

# 男女差を俯瞰するための階層型データ可視化： emd による男女間の分布差の導入

中井 祐希\* 伊藤 貴之\*

**概要.** データの偏りに起因する問題の解決には、人間による意思決定が必要である。言い換えれば、この問題を解決するには、データに潜む偏りを発見し、その偏りを人間が理解することが必要である。この観点から著者らは、多数の人物を対象としたデータからデータの分布の男女差を可視化する手法を開発している。この手法では、データ中の人物群を属性で階層的に分割し、帯グラフを搭載した階層型データ可視化手法を適用している。本報告では、emd による男女間の分布差算出を導入することで、男女差が大きな部分を強調表示する手法を提案するとともに、空調の温感に関する評価値の可視化例を紹介する。

## 1 はじめに

機械学習やデータサイエンスの普及にともない、データの偏りがもたらす問題が注目されるようになった。例えば機械学習に用いる訓練データに偏りがあることで、学習結果にも偏りが生じることがある。また、データの偏りはデータ全体に見られるとは限らず、むしろ特定の属性や要素を持つ部分集合に見られることが多い。そのため、データの偏りを発見するには、データが持つ属性ごとの数値分布の違いを観察することが重要である。ここでいう属性とは例えば、性別、地域、人種、世代などである。

著者らは、多数の人物を対象としたデータから、データの分布の男女差を可視化する一手法 [6] を開発している。データを構成する人物群を属性で階層的に分割し、「平安京ビュー」[4] という階層型データ可視化手法を用いて可視化する。ここで下位階層を表現する長方形領域を 2 列の帯グラフで表現し、その帯グラフによって所定の属性値に関する男性および女性の数値分布を表現する。本報告ではその拡張手法として、男性と女性の数値分布の分布間距離を earth mover's distance (emd)[7][8] により算出し、指定値よりも emd が高ければ男女の帯グラフをカラースケールで表示し、低ければグレスケールで表示することで、数値分布の男女差が大きな部分を強調表示する手法を報告する。本報告では空調の温感に関する評価値の可視化例を紹介する。

## 2 関連研究

Cabrera らによる FairVis[2] は、センシティブな属性を複合的にグループ化し、グループ間で発生する交差バイアスに注目する可視化解析システムであ

る。Munechika らによる VISUAL AUDITOR[5] は、モデルパフォーマンスの低いサブグループを可視化することで、モデルの偏りを監査し要約する。Ghai らによる D-BIAS[3] は、バイアスの特定と緩和のための HITL アプローチを具現化したビジュアルな対話型ツールである。

## 3 階層型データとしての偏りの可視化

3.1 節、3.2 節の処理は、我々が以前に報告した内容 [6] と同一である。

### 3.1 データの概要

提案手法では以下のデータを前提とする。A は人物集合によるデータ全体を表し、 $a_i$  は  $i$  番目の人物を表し、 $n$  はデータ中の人数を表す。

$$A = \{a_1, a_2, \dots, a_n\}$$

また、 $i$  番目の人物に相当する  $a_i$  は以下の変数を有するものとする。ここで  $e_i$  は可視化の対象となる実数値、 $g_i$  は  $i$  番目の人物の性別、 $r_{ij}$  は  $j$  番目の実数型変数の属性値、 $c_{ik}$  は  $k$  番目のカテゴリ型変数の属性値である。

$$a_i = \{e_i, g_i, r_{ij}, \dots, c_{ik}, \dots\}$$

### 3.2 木構造の生成

提案手法では、属性値  $r_{ij}$  または  $c_{ik}$  のうちユーザが選んだ複数の属性値を用いて人物群を階層的に分類し、木構造を構成する。この木構造の特定のノード配下の実数値  $e_i$  の偏りが見られるようであれば、ユーザが選択した複数の属性値がもたらす交差バイアスに起因する偏りが示唆される。

### 3.3 emd による分布間距離の算出

emd[7][8] は、ある分布をもう一方の分布に移動させるための最小コストとして定義される距離尺度である。本手法では、末端の葉ノード群に相当する

Copyright is held by the author(s). This paper is non-refereed and non-archival. Hence it may later appear in any journals, conferences, symposia, etc.

