

Data lineage and observability for ensemble Machine Learning serving in the edge

Dung Nguyen Anh

Advisor: Minh-Tri Nguyen

Data lineage and observability for ensemble Machine Learning serving in the edge

Dung Nguyen Anh

Thesis submitted in partial fulfillment of the requirements for the degree of
Bachelor of Science in Technology.
Otaniemi, 29 Jan 2023

Supervisor: professor Maarit Käpylä
Advisor: Minh-Tri Nguyen

Aalto University
School of Science
Bachelor's Programme in Science and Technology

Contents

Contents	ii
1. Abstract	iii
2. Introduction	iv
3. Background	vi
3.1 Data quality and Data dimensions	vi
3.2 Data lineage	vi
4. Data lineage for ensemble ML serving	viii
5. Experiment and Results	ix
6. Conclusion	x
Bibliography	xi

1. Abstract

2. Introduction

Recently, Machine Learning (ML) has been applied extensively to many problems, including customer journey optimization in marketing [7], particle identification in Physics [2], and predicting mental health problems in healthcare [6]. Then, with the increasing number of ML applications and their proven performance, several organizations are seeking third-party ML service providers to improve their operations, creating a new business model - ML as a Service (MLaaS). In MLaaS, there are two main business engagement models. The simpler version is the two-stakeholders engagement model, where the service provider will use the data from their customer to provide the ML service, including both model creation and running that model in production [8]. This paper focuses on the two-stakeholders model as it helps decrease the complexity of the research.

Like other kinds of service, MLaaS needs a contract, usually called a Service Level Agreement (SLA), that sets the expectations and describes the delivered service to preserve the benefits of both stakeholders. However, the current SLA has yet to fulfill its job because ML-specific attributes, such as data quality, inference accuracy, and explainability, have created new challenges. For example, the popular relationship between data quality and inference accuracy indicates that if the customer submits data unsuitable for the defined purpose, ML models can give false predictions, causing ML providers to be penalized based on the SLA. Moreover, in two-stakeholder MLaaS, the customer does not cooperate with the ML provider in the prediction task, so from their perspective, ML models are considered as a black box, which is not sufficient in human-life-related applications such as e-health or autonomous vehicles. One way to resolve these two challenges is by explicitly explaining the inference result to the customer, which would require tracing down the root cause (data quality) of the decreased performance.

To obtain that information, ML providers first need to manually identify how the data is processed from end to end, and in the current situation, it would demand a lot of work. It is because ML providers do not only use one model

but an ensemble of many models. Although the ensemble model increases the prediction accuracy and the robustness, it increases the pipeline complexity, e.g., the source data is processed differently for each base model and at different microservice. So, this research proposes the use of data lineage in MLaaS as one resolution to this problem. The rest of this paper is organized as follows: section 2 discusses the background of the research, while how data lineage will be implemented is explained in section 3; section 4 describes the experiment and its result; section 5 is the conclusion of the research.

3. Background

3.1 Data quality and Data dimensions

In [4], Liu et al. defined data quality as the suitability of the data for the application and data dimension as metrics that describe the level of data quality. As the definitions are pretty general, data dimensions depend highly on the application, and different organizations value data dimensions differently [1]. For example, although completeness is valuable when working with tabular data, it is not as relevant in computer vision. Instead, attributes like sharpness, noise, or dynamic range, are more significant and can affect the inference accuracy heavily. In edge computing and IoT applications, data characteristics are uncertain, erroneous, and noisy DQ in IoT, so it demands precise monitoring of the whole pipeline for fault tolerance and problem investigation.

3.2 Data lineage

Data lineage or provenance describes the origin of data, how it is derived, and how it changes over time [3]. Based on which questions it can answer, Ikeda and Widom [3] classify data lineage into two types: *where*-lineage and *how*-lineage. Then for each class, there are two *granularities*:

- Schema-level (coarse-grained) answers data lineage questions at the general level. It can be which datasets our current data originate from for where-lineage and which transformations have been used for how-lineage
- Instance-level (fine-grained), on the other hand, clarifies the origin of a specific data point and how different data points are combined to get such results.

The answers that data lineage can provide will shine in edge computing because the complicated data pipeline makes it very time-consuming and costly for ML providers to trace down the root cause of the problems manually. However, with data lineage, they can fastly identify the data input's origin and derivation and then confirm the data quality based on them, helping reduce the required work [5]. Although data lineage is precious in MLaaS, there are some challenges that we must solve before being able to apply it:

- Storing and querying lineage can be expensive when the number of IoT devices and the complexity of the pipeline increase[3]
- Current ML solution adopts many open source application and library that can be considered as a black box, so how can we accurately capture data provenance for black box operations and many others [5]

In progress

4. Data lineage for ensemble ML serving

5. Experiment and Results

6. Conclusion

Bibliography

- [1] Corinna Cichy and Stefan Rass. An Overview of Data Quality Frameworks. *IEEE Access*, 7:24634–24648, 2019.
- [2] Denis Derkach, Mikhail Hushchyn, Tatiana Likhomanenko, Alex Rogozhnikov, Nikita Kazeev, Victoria Chekalina, Radoslav Neychev, Stanislav Kirillov, Fedor Ratnikov, and on behalf of the LHCb collaboration. Machine-learning-based global particle-identification algorithms at the lhcb experiment. *Journal of Physics: Conference Series*, 1085(4):042038, sep 2018.
- [3] Robert Ikeda and Jennifer Widom. Data lineage: A survey. Technical report, Stanford InfoLab, 2009.
- [4] Aimad Karkouch, Hajar Mousannif, Hassan Al Moatassime, and Thomas Noel. Data quality in internet of things: A state-of-the-art survey. *Journal of Network and Computer Applications*, 73:57–81, September 2016.
- [5] Mingjie Tang, Saisai Shao, Weiqing Yang, Yanbo Liang, Yongyang Yu, Bikas Saha, and Dongjoon Hyun. SAC: A System for Big Data Lineage Tracking. In *2019 IEEE 35th International Conference on Data Engineering (ICDE)*, pages 1964–1967, April 2019. ISSN: 2375-026X.
- [6] Ashley E. Tate, Ryan C. McCabe, Henrik Larsson, Sebastian Lundström, Paul Lichtenstein, and Ralf Kuja-Halkola. Predicting mental health problems in adolescence using machine learning techniques. *PLOS ONE*, 15(4):1–13, 04 2020.
- [7] Alessandro Terragni and Marwan Hassani. Optimizing customer journey using process mining and sequence-aware recommendation. In *Proceedings of the 34th ACM / SIGAPP Symposium on Applied Computing, SAC '19*, page 57–65, New York, NY, USA, 2019. Association for Computing Machinery.
- [8] Linh Truong and Tri Nguyen. QoA4ML – A Framework for Supporting Contracts in Machine Learning Services. *2021 IEEE International Conference on Web Services (ICWS)*, page 11, 2021.