



コホートデータで 地域社会のみらいをつくる

- Making the Future of Communities
with Cohort Data -

Aug/09/2022



東京大学 未来ビジョン研究センター アクサ生命保険 株式会社

平松 雄司

- 保険数理, 統計
- 公衆衛生
- データサイエンス





- 終末期医療費に関する研究
- 生活習慣病に関する研究
 - 高血圧症・脂質異常・高血糖等
 - Metabolic syndrome
- 医療画像 AI研究のアドバイザー
- データサイエンス講座

大学・会社内に限定せず
データ分析を促進する活動もしています



Kaggle TensorFlow Help Protect
the Great Barrier Reef Object
Detection Competition
Feb. 2022, 15th / 2025 teams



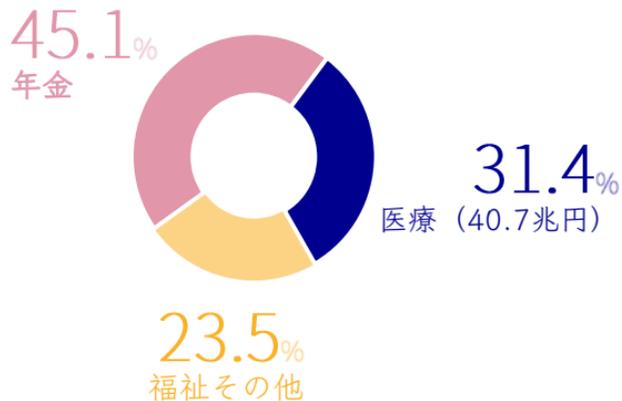


1. 日本の現状と課題
2. 国保データの活用例
3. コホートデータをより活かすために

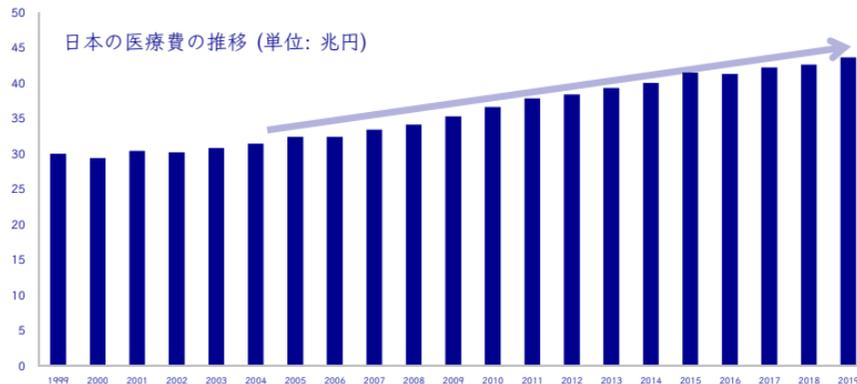


1 日本の現状と課題

高齢化にともない増大する医療費



社会保障給付費 (計129.6兆円, 2021年度予算ベース)
出典: 厚労省



OECD Health Statistics 2021, Health expenditure and financing
Financing scheme: Government/compulsory
<https://stats.oecd.org/>

健康増進および健康寿命の延伸による
持続可能な未来の構築は日本にとって喫緊の課題

東京大学とアクサ生命による研究（2019 - 現在）

● 終末期医療費

- 終末期医療費（死亡前医療費）は年間医療費の20%前後を占めており、社会的負担は決して軽くない（Hashimoto et al. 2010）
- 終末期医療費の要因に関する研究（後述）

● 生活習慣病

- 我が国の循環器系疾患の罹患率・死亡率・DALYsなどは改善が頭打ちとなっており、近年は悪化傾向にある（Hata et al. 2013）
- 特定健診の診断結果や保健指導に関する研究

