

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

Dung Nguyen Anh

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 Korpi-Lagg
Advisor: Minh-Tri Nguyen

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

AuthorDung Nguyen Anh

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

School School of Science

Degree programme Bachelor's Programme in Science and Technology

Major Data Science**Code** SCI3095

Supervisor professor Maarit Korpi-Lagg

Advisor Minh-Tri Nguyen

Level Bachelor's thesis**Date** 27 Nov 2017**Pages** 70**Language** English

Abstract

Abtract

Keywords keyword

urn <https://aaltodoc.aalto.fi>

Contents

Abstract	ii
Contents	iii
1. Introduction	iv
2. Background	vi
2.1 Provenance and Data quality	vi
2.2 Ensemble ML and its explainability challenges	vii
2.3 PROV Family	vii
2.4 Related work	ix
3. Data lineage for ensemble ML serving	x
4. Experiment and Results	xi
5. Conclusion	xii
Bibliography	xiii

1. Introduction

In recent years, Machine Learning (ML) has been applied extensively to many problems, including customer journey optimization in marketing [27], particle identification in Physics [9], and mental health problems prediction in healthcare [26]. With their proven performance and increasing number of ML applications, several organizations are seeking third-party ML service providers to improve their operations, thus leading to the creation of a new business model — ML as a Service (MLaaS). As with other kinds of service, MLaaS requires a contract which is 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 preservation job because ML-specific attributes, such as data quality, inference accuracy, and explainability, have created new challenges. For instance, in Internet of Things (IoT) applications, many issues of the ML models have their root causes in the quality of data, which is impacted by many factors such as erroneous measurement, environmental noise, and discrete observations [19]. Moreover, because of the relationship between data quality and inference accuracy, ML models can produce false predictions when data quality problems exist, causing ML providers to be penalized based on the SLA. Another issue is that in MLaaS, the customers only submit the data and does not cooperate with the ML provider in the prediction task. Thus, from their perspective, ML models are considered to be a black box, which is insufficient in human life-related applications, such as e-health and autonomous vehicles, where uncertain decisions cannot be tolerated.

One approach to resolve these two challenges is increasing model explainability by interpreting the inference result to the customer [13]. However, explainability aspects of ML models have been researched mostly in the training task [19], making it unclear about the appropriate utilization and implementation of ML-specific attributes and constraints. As an approach to support ML-specific service

contracts, Linh et al.[28] proposed the QoA4ML framework which outlines the essential components of such contracts, the definition of ML attributes and constraints, and the guidance on the process to monitor and assess these elements. Although QoA4ML has created a foundation for robust ML monitoring, current implementation only focuses on general ML metrics such as the final confidence of the ensemble model. Consequently, because base models are challenging to monitor and explain with the current state of QoA4ML, they remain like a black box to both stakeholders. Such a level of explainability is not sufficient for some sophisticated applications, such as digital assistants, where a complex chain of models is employed. Additionally, with the increasing amount of research in autoscaler for the inference serving system [20] and automated ensemble [3], it demands a new technique to capture and visualize underlying inference graph to ensure dynamic and robust inference capabilities of ML solutions.

Recognizing above problem can be formalized as capturing data lineage, this thesis performs an analysis of data lineage monitoring and its implementation in edge environment. From that, the research proposes an approach to improve QoA4ML framework with data lineage, which can be combined with other techniques to improve the current MLaaS. This paper provides a prototype with QoA4ML as the foundation so that ML developer and provider can integrate it with their deployed service.

The rest of this paper is organized as follows. Section 2 discusses the background of the research, and section 3 explains data lineage implementation in the QoA4ML framework. Section 4 describes the experiment and its result, while section 5 concludes the research.

2. Background

This section presents the background on provenance for AI applications, ensemble ML, and its explainability challenge. Lastly, Section 2.3 discusses one standard for data model of provenance which is utilized in section 3.

2.1 Provenance and Data quality

The terms lineage, provenance, and traceability can be used interchangeably to refer to the process of constructing a final product, whether it is a digital file or a physical item [12]. Regarding lineage for the digital file, Wang et al [29] was the first paper to discuss those issues, which helped formally define provenance as “data source tagging” and “intermediate source tagging” problems [24, p. 5]. From its original interpretation, provenance has been applied in many domains such as scientific databases [30] [4] [21] [11], data warehouses [6] [7], and recently big data platforms [17] [10], and IoT [23]. In different applications, the types of lineages, their generating techniques, and main issue that they can answer are diverse [12]. For instance, there are four main provenance classes, i.e., provenance meta-data, information systems provenance, workflow provenance, and data provenance which are arranged from the most general to most specific process.

