

SIMULATING FEDERATED LEARNING FOR CULVERT MANAGEMENT IN UTAH

Pouria Mohammadi¹

¹Dept. of Civil and Environmental Eng., University of Utah, Salt Lake City, UT, USA

ABSTRACT

Transportation agencies increasingly adopt machine learning (ML) to enhance infrastructure management strategies. However, traditional centralized ML approaches face significant challenges due to data scarcity, which limits its effectiveness. To address this issue, we proposed using Federated Learning (FL), a novel approach that enables multiple agencies to collaboratively train ML models on their local datasets without sharing raw data, thus preserving data privacy. In our study, we developed and compared centralized and FL-based Artificial Neural Network (ANN) models using culvert datasets from six states. The FL models, especially when augmented with synthetic data, showed considerable improvements in predictive accuracy (30%), nearly matching the performance of centralized models. Our findings suggest that FL can effectively overcome data limitations, offering transportation agencies a more collaborative and data-driven approach to infrastructure management. This approach particularly benefits agencies with smaller datasets, enhancing their predictive capabilities and maintenance strategies.

1 INTRODUCTION

Machine learning (ML) empowers transportation agencies to make smarter and more effective decisions regarding the maintenance and operation of infrastructure systems (Iyer, 2021). Despite the potential of ML for enhancing transportation infrastructure management, several significant obstacles prevent its widespread implementation. One of the primary challenges that agencies face is the scarcity of comprehensive data inventory. The lack of comprehensive datasets poses several challenges for ML models in infrastructure management, including not learning properly and hindering the use of powerful models (Theodorou et al., 2023). To address these challenges, we explored the application of Federated Learning (FL), a new paradigm of ML, to enhance culvert management in Utah. FL enables multiple entities, such as different transportation agencies, to collaboratively train ML models on their local data without sharing raw data, thereby overcoming data privacy and scarcity issues (Mohammadi et al., 2024). By simulating this approach within the context of culvert management, we aimed to demonstrate how FL can lead to more accurate predictions and proactive maintenance strategies, ultimately improving infrastructure reliability and reducing costs for transportation agencies in Utah.

2 METHODOLOGY

To achieve research goals, we collected data from six state inventories, each contributing a specific number of culvert records. These contributions were 272 records from Utah, 766 from Colorado, 3,884 from Vermont, 1,050 from New York, 417 from Massachusetts, and 1,851 from Ohio. Availability and similarities to Utah's culvert data were the primary criteria for selecting other inventories. Following the data collection phase, we performed a series of preprocessing steps to refine the datasets. After preprocessing the data, we developed and evaluated multiple ML models, comparing both centralized and federated learning approaches to highlight the distinct advantages of FL in the context of culvert management. We selected Artificial Neural Network (ANN) for our core ML algorithm due to its strong predictive capabilities, particularly in handling complex, multiclass classification tasks. Through this

approach, we aimed to simulate how FL, using ANN, can offer significant improvements over traditional centralized models, particularly in terms of data privacy and model accuracy. To achieve this, we first developed three centralized ANN models—Utah-CL, Utah-SMOTE, and ALL-CL—using Utah’s culvert dataset, Utah’s culvert dataset augmented with synthetic data generated through the SMOTE (Synthetic Minority Over-sampling Technique), and a combined dataset that fused culvert data from all six states involved in the study, respectively. Next, we developed five additional centralized models, each based on the culvert dataset from one of the other five states. Finally, we developed two FL models: ALL-FL-SMOTE, which used the augmented datasets from all six states, and ALL-FL, which used the original datasets from all six states without augmentation. This comprehensive approach enabled us to highlight the advantages of FL and assess the individual gains of participation within the FL network.

3 RESULTS

The performance of these models was assessed using various metrics, including accuracy, precision, recall, F1 score, and total loss. As depicted in Figure 1, the FL-based models, ALL-FL-SMOTE and ALL-FL, demonstrated superior performance to the Utah-CL model, exhibiting notable enhancements in precision and accuracy. The results from Utah-CL-SMOTE indicated that while data augmentation alone was insufficient to address the gaps in Utah’s dataset, the application of FL effectively bridged these gaps. In other words, employing the FL by the Utah Department of Transportation can enhance its productivity capacity and avoid culvert failures more quickly. Although the ALL-CL set a high benchmark, the ALL-FL-SMOTE’s accuracy approached closely. The importance of this result lies in the fact that it demonstrates that the FL approach can still produce results that are comparable with those obtained by centralized learning, even without direct access to the data.

A comparison of the centralized models developed on local datasets with the FL models revealed that states like Utah, Colorado, Massachusetts, and New York achieved more accurate predictive outcomes by participating in an FL network. However, the FL model exhibited comparatively modest performance in recall and F1 score metrics, likely due to the class distribution disparities among the local datasets. Even the ALL-CL model, which utilized fused data from all six states, failed to outperform the Vermont-CL and Ohio-CL models. The results indicate that the benefits an agency gains from FL are closely tied to the size of its own dataset, with agencies possessing smaller datasets deriving the most significant advantages. Consequently, agencies with limited data are more inclined to adopt FL but should also be prepared to shoulder a larger share of the associated costs.

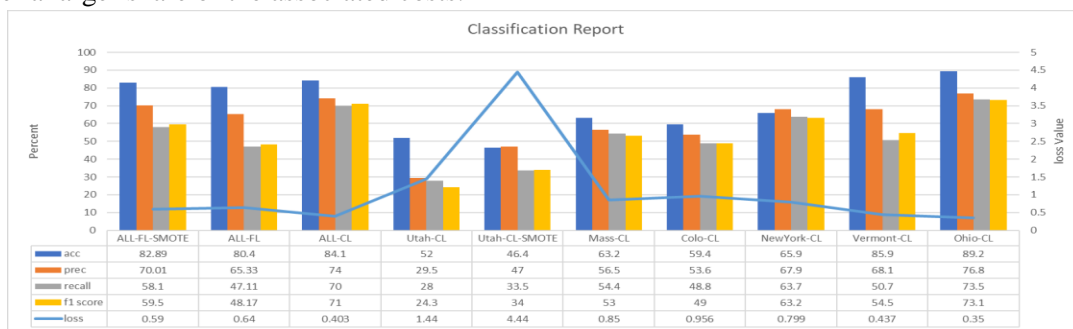


Figure 1: Comparison of classification reports for eight CL models and two FL models.

REFERENCES

Iyer, L. S. 2021. “AI enabled applications towards intelligent transportation”. *Transp. Eng.*, 5: 100083. Elsevier.

Mohammadi, P., A. Rashidi, and S. Asgari. 2024. “Privacy-preserving culvert predictive models: A federated learning approach”. *Advanced Engineering Informatics* 61:102483.

Theodorou, E., E. Spiliotis, and V. Assimakopoulos. 2023. “Optimizing inventory control through a data-driven and model-independent framework”. *EURO J. Transp. Logist.* 12: 100103.