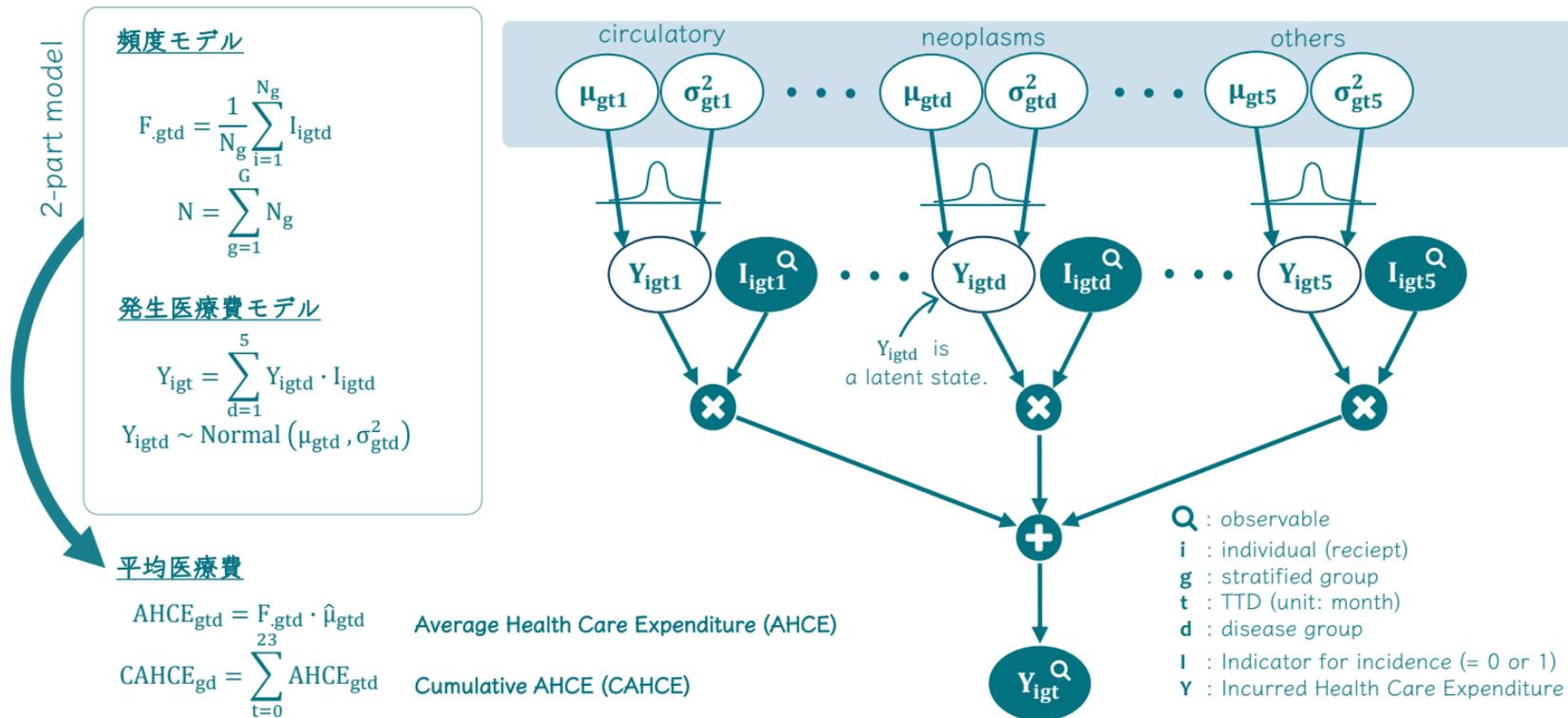


国保データの活用例

終末期医療費の要因研究

- 分析対象とした要因
 - 性・年齢
 - 死への近接性 (Proximity To Death: PTD)
 - 傷病
- 静岡県の国民健康保険加入者のうち、65 - 95 歳、2012 - 2018 年内に死亡した人を対象
- 死亡前二年間における各月の医療費を性・年齢・傷病グループ別に集計し分析
- 傷病グループ
 - I00 - I99 : 循環器系 (circulatory)
 - N18 : 慢性腎臓病 (CKD)
 - C00 - D48 : 悪性新生物 (neoplasms)
 - J00-J99 : 呼吸器系 (respiratory)
 - Others : その他の ICD10 全て
- ベイズ統計を利用した傷病グループ別医療費への分解 (**cost allocation**)
レセプトの医療費は月単位でまとまっており、傷病別に割り当てられていない

