性差に基づく科学技術イノベーション検討 のための話題提供

- AI における公平性 -

NTTコミュニケーション科学基礎研究所 (兼) 理化学研究所革新知能統合研究センター

上田 修功

本日の内容

- 1. AI(機械学習) 技術のおさらい
- 2. AIにおける(性)差別の話題

本日の内容

- 1. AI(機械学習) 技術のおさらい
- 2. AIにおける(性)差別の話題

AI (=機械学習)の定義

明示的にプログラミングすることなく, コンピュータに学ぶ能力 を 与えようとする研究分野 (A. L. Samuel, 1959)

機械学習分野では、経験から自動的に改善を図れるような コンピュータプログラムを構築する方法について議論している (T. M. Mitchell, 1999)

人工知能(AI)ブームの変遷

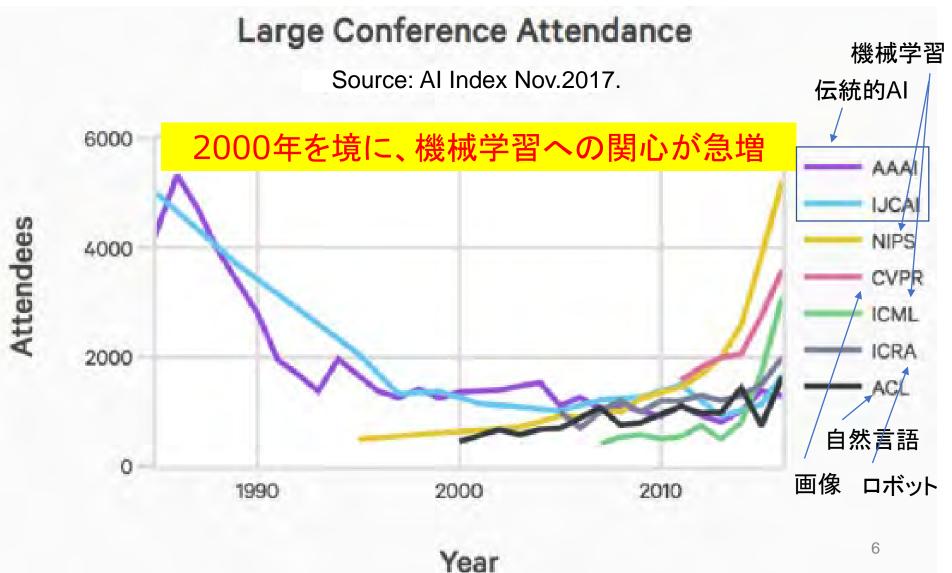
第1次AIブーム(1950~1960年代) 探索&推論

第2次AIブーム(1980年代) 知識表現, 学習

第3次AIブーム(2000年代) ビッグデータ+深層学習

第3次AIブーム=機械学習ブーム

現在の機械学習ブーム=深層学習ブーム



学習のタイプ

人の学習と同じ





教師なし学習

(自習する)

ex)特徴選択、特徴抽出、 次元圧縮、クラスタリング

現在、深層学習はこれら全てに貢献

教師あり学習

(先生に習う)

ex)画像認識, 音声認識, 言語翻訳などのパターン認識全般



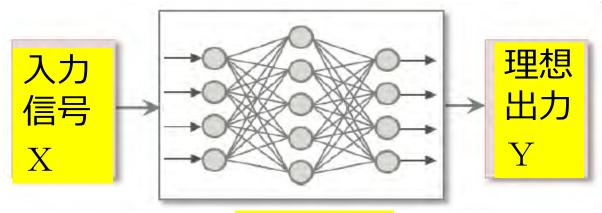
強化学習

(試行錯誤して報酬を最大化する)

ex)歩行ロボット, ゲーム(アルファ碁)

教師あり学習(主流)

