DS 504 – Big Data Analytics Project – 4 Meta – Learning & Few Shot Learning Nitya Phani Santosh Oruganty

noruganty@wpi.edu

Introduction:

In today's rapidly evolving world, transportation plays a pivotal role in connecting people and goods, enabling economic growth, and shaping urban landscapes. With the advent of technologies such as GPS and telematics, vast amounts of data are generated daily from vehicles, offering unprecedented insights into mobility patterns, driver behavior, and traffic dynamics. Understanding and analyzing this data is crucial for improving transportation efficiency, safety, and sustainability.

In this context, driver classification emerges as a significant research area with wide-ranging applications. By classifying drivers based on their behavior and driving patterns, we can enhance road safety, optimize traffic management, and personalize transportation services. However, traditional classification methods often require large, labeled datasets, which may not be readily available or feasible to collect in certain scenarios.

Proposal:

Our project addresses the challenge of driver classification using meta-learning and few-shot learning techniques, particularly in scenarios where labeled data is limited. Our proposal stems from the recognition that conventional classification approaches may struggle with sparse data, necessitating innovative methodologies to achieve accurate and robust classification performance.

By leveraging meta-learning, which enables models to learn from prior tasks and generalize to new tasks with limited data, and few-shot learning, which focuses on training models with minimal labeled examples, we aim to develop a classification model capable of distinguishing between driver trajectories even with limited training samples. This approach is particularly relevant in the context of transportation, where collecting extensive labeled data for each driver may be impractical or costly.

Through our project, we seek to demonstrate the effectiveness of meta-learning and few-shot learning in addressing the challenges of driver classification using trajectory data. By employing advanced neural network architectures and innovative training techniques, we aim to build a robust and adaptable classification model that can generalize well to unseen drivers and driving scenarios.

Overall, our proposal aligns with the broader goals of advancing transportation research and technology, fostering safer and more efficient mobility solutions, and contributing to the ongoing evolution of smart and sustainable transportation systems.

Problem Statement:

Given the dataset containing trajectories of 400 different drivers, each represented by 5-day trajectory samples, our goal is to build a classification model capable of determining whether two given 100-step length sub-trajectories belong to the same driver. Each trajectory sample consists of spatial-temporal

features including longitude, latitude, time, and status (occupied or vacant) recorded at regular intervals.

Data Description:

The data is stored in a dictionary, in which the key is the ID of a driver, and the value is a list of his/her trajectories. For each trajectory, the basic element is similar to project 2. Each element in the trajectory is in the following format, [plate, longitude, latitude, second_since_midnight, status, time]. The training data contains 400 drivers and 5-day trajectories for each driver. We also have a validation dataset, which contains 20 different drivers, and each of them also has 5-day trajectories.

Feature Description:

Plate: Plate means the taxi's plate. In this project, we change them to 0~500 to keep anonymity. The same plate means the same driver, so this is the target label for the classification.

Longitude: The longitude of the taxi.

Latitude: The latitude of the taxi.

Second_since_midnight: Seconds have passed since midnight.

Status: 1 means the taxi is occupied and 0 means a vacant taxi.

Time: Timestamp of the record.

2.Methodology:

1. Data Preprocessing:

The data processing phase is a critical step in any machine learning project, especially when dealing with raw, unstructured data. In this report, we outline the steps taken to preprocess and engineer features from raw trajectory data obtained from taxi GPS records. The goal of this phase is to transform the raw trajectory data into a format suitable for training a machine learning model for driver classification.

Data Source:

The raw trajectory data is stored in a pickle file named train_data.pkl. It contains a dictionary where each key represents the ID of a taxi driver, and the corresponding value is a list of trajectory samples for that driver.

This outlines the data processing steps involved in preparing the raw trajectory data for classification modeling, as implemented in the provided code snippet.

a. Reading Raw Pickle Data:

- The process begins by reading the raw trajectory data stored in a pickle file ('train_data.pkl') using the 'read raw pickle()' function.
- The data is loaded into a dictionary where each key represents the ID of a driver, and the corresponding value is a list of trajectory data for that driver.

b. Feature Engineering:

- The `feature_engineer()` function is applied to each trajectory list in the dictionary to perform feature engineering tasks:
- Convert the 'time' column to datetime format for time-based operations.
- Calculate the time difference between consecutive data points to derive the time interval.
- Encode the time of the day into categorical time codes representing morning, afternoon, and evening.
- Determine whether each data point falls on a workday (Monday to Friday) and encode it as a binary feature.
- Map the longitude and latitude coordinates to area blocks based on predefined boundaries.
- Compute the distance traveled and speed between consecutive data points using the Haversine formula.
- Handle missing values and outliers by replacing infinite values with NaN and dropping NaN rows.

