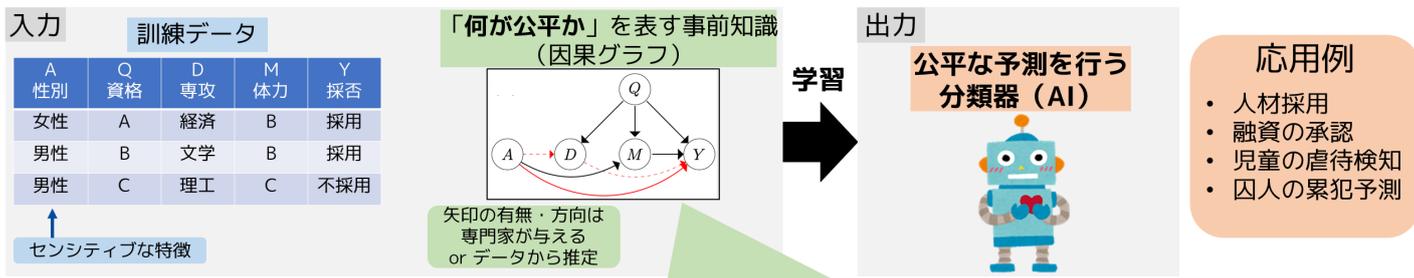


### みんなにとって公平な決め方、教えます

<b>どんな研究</b>	融資承認や人材採用、囚人の累犯予測など、人間に対する重大な意思決定を機械学習予測によって行うための研究です。この実現に向け、性別・人種・宗教・性的指向など、人間が持つセンシティブな特徴に関して公平で、かつ高精度な予測を行うための技術を考案しました。
<b>どこが凄い</b>	例えば「体力を要する職種なので体力が理由の不採用は差別的でない」といった要請に基づき、不要な制約なしに学習することで予測精度を高めます。従来法ではデータが特定のモデルから生じている場合のみ公平性を保証可能でしたが、提案法では様々なデータで公平性を保証できます。
<b>めざす未来</b>	「何が公平か」は人や社会の価値観とともに移り変わるものです。社会的な要請と技術的な限界に向き合いながら、より柔軟に社会的要請に応じられる機械学習技術を実現し、一人一人が不利益を被ることなく、同時に、効率的な意思決定を行える社会に貢献できればと考えています。

### 問題：人間に対して公平な予測（意思決定）を行うAIを学習



### ポイント1：高精度な予測を実現

#### 例：体力を要する職種における人材採用の場合

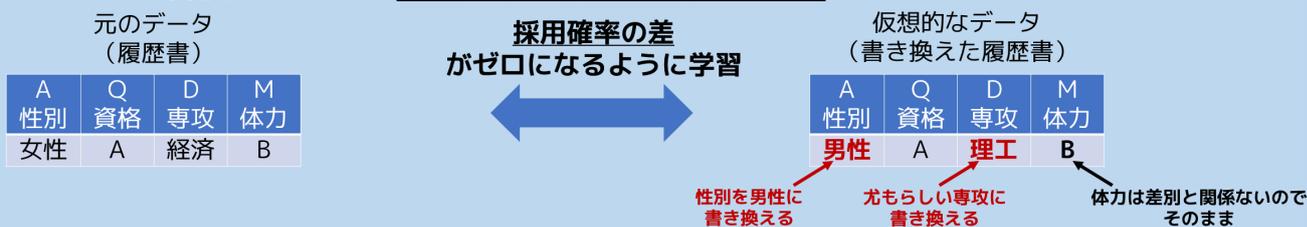
個々の意思決定に求められる要請を事前知識として活用することで不要な制約を課すことなく学習し高精度な予測を実現

- 性別Aによって採否Yを決定するのは**差別的** (A→Y)
- 専攻Dによって採否Yを決定して採用率に性差が生じるのは**差別的** (A→D→Y)
- (体力は必要なので) 体力Mによって採否Yを決定するのは**差別的でない** (A→M→Y)
  - 1., 2.によって生じる採用率の性差には制約を課す必要がある
  - 3.によって生じる採用率の性差に制約を課すのは**不要** (予測精度を下げてもいい)

### ポイント2：様々なデータに対し個人公平性を保証可能

どのような関数形（モデル）から生じたデータでも各個人一人一人に対し**予測が公平**になるように学習できる

#### 各個人ごとに、差別による決定結果のちがいをゼロにする



### 難しさ

履歴書を適切に書き換えるためにはデータの生成式を正しく表す必要があり、扱いやすい関数形（モデル）から生じたデータでなければ近似不可能◎

### 提案法

- データの生成式によらず推定できる差別度を新たに提案
- この差別度をゼロにすれば**差別による決定結果のちがいをゼロ**にできる

### 関連文献

[1] Y. Chikahara, S. Sakaue, A. Fujino, H. Kashima, "Learning Individually Fair Classifier with Path-Specific Causal-Effect Constraint," in Proc. the 24-th International Conference on Artificial Intelligence and Statistics (AISTATS), 2021.

### 連絡先

近原 鷹一 (Yoichi Chikahara) 協創情報研究部 知能創発環境研究グループ  
Email: cs-openhouse-ml@hco.ntt.co.jp