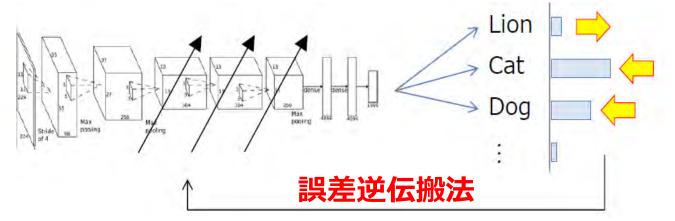
 $Y=f(X; \Theta)$



学習機械

・学習データ: D = {x,y} y: 教師データ





ニューラルネットワークの変遷



1958 Perceptron F. Rosenblatt

1986: Nature論文: 逆誤差伝搬法 Rumelhart, Hinton, Williams

1974 Backpropagation
P. Werbos

CNN Y. LeCun





Convolution Neural Networks for Handwritten Recognition

1998



Google Brain Project on 16k Cores 2012

第1次ニューロブーム

awkward silence (Al Winter)

第3次ニューロブーム

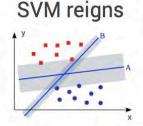
L.Faifei

1969 Perceptron criticized



線形分離不能なケースでは限界 (by M. Minsky: Alの父が 第1次ニューロブームを終焉させた)

第2次ニューロブーム 1995



V. Vapnik

2006 Restricted Boltzmann Machine



2012 AlexNet wins ImageNet

IM. GENET

1400万画像、2万 クラスの良質な 画像データベース

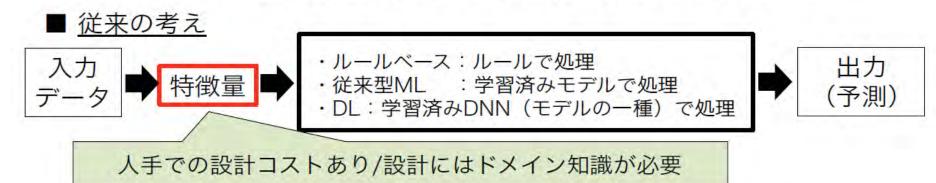
G. Hinton

2009年、DNNが音声認識において従来の state-of-the-art手法の性能を超えた!

Source: https://www.slideshare.net/LuMa921/deep-learning-a-visual-introduction

深層学習の特徴

- ・入力と教師データの対である学習データのみを必要とし、 モデルの構築は不要(データ任せ)
- ・アルゴリズム開発のようなプロフェッショナルな知識は不要 (フリーのライブラリーも整備!)
- 深層学習は、特徴量の抽出ステップを省略しても高い精度が出せる.



■ 深層学習は下記フローでも精度が出せる (パラダイムシフト)



ImageNet Large Scale Visual Recognition Challenge (ILSVRC)

ImageNetのデータセットを利用して1000クラスの物体認識を行う



Stanford大により作成された、深層学習躍進の一躍を担っ ■ Company of the Land Stanford大により作成された、深層字習躍進の一躍を担った大規模(1400万画像、2万クラス)かつ良質なデータセット



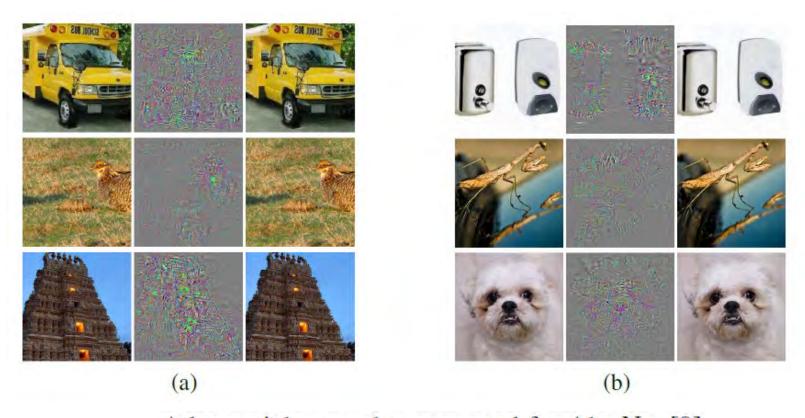
Prof. Li Feifei

http://image-net.org

深層学習は人間を超える?