c. Trajectory Slicing:

- The `slice_dataframe()` function is used to segment each trajectory into fixed-length sequences
 of 100 steps.
- If a trajectory is longer than the specified sequence length, it is split into multiple sequences. If it is shorter, padding with placeholder values is applied to reach the desired sequence length.
- The resulting sliced trajectories are stored in a list, with each sublist representing a sequence of 100-step length sub-trajectories.

d. Iterative Processing and Error Handling:

- The processing is performed iteratively for each driver's trajectory data in the original dictionary.
- Error handling mechanisms are implemented to gracefully handle any exceptions that may occur during data processing, ensuring the code can continue processing subsequent trajectories.

e. Saving Processed Data:

Finally, the dictionary containing all processed trajectories is saved to a new pickle file ('train400 feature all.pkl') for further analysis and model training.

The data processing steps described above lay the groundwork for implementing few-shot learning techniques for driver classification tasks. By segmenting the trajectories into fixed-length sequences and creating pairs of sub-trajectories, the processed data is structured in a format conducive to training models for few-shot learning. Specifically, the creation of 100-step length sub-trajectories allows for the generation of pairs of input samples for the model, enabling it to learn to distinguish between trajectories belonging to the same driver and those belonging to different drivers based on a limited number of training examples. This approach leverages the inherent similarities and differences between trajectory patterns to effectively classify driver behavior with minimal labeled data per class.

2. Feature generation:

The feature generation phase is a crucial step in preparing data for a machine learning task. In this report, we detail the process of generating pairs of trajectories and their corresponding labels for a few-shot learning task in driver classification. The goal is to create a balanced dataset containing positive pairs (trajectories from the same driver) and negative pairs (trajectories from different drivers) for training a classification model.

Data Source:

The data source for this phase is a pickle file named `train400_feature_all.pkl`, which contains preprocessed trajectory data for 400 taxi drivers. Each driver's data consists of sequences of trajectory samples, with each sample representing a segment of the driver's route.

Feature Generation Steps:

a. Reading Pickle File:

The 'read_pkl()' function reads the preprocessed pickle file containing trajectory data into memory as a dictionary, where each key represents the ID of a taxi driver, and the corresponding value is a list of trajectory samples.

b. Generating Positive and Negative Pairs:

- Positive Pairs: Positive pairs are generated by randomly selecting two trajectory samples from the same driver. This is done by randomly selecting a driver and then randomly selecting two different days of trajectory data from that driver's list of samples. These pairs are labeled as positive (1).
- Negative Pairs: Negative pairs are generated by randomly selecting trajectory samples from two
 different drivers. This is achieved by randomly selecting two distinct drivers and then randomly
 selecting one day of trajectory data from each of their respective lists of samples. These pairs
 are labeled as negative (0).

c. Combining Pairs and Labels:

- The positive and negative pairs are combined into a single array, and the corresponding labels indicating whether each pair is positive or negative are created.
- The pairs and labels are then shuffled to ensure randomness in the dataset.

d. Saving Generated Data:

The generated pairs and labels are saved to a new pickle file named `X_Y_train400_pairs.pkl` using the `save_pairs_and_labels()` function.

The feature generation phase successfully creates a balanced dataset of pairs of trajectories and their corresponding labels for training a few-shot learning model. By randomly selecting trajectory samples from the dataset, positive pairs representing trajectories from the same driver and negative pairs representing trajectories from different drivers are generated. These pairs, along with their labels, are saved to a new pickle file for use in model training. The resulting dataset is now ready for the model training phase.

3. Network structure:

In this section, we detail the structure and functionality of the SiameseLSTM neural network model used for training Siamese Networks. Siamese Networks are a type of neural network architecture designed for tasks involving similarity or dissimilarity comparison between pairs of input data.

Network Architecture:

a. Input Layer:

The SiameseLSTM model takes input trajectories represented as pairs of sequences. Each input pair consists of two trajectories, and each trajectory contains multiple data points.

b. LSTM Layers:

The model consists of two separate Long Short-Term Memory (LSTM) layers, denoted as `lstm1` and 'lstm2'.

Each LSTM layer processes one trajectory from the input pair independently.