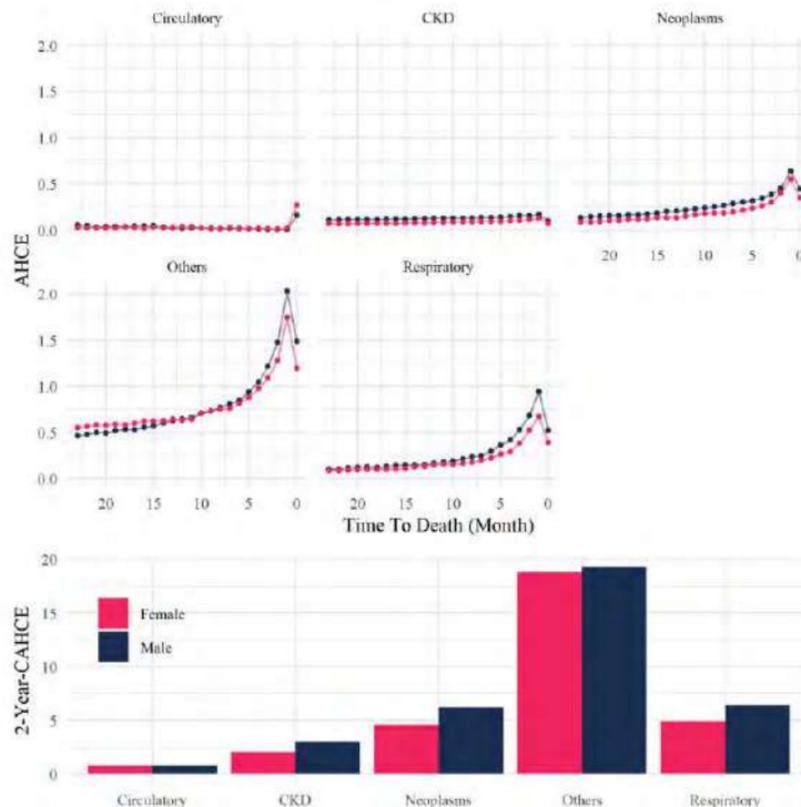
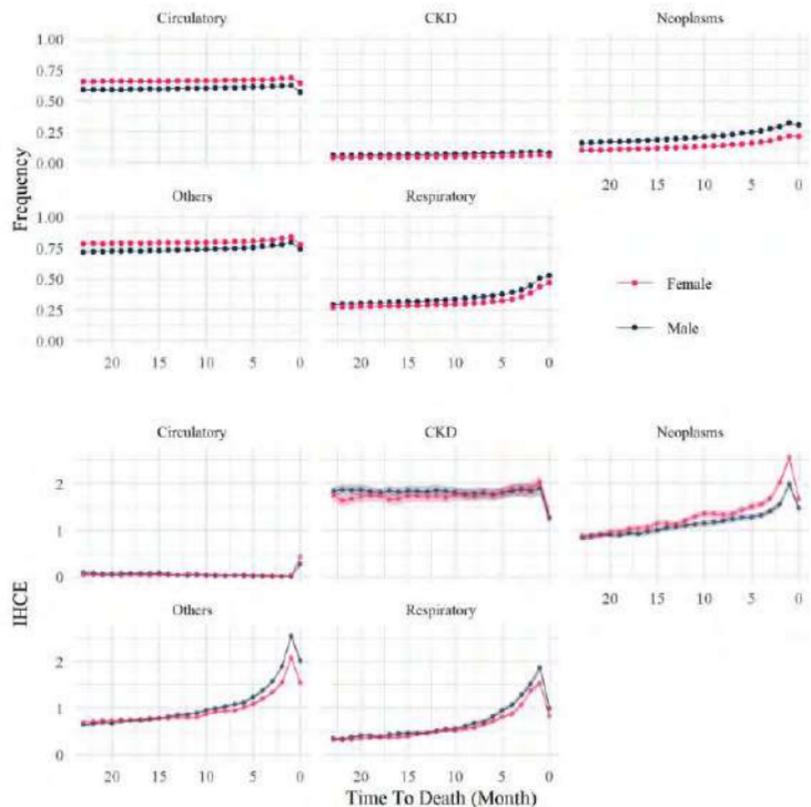
本研究提案の傷病グループ別の医療費割り当て手法



