

令和 5 年 5 月 5 日現在

機関番号：12611

研究種目：基盤研究(C)（一般）

研究期間：2020～2022

課題番号：20K11917

研究課題名（和文）機械学習の訓練データ検証のための対話的可視化手法の研究

研究課題名（英文）Interactive visualization for verification of training data for machine learning

研究代表者

伊藤 貴之（Takayuki, Itoh）

お茶の水女子大学・基幹研究院・教授

研究者番号：80401595

交付決定額（研究期間全体）：（直接経費） 3,400,000円

研究成果の概要（和文）：機械学習に用いる訓練データの分布や作成過程を可視化することで、高品質な訓練データの構築を支援する諸手法の研究に着手した。具体的な研究成果として、以下の3種類の可視化手法を提案した。1)訓練データのアノテーション作業の半自動化と、その根拠となる決定木の構築結果とその動作結果の可視化。2)複数の訓練データの特徴量分布とラベル分布の比較可視化。3)複数の作業員によるアノテーションの作業工程の可視化と、アノテーションの信頼性に関する検証。

研究成果の学術的意義や社会的意義

大規模で複合的な訓練データの分布や制作過程を視認性の高い形で情報提示する手法の開発は、可視化の研究における学術面での本質的な課題であり、これを解くことに学術的意義があった。一方で、機械学習の普及により訓練データの品質は社会的に大きな課題となっている。訓練データ制作の半自動化による信頼性の向上、複数の訓練データ間での特徴量やラベルの分布の検証、訓練データの制作過程での各作業員による工程の信頼性の検証、といった各課題は機械学習の品質を向上するために重要な課題であり、これらの解決には大きな社会的意義があった。

研究成果の概要（英文）：This research focused on the development of various methods to support the construction of high-quality training data for machine learning, by visualizing the distribution and creation process of the training data. As the results, we proposed the following three types of visualization methods: 1) Semi-automation of the annotation process of training data and visualization of the results of the construction and operation of the decision tree that serves as the basis for the annotation process; 2) Comparative visualization of distributions of features and labels of multiple training data; and 3) Visualization of the annotation process by multiple workers and verification of the reliability of the annotations.

研究分野：可視化

キーワード：可視化 訓練データ アノテーション

### 1 . 研究開始当初の背景

機械学習の普及にともない、数千個・数万個といった大量のマルチメディアデータ(画像や音声など)を訓練データとして扱う機会が増えた。これらの訓練データには複数の作業者の手作業によって注釈(アノテーション)と呼ばれる教師信号がつけられることが多い。このようにして形成された訓練データの品質は機械学習結果の品質に大きな影響を及ぼすことがある。

機械学習の訓練データには大規模なものが多く、利用者がその全貌を理解しているとは限らない。一方で、機械学習を実用する現場では機械学習の動作結果に対する説明責任が求められる場面も多く、運用者が訓練データの中身を理解していないでは済まされない状況も起こりえる。AI 関連技術の法令化によっては運用者の説明責任が強化される可能性もある。この状況に対する解決手段として、機械学習のための訓練データの内容を利用者が検証し理解するための可視化手法の確立が重要であると考えた。

### 2 . 研究の目的

本研究の目的は、「訓練データの構成を可視化することで、機械学習の利用者が高い品質の訓練データを構築できるようにする」というものである。機械学習のための可視化手法は2017年頃から急激に多く発表されている。その多くは機械学習の動作を理解するためにその内部構造(例えばニューラルネットワークの中間層)の振る舞いを可視化しているのに対して、訓練データの構成に特化した可視化手法は少ない。特に、複数の作業者による訓練データへの注釈の揺れに着目した可視化手法、テストデータにおける誤動作と訓練データの対応に着目した可視化手法は、申請時点ではほとんど見当たらなかった。この点が申請時点での本研究の学術的独自性であったと考えられる。また、本研究が対象とする訓練データは複合的なものであり、これを視認性・可読性の高い形で可視化するためのビジュアルデザインを追求することで、研究の創造性を追求するものであった。

本研究は機械学習の利用者に対して

- ・ 数千個～数万個という大量のデータ要素の特徴分布
- ・ 複数の作業者による訓練データの注釈の揺れ
- ・ テストデータでの実行結果のエラーと訓練データの対応関係

を含む複雑な構造をもった情報を提示することを目標とした。このような複合的なデータから傾向や知識を発見するための視認性・可読性・対話操作性の高い情報デザインを実現することは可視化の学術研究における本質的な課題であった。