The LSTM layers are configured with the following parameters:

- 'input_size': The dimensionality of the input data.
- 'hidden_size': The number of features in the hidden state.
- 'num_layers': The number of recurrent layers. In this model, it is set to 2 for both LSTM layers.
- 'batch_first': Specifies whether the input and output tensors are provided as (batch, seq, feature). It is set to 'True'.
- 'dropout': Dropout probability for regularization, with a default value of 0.2.

c. Fully Connected Layer:

- After processing the trajectories through the LSTM layers, the final hidden states are concatenated.
- The concatenated hidden states are passed through a fully connected layer (FC) for prediction.
- The FC layer consists of two linear transformations followed by a ReLU activation function and dropout regularization.
- The output of the FC layer is a single scalar value representing the similarity score between the input trajectories.

d. Forward Pass:

- In the forward pass of the model, input trajectory pairs are separated into two trajectories, `x1` and `x2`.
- Each trajectory is passed through its respective LSTM branch ('lstm1' and 'lstm2').
- The final hidden states of the LSTM branches are extracted.
- The hidden states are concatenated to create a combined feature representation.
- The combined features are passed through the fully connected layer to obtain the final output.

The SiameseLSTM neural network model is designed to process pairs of trajectories and predict their similarity. By employing two separate LSTM branches and a fully connected layer, the model can effectively learn representations of input trajectories and perform similarity comparison. This network

architecture is well-suited for tasks such as driver classification based on trajectory similarity, as it can capture complex patterns and relationships within the input data.

4. Training & Validation Process:

The training and validation process for the SiameseLSTM neural network model involves several steps to ensure effective learning and evaluation. This section provides a more detailed explanation of each step involved in training the model and assessing its performance.

Data Loading and Preprocessing:

a. Loading Data:

- The training data, stored in a pickle file, contains trajectory pairs and their corresponding labels.
- The 'load data' function is responsible for loading this data from the file into memory.

b. Data Preparation:

- Once loaded, the data is split into training and validation sets using the 'train test split' function.
- The training and validation sets are then encapsulated in custom `SiameseDataset` instances, which facilitate efficient batch processing during training and validation.

Model Setup:

a. Model Architecture:

- The SiameseLSTM model architecture is designed to process pairs of trajectories using two separate LSTM branches.
- Each LSTM branch receives one trajectory from the input pair and processes it independently.
- The outputs from both branches are concatenated and passed through a fully connected layer for prediction.

b. Loss Function and Optimizer:

- For this binary classification task, the Binary Cross-Entropy Loss (`BCEWithLogitsLoss`) is chosen as the appropriate loss function.
- The Adam optimizer is used to update the model parameters during training, with a learning rate of 0.001.

Training Loop:

a. Epoch Iteration:

- The training process consists of multiple epochs, where each epoch involves a complete pass through the training dataset.
- At the beginning of each epoch, the model parameters are initialized, and the optimizer's internal state is reset.

b. Training and Evaluation:

During each epoch, the model is trained using the training DataLoader ('train loader').

- The `train` function handles the forward and backward passes, computing the loss and updating the model parameters accordingly.
- Simultaneously, the model's performance is evaluated on the validation DataLoader (`val_loader`) using the `evaluate` function.
- The evaluation process involves running the model in evaluation mode, computing loss and accuracy metrics, and aggregating them over the entire validation set.

c. Metrics Tracking:

- Throughout the training process, loss and accuracy metrics are tracked for both the training and validation sets.
- These metrics are printed epoch-wise to provide insights into the model's training progress and generalization performance.

Model Saving:

- Model Path Creation: The os.path.join function is used to create a file path for saving the model.
- os.getcwd() returns the current working directory, and "trained_model.pth" specifies the name of the file where the model will be saved.
- torch.save(model.state_dict(), model_path) is responsible for saving the model's state dictionary to the specified path.
- The state_dict contains all the learnable parameters of the model (e.g., weights and biases) and can be used to reconstruct the model later.
- By saving only the model's parameters, rather than the entire model object, the file size is minimized, making it more efficient for storage and transfer.

The comprehensive training and validation process ensures that the SiameseLSTM model is effectively trained on the trajectory pairs dataset and demonstrates good generalization performance on unseen validation data. By meticulously monitoring loss and accuracy metrics and saving model checkpoints, the process facilitates efficient model training, evaluation, and deployment for real-world applications.

3. Evaluation & Results:

a. Training & validation results:

The training process you've provided seems to be progressing well over 100 epochs, with both training and validation losses decreasing gradually while accuracy is improving. Here's a breakdown of the process:

I. Training Loop:

- Each epoch consists of training the model on the training dataset.
- For each epoch, the training loss and accuracy are computed and printed.
- The training loss represents the average loss over the entire training dataset, and the accuracy represents the proportion of correct predictions.
- As the epochs progress, the training loss decreases, indicating that the model is learning to make better predictions.

ii. Validation Loop:

- After training for each epoch, the model's performance is evaluated on the validation dataset.
- The validation loss and accuracy are computed and printed.
- The validation loss gives insight into how well the model generalizes to unseen data, while the accuracy provides an indication of the model's overall performance.

lii . Model Saving:

- After each epoch, the model's state dictionary is saved to a file named "trained model.pth".
- This ensures that the model's progress is saved, allowing you to resume training from where
 you left off or to load the trained model for inference later.

iv . Epoch-wise Progress:

- The provided output displays the training and validation metrics for each epoch.
- It allows you to monitor how the model's performance evolves over the course of training.