各傷病グループにかかる医療費（平均値： μ_{gtd} ）を統計的に推定する手法
損害保険の数理領域でよく使用される2-partモデルを元に着想

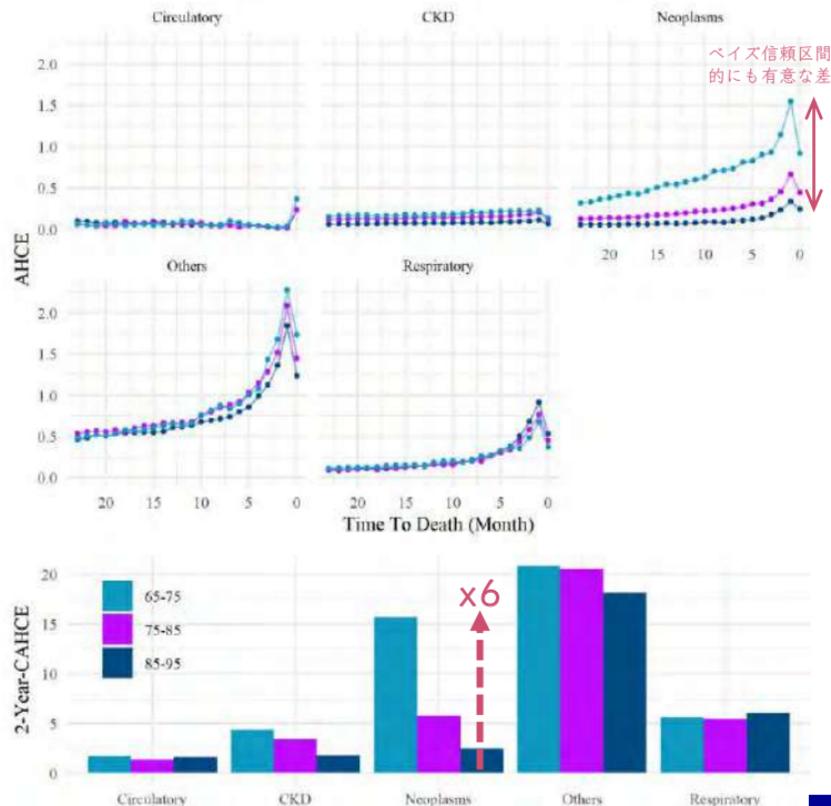
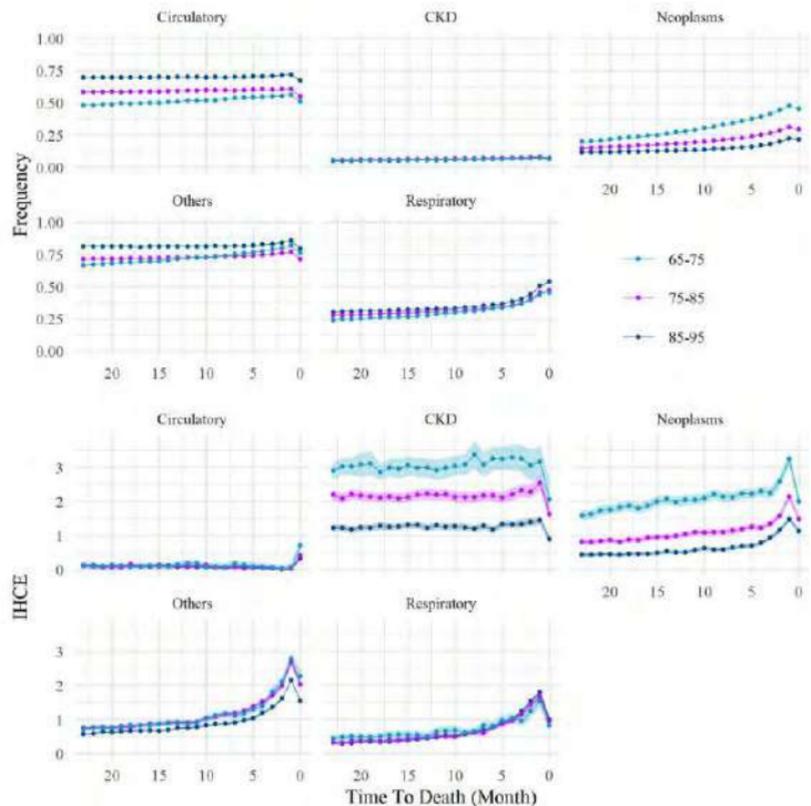
男女別：各月の死亡前医療費

単位 (10万円)



年齢グループ別: 各月の死亡前医療費

単位 (10万円)



本研究で使用的是のは
population-based cohort
であり，日本の都道府県
のうちの一つのみ



静岡県



ビッグデータと産学連携のもつ可能性

全国の都道府県から集めたデータを突合すれば，
より精緻で細かな傷病グループ別の分析が可能に



BMC Health Economics Review
in Springer Nature Group

[Examining proximity to death and health care expenditure by disease: a Bayesian-based descriptive statistical analysis from the National Health Insurance database in Japan. Health Econ Rev 12, 6 \(2022\). https://doi.org/10.1186/s13561-021-00353-9](https://doi.org/10.1186/s13561-021-00353-9)

Hiramatsu et al. Health Economics Review (2022) 12:6
<https://doi.org/10.1186/s13561-021-00353-9>

Health Economics Review

RESEARCH

Open Access



Examining proximity to death and health care expenditure by disease: a Bayesian-based descriptive statistical analysis from the National Health Insurance database in Japan

Yuji Hiramatsu^{1,2*}, Hiroo Ide¹, Atsuko Tsuchiya¹ and Yuji Furuji¹

1: UTokyo
2: AXA Life Japan
3: Shizuoka prefecture

Abstract

Background: Japan is one of the Organization for Economic Co-operation and Development (OECD) countries where population aging and increasing health care expenditures (HCE) are urgent issues. Recent studies have identified factors other than age, such as proximity to death and morbidity, as contributing factors to the increase in medical costs. It is important to assess HCE by disease and analyze their factors to estimate and improve future HCE.