Intriguing properties of neural networks

Christian Szegedy et al.



Adversarial examples generated for AlexNet [9].



Minimize $||r||_2$ subject to:

1.
$$f(x+r) = l$$

2.
$$x + r \in [0, 1]^m$$

深層学習の限界

データ駆動型アプローチの課題

・観測データと教師信号が大量に与えられた教師有り学習では 最強のツール(画像、音声、自然言語処理では、当深層学習は 必須要素技術)

・ただし、学習データが容易に準備できない応用(異常現象 (例:故障、災害などの希少データからの推定、予測)や、 結果に対する説明が必要な応用ではDNNの適用は困難

本日の内容

- 1. AI(機械学習) 技術のおさらい
- 2. AIにおける(性)差別の話題

AIにおける(性)差別関連事例

■機械翻訳で代名詞に性別の差がないトルコ語をhe/sheの差のある英語に翻訳すると、職業との共起に差が生じる。

https://medium.com/coinmonks/ai-doesnt-have-to-be-conscious-to-be-harmful-385d143bd311

■AmazonのAI(人工知能)を活用した人材採用システム(履歴書に書かれた約5万個のキーワードを抽出・分析)は、女性を差別するという機械学習面の欠陥が判明し、運用を取りやめる結果になった。

https://jp.reuters.com/article/amazon-jobs-ai-analysis-idJPKCN1ML0DN

- 2019年UNESCO(ユネスコ:国連教育科学文化機関)は、AI音声アシスタントのデフォルトの声が女性となっていることはジェンダーの偏りを強めると主張する報告書を発表した。
- ■顔認識ソフトウェアにおける人種差別: Googleフォトのソフトウェアは、写真に 写っている物体や顔を認識して、人や動物などに基づいて分類し、写真を保存・整理 するが、ある男性、ジャッキー・アルシネは、アフリカ系アメリカ人の友人の一人が写真 の中でゴリラと表示されていた

https://medium.com/coinmonks/ai-doesnt-have-to-be-conscious-to-be-harmful-385d143bd311

■アフリカ系アメリカ人の人名でGoogle検索すると、逮捕記録を示唆する広告が表示される。 その広告には個人の逮捕記録をチェックできるサービスを提供するリンク(Instant Checkmate)が設定されていた。 https://queue.acm.org/detail.cfm?id=2460278

何故AIで(性)差別が生じるのか

機械学習タスクは、学習データ作成(特徴量の選定、 データの収集)および、目的関数の設定からなるが、 これらに(暗に)差別が混入する

→ AIが差別を生み出すのではなく、むしろ社会の差別や 偏見が差別的なAIを生み出している

例:SNSのAIチャボットが差別的な発言をする

→ データ駆動型アプローチであるAI技術においては、 個々人の(性)差別を解消しない限り、差別の根絶 は困難

人間中心のAI社会原則 内閣府AI戦略

(6) 公平性、説明責任及び透明性の原則

「AI-Ready な社会」においては、AI の利用によって、人々が、その人の持つ背景によって不当な差別を受けたり、人間の尊厳に照らして不当な扱いを受けたりすることがないように、公平性及び透明性のある意思決定とその結果に対する説明責任(アカウンタビリティ)が適切に確保されると共に、技術に対する信頼性(Trust)が担保される必要がある。