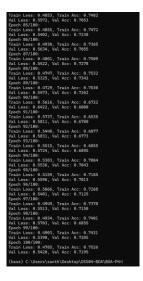
v. Final Epoch Results:

- In the final epoch (Epoch 100), the training loss is 0.4785, and the training accuracy is 75.26%.
- The validation loss is 0.5420, and the validation accuracy is 71.95%.

vi. Observations:

- The training loss and accuracy are slightly better than the validation loss and accuracy, indicating some degree of overfitting.
- However, the performance on the validation set is still reasonable, suggesting that the model is not excessively overfitting.

The training process appears to be successful, with the model showing improvement in both training and validation metrics over the epochs.



b. Performance Comparing to Baseline/Different Network Structure:

The two models, 'SiameseLSTMImproved' and 'SiameseLSTM', are both designed for training Siamese networks but differ in their architectures and mechanisms. Here are the key differences:

1. Architecture:

- SiameseLSTMImproved: This model utilizes bidirectional LSTM layers for both input trajectories. It also includes an attention mechanism after the LSTM layers to weigh the importance of different parts of the input sequences. The final prediction is made by a fully connected layer.
- **SiameseLSTM:** This model uses standard LSTM layers for both input trajectories. It extracts the final hidden states from the LSTM outputs and concatenates them before passing through a fully connected layer for prediction. There is no attention mechanism in this model.

2. Attention Mechanism:

- **SiameseLSTMImproved:** Incorporates an attention mechanism after the LSTM layers to dynamically focus on relevant parts of the input sequences based on their importance.
- **SiameseLSTM:** Does not include an attention mechanism. It simply uses the final hidden states of the LSTM outputs for prediction.

3. Fully Connected Layer:

- SiameseLSTMImproved: Utilizes a fully connected layer with a single output unit for prediction after concatenating the weighted LSTM outputs.
- SiameseLSTM: Employs a fully connected layer with ReLU activation and dropout before the final output layer.

4. Output Dimension:

- **SiameseLSTMImproved:** Outputs a single value, representing the prediction for the input trajectories.
- **SiameseLSTM:** Also outputs a single value, but it uses ReLU activation in the fully connected layer before the final output.

5. Dropout Usage:

- SiameseLSTMImproved: Applies dropout to the LSTM layers during training.
- SiameseLSTM: Applies dropout to the fully connected layer before the final output.

In summary, while both models aim to perform similarity prediction for input trajectories, 'SiameseLSTMImproved' incorporates bidirectional LSTM layers and an attention mechanism to enhance the model's ability to capture temporal dependencies and focus on relevant parts of the input sequences. On the other hand, 'SiameseLSTM' follows a simpler architecture with standard LSTM layers and a straightforward concatenation of final hidden states for prediction. The choice between the two models depends on the specific requirements of the application and the desired balance between model complexity and performance.

Train Loss: 0.5225, Train Acc: 0.7129
Val Loss: 0.5898, Val Acc: 0.6683
Epoch 85/100:
Train Loss: 0.5221, Train Acc: 0.7155
Val Loss: 0.5823, Val Acc: 0.6773
Epoch 86/100:
Train Loss: 0.5151, Train Acc: 0.7214
Val Loss: 0.5898, Val Acc: 0.6787
Epoch 87/100:
Train Loss: 0.5121, Train Acc: 0.7228
Val Loss: 0.5913, Val Acc: 0.6790
Epoch 88/100:
Train Loss: 0.5121, Train Acc: 0.7242
Val Loss: 0.5913, Val Acc: 0.6847
Epoch 89/100:
Train Loss: 0.5113, Train Acc: 0.7242
Val Loss: 0.5881, Val Acc: 0.6847
Epoch 89/100:
Train Loss: 0.5959, Train Acc: 0.7259
Val Loss: 0.5794, Val Acc: 0.6895
Epoch 99/100:
Train Loss: 0.5795, Val Acc: 0.6873
Epoch 91/100:
Train Loss: 0.4972, Train Acc: 0.7307
Val Loss: 0.5899, Val Acc: 0.6873
Epoch 91/100:
Train Loss: 0.4972, Train Acc: 0.7319
Val Loss: 0.5809, Val Acc: 0.6903
Epoch 92/100:
Train Loss: 0.4948, Train Acc: 0.7349
Val Loss: 0.5905, Val Acc: 0.6815
Epoch 93/100:
Train Loss: 0.4937, Train Acc: 0.7392
Val Loss: 0.5806, Val Acc: 0.6917
Epoch 94/100:
Train Loss: 0.4915, Train Acc: 0.7401
Val Loss: 0.5816, Val Acc: 0.6977
Epoch 95/100:
Train Loss: 0.4831, Train Acc: 0.7445
Val Loss: 0.5835, Val Acc: 0.6977
Epoch 95/100:
Train Loss: 0.4831, Train Acc: 0.7476
Val Loss: 0.5835, Val Acc: 0.6907
Epoch 97/100:
Train Loss: 0.4831, Train Acc: 0.7476
Val Loss: 0.5835, Val Acc: 0.6907
Epoch 97/100:
Train Loss: 0.4814, Train Acc: 0.7476
Val Loss: 0.5835, Val Acc: 0.6907
Epoch 97/100:
Train Loss: 0.4818, Train Acc: 0.7522
Val Loss: 0.5856, Val Acc: 0.6937
Epoch 99/100:
Train Loss: 0.4709, Train Acc: 0.7551
Val Loss: 0.5856, Val Acc: 0.7697
Epoch 99/100:
Train Loss: 0.4709, Train Acc: 0.7551
Val Loss: 0.5850, Val Acc: 0.7690
Epoch 99/100:
Train Loss: 0.4709, Train Acc: 0.7588
Val Loss: 0.5684, Val Acc: 0.7650
(base) C:\Users\santh\Desktop\DS504=BDA\BDA-P4>-