Methods: We extracted individual records spanning approximately 2 years prior to the death of persons aged 65 to 95 years from the National Health Insurance data in Japan, and used a Bayesian approach to decompose monthly HCE into five disease groups (circulatory, chronic kidney disease, neoplasms, respiratory, and others). The relationship between the proximity to death and the average HCE in each disease group was stratified by sex and age and analyzed using a descriptive statistical method similar to the two-part model.

Results: The average HCE increased rapidly as death approached in most disease groups, but the increase-pattern differed greatly among disease groups, sex, and age groups. The effect of proximity to death on average HCE was small for chronic diseases, but large for lethal diseases. When stratified by age and sex, younger and male decedents tended to have higher average HCE, but the extent of this varied by disease group. The two-year cumulative average HCE for neoplasms in the 65–75-years age group was about six times larger than those in the 85–95-years age group.





3
コホートデータをより活かすために

データサイエンスにおけるデータの価値

画像認識 (Computer Vision)

2012年のAlexNet (Krizhevsky et al. 2012) を皮切りに深層学習による画像認識は今日までに大きな進歩を遂げた。近年は、Google researchによるattention(Transformer)-basedのモデルも提案されている(Dosovitskiy et al. 2020)。

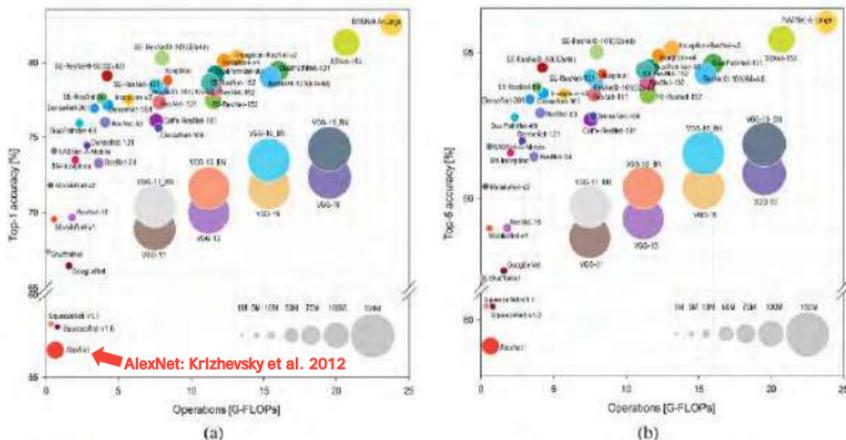


FIGURE 1. Ball chart reporting the Top-1 and Top-5 accuracy vs. computational complexity. Top-1 and Top-5 accuracy using only the center crop versus floating-point operations (FLOPs) required for a single forward pass are reported. The size of each ball corresponds to the model complexity. (a) Top-1; (b) Top-5.

Bianco et al. 2018

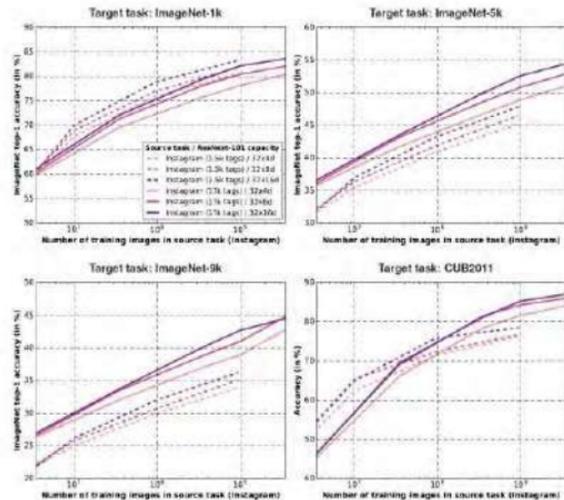


Fig. 2: Classification accuracies on IN-1k, 5k, 9k and CUB2011 target tasks as a function of the number of Instagram images used for pretraining for three network architectures (colors) and two hashtag vocabularies (dashed / solid lines). Only the linear classifier is trained on the target task. Higher is better. Mahajan et al. 2018

深層学習では、おおよそ、データ数の対数に対して線形にモデル性能が向上する。上図のx軸の右端は35億枚の画像データを学習に使用したことを示しているが、まだモデルの性能が伸びる途中にある。

良い品質のデータはあればある程良い

自然言語処理 (Natural Language Processing)

Google Brainによる2017年のTransformer (Vaswani et al. 2017), そして, Google AIによる2018年のBERT (Devlin et al. 2018) の提案以降, 爆発的な進歩を遂げている。

