣モデルは 4.1 節で説明した Model1 を構築する。なお、4.2.1 項で説明した手法は、逐次的にワーカの作業の質を確認する手法であるため、データセットの中から任意の数のデータを一つのまとまりとして実験を行った。

ワーカ w_i のデータセット D_i から抽出するデータを 200 件とし、このデータを訓練データとして模倣モデルを構築する。訓練データは、各評価ラベルのデータ数が等しくなるように、データ数が最も多い評価ラベルのデータ数に合わせてデータの複製を行う。模倣モデルを構築するとき、300 エポックの学習を行う。また、データセット D'_i から抽出するデータの数を 6,000 件とし、このデータセットを D''_i とする。構築した模倣モデルを使用してデータセット D''_i の評価ラベルを予測する。予測によって付与された評価ラベルと複数人の作業に対する多数決にて付与された

評価ラベルとの一致率を式 (4.2.1) を用いて算出する。算出した一致率 A_i が定めた閾値を上回っているかどうかを判定する。このときの閾値は、 $k = 8, \theta_{max} = 0.35$ として式 (4.2.3) に従って定める。一致率 A_i が閾値を上回っていた場合、訓練データを 200 件増やして、模倣モデルを構築する作業から一致率の判定をする作業までの流れを再び行う。一致率 A_i が閾値を下回る、あるいは、データセット D_i に含まれるすべてのデータを使用して模倣モデルを構築した場合は作業を終える。一致率 A_i が閾値を下回ったワーカーの作業データは取り除き、取り除かれなかったワーカーの作業データを使用してデータセット D_{eva} に含まれるデータの評価ラベルを多数決にて付与し、*Accuracy* を算出する。

5.2.2 結果・考察

Model1 を用いてワーカーの品質を測定することによってワーカーの選定を行い、多数決をとった結果を表 5.2 に示す。複数人の作業に対する多数決にて集約された結果の *Accuracy* と比較すると、実験 1 の手法を用いて算出した *Accuracy* の方が低くなっていることがわかる。

Step1~3 を繰り返し行うことによって取り除かれたワーカーは 127 人中 4 人のみであった。取り除かれたワーカーが実際に行った作業結果と複数人の作業に対する多数決にて付与された評価ラベルとの一致率を確認したところ、他のワーカーと比べて少し低い値になっていた。しかし、取り除かれたワーカーよりも一致率の低いワーカーが存在していた。そのため、必ずしもワーカーが実際に行った作業結果と複数人の作業に対する多数決にて付与された評価ラベルとの一致率の低いワーカーが取り除かれるわけではないと考えられる。また、訓練データとして使用するデータ数が多くなるほど閾値が大きくなるため、作業データ数の多いワーカーの方が取り除かれやすい

表 5.2 ワーカーを選定して多数決を取り直した時の *Accuracy*

	<i>Accuracy</i>
all	0.6625
Model1	0.5845

傾向にあると考えられる。

表 5.3 は、複数人の作業に対する多数決にて付与された評価ラベルと手法を適用することによって付与された評価ラベルに変化があったツイートの例である。一つ目のツイートは、複数人の作業に対する多数決にて付与された評価ラベルは不正解であったが、手法を用いて付与された評価ラベルは正解であった例である。ワーカによって付与された評価ラベルはネガティブのラベルが多数であり、投票作業によって「攻撃性なし」の評価ラベルが付与された。しかし、手法を用いて付与された評価ラベルは「ニュートラル」であった。そのため、ネガティブの評価を付与したワーカは、一致率が閾値を下回ってしまって取り除かれた品質の低いワーカであると考えられる。また、複数人の作業に対する多数決にて付与された評価ラベルが不正解であり、手法を用いて付与された評価ラベルが正解であったデータは 2 件のみであった。このことから、本実験にて取り除かれたワーカ以外にも品質の低いワーカが存在し、その全てを取り除くことができなかったと考えられる。

二つ目のツイートは、複数人の作業に対する多数決にて付与された評価ラベルは正解であったが、手法によって付与された評価ラベルは不正解であった例である。手法を用いることによって、評価ラベルが「攻撃性なし」から「ニュートラル」に

