

Report Summary: "Federated GNNs Learning from Dynamic Graphs for Traffic Forecasting"

Introduction:

This project, conducted by Muhammad Usman, focuses on the critical role of traffic forecasting in optimizing transportation systems and urban planning. Traditional forecasting methods often struggle with the complexity and dynamism of transportation networks. The study explores the application of Federated Graph Neural Networks (GNNs), a cutting-edge machine learning technique, to enhance traffic forecasting using dynamic graphs.

Objectives:

The primary objective of this project is to develop an advanced methodology using Federated GNNs for predicting traffic patterns in scenarios with evolving network structures.

Challenges:

Training distributed dynamic GNNs poses several challenges, including massive parameter communication, model convergence issues, and workload imbalance among workers. Addressing these challenges is crucial for achieving efficient and accurate traffic forecasting.

Data Sources:

The project utilizes openly available dynamic graph benchmark datasets, including METR_LA and PEMS_Bay, as primary data sources.

Methodologies:

The project's approach centers on Federated GNNs as the core machine learning framework for traffic forecasting. The methodology comprises data collection from real-world traffic historical speed data, dynamic graph embedding, unsupervised data partitioning, the implementation of a standalone Federated Learning Framework, and model evaluation using metrics like MAE, RMSE, and MAPE.

Timeline:

The project is expected to span a duration of 4-6 months to complete the outlined tasks and achieve the defined objectives.

Project Summarization.

In "Federated GNNs Learning from Dynamic Graphs for Traffic Forecasting" project the focus is on enhancing traffic speed forecasting in the context of complex and ever-changing transportation networks. The project's objectives include developing an advanced methodology that utilizes Federated Graph Neural Networks (GNNs) to predict traffic patterns. Challenges related to distributed GNN training are addressed, and openly available benchmark datasets like METR_LA and PEMS_Bay are used as primary data sources. The proposed methodology encompasses data collection, dynamic graph embedding, unsupervised data partitioning, grouping relevant clients based on their previous performance and its data distribution, the implementation of a standalone Federated Learning algorithm, and model evaluation. This project is anticipated to span 4-6

months, aiming to provide valuable insights into traffic forecasting techniques using state-of-the-art machine learning approaches.

Technical Report:

1. Graph Normalization:

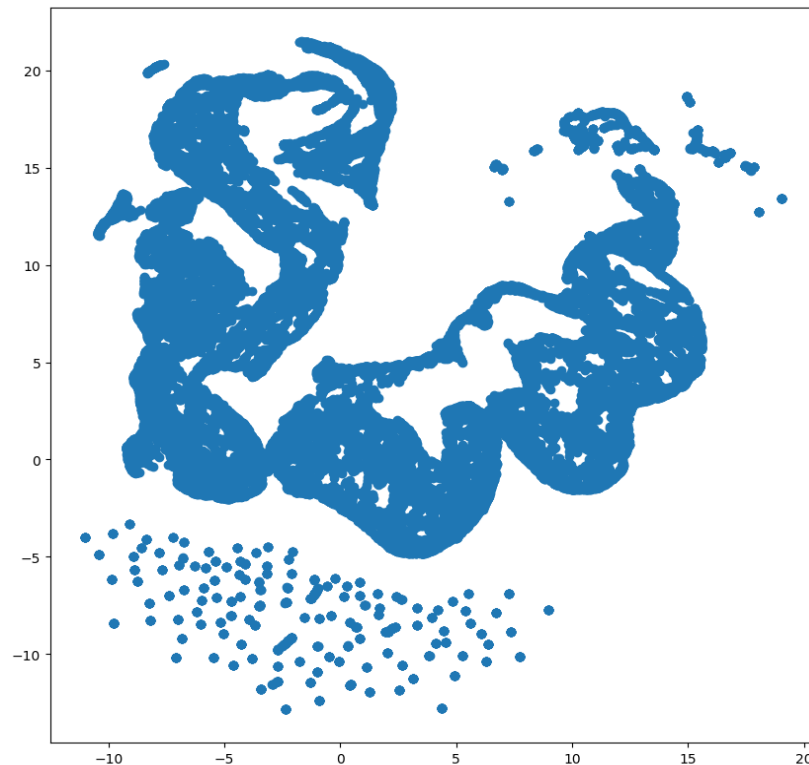
We first normalize all the Graphs features(X) by zero mean and unit standard deviation. After that the overall features are again normalized with the maximum value of the entire dataset.

```
means = np.mean(X, axis=(0, 2))
self.means = means
X = X - means.reshape(1, -1, 1)
stds = np.std(X, axis=(0, 2))
self.stds = stds
X = X / stds.reshape(1, -1, 1)

self.max_value = np.max(X)
X = X / np.max(X)
```