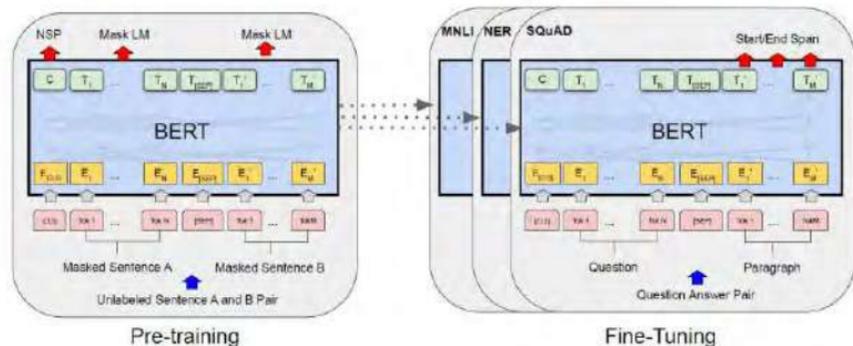


Figure 1: Overall pre-training and fine-tuning procedures for BERT. Apart from output layers, the same architectures are used in both pre-training and fine-tuning. The same pre-trained model parameters are used to initialize models for different down-stream tasks. During fine-tuning, all parameters are fine-tuned. [CLS] is a special symbol added in front of every input example, and [SEP] is a special separator token (e.g. separating questions/answers).

Devlin et al. 2018

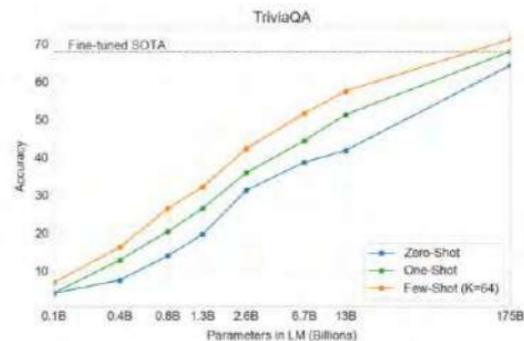


Figure 3.3: On TriviaQA GPT3's performance grows smoothly with model size, suggesting that language models continue to absorb knowledge as their capacity increases. One-shot and few-shot performance make significant gains over zero-shot behavior, matching and exceeding the performance of the SOTA fine-tuned open-domain model, RAG [LPP⁺20]

近年は, 事前学習として, 巨大なパラメータ数 (GPT-3 175B: 1750億) をもつモデルを, 巨大なデータ (Common Crawl: 570 GB, 元は50 TB程度) で学習させるのがトレンドになっている (Brown et al. 2020).

その結果, Big Tech以外の組織にとって, 研究のハードルが上がってきている

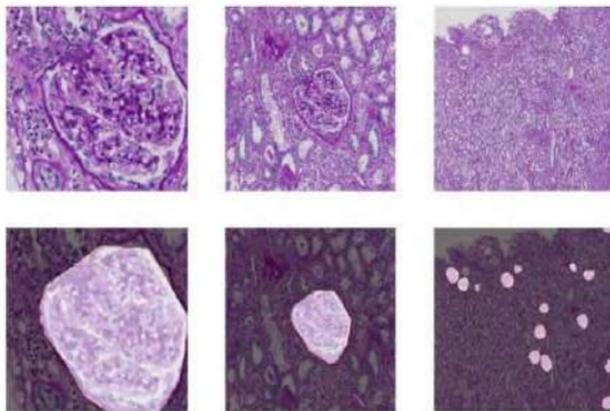
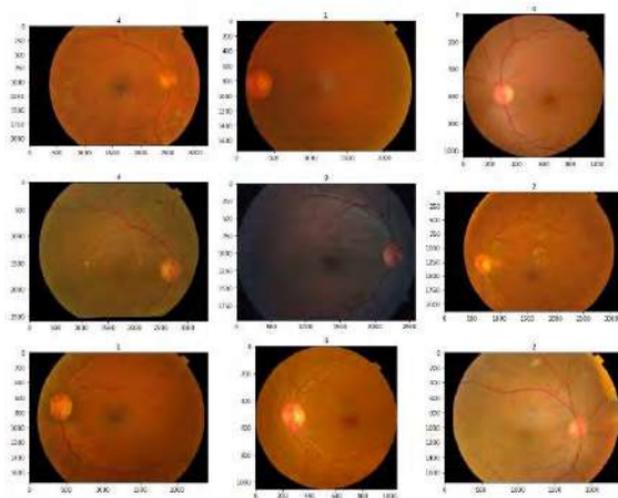
医療領域における深層学習

Computer Aided Diagnosis (CAD)