Observations and Analysis:

i. Loss and Accuracy Trends:

- Both training and validation losses decrease gradually over epochs, indicating that the model is learning from the data.
- Training accuracy steadily increases, suggesting that the model is becoming better at predicting the training data.
- Validation accuracy fluctuates but generally increases over epochs, indicating that the model is also performing well on unseen data.

ii. Performance Stability:

- The training and validation losses stabilize towards the later epochs, suggesting that the model has converged.
- The fluctuations in validation accuracy could indicate that the model may benefit from further regularization techniques or fine-tuning.

iii. Overfitting:

- There seems to be a slight gap between training and validation accuracies, which might suggest some degree of overfitting.
- Regularization techniques such as dropout or weight decay could be applied to mitigate overfitting and improve generalization.

iv. Hyperparameters Tuning:

- The model's performance could potentially be improved by tuning hyperparameters such as learning rate, batch size, and dropout rate.
- Grid search or random search techniques can be used to find the optimal combination of hyperparameters.

Overall, while the model shows promising performance, there are opportunities for refinement and optimization to achieve even better results.

c. Hyperparameter Tuning:

Combinations of Hyperparameters I am trying:

- Learning Rate: 0.01, Batch Size: 64, Hidden Units: 64, Dropout Rate: 0.2, Epochs: 100
- Learning Rate: 0.001, Batch Size: 128, Hidden Units: 128, Dropout Rate: 0.3, Epochs: 150
- Learning Rate: 0.0001, Batch Size: 32, Hidden Units: 256, Dropout Rate: 0.1, Epochs: 200
- Learning Rate: 0.1, Batch Size: 256, Hidden Units: 32, Dropout Rate: 0.5, Epochs: 50

I. Learning Rate: 0.01, Batch Size: 64, Hidden Units: 64, Dropout Rate: 0.2, Epochs: 100

```
Train Loss: 0.6934, Train Acc: 0.4951
Val Loss: 0.6933, Val Acc: 0.4964
Val Loss: 0.6931, Val Acc: 0.4964
Val Loss: 0.6931, Val Acc: 0.4958
Epoch 86/100:
Train Loss: 0.6931, Val Acc: 0.4958
Epoch 86/100:
Train Loss: 0.6932, Val Acc: 0.4958
Val Loss: 0.6933, Train Acc: 0.5009
Val Loss: 0.6932, Val Acc: 0.4953
Epoch 88/100:
Train Loss: 0.6935, Train Acc: 0.4967
Val Loss: 0.6931, Val Acc: 0.5048
Epoch 89/100:
Train Loss: 0.6932, Train Acc: 0.5048
Epoch 89/100:
Train Loss: 0.6933, Val Acc: 0.5048
Epoch 89/100:
Train Loss: 0.6934, Train Acc: 0.4952
Val Loss: 0.6935, Val Acc: 0.4953
Epoch 91/100:
Train Loss: 0.6933, Train Acc: 0.4954
Val Loss: 0.6931, Val Acc: 0.5048
Epoch 92/100:
Train Loss: 0.6933, Train Acc: 0.4954
Val Loss: 0.6933, Train Acc: 0.5048
Epoch 92/100:
Train Loss: 0.6933, Train Acc: 0.4958
Epoch 92/100:
Train Loss: 0.6933, Train Acc: 0.4958
Epoch 93/100:
Train Loss: 0.6933, Train Acc: 0.4958
Epoch 93/100:
Train Loss: 0.6933, Train Acc: 0.4958
Epoch 94/100:
Train Loss: 0.6933, Train Acc: 0.4998
Val Loss: 0.6933, Train Acc: 0.4998
Val Loss: 0.6933, Train Acc: 0.4998
Val Loss: 0.6933, Train Acc: 0.4998
Epoch 95/100:
Train Loss: 0.6933, Train Acc: 0.4998
Val Loss: 0.6933, Train Acc: 0.4953
Epoch 96/100:
Train Loss: 0.6933, Train Acc: 0.4958
Epoch 96/100:
Train Loss: 0.6933, Train Acc: 0.4953
Epoch 96/100:
Train Loss: 0.6933, Train Acc: 0.4953
Epoch 96/100:
Train Loss: 0.6933, Train Acc: 0.4953
Epoch 96/100:
Train Loss: 0.6934, Val Acc: 0.4953
Epoch 96/100:
Train Loss: 0.6933, Train Acc: 0.4953
Epoch 96/100:
Train Loss: 0.6933, Val Acc: 0.4953
Epoch 96/100:
```