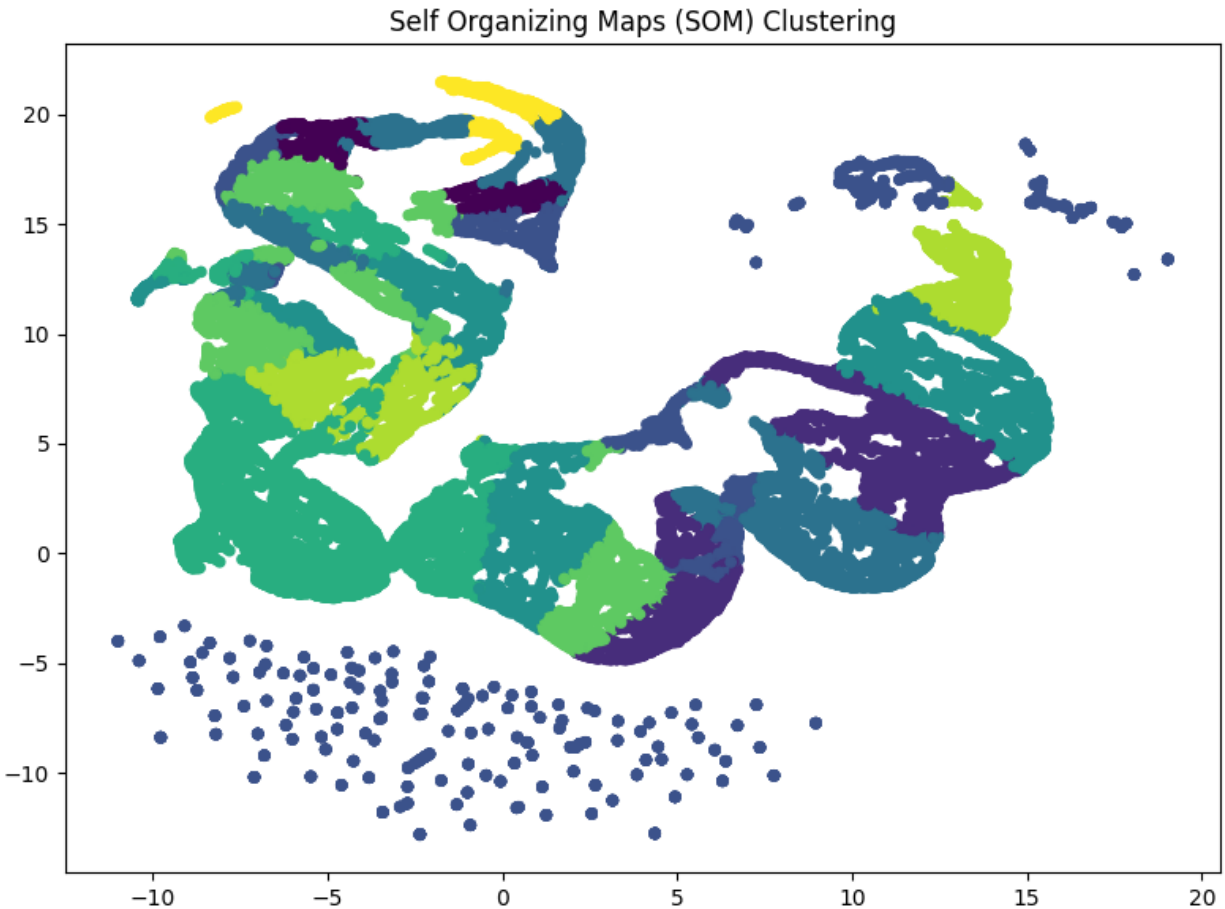
2. Converting all the graphs to its Graph Level Embeddings using A3TGCN:

We then convert all the graphs to Graph Level Embeddings for low level vector representation. All the graphs' embeddings have been plotted and shown in the following image.