- 人力では困難な生検などの自動化
- 医師の診断の補助ツール

臨床の現場に、深層学習を利用した装置やツールが導入され始めている

Kaggle
APTOS 2019 Blindness Detection
糖尿病網膜症の重症度検出



Kaggle
HuBMAP - Hacking the Kidney

PAS染色された腎臓の顕微鏡画像における系
球体のセグメンテーション

医療データにおける問題点

- 医療（画像）データは手に入りやすく、希少疾患であれば尚更（Pati et al. 2022）
- 欧州のGDPRや米国のHIPAAにみられるように、個人情報保護に関する法律が各国で整備されており、それらの内容は異なっている。そのため、各国をまたいで世界中からデータを集める（Centralised Data Lake）にはハードルがある。
- うまく匿名化できているように思えても、わずかなデータ要素で識別が可能な場合がある（Rocher et al. 2019）
- 学習に使用するデータは、課題の背景にあるデータの分布を再現していることを前提としている。学習データが十分ではなく、学習データ内にない分布をもつデータに対して、良好な汎化性能を維持することは原理上難しいことが多い。

Federated Learning (連合学習)

Rieke et al. 2020.

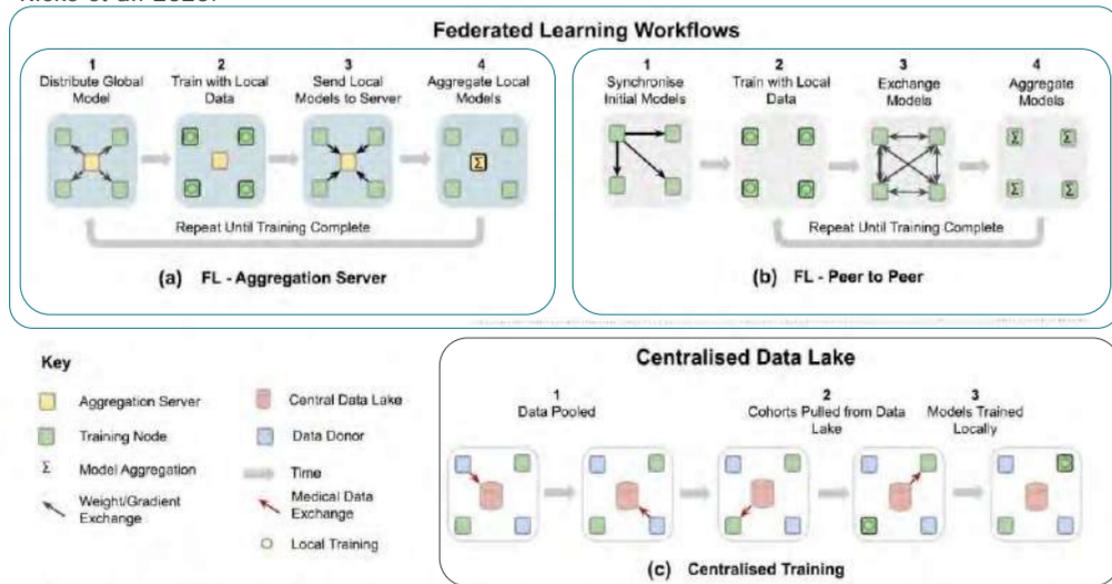


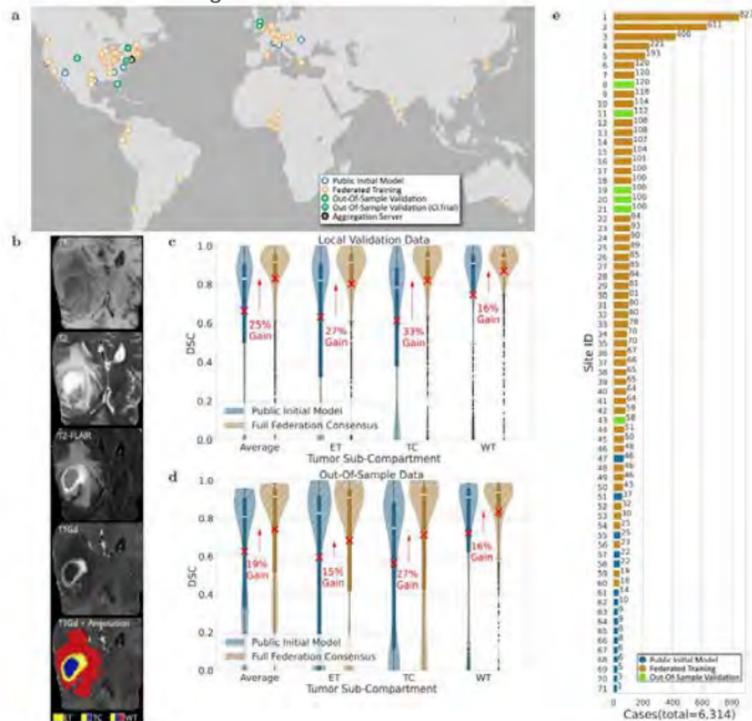
Fig. 1 Example federated learning (FL) workflows and difference to learning on a Centralised Data Lake. **a** FL aggregation server—the typical FL workflow in which a federation of training nodes receive the global model, resubmit their partially trained models to a central server intermittently for aggregation and then continue training on the consensus model that the server returns. **b** FL peer to peer—alternative formulation of FL in which each training node exchanges its partially trained models with some or all of its peers and each does its own aggregation. **c** Centralised training—the general non-FL training workflow in which data acquiring sites donate their data to a central Data Lake from which they and others are able to extract data for local, independent training.