- AI の設計思想の下において、人々がその人種、性別、国籍、年齢、政治的信念、宗教等の多様なバックグラウンドを理由に不当な差別をされることなく、全ての人々が公平に扱われなければならない。
- ▶ AI を利用しているという事実、AI に利用されるデータの取得方法や使用方法、 AI の動作結果の適切性を担保する仕組みなど、用途や状況に応じた適切な説明が得られなければならない。
- ▶ 人々が AI の提案を理解して判断するために、AI の利用・採用・運用について、 必要に応じて開かれた対話の場が適切に持たれなければならない。
- ▶ 上記の観点を担保し、A! を安心して社会で利活用するため、A! とそれを支える データないしアルゴリズムの信頼性(Trust)を確保する仕組みが構築されなければならない。

https://www8.cao.go.jp/cstp/ai/aigensoku.pdf

機械学習と公平性に関する声明

2019年12月10日

人工知能学会 倫理委員会 日本ソフトウェア科学会 機械学習高額研究会 電子情報通信学会 情報論的学習理論と機械学習研究会

- 1. 機械学習は道具にすぎず人間の意思決定を補助するものであること
- 2. 私たちは、公平性に寄与できる機械学習を研究し、社会に 貢献できるよう取り組んでいること

関連シンポジウム:機械学習と公平性に関するシンポジウム (2020年1月9日)

1. 機械学習は道具にすぎません

機械学習はあくまでも道具にすぎず、その使い方を定めるのは人間です。機械学習は人類社会の繁栄に大きく貢献できる可能性を秘めているとともに、不適切な利用をすれば人類社会の利益に反する可能性もあります。機械学習は過去の事例に基づいて未来を予測しますから、偏りのある過去に基づいて予測する未来は、やはり偏りのあるものになりかねません。もし、過去と異なる「あるべき未来」を求めるのであれば、機械学習による予測や判断が公平性を欠くことがないように人間が機械学習に注意深く介入する必要があります。

同時に、「何が公平か」については、科学技術や工学だけの問題ではなく、現在の人類社会が何を求めているか、という価値観の問題抜きには語れません。機械学習という「道具」を正しく使うためには、それが「公平性」という私たち人類社会の価値観に対して、どのような影響を与えるかを正しく理解し、そのリスクを評価し、方策について合意しなければならないのです。この点は、私たちだけではなく、機械学習に携わる技術者や利用者、経営者、そして組織や社会の全体が把握し向き合っていく必要があります。

引用元: http://ai-elsi.org/wp-content/uploads/2019/12/20191210MLFairness.pdf

2. 私たちは機械学習で公平性に寄与します

私たちは、機械学習の利用が社会の不利益になってはならないと考え、この問題を解決するために、行動指針と技術開発の双方から真摯に取り組んでいます。IEEE Ethically Aligned Design では機械学習の不適切な利用ないしは誤用、悪用を戒め、その対策を具体的に記述しています[3]。人工知能学会では、自らの社会における責任を自覚し、社会と対話するために、学会会員の倫理的な価値判断の基礎となる倫理指針を2017年に定めました[4]。我が国社会の様々なステークホルダ(その一部は、私たちでした)が集まって、高度な情報技術を社会でどのように使っていくべきかを議論し、その結果が、内閣府「人間中心の AI 社会原則」として2019年3月に公開されました[5]。その基本理念の1つは多様性と包摂であり、高度な情報技術の利用にあたっては「公平性のある意思決定とその結果に対する説明責任」を担保するように求めています。

これらに呼応して、私たちも公平性の様々な側面をいかに定量的に評価し、実現していくかについての研究を進めています。最近の主要な研究集会では必ず機械学習の公平性に関する研究発表がありますし、世界的にも公平性に関する研究論文の数は増えています。実は「公平性とは何か」を機械学習の言葉で数理的に突き詰めていくと、多数のバリエーションがあることがわかります。人々が何を公平と考えるか、様々な基準を機械学習の言葉で表現しなおすことによって、「公平」という概念をより明確なものにしていくこともできるのです。このように、私たちは、機械学習によって公平性に起きうる問題を防ぐだけでなく、機械学習をきっかけとして公平性のあり方を定義、議論することにも真摯に取り組んでいます。

引用元: http://ai-elsi.org/wp-content/uploads/2019/12/20191210MLFairness.pdf

Fairness-aware Machine Learning

公平性を配慮した機械学習技術 2017年ごろより研究が加速(EUのGDPRの影響?)