\* お茶の水女子大学

人物群のうち、男性と女性の  $e_i$  の分布間距離を emd を用いて算出し、この結果を男女間の分布の非類似度とする。emd の値が大きいほど男女の分布が異なり、反対に emd の値が小さいほど男女の分布が類似していることが示唆される。

### 3.4 「平安京ビュー」を用いた可視化

「平安京ビュー」[4] では葉ノードを正方形のアイコンで表現したのに対し、提案手法では葉ノード群に相当する人物群が有する  $e_i$  の分布を男女別に2列の帯グラフで表現する。帯グラフの各領域の色は HSI 表色系を採用し、以下の原則に沿って算出する。

**色相 (H):** emd が指定値より低い場合は男女共にグレースケールで、emd が指定値より高い場合は男性を青、女性を赤にする。

**彩度 (S):** 平均値に近いほど低く、最大値/最小値に近いほど高くする。

**明度 (I):** 値が大きいほど高くする。

## 4 空調温感データでの適用事例

本報告では空調の温感に関するオープンデータ [1] を適用した事例を示す。このデータから著者らは 32, 373 人を対象として以下の属性値を抽出した。

**TS:** 温感に対する 7 段階の評価値。

**Sex:** 生物学的な意味での性別。「男性」「女性」以外で回答した人物は本事例では対象外とした。

**Age:** 年齢。4 段階に分類。

**Cloth:** 服装の厚さの実数値。大きいほど厚い。4 段階に分類。

**Metab:** 代謝量に関する実数値。4 段階に分類。

**Season:** 春夏秋冬の 4 種類のカテゴリ値。

**Building:** オフィス・教室・住居・高齢者施設・その他の 5 種類のカテゴリ値。

**Strategy:** エアコン・換気・混合の 3 種類のカテゴリ値。

図 1 は Cloth, Season, Age の順に属性値を参照して人物を分類した可視化結果である。可視化画面左下の服装が最も厚着である A の枠の内側に注目する。秋に対応する (a) の枠において、全ての男女の帯グラフがカラースケールで表示されていることから、服装が最も厚着で季節が秋であると、年齢を問わず男女の温感の判断が食い違う傾向にあることがわかる。冬に対応する (b) の枠では、年齢が最も若い人物群を表す帯グラフのみカラースケールで表示されており、他はグレースケールである。このことから、服装が最も厚着で季節が冬であると、年齢が最も若い人物群では温感に男女差が生じるが、他

の年齢層では大きな男女差は生じていないことがわかる。冬に対応する (b) の枠だけでなく、水色の枠で囲まれた各部分でも同様な議論が成立する。続いて、可視化画面右上の服装が最も薄着である B の枠の内側に注目する。冬に対応する (c) の枠において、全ての男女の帯グラフがグレースケールで表示されている。このことから、服装が最も薄着で季節が冬であると、年齢を問わず、温感に大きな男女差が生じないことがわかる。オレンジ色の枠で囲まれた部分でも同様な議論が成立する。総じて、男女間の分布の違いは、データ全体ではなく、特定の属性や要素を持つ部分集合に見られることがこの可視化結果からも示唆される。

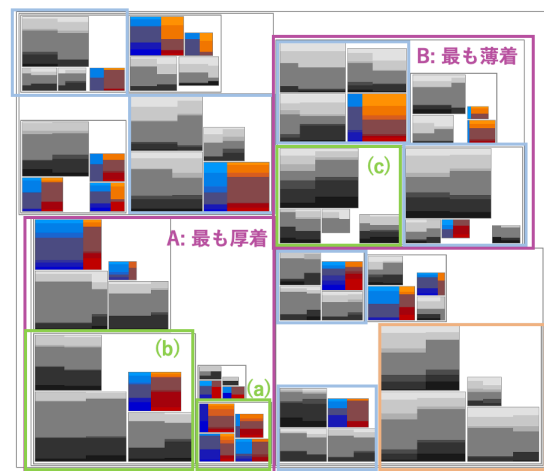


図 1. Cloth, Season, Age の順に人物を分類した例

## 5 まとめ・今後の課題

本研究では、多数の人物を対象としたデータ中に潜む男女差の可視化手法の拡張として、emd による男女間の分布差算出を導入した手法を提案した。データを構成する人物群を属性に沿って階層的に分類し、「平安京ビュー」を用いて可視化する。この際に、階層構造の特定のノード配下に属する人物群が有する数値を男女別に2列の帯グラフで表現し、男女の数値分布の分布間距離を emd を用いて算出する。emd が平均より低い場合はグレースケール、高い場合はカラースケールで帯グラフを表現することで、男女差の大きい部分に焦点を当てた可視化結果を表示する。本報告では空調の温感データを題材として可視化の実例を示した。

