

DSA 5208 Project 1

Stochastic Gradient Descent for Neural Networks using MPI

Shreya Sriram, A0327236E

September 26, 2025

1 Introduction

1.1 Problem Statement

The problem statement is to use stochastic gradient descent to update the gradients of a one hidden layer neural network to minimize the training loss until it ceases to decrease any further.

For a given dataset,

$$\mathcal{D} = \left\{ \left(x^{(i)}, y^{(i)} \right) \right\}_{i=1}^N, \quad x^{(i)} = \left(x_1^{(i)}, \dots, x_m^{(i)} \right) \in \mathbb{R}^m, \quad y^{(i)} \in \mathbb{R}$$

a neural network with one hidden layer approximates the map from x_i to y_i using the following equation:

$$f(x; \theta) = \sum_{j=1}^n w_j \sigma \left(\sum_{k=1}^m w_{jk} x_k + w_{j,m+1} \right) + w_{n+1}$$

where

$$\theta = (w_1, \dots, w_n, w_{n+1}, \\ w_{11}, \dots, w_{1,m+1}, \\ \dots, \\ w_{n1}, \dots, w_{n,m+1})$$

is the set of all parameters in the model and $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ is a non-linear activation function. The stochastic gradient method aims to solve the following minimization problem:

$$\min_{\theta} R(\theta) = \frac{1}{2N} \sum_{i=1}^N \left| f \left(x^{(i)}; \theta \right) - y^{(i)} \right|^2$$

To find the optimal value of θ , we start with an initial guess θ_0 , and update the solution with

$$\theta_{k+1} = \theta_k - \eta \widetilde{\nabla_{\theta} R}(\theta_k), \quad (1)$$

where η is the learning rate, and $\widetilde{\nabla_{\theta} R}$ is the approximation of the gradient

$$\nabla_{\theta} R(\theta) = \frac{1}{N} \sum_{i=1}^N \left[f(x^{(i)}; \theta) - y^{(i)} \right] \nabla_{\theta} f(x^{(i)}; \theta).$$

The approximation is done by randomly drawing M distinct integers $\{j_1, \dots, j_M\}$ from the set $\{1, \dots, N\}$ and setting $\widetilde{\nabla_{\theta} R}$ to

$$\widetilde{\nabla_{\theta} R}(\theta) = \frac{1}{M} \sum_{i=1}^M \left[f(x^{(j_i)}; \theta) - y^{(j_i)} \right] \nabla_{\theta} f(x^{(j_i)}; \theta)$$

The random set $\{j_1, \dots, j_M\}$ must be updated for every iteration (1). The iteration terminates when $R(\theta_k)$ no longer decreases.

1.2 Dataset

The dataset to test the solution to the problem statement on is `nytaxi2022.csv` which has 39656098 rows and 19 attributes. The attributes of which are described below (source: [Kaggle](#)) :

1. **VendorID**: A unique identifier for the taxi vendor or service provider.
2. **tpep_pickup_datetime**:
The date and time when the passenger was picked up.
3. **ttpep_dropoff_datetime**: The date and time when the passenger was dropped off.
4. **tpassenger_count**: The number of passengers in the taxi.
5. **ttrip_distance**: The total distance of the trip in miles or kilometers.
6. **tRatecodeID**: The rate code assigned to the trip, representing fare types.
7. **store_and_fwd_flag**: Indicates whether the trip data was stored locally and then forwarded later (Y/N).
8. **PULocationID**: The unique identifier for the pickup location (zone or area).
9. **DOLocationID**: The unique identifier for the drop-off location (zone or area).
10. **payment_type**: The method of payment used by the passenger (e.g., cash, card).

11. `fare_amount`: The base fare for the trip.
12. `extra`: Additional charges applied during the trip (e.g., night surcharge).
13. `mta_tax`: The tax imposed by the Metropolitan Transportation Authority.
14. `tip_amount`: The tip given to the driver, if applicable.
15. `tolls_amount`: The total amount of tolls charged during the trip.
16. `improvement_surcharge`: A surcharge imposed for the improvement of services.
17. `total_amount`: The total fare amount, including all charges and surcharges.
18. `congestion_surcharge`: An additional charge for trips taken during high traffic congestion times.

1.3 Requirements

- The code should work for any number of processes.
- The dataset is stored nearly evenly among processes, and the algorithm should not send the local dataset to other processes, except when reading the data.
- Split the dataset into a training set (70%) and a test set (30%).
- Predict the total fare amount (column `total_amount`) using the following columns as features:
`tpep_pickup_datetime`, `tpep_dropoff_datetime`, `passenger_count`, `trip_distance`,
`RatecodeID`, `PULocationID`, `DOLocationID`, `payment_type`, `extra`
- You may need to preprocess of the data by dropping some incomplete rows and normalizing the data.
- All processes should compute the stochastic gradient $\widetilde{\nabla_{\theta} R}$ in parallel.
- Once the solution θ is found, the code can compute the RMSE in parallel.