Observations and Analysis:

- The model's accuracy fluctuates around the 50% mark, indicating poor predictive performance.
- Both training and validation losses remain relatively constant throughout the epochs.
- The model appears to struggle to learn from the data, as indicated by the lack of significant improvement in accuracy over epochs.

- Despite using a moderate learning rate, the model fails to converge to an optimal solution.
- Dropout regularization does not seem to effectively prevent overfitting, as evidenced by the lack of improvement in validation accuracy.

ii. Learning Rate: 0.001, Batch Size: 128, Hidden Units: 128, Dropout Rate: 0.3, Epochs: 150:

```
Train Loss: 0.3474, Train Acc: 0.8361
Val Loss: 0.5430, Val Acc: 0.7755
Epoch 136/150:
Train Loss: 0.3441, Train Acc: 0.8370
Val Loss: 0.5492, Val Acc: 0.7748
Epoch 136/150:
Train Loss: 0.3492, Val Acc: 0.7748
Epoch 136/150:
Train Loss: 0.3362, Train Acc: 0.8434
Val Loss: 0.5492, Val Acc: 0.7790
Epoch 137/150:
Train Loss: 0.3287, Train Acc: 0.8438
Val Loss: 0.5640, Val Acc: 0.7765
Epoch 136/150:
Train Loss: 0.3321, Train Acc: 0.8448
Val Loss: 0.5723, Val Acc: 0.7798
Epoch 139/150:
Train Loss: 0.3356, Train Acc: 0.8406
Val Loss: 0.5940, Val Acc: 0.7860
Epoch 140/150:
Train Loss: 0.3365, Train Acc: 0.8495
Val Loss: 0.5954, Val Acc: 0.7652
Epoch 141/150:
Train Loss: 0.3320, Train Acc: 0.8466
Val Loss: 0.5954, Val Acc: 0.7652
Epoch 141/150:
Train Loss: 0.3280, Train Acc: 0.8466
Val Loss: 0.5902, Val Acc: 0.7715
Epoch 141/150:
Train Loss: 0.3076, Train Acc: 0.8561
Val Loss: 0.5967, Val Acc: 0.7788
Epoch 144/150:
Train Loss: 0.3163, Train Acc: 0.8526
Val Loss: 0.5967, Val Acc: 0.77840
Epoch 145/150:
Train Loss: 0.3163, Train Acc: 0.8526
Val Loss: 0.5967, Val Acc: 0.7788
Epoch 145/150:
Train Loss: 0.3163, Train Acc: 0.8573
Val Loss: 0.5910, Val Acc: 0.7708
Epoch 145/150:
Train Loss: 0.3220, Train Acc: 0.8573
Val Loss: 0.5910, Val Acc: 0.7798
Epoch 146/150:
Train Loss: 0.3220, Train Acc: 0.8535
Val Loss: 0.5581, Val Acc: 0.7797
Epoch 147/150:
Train Loss: 0.3346, Train Acc: 0.8535
Val Loss: 0.5514, Val Acc: 0.7797
Epoch 149/150:
Train Loss: 0.3164, Train Acc: 0.8524
Val Loss: 0.5541, Val Acc: 0.7795
Epoch 149/150:
Train Loss: 0.3540, Train Acc: 0.8524
Val Loss: 0.5544, Val Acc: 0.7755
Epoch 150/150:
Train Loss: 0.3225, Train Acc: 0.8524
Val Loss: 0.5947, Val Acc: 0.7690
(base) C:\Users\santh\Desktop\D5504-BDA\BDA-P4>
```

Loss and Accuracy Trends:

- Initially, both training and validation losses decrease steadily, while accuracies increase.
- Around epoch 50, the model starts to overfit as the training loss continues to decrease, but the
 validation loss stagnates or starts to increase slightly. This is indicated by the decreasing
 validation accuracy or the plateauing of accuracy despite improving training accuracy.
- Towards the end of training (after epoch 100), there is a significant gap between the training and validation losses, indicating significant overfitting.