表 5.3 複数人の作業に対する多数決にて付与された評価ラベルと実験 1 で付与された評価ラベルに変化があったツイート例

「笑」を含むツイート	vote	true	all	Model1
最寄りのゲーセン 10km 以上離れてんの笑!?	攻撃性なし 自虐 ニュートラル 攻撃性なし ニュートラル	ニュートラル	攻撃性なし	ニュートラル
え、ちょっと確認!今週の金曜日に虹会 20 時から??えっ、待てよ、!!部活とか、帰りの時間で、見れんかもじゃん、やばいやばい!!え、終わった…笑	ニュートラル 自虐 攻撃性なし ニュートラル 攻撃性なし	攻撃性なし	攻撃性なし	ニュートラル

vote ワーカが付与した評価ラベル
true 自身で付与した正解ラベル
all 複数人の作業に対する多数決にて付与された評価ラベル
Model1 Model1 と手法を適用して付与された評価ラベル

変化した。このことから、ネガティブの評価ラベルを付与したワーカは、手法を用いたことによって取り除かれた品質の低いワーカであると考えられる。しかし、手法を用いて付与された評価ラベルが正解ラベルと一致しなかったことから、取り除かれたワーカの作業結果の中にも有用なデータがあったと考えられる。また、複数人の作業に対する多数決にて付与された評価ラベルが正解であり、手法を用いて付与された評価ラベルが不正解のデータは 210 件であった。このことから、取り除かれたワーカの作業データを全て取り除いてしまうと、有用なデータまで取り除かれてしまうということが考えられる。

本実験では、品質の低いワーカを全て取り除くことはできず、質の悪い作業データが残ってしまった。そのため、正解ラベルと同じ評価を得ることができなかったと考えられる。また、必要なデータが取り除かれてしまうという結果になったことから、取り除かれたワーカの作業データの中にも有用な作業データがあったと考えられる。

5.3 実験 2：一致率による票の重み設定

本節では、4.2.2 項で説明した、一致率をもとにワーカごとに票の重みを設定することによって、多数決をとる手法の有効性を確かめるための実験について述べる。

5.3.1 実験手順

5.1 節で説明したクラウドワーカ w_i のデータセット D_i を使用して模倣モデルを構築する。このとき、5.2 節と同様にポジティブかネガティブかなどを判定するモデルと、ネガティブの中でさらに細かい評価ラベルを判定するモデルの二つの模倣モデルを構築する。モデルの構築には 4.1 節で説明した二種類の方法を用いて行い、合計で四つのモデルを構築する。また、模倣モデルを構築するとき 5.2 節と同様に 300 エポックの学習を行う。データセット D_i を訓練データ 6 割、検証データ 2 割、テストデータ 2 割に分けて学習を行う。このとき、訓練データは 5.2 節と同様に最もデータ数が多い評価ラベルのデータ数に合わせてデータの複製を行う。構築した二種類の模倣モデルを使用して評価ラベルの予測を行い、4.2 節の式 (4.2.1)

を用いてワーカーごとに一致率 A_i を算出する。算出した全ワーカーの一致率と 8 クラス分類のチャンスレート 0.125 を用いて最大値 1, 最小値 0 となるように正規化を行い, ワーカー w_i の品質 q_i を求める。この品質 q_i をワーカー w_i の票の重みとしてデータセット D_{eva} に含まれるデータの評価ラベルを多数決にて付与し, *Accuracy* を算出する。

5.3.2 結果・考察

Model1 と Model2 のそれぞれを用いて算出した一致率をもとにワーカーの品質を求め, その品質を用いて各ワーカーの票に重み付けして多数決をとった結果を表 5.4 に示す。表中の太字は, 複数人の作業に対する多数決にて集約された結果の *Accuracy* よりも高くなった結果である。Model1 と手法を適用した場合と Model2 と手法を適用した場合の *Accuracy* を比べると, Model1 を用いてワーカーの品質 q_i を定めた場合の *Accuracy* が高くなっている。また, Model1 と手法を適用した場合の *Accuracy* は, 複数人の作業に対する多数決にて集約された結果の *Accuracy* より高くなっている。

Model2 では, 各ワーカー用にファインチューニングする前のモデルが, ワーカーの一致率を出す際の正解のデータを使用して学習している。そのため, データ数が少ないワーカーほどファインチューニングによるパラメータの変化量が少なく, 各ワーカーの特徴が捉えきれしていない。それにより, ファインチューニング前とほとんど同じ予測をするため, ワーカーの質に関係なく一致率が高くなりやすい。ワーカーの質が関係なくなるため, 全体を通して一致率の差が現れにくい。そのため, ワーカー w_i の品質 q_i を票の重みとして設定したとしても, 複数人の作業に対する多数決にて付