With the development of many AI applications, researchers have tried to incorporate provenance into the AI systems [2] by linking the input and output of the model, which can be a valuable source in interpreting the inference result. However, in this domain, provenance is commonly utilized in a rather general way, where entire algorithms or data transformations are merely represented by semantic relationships [14]. Consequently, while entire pipelines can be documented with provenance, the specific inner workings of individual models remain opaque.

Additionally, provenance is often considered when assessing data quality tasks,

e.g., evaluating integrity, trust, and accuracy. By analyzing provenance, the AI systems can detect errors in data generation and processing, which is valuable for IoT applications where data is uncertain, erroneous, and noisy [16].

2.2 Ensemble ML and its explainability challenges

Ensemble ML is a conventional technique that involves combining multiple models, which can reduce the variance of prediction, to improve the accuracy and robustness of predictions. This is achieved by utilizing multiple base models with different algorithms, hyperparameters, or subsets of features on the same input dataset, and then combining their results by applying various aggregation methods such as averaging, voting, or stacking to produce the final decision.

Such a technique has been shown to be highly effective in many applications including financial forecasting [25], image recognition [5], and natural language processing [8]. Then, as the demand for more precise forecasts and advanced applications, such as AI assistants or autonomous vehicles, increases, it requires a complex dataflow graph, usually represented by a direct acyclic graph (DAG), comprised of multiple base models that are interconnected and orchestrated to handle the input data sequentially or in parallel [20]. For each prediction request, the whole DAGs or only a subset of it can be involved, requiring ML providers to record the data flow of the inference task for further analysis and explanation. Additionally, in resource-constrained environments like the edge, ensemble selection rules can be applied so that only models that provide the best inference performance are utilized [3]. Although this approach can help resolve the challenges of high computation complexity, high resource occupation, and the moderate inference time of ensemble ML, it introduces uncertainty to the inference flow, making the data pipeline less transparent and explainable. Consequently, it is very challenging for current interpreting techniques to produce an explainable and transparent ML solutions without the data lineage of the ensemble model.

2.3 PROV Family

Researchers have proposed various models, languages, and tools to facilitate the documentation of provenance, including those tailored for AI/ML models [1]. However, it is still crucial to have a standard reusable ontology to capture provenance in heterogeneous environments, such as the web, for further development of provenance documentation. As an approach to resolve this problem, the PROV family of

PROV Concepts	Classes	Name
Entity Activity Agent	PROV-DM types	Entity Activity Agent
Generation Usage Communication Derivation Attribution Association Delegation	PROV-DM relations	WasGeneratedBy Used WasInformedBy WasDerivedFrom WasAttributedTo WasAssociatedWith ActedOnBehalfOf

Table 2.1. PROV-DM core concepts and types [18]

documents (PROV-OVERALL) was introduced by The World Wide Web Consortium (W3C) in 2013 [18]. In those documents, the most significant one is the PROV data model (PROV-DM) which provides a universal data model for provenance that can be applied to translate domain-specific and application-specific representations of provenance into a standardized format that can be shared between different systems. PROV-DM has ten core concepts which can be divided into types and relations like in Table 2.1. Then, as seen in Figure 2.1, PROV-DM describes the employment and production of *entities* by *activities*, which may be influenced in several manner by *agent*. In the following paragraphs, the proper definition and application of those data models will be discussed

Firstly, entities are the fundamental building blocks of the PROV-DM data model. They represent physical or digital objects involved in a process, such as a document, dataset, or prediction result. Entities can have fixed attributes that describe their characteristics, such as their size, quality, or creation time. Then, when the value of an attribute is adjusted, a new entity is produced and related to the previous version (“wasDerivedFrom”).

Secondly, activities are actions or processes that transform or manipulate entities. For instance, activities include data processing, predicting, and request serving. Activities can have attributes describing their properties, such as the start and end times of the task or the ML models employed to generate the prediction. Then, this data model can be “used” by entities as their inputs, and it can output a new entity which “wasGeneratedBy” the activities. Moreover, an activity can be informed by other activities with “wasInformedBy” relations.