- 各機関同士がやり取りするのはモデルのパラメータのみ (e.g. 画像認識モデルの重み係数など)
- 中央集権的なデータのストア (Centralised Learning) を行わないため、各医療機関のもつデータの秘匿性が保たれる
- 大別して
(a) aggregation server
(b) peer to peer
の二通りの方式がある

Federated Learning 応用例

The Federated Tumor Segmentation initiative (FetS)

Pati et al. 2022. Fig. 1



希少疾患である脳内悪性腫瘍（Glioblastoma）
の領域検出（Segmentation Task）

- 六カ国内の71施設が参加する巨大プロジェクト
- mpMRI: T1/T2/T1Gd/T2-FLAIR
- 25,256 MRI images from 6,314 patients
- 腫瘍の領域検出において、大幅な精度向上を達成（~ 25%）

大学・企業・国がデータ・人材で協力し、
持続可能な未来へ

References

1. Hashimoto H, Horiguchi H, Matsuda S. Micro data analysis of medical and long-term care utilization among the elderly in Japan. *International Journal of Environmental Research and Public Health*. 2010;7:3022–37.
2. Hata J, Ninomiya T, Hirakawa Y, Nagata M, Mukai N, Gotoh S, et al. Secular Trends in Cardiovascular Disease and Its Risk Factors in Japanese. *Circulation* 2013;128:1198–205. <https://doi.org/10.1161/CIRCULATIONAHA.113.002424>.
3. Hiramatsu Y, Ide H, Tsuchiya A, Furui Y. Examining proximity to death and health care expenditure by disease: a Bayesian-based descriptive statistical analysis from the National Health Insurance database in Japan. *Health Econ Rev*. 2022;12:6.
4. Krizhevsky A, Sutskever I, Hinton GE. ImageNet Classification with Deep Convolutional Neural Networks. In: *Advances in Neural Information Processing Systems*. Curran Associates, Inc.; 2012.
5. Dosovitskiy A, Beyer L, Kolesnikov A, Weissenborn D, Zhai X, Unterthiner T, et al. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. 2021.
6. Bianco S, Cadene R, Celona L, Napolitano P. Benchmark Analysis of Representative Deep Neural Network Architectures. *IEEE Access*. 2018;6:64270–7.
7. Mahajan D, Girshick R, Ramanathan V, He K, Paluri M, Li Y, et al. Exploring the Limits of Weakly Supervised Pretraining. 2018.
8. Vaswani A, Shazeer N, Parmar N, Uszkoreit J, Jones L, Gomez AN, et al. Attention is All You Need. 2017.
9. Devlin J, Chang M-W, Lee K, Toutanova K. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. 2019.
10. Brown TB, Mann B, Ryder N, Subbiah M, Kaplan J, Dhariwal P, et al. Language Models are Few-Shot Learners. 2020.
11. Common Crawl. <https://commoncrawl.org/>.
12. APTOS 2019 Blindness Detection | Kaggle. <https://www.kaggle.com/competitions/aptos2019-blindness-detection>.
13. HuBMAP - Hacking the Kidney | Kaggle. <https://www.kaggle.com/competitions/hubmap-kidney-segmentation/overview>.
14. Pati S, Baid U, Edwards B, Sheller M, Wang S-H, Reina GA, et al. Federated Learning Enables Big Data for Rare Cancer Boundary Detection. 2022.
15. Rocher L, Hendrickx JM, de Montjoye Y-A. Estimating the success of re-identifications in incomplete datasets using generative models. *Nat Commun*. 2019;10:3069.
16. Rieke N, Hancox J, Li W, Milletari F, Roth HR, Albarqouni S, et al. The future of digital health with federated learning. *npj Digit Med*. 2020 Sep 14;3(1):1–7.
17. The Federated Tumor Segmentation (FeTS) initiative | CBICA | Perelman School of Medicine at the University of Pennsylvania. <https://www.med.upenn.edu/cbica/fets/>.