それと同時に本研究では、同一作業者による注釈の一貫性、作業者間の注釈の不一致の傾向、あるいは注釈の不一致を生じやすいデータ要素の傾向、といったヒューマンファクタを発見することも目標とした。作業者による訓練データの注釈は必ずしも品質を保証できる作業ではない。作業者の熟達度・感性・知識などが影響を及ぼすこともあれば、作業日の疲労度や健康状態が影響を及ぼすこともある。機械学習の品質向上の指針の一つとして、訓練データ作成にかかわる作業者に関する知見を得ることも本研究課題の目標とした。

### 3 . 研究の方法

本研究では3年間で以下の3種類の方法に着手した。

**[方法1]** 訓練データの「アノテーションの半自動化」のための可視化に取り組んだ。画像の印象を推測する機械学習の構築において一般的に、学習に使用する訓練データ作成は多数の画像に

アノテーションを付与する作業が必要である。しかし作業者が有する個人の印象回答が学習結果に依存する問題がある。本研究では画像から受ける印象の個人差を低減し、かつ作業者の負担を低減するために、多人数の印象回答値を機械学習した結果にそって画像の印象のアノテーション付与を半自動化するシステムと、その過程を可視化する手法を開発した。本手法ではまずSD法を採用した印象評価を実施し、続いて各画像の印象値を用いてファジィ決定木を生成する。このファジィ決定木によって画像を自動分類したのち、その結果と過程を可視化することで、ユーザによる画像再分類を支援する。決定木の表示と類似画像の一覧表示を連動させることで、決定木の可読性を向上させ、アノテーション付与の傾向理解を促す。

**[方法 2]** 複数のラベル付き多次元データの差異を可視化する手法を開発した。本手法では与えられた複数の多次元データに対して同一の次元削減を施して同一画面空間上に配置するとともに、標本ごとに付与されているラベルを用いて標本群を半透明表示する。この表現により、複数の多次元データの間でどのラベルを有する標本に大きな共通点または差異が見られるか、またどのラベルを有する標本に外れ値が生じやすいか、といった点を観察できるようになる。この可視化手法を用いることで、機械学習の複数のラベル付き訓練データにおける特徴量とラベルの差異を比較できるようになる。

**[方法 3]** アノテーションの付与作業の過程を観察するための可視化手法を開発した。アノテーションはデータに人手で注釈を付与する作業であり、機械学習の訓練データ作成にも用いられる。アノテーションの信頼性は機械学習の信頼性の観点からも極めて重要である。アノテーションには作業者(ワーカー)ごとの傾向があり、これらの傾向差がデータの信頼性を損ねる可能性がある。特にワーカーの主観に依存するタスクにおいて、ワーカーごとの傾向差は顕著に現れる。そこでワーカーのアノテーション結果を観察することで、信頼性の高いアノテーションの実現を目指した。

#### 4. 研究成果

本研究で着手した3種類の可視化手法の各々について、以下の通り成果を示す。

**[方法 1]** 本方法に沿って開発した可視化ソフトウェアのスナップショットを図1に示す。図1(上)は本手法に沿って自動生成されたファジィ決定木の可視化機能であり、図1(下)はファジィ決定木によって自動分類された画像群の一覧表示のための可視化機能である。図1は尺度「暗い 明るい」に沿って画像を自動分類した結果を示しているが、他にも尺度「フィット ルーズ」「フォーマル カジュアル」「シンプル ゴージャス」「日常的な 非日常的な」についても同様に可視化結果を考察し、各々の尺度に対して特徴的なファジィ決定木と画像分類結果が得られたことを検証している。

本方法によるシステムの操作性と有効性を示すための評価実験を実施した。評価者には1)機械学習の研究に従事する学生、2)可視化の研究に従事する学生、3)いずれの研究にも従事しない学生、の3グループの学生に依頼した。評価実験の画面録画を観察することで、3グループの各々には特徴的な操作手順があり、機械学習の経験がある人ならではの操作手順、可視化の経験がある人ならではの操作手順が観察された。また本システムによる画像自動分類結果を評価者に再分類させたところ、各グループとも再分類操作をした画像枚数は少なく、自動分類がユーザの手間を削減していることがわかった。ベースライン手法と比較した実験においても、画像再分類の所要時間が大幅に削減されており、作業の短時間化に貢献可能であることが示された。