表 5.4 票の重みを設定して多数決を取り直した時の *Accuracy*

	<i>Accuracy</i>
all	0.6625
Model1	0.6685
Model2	0.6625

与された評価ラベルと変わらない評価ラベルが付与されやすくなる。実際に、複数人の作業に対する多数決にて付与された評価ラベルと、Model2 を用いて評価ラベルの予測をした場合に付与された評価ラベルが変化していたデータは 2,000 件のうち 9 件のみであった。

一方、Model1 を用いて模倣モデルを構築すると、ワーカ w_i の特徴をそのまま反映させたモデルが構築できる。そのため、ワーカ w_i の本来の一致率 A_i を算出することができると考えられ、複数人の作業に対する多数決にて付与された評価ラベルに依存しない品質 q_i を求めることができる。そして、より品質の高いワーカの作業結果が反映されやすくなると考えられる。また、複数人の作業に対する多数決にて付与された評価ラベルと、Model1 を用いて評価ラベルの予測をした場合に付与された評価ラベルが変化していたデータは 2,000 件のうち 85 件であった。

表 5.5 は、複数人の作業に対する多数決にて付与された評価ラベルと、4.2.2 項にて説明した手法を用いて付与された評価ラベルが変化していたツイートの例である。一つ目のツイートは、複数人の作業に対する多数決にて付与された評価ラベルは不正解であったが、Model1 と手法を適用して付与された評価ラベルは正解で

表 5.5 複数人の作業に対する多数決にて付与された評価ラベルと実験 2 で付与された評価ラベルに変化があったツイート例

「笑」を含むツイート	vote	true	all	Model1	Model2
今週末の楽しみ無くなり そう一さすがにやさぐれ そうだわ笑	攻撃性なし 自虐 攻撃性なし 自虐 自虐	攻撃性なし	自虐	攻撃性なし	自虐
そやゆとって何でこんな にいいの？ふたりの全部 がいい笑	ニュートラル ポジティブ ポジティブ ニュートラル ポジティブ	ポジティブ	ポジティブ	ニュートラル	ポジティブ

- vote ワーカが付与した評価ラベル
- true 自身で付与した正解ラベル
- all 複数人の作業に対する多数決にて付与された評価ラベル
- Model1 Model1 と手法を適用して付与された評価ラベル
- Model2 Model2 と手法を適用して付与された評価ラベル

あった例である。ワーカによって付与された評価ラベルは「攻撃性なし」が2票、「自虐」が3票であったため、複数人の作業に対する多数決にて付与された評価ラベルは「自虐」となっている。しかし、Model1 と手法を適用して付与された評価ラベルは「攻撃性なし」となっており、正解ラベルと一致している。このことから、「攻撃性なし」の評価ラベルを付与したワーカの品質が高いということが予想できる。また、品質の低いワーカの影響力を抑えることができたため、正しい結果を得ることができたと考えられる。

二つ目のツイートは、複数人の作業に対する多数決にて付与された評価ラベルは正解であったが、Model1 と手法を適用して付与された評価ラベルは不正解であった例である。一つ目のツイートと同様に、ワーカによって付与された評価ラベルが2票と3票に割れている。しかし、Model1 と手法を適用して付与された評価ラベルは2票しか付与されなかった評価ラベルである。品質の低いワーカによって付与された多数派の評価ラベルではなく、品質の高いワーカによって付与された少数派の評価ラベルが集約結果として得られたと考えられる。このことから、提案手法を用いることによって、品質の低いワーカの影響力を抑えることができたと考えられる。しかし、その評価ラベルは正解ラベルと一致していないことから、品質の高いワーカの評価が常に正しいとは限らないということが考えられる。また、品質の低いワーカが行った作業データの中にも有用なデータが存在すると考えることができ、全ての作業データに対して影響力を抑えることが良いわけではないと考えられる。

5.4 実験3：一致率による無効票の決定

本節では、4.2.3 項で説明した、一致率の値をワーカがデータに対して投票できる確率として、多数決をとる手法の有効性を確かめるための実験について述べる。

5.4.1 実験手順

5.3 節の実験と同じようにワーカ w_i の模倣モデルを構築し、一致率 A_i を算出する。本節でも模倣モデルは5.2 節と同様に二つ構築する。モデルの構築には4.1 節