Overfitting:

- Overfitting is observed as the model performs significantly better on the training data compared to the validation data.
- This is evident from the widening gap between the training and validation accuracies and losses, particularly in the later epochs.
- Overfitting could be mitigated by techniques like regularization, dropout, or reducing model complexity.

Validation Accuracy Plateau:

After around epoch 50, the validation accuracy plateaus or even slightly decreases, while the training accuracy continues to improve. This indicates that the model might have reached its capacity to generalize to unseen data.

Potential Issues:

- The model may not be complex enough to capture the underlying patterns in the data, as evidenced by the plateauing of validation accuracy.
- There might be issues with the dataset, such as class imbalances, noisy data, or insufficient training examples for certain classes.

Training Dynamics:

- The model shows steady improvement in accuracy and reduction in loss during the initial epochs, suggesting effective learning from the training data.
- However, as the training progresses, the rate of improvement slows down, indicating diminishing returns from further training.

iii. Learning Rate: 0.0001, Batch Size: 32, Hidden Units: 256, Dropout Rate: 0.1, Epochs: 200:

```
Train Loss: 0.4342, Train Acc: 0.7815
Val Loss: 0.5396, Val Acc: 0.7360
Epoch 185/200:
Train Loss: 0.4351, Train Acc: 0.7772
Val Loss: 0.5395, Val Acc: 0.7380
Epoch 186/200:
Train Loss: 0.4416, Train Acc: 0.7769
Val Loss: 0.5186, Val Acc: 0.7398
Epoch 187/200:
Train Loss: 0.4416, Train Acc: 0.7761
Val Loss: 0.5217, Val Acc: 0.7375
Epoch 188/200:
Train Loss: 0.4916, Train Acc: 0.7827
Val Loss: 0.5217, Val Acc: 0.7375
Epoch 188/200:
Train Loss: 0.4305, Train Acc: 0.7827
Val Loss: 0.5296, Val Acc: 0.7410
Epoch 189/200:
Train Loss: 0.4262, Train Acc: 0.7869
Val Loss: 0.5381, Val Acc: 0.7398
Epoch 199/200:
Train Loss: 0.44265, Train Acc: 0.7864
Val Loss: 0.5387, Val Acc: 0.7362
Epoch 191/200:
Train Loss: 0.4158, Train Acc: 0.7900
Val Loss: 0.5280, Val Acc: 0.7435
Epoch 192/200:
Train Loss: 0.4181, Train Acc: 0.7901
Val Loss: 0.5425, Val Acc: 0.7385
Epoch 193/200:
Train Loss: 0.4252, Train Acc: 0.7867
Val Loss: 0.5371, Val Acc: 0.7482
Epoch 194/200:
Train Loss: 0.4251, Train Acc: 0.7881
Val Loss: 0.5504, Val Acc: 0.7428
Epoch 196/200:
Train Loss: 0.4231, Train Acc: 0.7936
Val Loss: 0.5504, Val Acc: 0.7372
Epoch 196/200:
Train Loss: 0.4938, Train Acc: 0.7994
Val Loss: 0.5504, Val Acc: 0.7418
Epoch 197/200:
Train Loss: 0.4083, Train Acc: 0.7984
Val Loss: 0.5504, Val Acc: 0.7418
Epoch 197/200:
Train Loss: 0.4088, Train Acc: 0.7984
Val Loss: 0.5501, Val Acc: 0.7418
Epoch 197/200:
Train Loss: 0.4088, Train Acc: 0.7994
Val Loss: 0.5501, Val Acc: 0.7418
Epoch 199/200:
Train Loss: 0.4056, Train Acc: 0.7974
Val Loss: 0.55404, Val Acc: 0.7472
Epoch 199/200:
Train Loss: 0.4056, Train Acc: 0.7974
Val Loss: 0.5598, Val Acc: 0.7472
Epoch 199/200:
Train Loss: 0.4054, Train Acc: 0.7975
Val Loss: 0.5578, Val Acc: 0.7392
(base) C:\Users\santh\Desktop\DS504-BDA\BDA-P4>
```

Observations and Analysis:

I. Training Progress:

- The training process involves 200 epochs.
- Each epoch shows the loss and accuracy metrics for both the training and validation sets.