本方法は国際会議で発表され、さらに国内会議にて最優秀論文賞を受賞して論文誌に掲載されている。さらに協力学生は企業の懸賞論文企画にて最優秀賞を受賞している。

**[方法 2]** 本手法による可視化結果のスナップショットを図 2,3 に示す .

図 2 は手書き文字画像データセット MNIST を題材として . 数字の「1」「7」の画像特徴量を投影した結果を示している . 黄土色の三角形群が「1」の画像群の特徴量を . 緑色の三角形群が「7」の画像群の特徴量を示している . 両者には重なりが見られることから . 両画像の中には「1」と「7」の判別が難しい画像が混在していることが示唆される .

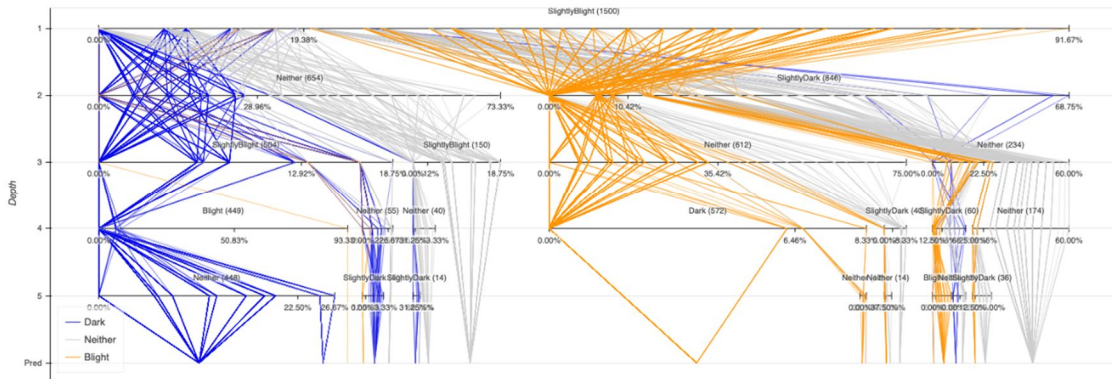
本方法は国際会議で発表され . 当該国際会議の代表論文として海外出版社による書籍への掲載に招待された . 書籍原稿として国際会議原稿の拡張版を提出済みであり . 2023 年末までに出版される予定である .

**[方法 3]** 本手法によるワーカの観察結果を示す .

今回採用した実験では 3 人のワーカに 977 枚の顔画像を提示してアノテーション作業を依頼した . その結果として . アノテーションの信頼性は作業開始時よりも終盤の方が高く , そして所要時間の長い作業よりも所要時間の短い作業の方が高い信頼性を得られる傾向が見られた . この結果から , 作業の序盤ではアノテーションの判断基準が定まっていないためにブレが生じること , アノテーションに時間がかかるタスクは判断が難しいタスクであり信頼性も低くなることが示唆される . 一方で , 作業の終盤になると疲労感によって信頼性が下がる可能性を指摘したワーカもいた . このことから . ワーカによっては休憩などの配慮が訓練データの信頼性を向上できる可能性が考えられる .

また . アノテーションの項目ごとの信頼度を算出した結果 . 3 人のワーカともに信頼性の高い項目・信頼性の低い項目が一致していた . このことから . アノテーションの品質を上げるには特定の項目への着目が重要であることが示唆された . さらに . アノテーションの信頼性が低いとされる画像にはワーカごとに異なる傾向があることから . ワーカごとの傾向を把握することも重要であることが示唆された .

本方法は国内会議および国際会議で発表され . 国内会議では共同研究学生が発表して学生プレゼンテーション賞を受賞した .