で説明した二種類の方法を用いて行い、合計で四つのモデルを構築する。また、模倣モデルの構築については 5.3 節と同様の手順で行う。一致率 A_i の値をワーカー w_i のデータ d_j に対する票が有効になる確率とする。例えば、一致率 A_i が 0.6 のワーカー w_i の票は 6 割の確率で有効になり、4 割の確率で無効となる。4.2.3 項で説明した手法をデータセット D_{eva} に含まれるデータに対して行い多数決をとって評価ラベルを付与し、*Accuracy* を算出する。

5.4.2 結果・考察

Model1 と Model2 のそれぞれを用いて算出した一致率の値をワーカーが一つのデータに対して投票できる確率として、多数決をとった結果を表 5.6 に示す。Model1 と手法を適用した場合と Model2 と手法を適用した場合の *Accuracy* を比べると、Model2 と手法を適用した場合の *Accuracy* が高くなっている。しかし、どちらの方法を用いた場合でも、複数人の作業に対する多数決にて集約した結果の *Accuracy* より低くなっている。

表 5.7 には作業数が多い順に 5 人のワーカーの各模倣モデルによる評価ラベルの予測と、複数人の作業に対する多数決にて付与した評価ラベルの一致率が示してある。ワーカーごとに Model1 と Model2 の模倣モデルを比べると、Model2 の方が一致率の値が高くなっている。しかし、この結果はモデルの構築方法に由来するものであると考えられるため、ワーカーの特徴を捉えきれていないわけではない。一致率が高いほど票が有効になる確率が高くなるため、Model2 を用いて実験を行った場合の *Accuracy* が高くなったと考えられる。すべての票が無効票となったデータの数を調べたところ、Model1 を用いた場合は 2,000 件のうち 444 件、Model2 を用い

表 5.6 一致率を確率として多数決を取り直した時の *Accuracy*

	<i>Accuracy</i>
all	0.6625
Model1	0.3460
Model2	0.4335

た場合は 2,000 件のうち 74 件となっており、Model1 の方がすべての票が無効票となったデータの数が多くなっている。つまり、評価ラベルをつけられないデータが多く存在したために、*Accuracy* が低くなってしまったと考えられる。また、一致率の値をそのまま確率としているため、作業結果が使用される確率が低い。そのため、有用なデータも取り除かれてしまい、正しい評価ラベルを得られなかったと考えられる。

表 5.8 は複数人の作業に対する多数決にて付与された評価ラベルと、4.2.3 項にて説明した手法を用いて付与された評価ラベルが変化していたツイートの例である。一つ目のツイートは、複数人の作業に対する多数決にて付与された評価ラベルは不正解であったが、Model1 と手法を適用して付与された評価ラベルは正解であった例である。複数人の作業に対する多数決にて付与された評価ラベルは多数派の評価ラベルであったのに対して、Model1 と手法を適用して付与された評価ラベルは少数派の評価ラベルであった。このことから、少数派の評価ラベルを付与したワーカの品質が高く、票が有効となる確率が高くなったのではないかと考えられる。しかし、Model2 と手法を適用して付与された評価ラベルを見ると、複数人の作業に対する多数決にて付与された評価ラベルと同じであることがわかる。ワーカ w_i の一致率 A_i は、Model2 を用いた場合の方が Model1 を用いた場合よりも高い。そのため、Model2 を用いた場合の方が票が有効となる確率が高くなりやすい。このことから、Model1 を用いた場合に正解ラベルと同じ評価ラベルが付与されたことは、必ずしも結果が改善されたとは言えないと考えられる。

二つ目のツイートは、複数人の作業に対する多数決にて付与された評価ラベルは正解であったが、Model1 と手法を適用して付与された評価ラベルは不正解であっ

表 5.7 作業数が多い上位 5 人の各模倣モデルによる評価予測の一致率

ワーカ ID	Model1 の一致率	Model2 の一致率
836	0.4059	0.5579
20	0.3648	0.4948
944	0.3801	0.5126
913	0.2801	0.4172
946	0.4148	0.5519