今後の課題として、帯グラフをカラースケールで表示する場合とグレースケールで表示する場合を分類するための適切な閾値の議論や、非常に多くの属性を有するデータにおいて、可視化する価値のある属性の組み合わせを自動選出する手法の開発が挙げられる。

## 参考文献

- [1] ASHRAE Global Thermal Comfort Database II. <https://www.kaggle.com/datasets/claytonmiller/ashrae-global-thermal-comfort-database-ii>.
- [2] Á. A. Cabrera, W. Epperson, F. Hohman, M. Kahng, J. Morgenstern, and D. H. Chau. FairVis: Visual analytics for discovering intersectional bias in machine learning. In *2019 IEEE Conference on Visual Analytics Science and Technology (VAST)*, pp. 46–56. IEEE, 2019.
- [3] B. Ghai and K. Mueller. D-BIAS: a causality-based human-in-the-loop system for tackling algorithmic bias. *IEEE Transactions on Visualization and Computer Graphics*, 29(1):473–482, 2022.
- [4] T. Itoh, Y. Yamaguchi, Y. Ikehata, and Y. Kajinaga. Hierarchical data visualization using a fast rectangle-packing algorithm. *IEEE Transactions on Visualization and Computer Graphics*, 10(3):302–313, 2004.
- [5] D. Munechika, Z. J. Wang, J. Reidy, J. Rubin, K. Gade, K. Kenthapadi, and D. H. Chau. Visual Auditor: Interactive Visualization for Detection and Summarization of Model Biases. In *2022 IEEE Visualization and Visual Analytics (VIS)*, pp. 45–49. IEEE, 2022.
- [6] Y. Nakai and T. Itoh. Hierarchical Data Visualization of Gender Difference: Application to Feeling of Temperature. In *27th International Conference on Information Visualisation (IV2023)*, pp. 178–183, 2023.
- [7] O. Pele and M. Werman. A linear time histogram metric for improved sift matching. In *Computer Vision–ECCV 2008: 10th European Conference on Computer Vision, Marseille, France, October 12–18, 2008, Proceedings, Part III 10*, pp. 495–508. Springer, 2008.
- [8] O. Pele and M. Werman. Fast and robust earth mover’s distances. In *2009 IEEE 12th international conference on computer vision*, pp. 460–467. IEEE, 2009.

## 未来ビジョン

本研究は可視化技術の発展と社会的問題の解決という二つの側面で未来ビジョンを持つ。

可視化技術の発展という側面では、「俯瞰性の高い可視化の追求」を目指す。様々な可視化表現を組み合わせ、反復的な操作繰り返すことによって可視化結果から知見を得たり、結論を導いたりするのではなく、最初の一画面で一定の知見を見出し、最小限の操作によって結論を導くことができるようなシステムを構築することは、短時間でのデータ分析を可能にする。この観点からデータを俯瞰する最初の一画面でデータ中の特定の部分に潜む偏りを発見させることに重点を置いており、これからも俯瞰性の高い可視化手法の開発を続けたい。また、データの偏りに限らず多くのデータ分析業務

分野において、俯瞰性の高い可視化手法の有用性を実証したい。

もう一つの社会的問題の解決という側面では、「社会的問題解決のための意思決定支援」を目指す。本研究の題材となったデータの偏りを発見することは、可視化を用いなくても数理的手法によって実現できる。しかし、データの偏りは文脈に左右されるため、必ずしも定量的に判断できるとは限らない。また、発見された偏りがどのくらい是正されることが望ましいのか、どのような手段で是正されるべきなのか、といった点を模索するには、人間による意思決定が必要である。これらの観点から、人間がデータを理解することを容易にし、人間の意思決定を支援する可視化技術の開発に努める。社会的問題解決に貢献できる可視化技術・ツールの構築を目指したい。