Lastly, agents are entities that are responsible for carrying out activities (“wasAssociatedWith”), and when one activity is finished, an entity will be produced which

is attributed back to the agent (“wasAttributedTo”). Agents can be people, software tools, ML models, or other entities that can initiate or control activities. Additionally, an agent can “actedOnBehalfOf” other agents, describing the delegation relationship among them.

Section 3 will discuss the implementation of PROV-DM to capture data lineage of ensemble ML models.

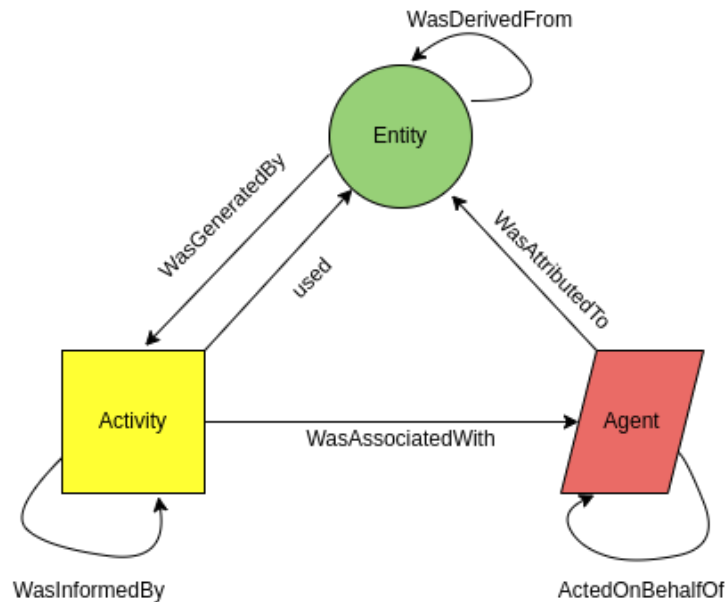


Figure 2.1. W3C PROV-DM structure

2.4 Related work

Data lineage/provenance in AI systems: Data provenance is getting more attention when the demand for Explainable AI (XAI) increases. For instance, some researchers have proposed how provenance should be implemented: [15] introduced ‘Six Ws’ framework for provenance graph-based XAI; [22] combined abstraction and reasoning support offered by models with provenance graph which allows tracing the process of producing the current state.

3. Data lineage for ensemble ML serving

4. Experiment and Results

5. Conclusion

Bibliography

- [1]
- [2] Leonardo Guerreiro Azevedo, Renan Souza, Raphael Melo Thiago, Elton Soares, and Marcio Moreno. Experiencing ProvLake to Manage the Data Lineage of AI Workflows. In *Anais Estendidos do XVI Simpósio Brasileiro de Sistemas de Informação (Anais Estendidos do SBSI 2020)*, pages 206–209, Brasil, November 2020. Sociedade Brasileira de Computação (SBC).
- [3] Yang Bai, Lixing Chen, Mohamed Abdel-Mottaleb, and Jie Xu. Automated Ensemble for Deep Learning Inference on Edge Computing Platforms. *IEEE Internet of Things Journal*, 9(6):4202–4213, March 2022.
- [4] Peter Buneman, Adriane Chapman, and James Cheney. Provenance management in curated databases. In *Proceedings of the 2006 ACM SIGMOD International Conference on Management of Data*, SIGMOD '06, page 539–550, New York, NY, USA, 2006. Association for Computing Machinery.
- [5] Yushi Chen, Ying Wang, Yanfeng Gu, Xin He, Pedram Ghamisi, and Xiuping Jia. Deep learning ensemble for hyperspectral image classification. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 12(6):1882–1897, 2019.
- [6] Yingwei Cui, Jennifer Widom, and Janet L. Wiener. Tracing the lineage of view data in a warehousing environment. *ACM Trans. Database Syst.*, 25(2):179–227, jun 2000.
- [7] YW Cui and J. Widom. Lineage tracing for general data warehouse transformations. *The VLDB Journal — The International Journal on Very Large Data Bases*, 12:41–58, 09 2001.
- [8] Saikat Das, Mohammad Ashrafuzzaman, Frederick T. Sheldon, and Sajjan Shiva. Network intrusion detection using natural language processing and ensemble machine learning. In *2020 IEEE Symposium Series on Computational Intelligence (SSCI)*, pages 829–835, 2020.
- [9] 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.
- [10] Boris Glavic. *Big Data Provenance: Challenges and Implications for Benchmarking*, volume 8163, pages 72–80. 01 2014.