た例である。このツイートでは、ワーカが付与した評価ラベルにばらつきが見られる。Model1 と手法を適用して付与された評価ラベルは元々 1 票しか付与されていない評価ラベルである。このことから、その票以外の票が無効票となってしまったと考えられる。つまり、残りの 4 票の中に有効なデータが含まれていたとしても、そのデータが取り除かれてしまったと考えられる。

表 5.8 複数人の作業に対する多数決にて付与された評価ラベルと実験 3 によって付与された評価ラベルに変化があったツイート例

「笑」を含むツイート	vote	true	all	Model1	Model2
架純ちゃんの髪型も良い なぁ前髪伸ばそかな笑	ニュートラル ポジティブ ポジティブ ニュートラル	ニュートラル	ポジティブ	ニュートラル	ポジティブ
じゃ、途中だけど、また明日ね！笑	ニュートラル ポジ + ネガ ニュートラル ポジティブ その他	ニュートラル	ニュートラル	ポジティブ	ニュートラル

vote	ワーカが付与した評価ラベル
true	自身で付与した正解ラベル
all	複数人の作業に対する多数決にて付与された評価ラベル
Model1	Model1 と手法を適用して付与された評価ラベル
Model2	Model2 と手法を適用して付与された評価ラベル

第6章 おわりに

本研究では、クラウドワーカの品質を考慮した結果集約方法を行うことによって、クラウドソーシングの質を向上させることを目的としている。クラウドワーカの中にはスパムワーカと呼ばれるワーカが一定数存在しており、そのスパムワーカの影響を受けることによって作業依頼者の求める結果を得られないことがある。しかし、スパムワーカは他のワーカとは違う作業結果であるため、正しい作業結果と比較することによってスパムワーカを見つけ出し排除することが可能である。そのとき使用する作業結果は作業依頼者の求める結果に近い作業結果である必要がある。そして、多くの意見を集めることによって作業依頼者の求める結果に近い作業結果を得ることができると考えられる。しかし、ワーカを集めるには時間と費用がかかってしまうため、ワーカの模倣モデルを構築して擬似的に作業結果を集めることができないかと考えた。そこで本論文では、クラウドワーカの模倣モデルを構築して算出する一致率を、ワーカの品質として利用する異なる三つの結果集約方法を提案した。構築した模倣モデルを使用して得られた予測結果をもとに算出した一致率を用いて結果集約を行うことによって、スパムワーカによる影響を抑えられ、作業依頼者の求める結果に近い作業結果を得られることを考えた。提案手法の有効性を確かめるために実験を行った。

一つ目の実験は、逐次的にワーカの品質を測定することによってワーカの選定を行い、取り除かれなかったワーカの作業結果のみを用いて多数決をとる手法を用いて行った。その結果、手法を適用して集約した結果の *Accuracy* は 0.5845 であった。一方、複数人の作業に対する多数決にて集約した結果の *Accuracy* は 0.6625 であり、手法を適用して付与された結果の方が低くなったことがわかった。この結果から、品質の低いワーカとして取り除かれたワーカの作業データの中に含まれる有用な作業データまで取り除かれてしまったと考えられる。

二つ目の実験は、ワーカの模倣モデルを用いて算出した一致率をもとにワーカの品質を求め、その品質を用いて各ワーカの票に重み付けして多数決をとる手法を用いて行った。その結果、Model1 と手法を適用して集約した結果の *Accuracy* は 0.6685 であり、Model2 と手法を適用して集約した結果の *Accuracy* は 0.6625 であった。Model1 を用いた場合の *Accuracy* が高くなった理由は、より各ワーカの

特徴を捉えたモデルを構築できたことによって、ワーカ本来の品質を求めることができるためであると考えられる。そのため、品質の低いワーカの作業結果の重みが小さくなり、評価ラベルを付与する際に及ぼす影響力が抑えられた。

三つ目の実験は、一致率の値をワーカが一つのデータに対して投票できる確率として多数決をとる手法を用いて行った。その結果、Model1 と手法を適用して集約した結果の *Accuracy* は 0.3460 であり、Model2 と手法を適用して集約した結果の *Accuracy* は 0.4335 であった。どちらのモデルにおいてもワーカごとの一致率が低く、その値をそのまま確率として用いているためほとんどの票が無効票となってしまっている。そのため、有用な作業データも取り除かれてしまったと考えられる。