3. Unsupervised Client Construction:

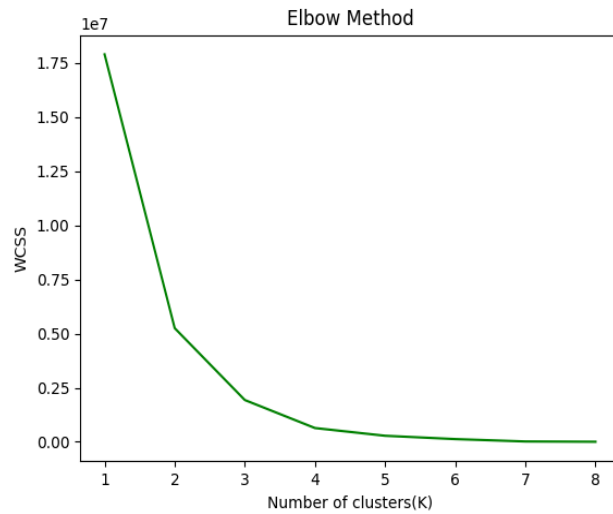
We then create clusters of embeddings using unsupervised learning and assign each cluster to a client. This means that all the graph embeddings within a specific cluster serve as an index to retrieve graphs from the main dataset and assign them to that specific client. Here is a detailed result of clustering the graph-level embeddings using the SOM method, for better understanding.



Each color on the map represents a cluster. The above image has nine colors, so there are nine clusters. Since we are dealing with temporal graphs, similar patterns that occur at a specific time stamp and belong to a cluster have been colored the same. Same colors at different places on the map represent which cluster these dynamic patterns are associated with.

4. Grouping the Clients:

After the initial round of communication between the server and clients, the server proceeds to categorize clients based on their past performances and data distribution. It extracts fourteen statistical features from each client and utilizes a k-means algorithm to assign labels to them, effectively grouping clients with similar labels. To determine the optimal number of clusters, we employ the Elbow method. Here's an example that suggests that having "K=2" groups may yield better results. A snapshot of the output is also provided for better understanding.



```

Training STARTED. Trainer GPU: cuda:0
14:22:55 - root - INFO - Local training with client id list: [1, 0, 5, 2, 8, 4, 7, 6, 3]

Testing--Round:0, MSE:83.264, MAE:3.733, RMSE:8.983, MAPE:14.190, ACC:0.776, R2:0.906, VAR:0.906

Grouping clients, communication round: 1
14:39:49 - root - INFO - Local training with client id list: [5, 6, 2]
14:47:14 - root - INFO - Local training with client id list: [7, 4, 3, 1, 0, 8]

Grouping clients, communication round: 2
14:57:05 - root - INFO - Local training with client id list: [5, 6, 2]
15:04:31 - root - INFO - Local training with client id list: [7, 4, 3, 1, 0, 8]

```

5. Training Federated Global Attention based GNN Model:

In this example, we trained a federated GNN model with 9 clients using the Metr_LA dataset. Initially, the dataset was partitioned among these 9 clients based on their shared characteristics. Each client possessed its own local dynamic graph dataset and employed a GNN-based A3TGCN topology for modeling. Clients trained their models locally and periodically updated the server (global model) with their local parameters. The server, in turn, incorporated these updates from each client. In this project, all clients utilized the A3TGCN machine learning topology, which combines GCN (Graph Convolutional Network) with an attention mechanism to capture spatial information in dynamic graphs and GRU (Graph Recurrent Unit) to capture temporal information.

The server employed the FedAvg mechanism to aggregate parameters received from the client groups. The resulting figure illustrates the algorithm's performance on unseen data, with an MAE score of 3.73 achieved so far.