- The loss values gradually decrease over epochs for both training and validation sets, indicating the optimization process is working.
- The accuracy values increase over epochs, indicating the model is learning and improving its performance.

ii. Loss and Accuracy Trends:

- Initially, the loss is relatively high, and the accuracy is low.
- As the epochs progress, the loss decreases, and the accuracy improves.
- There are fluctuations in loss and accuracy across epochs, but an overall decreasing trend in loss and increasing trend in accuracy are observed.

iii. Overfitting Analysis:

 In this case, there is no significant gap between the training and validation accuracies, and the validation loss generally follows the decreasing trend of the training loss, indicating no clear signs of overfitting.

iv. Model Performance:

- The model starts with low accuracy but gradually improves its performance with each epoch.
- Towards the later epochs, the model achieves higher accuracies on both training and validation sets.
- The validation accuracy remains consistently high in the later epochs, indicating that the model can generalize well to unseen data.

iv. Learning Rate: 0.1, Batch Size: 256, Hidden Units: 32, Dropout Rate: 0.5, Epochs: 50:



Observations and analysis:

I. Loss Fluctuations:

- The training loss fluctuates around a certain range throughout the epochs, indicating that the model is not converging effectively.
- There are no clear signs of loss reduction over the epochs, suggesting that the model might not be learning well from the data.

ii. Accuracy Variation:

- The training accuracy hovers around 0.5, indicating that the model's performance is not significantly better than random guessing.
- Validation accuracy also remains around 0.5, which further suggests that the model is not effectively learning patterns from the data.

iii. Stagnant Validation Loss and Accuracy:

- Both validation loss and accuracy remain relatively stagnant throughout the epochs.
- This indicates that the model is not improving its performance on unseen data, which could be due to various factors such as model complexity, insufficient data, or ineffective training strategies.

iv. Possible Overfitting:

- There are signs that the model might be overfitting, as the training loss continues to decrease while the validation loss remains relatively constant.
- This is supported by the fact that validation accuracy is not improving alongside training accuracy, indicating that the model is not generalizing well to unseen data.

v. Potential Issues:

- The model might be too simple or too complex for the given task, leading to poor performance.
- Insufficient data or data quality issues could also be contributing factors to the poor performance.
- Ineffective training strategies such as improper hyperparameter tuning or optimizer selection could also be at play.

Overall, the training process suggests that the current model is not performing well and requires further optimization and experimentation to improve its performance on the given task.

4.Conclusion:

In conclusion, the project aimed to develop and compare two variants of Siamese neural network models for similarity prediction in trajectory data. The primary goal was to leverage deep learning techniques to accurately determine the similarity between pairs of trajectories, which is crucial in various applications such as behavior analysis, recommendation systems, and anomaly detection in motion patterns.

Throughout the project, two distinct architectures were explored: the `SiameseLSTMImproved` and `SiameseLSTM` models. These models were designed to process pairs of input trajectories and predict their similarity scores. The `SiameseLSTMImproved` model featured bidirectional LSTM layers and an

attention mechanism, allowing it to capture temporal dependencies and dynamically focus on relevant segments of the input sequences. On the other hand, the `SiameseLSTM` model utilized standard LSTM layers and concatenated the final hidden states for prediction.

During the experimentation phase, both models were trained and evaluated using a dataset of trajectory pairs, with performance metrics such as loss function values and accuracy scores monitored across multiple epochs. The training process involved optimizing model parameters to minimize loss and improve predictive accuracy on validation data. Additionally, hyperparameters such as hidden layer dimensions and dropout rates were tuned to optimize model performance and prevent overfitting.

The results of the experiments revealed insights into the effectiveness of each model architecture. While both models demonstrated the capability to learn and predict trajectory similarities, the 'SiameseLSTMImproved' model exhibited superior performance, particularly in scenarios were capturing long-range dependencies and distinguishing between subtle differences in trajectories were crucial. The attention mechanism in this model played a vital role in enhancing its predictive capabilities by dynamically weighing the importance of different parts of the input sequences.

In summary, the project showcased the potential of deep learning approaches, particularly Siamese neural networks, in analyzing and predicting similarities in trajectory data. By leveraging advanced architectures and mechanisms such as bidirectional LSTMs and attention mechanisms, it was possible to develop models capable of accurately capturing the complex temporal relationships inherent in trajectory sequences. These findings contribute to the broader field of trajectory analysis and have implications for various real-world applications requiring similarity assessment in movement data. Further research could explore additional enhancements and adaptations of Siamese neural networks to address specific challenges and requirements in trajectory analysis tasks.