暗い

Depth 5

明るい

Depth 4

どちらでもない

Depth 3

**暗い**  
黒に近い衣服が多い

**明るい**  
色鮮やかな衣服が多い

**どちらでもない**  
青や緑または複数混合色の衣服が多い

図 1. [方法 1]におけるファジィ決定木の可視化機能と画像分類結果の一覧表示機能

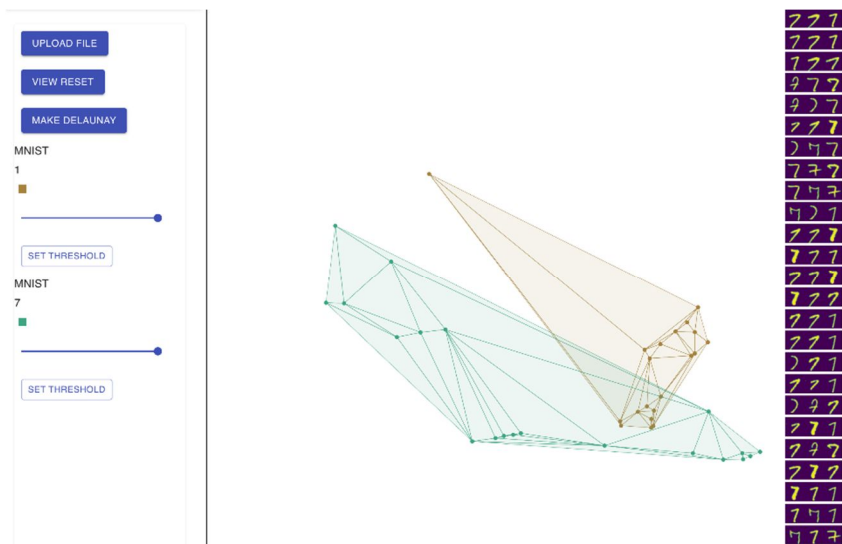


図 2. 手書き文字画像データセット MNIST の「1」「7」に関する可視化結果

5. 主な発表論文等

〔雑誌論文〕 計1件（うち査読付論文 1件 / うち国際共著 0件 / うちオープンアクセス 0件）

|  |                       |
|--|-----------------------|
| 1. 著者名<br>飯島絋理, 伊藤貴之                           | 4. 巻<br>21            |
| 2. 論文標題<br>印象評価にもとづくアノテーション作業の半自動化を支援する可視化システム | 5. 発行年<br>2022年       |
| 3. 雑誌名<br>芸術科学会論文誌                             | 6. 最初と最後の頁<br>186-198 |
| 掲載論文のDOI (デジタルオブジェクト識別子)<br>なし                 | 査読の有無<br>有            |
| オープンアクセス<br>オープンアクセスではない、又はオープンアクセスが困難         | 国際共著<br>-             |

〔学会発表〕 計18件（うち招待講演 0件 / うち国際学会 5件）

|  |
|--|
| 1. 発表者名<br>Karen Kosaka, Takayuki Itoh   |
| 2. 発表標題<br>A Visualization Method for Training Data Comparison                       |
| 3. 学会等名<br>5th International Conference on Information Visualisation (IV2021) (国際学会) |
| 4. 発表年<br>2021年  |

|   |
|---|
| 1. 発表者名<br>Akari Iijima, Takayuki Itoh  |
| 2. 発表標題<br>Visualization for Image Annotations based on Semantic Differential |
| 3. 学会等名<br>IEEE VIS, Poster Session (国際学会)                                    |
| 4. 発表年<br>2021年   |

|                               |
|-------------------------------|
| 1. 発表者名<br>高坂夏怜, 伊藤貴之         |
| 2. 発表標題<br>訓練データ比較のための可視化の一手法 |
| 3. 学会等名<br>可視化情報シンポジウム        |
| 4. 発表年<br>2021年               |

|                                   |
|-----------------------------------|
| 1. 発表者名<br>伊藤貴之                   |
| 2. 発表標題<br>説明性の高い機械学習の実現のための可視化技術 |
| 3. 学会等名<br>2021年電子情報通信学会ソサエティ大会   |
| 4. 発表年<br>2021年                   |

|   |
|---|
| 1. 発表者名<br>飯島紆理, 伊藤貴之                         |
| 2. 発表標題<br>SD法による大規模印象評価に基づくアノテーションを支援する可視化   |
| 3. 学会等名<br>第29回インタラクティブシステムとソフトウェアに関するワークショップ |
| 4. 発表年<br>2022年                               |

|  |
|--|
| 1. 発表者名<br>高坂夏怜, 伊藤貴之                  |
| 2. 発表標題<br>訓練データの比較可視化のための閾値自動設定       |
| 3. 学会等名<br>第14回データ工学と情報マネジメントに関するフォーラム |
| 4. 発表年<br>2022年                        |

|  |
|--|
| 1. 発表者名<br>三浦梨花, 伊藤貴之                  |
| 2. 発表標題<br>主観を要するアノテーションタスクの観察と可視化     |
| 3. 学会等名<br>第14回データ工学と情報マネジメントに関するフォーラム |
| 4. 発表年<br>2022年                        |