一つ目と三つ目の手法では、作業データを取り除くというアプローチを行ったため、有用な作業データまで取り除かれてしまうという結果となった。しかし、二つ目の手法では、品質の低いワーカの影響力を抑えるだけで、そのワーカが行った有用な作業データを活用することができる。そのため、他の手法と比べて良い結果を得ることができた。これらの手法は作業の依頼が終了してデータが揃った状態に適用しているため、依頼終了後にしか品質の低いワーカの特定ができない。そのため、品質の低いワーカに割り振られる作業数が増えてしまう。

今後の展望として作業依頼中に適用できるように拡張することを考えている。それによって、品質の低いワーカの特定が早まり、途中でそのワーカを取り除くことができるようになる。そして、品質の低いワーカに割り振られる予定だった作業を品質の高いワーカに割り振ることができ、クラウドソーシングの質が向上すると考えられる。

謝辞

本研究を進めるにあたって、指導教員である鈴木優准教授にはたくさん指導していただきました。また、女一人で周りに馴染めないことを気にかけ、研究室の皆さんと仲良くなる場を設けてくださいました。学会やアルバイトの手続きにあたって、秘書の佐野さん、井尾さんに大変お世話になりました。研究に行き詰まったとき、研究室の先輩方、同期の皆さんに相談に乗っていただきました。また、9月に行われた学会発表を機に研究室内の雰囲気良くなり、楽しく過ごすことができました。研究がどうしても嫌になったときやつらいとき、家族には話を聞いてもらったり励ましてもらったり、時にやさしく時にきびしく支えていただきました。皆様のご助力により一応の終わりを迎えられること、心より深く感謝申し上げます。

参考文献

- [1] Jeff Howe, et al. The rise of crowdsourcing. *Wired magazine*, Vol. 14, No. 6, pp. 1–4, 2006.
- [2] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics.
- [3] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, Vol. 30, pp. 5998–6008, 2017.
- [4] 西智樹, 小出智士, 大野宏司, 長屋隆之. ソーシャルネットワークを用いたクラウドソーシングの品質向上. 人工知能学会全国大会論文集 第 27 回 (2013), pp. 3M3OS07d4–3M3OS07d4. 一般社団法人 人工知能学会, 2013.
- [5] 芦川将之, 川村隆浩, 大須賀昭彦. プライベートクラウドソーシングにおける精度向上手法. 人工知能学会全国大会論文集 第 28 回 (2014), pp. 1J5OS18b4–1J5OS18b4. 一般社団法人 人工知能学会, 2014.
- [6] 芦川将之, 川村隆浩, 大須賀昭彦. マイクロタスク型クラウドソーシングプラットフォーム環境における精度向上手法の導入と評価. 人工知能学会論文誌, Vol. 29, No. 6, pp. 503–515, 2014.
- [7] Harry Halpin and Roi Blanco. Machine-learning for spammer detection in crowd-sourcing. In *Workshops at the Twenty-Sixth AAAI Conference on Artificial Intelligence*, pp. 85–86, 2012.
- [8] 松原繁夫, 水島拓也. クラウドソーシングにおける複数タスク割当て. 人工知能学会全国大会論文集 第 27 回 (2013), pp. 3M4OS07e3–3M4OS07e3. 一般社団法人 人工知能学会, 2013.

- [9] Alexander Philip Dawid and Allan M Skene. Maximum likelihood estimation of observer error-rates using the em algorithm. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, Vol. 28, No. 1, pp. 20–28, 1979.
- [10] 小山聡, 馬場雪乃, 櫻井祐子, 鹿島久嗣. クラウドソーシングにおけるワーカーの確信度を用いた高精度なラベル統合. 人工知能学会全国大会論文集 第 27 回 (2013), pp. 2M5OS07b2–2M5OS07b2. 一般社団法人 人工知能学会, 2013.

発表リスト

- [1] 太田奈那, 鈴木優『BERT を用いた分類器によるクラウドソーシングの質の向上』第 21 回情報科学技術フォーラム, 2022.
- [2] 太田奈那, 鈴木優『BERT を用いた分類器によるクラウドソーシングの質の向上』東海関西データベースワークショップ, 2022.
- [3] 太田奈那, 鈴木優『クラウドワーカの模倣モデルと投票作業の一致率を用いた結果集約方法』第 15 回データ工学と情報マネジメントに関するフォーラム, 2023.