- [11] Thomas Heinis and Gustavo Alonso. Efficient lineage tracking for scientific workflows. In *Proceedings of the 2008 ACM SIGMOD International Conference on Management of Data*, SIGMOD '08, page 1007–1018, New York, NY, USA, 2008. Association for Computing Machinery.
- [12] Melanie Herschel, Ralf Diestelkämper, and Housseem Ben Lahmar. A survey on provenance: What for? What form? What from? *The VLDB Journal*, 26(6):881–906, December 2017.
- [13] Senthil Kumar Jagatheesaperumal, Quoc-Viet Pham, Rukhsana Ruby, Zhaohui Yang, Chunmei Xu, and Zhaoyang Zhang. Explainable AI Over the Internet of Things (IoT): Overview, State-of-the-Art and Future Directions. *IEEE Open Journal of the Communications Society*, 3:2106–2136, 2022. Conference Name: IEEE Open Journal of the Communications Society.
- [14] Fariha Tasmin Jaigirdar, Carsten Rudolph, Gillian Oliver, David Watts, and Chris Bain. What Information is Required for Explainable AI? : A Provenance-based Research Agenda and Future Challenges. In *2020 IEEE 6th International Conference on Collaboration and Internet Computing (CIC)*, pages 177–183, Atlanta, GA, USA, December 2020. IEEE.
- [15] Fariha Tasmin Jaigirdar, Carsten Rudolph, Gillian Oliver, David Watts, and Chris Bain. What Information is Required for Explainable AI? : A Provenance-based Research Agenda and Future Challenges. In *2020 IEEE 6th International Conference on Collaboration and Internet Computing (CIC)*, pages 177–183, Atlanta, GA, USA, December 2020. IEEE.
- [16] 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.
- [17] Dionysios Logothetis, Soumyarupa De, and Kenneth Yocum. Scalable lineage capture for debugging disc analytics. In *Proceedings of the 4th Annual Symposium on Cloud Computing*, SOCC '13, New York, NY, USA, 2013. Association for Computing Machinery.
- [18] Paolo Missier Luc Moreau. Prov-dm: The prov data model.
- [19] My-Linh Nguyen, Thao Phung, Duong-Hai Ly, and Hong-Linh Truong. Holistic Explainability Requirements for End-to-End Machine Learning in IoT Cloud Systems. In *2021 IEEE 29th International Requirements Engineering Conference Workshops (REW)*, pages 188–194, Notre Dame, IN, USA, September 2021. IEEE.
- [20] Kamran Razavi, Manisha Luthra, Boris Koldehofe, Max Mühlhäuser, and Lin Wang. FA2: Fast, Accurate Autoscaling for Serving Deep Learning Inference with SLA Guarantees. In *2022 IEEE 28th Real-Time and Embedded Technology and Applications Symposium (RTAS)*, pages 146–159, May 2022. ISSN: 2642-7346.
- [21] Christopher Ré and Dan Suciu. Approximate lineage for probabilistic databases. *Proc. VLDB Endow.*, 1(1):797–808, aug 2008.
- [22] Owen Reynolds, Antonio García-Domínguez, and Nelly Bencomo. Automated provenance graphs for models@run.time. In *Proceedings of the 23rd ACM/IEEE International Conference on Model Driven Engineering Languages and Systems: Companion Proceedings*, pages 1–10, Virtual Event Canada, October 2020. ACM.

- [23] Xie Rong-na, Li Hui, Shi Guo-zhen, Guo Yun-chuan, Niu Ben, and Su Mang. Provenance-based data flow control mechanism for internet of things. *Transactions on Emerging Telecommunications Technologies*, 32(5):e3934, 2021.
- [24] Martin Schmitz. Survey on Data Quality and Provenance.
- [25] Shaolong Sun, Shouyang Wang, and Yunjie Wei. A new ensemble deep learning approach for exchange rates forecasting and trading. *Advanced Engineering Informatics*, 46:101160, 2020.
- [26] 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.
- [27] 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.
- [28] 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.
- [29] Y Richard Wang and Stuart E Madnick. A Polygen Model for Heterogeneous Database Systems: The Source Tagging Perspective.
- [30] A. Woodruff and M. Stonebraker. Supporting fine-grained data lineage in a database visualization environment. In *Proceedings 13th International Conference on Data Engineering*, pages 91–102, 1997.