|  |
|--|
| 1. 発表者名<br>飯島緋理, 伊藤貴之                        |
| 2. 発表標題<br>SD法による大規模印象評価にもとづくアノテーションを支援する可視化 |
| 3. 学会等名<br>第14回データ工学と情報マネジメントに関するフォーラム       |
| 4. 発表年<br>2022年                              |

|   |
|---|
| 1. 発表者名<br>飯島緋理, 伊藤貴之                       |
| 2. 発表標題<br>SD法による大規模印象評価に基づくアノテーションを支援する可視化 |
| 3. 学会等名<br>情報処理学会第84回全国大会                   |
| 4. 発表年<br>2022年                             |

|  |
|--|
| 1. 発表者名<br>A. Iijima, T. Itoh, N. Grossmann, H.-Y. Wu  |
| 2. 発表標題<br>Visualization of semantic differential studies with a large number of images, participants and attributes |
| 3. 学会等名<br>24th International Conference on Information Visualisation (IV2020) (国際学会)                                |
| 4. 発表年<br>2020年  |

|  |
|--|
| 1. 発表者名<br>飯島緋理, 伊藤貴之                        |
| 2. 発表標題<br>SD法による画像印象のタグ付けを支援する可視化           |
| 3. 学会等名<br>第13回データ工学と情報マネジメントに関するフォーラム(DEIM) |
| 4. 発表年<br>2021年                              |



|  |
|--|
| 1. 発表者名<br>村上綾菜, 伊藤貴之                        |
| 2. 発表標題<br>高校生向けデータサイエンス教材の提案と操作ログの解析        |
| 3. 学会等名<br>第13回データ工学と情報マネジメントに関するフォーラム(DEIM) |
| 4. 発表年<br>2021年                              |

|  |
|--|
| 1. 発表者名<br>高坂夏怜, 伊藤貴之                        |
| 2. 発表標題<br>訓練データ比較のための可視化の一手法                |
| 3. 学会等名<br>第13回データ工学と情報マネジメントに関するフォーラム(DEIM) |
| 4. 発表年<br>2021年                              |

|                               |
|-------------------------------|
| 1. 発表者名<br>村上綾菜, 伊藤貴之         |
| 2. 発表標題<br>高校生向けデータサイエンス教材の開発 |
| 3. 学会等名<br>情報処理学会インタラクシヨン2021 |
| 4. 発表年<br>2021年               |

|                                    |
|------------------------------------|
| 1. 発表者名<br>飯島絢理, 伊藤貴之              |
| 2. 発表標題<br>SD法による画像印象のタグ付けを支援する可視化 |
| 3. 学会等名<br>情報処理学会第83回全国大会          |
| 4. 発表年<br>2021年                    |

|   |
|---|
| 1. 発表者名<br>Takayuki Itoh, Ayana Murakami  |
| 2. 発表標題<br>Visualization of Individual Variation of Multiple Annotators Working on Training Datasets for Machine Learning |
| 3. 学会等名<br>NICOGRAPH International 2020 (国際学会)  |
| 4. 発表年<br>2020年   |

|   |
|---|
| 1. 発表者名<br>Rika Miura, Takayuki Itoh  |
| 2. 発表標題<br>Observation and Visualization of Subjectivity-based Annotation Tasks |
| 3. 学会等名<br>26th International Conference on Information Visualisation (国際学会)    |
| 4. 発表年<br>2022年   |

|  |
|--|
| 1. 発表者名<br>高坂夏怜, 伊藤貴之                        |
| 2. 発表標題<br>ラベル付き多次元データの比較可視化の一手法と機械学習データへの応用 |
| 3. 学会等名<br>NICOGRAPH 2022                    |
| 4. 発表年<br>2022年                              |

〔図書〕 計0件

〔産業財産権〕

〔その他〕

-

|         |                           |                       |    |
|---------|---------------------------|-----------------------|----|
| 6. 研究組織 | 氏名<br>(ローマ字氏名)<br>(研究者番号) | 所属研究機関・部局・職<br>(機関番号) | 備考 |
|---------|---------------------------|-----------------------|----|

7. 科研費を使用して開催した国際研究集会

〔国際研究集会〕 計0件

8 . 本研究に関連して実施した国際共同研究の実施状況

| 共同研究相手国 | 相手方研究機関 |
|---------|---------|
|---------|---------|