各種のバイアス Barocas, 2016

- データバイアス: データ作成者の偏見や認知バイアスなどにより学習データに偏り(bias)が生じる
- 標本選択バイアス: 予測対象の集団が学習データに含まれていないことによるバイアス
- 帰納バイアス: 少数事例を例外、外れ値として扱うことによるバイアス(例:データの大半が男性データ)

Biased algorithms are easier to fix than biased people

Mullainathan 2019

Note: Bugbears or Legitimate Threats? (Social) Scientists' Criticisms of Machine Learning

公平性の形式的定義

S: モデル中のセンシティブ特徴 S=1(0): 配慮不要(要)な特徴

X: その他の特徴

Y: Y=1(0): 有利(不利)な判定、 $\widehat{Y}:$ 予測値

 $ightharpoonup データバイアスの解消 <math>\hat{Y} \perp \!\!\! \perp S$ (\hat{Y} とSが統計的に独立)

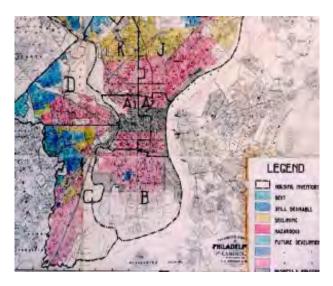
不利な判定と有利な判定がなされる割合をS=1,0の両方のケースで同じとすることでバイアスを解消する

Note: Red-Lining Effect Calders, 2010

ローン貸付の人種差別

住居情報が間接的に人種情報が 使われている

XとSの統計的独立性が重要



Wikipedia

■ 選択バイアスの解消

観測データセットが母集団の代表となっていない

- → ランダムサンプリング、層別化
- 帰納バイアスの解消: (Yの値は正しいと仮定)
 - $\Rightarrow \hat{Y} \perp S \mid Y$ Zafar, 2017

$$P(\hat{Y}=1|Y=0,X,S=1) = P(\hat{Y}=1|Y=0,X,S=0)$$
 かつ $P(\hat{Y}=1|Y=1,X,S=1) = P(\hat{Y}=1|Y=1,X,S=0)$

偽陽性率($Y=0を\hat{Y}=1$ と誤る率)と真陽性率(Y=1を $\hat{Y}=1$ と 正しく判定する率)をS=0,1の両者のケースで等しくする

集団公平性と個人公平性

集団公平性: 個々人は集団として公平に扱われる

 \rightarrow P(Y | S=s) = P(Y) for all s in S

例: 男性と女性の各グループで平均的に取り扱いが 等しければOK

個人公平性: 個々人が公平に扱われる

P($Y \mid S, X=x$) = P($Y \mid X=x$) for all x in X Xが所与の下で、YとSが統計的に独立

例: 性別以外の特徴が全て同じであれば、同一に扱われるべき

データマイニング,機械学習における公平性関連国際会議

- ICDM2012併設WS:
 Discriminant and Privacy-aware Data Mining
- NIPS2013, ICML2014 併設WS: Fairness, Accountability, and Transparency in Machine Learning
- NIPS2016シンポジウム:
 Machine Leaning and Law
- ICDM2016併設WS: Privacy and Discriminant in Data Mining
- KDD2019 チュートリアル:
 Fairness-aware Machine Learning: Practical Challenges and Lessons Learned

まとめ

- AI(機械学習) は道具に過ぎず、データ収集、アルゴリズム 設計、運用面で人間を介して差別が生まれ得る
- AIにおける公平性の議論は2017年代より活発化し、機械学習の国際会議でもWSやチュートリアル等で不公平性を防ぐ技術の研究が進行している
- ■公平性とは何か、公平性の在り方についての議論が必要