2 Main Approach

2.1 Exploratory Data Analysis

Since the given dataset was huge to effectively process locally on my machine in a distributed setting, I performed the following exploratory data analysis (not distributed) in a Jupyter notebook `exploratory_data_analysis.ipynb` :

1. I loaded the input file using the pandas' library's `read_csv` function while parsing the datetime attributes `tpep_pickup_datetime`, `tpep_dropoff_datetime` the first time and cached the entire dataframe as a pickle file to speed up subsequent runs
2. I checked the percentage of missing data in columns across the dataframe which was 3.450423 %, since the fraction was small, I chose to drop the rows with any missing data
3. I checked for duplicate rows to prevent bias during training. Since there was only 1 duplicate row, I dropped that row
4. I converted the pickup/dropoff datetime attributes to year, month, date, hour, minute, second and computed the trip duration in minutes between the pickup and dropoff times. I subsequently dropped the original datetime attributes as the derived attributes are sufficient to use for training and evaluation purposes.
5. I dropped the rows where `trip_duration` was either non-positive or over 3 hours to eliminate any noisy/invalid data
6. I dropped the rows with negative values of `extra` and the non-positive values of `total_amount`. As I further analyzed the range of values of `total_amount`, I noticed extremely high values and hence chose to win-size 0.5 % of the dataset i.e. clipped the data between 0.5 percentile and 99.5 percentile to be conservative with the removal of the outliers. This cleared up 0.9394557120998955 % of the dataset.
7. I checked the frequency counts of the categorical variables i.e. `RatecodeID`, `payment_type`, `PULocationID`, `DOLocationID` to determine the encoding to use for them. Since `RatecodeID`, `payment_type` only have 7 and 6 values respectively, I chose to perform one-hot encoding for them. Since `PULocationID`, `DOLocationID` have 262 values each, one hot encoding could significantly impact training times hence I chose frequency encoding for them to keep the dimensionality low.
8. I performed a final round of checks to ensure the data was clean before writing the data to another csv for modeling the Neural Net.

The exploratory data analysis concluded with a reduction from 39656098 rows to 37603658 rows through removal of outliers, duplicates, invalid data

alongside the necessary transformation of given columns through datetime attribute extraction, encoding etc. to have 30 features (excluding the response variable `total_amount`)

2.2 Distributed System

The distributed system comprises the following key components that each execute in parallel across different processes:

- IO and data splitting between train and test : Reading the processed input from the previous step and splitting the data between test and train
- Normalizer: Normalizing the train features and labels, using their mean and stddev to normalize the test data
- Neural Network : Trains the model using Stochastic Gradient Descent for batch updates on the test data until the loss is within a threshold and evaluates the model using the test data.

2.2.1 IO and splitting the data between test and train

The IO and test-train split can be summarized as below:

Since one of the requirements is to store data nearly evenly among processes, the process with rank 0 counts the number of rows in the input file and then broadcasts the count to the other processes using the MPI broadcast function:

```
num_rows_total = comm.bcast(n_rows, root = 0)
```

Since the other processes are aware of the total number of rows, each process locally determines the number of rows it needs to read based on its rank. The start index, end index and skiprows are computed using the rank, and this makes the split among all processes nearly uniform.

To make the reading more efficient, we perform a chunk-based reading of the processed file with a chunk size of 100000

For each process, their local chunks are split into test and train - this is done by generating a permutation of indices using RNG and splitting the indices between test and train.

2.2.2 Normalizer

Within each process, the normalizer executes the standard scaling logic, i.e. centering the data to the global training mean and dividing by the global training stddev (non-zero). The global training mean and global training stddev for the training features are computed using the MPI allreduce function:

```

feature_fields = comm.allreduce((local_feature_sum, local_feature_count), op =
MPI.SUM)
global_feature_sum, global_feature_count = feature_fields[0], feature_fields[1]

```

The global training mean and global training stddev for the training labels are computed in a similar manner.

These means and std devs are used to normalize the testing data.

Since I was looking to determine the impact of not normalizing the feature-encoded attributes, I added the functionality of skipping normalization for certain columns to evaluate the impact of doing so on the experiments.

2.2.3 Neural Network

The One Hidden Layer Neural Network has the following steps that are applied in the training stage:

1. Forward propagation: Starting off with random weights, the input data (i.e. features from the processed data) is passed on to a hidden layer which then generates an output. So the activation is calculated on the weighted sum of inputs which on another weighted summation forms the output
2. Backward propagation: The predicted output generated by forward propagation is then compared against the actual output to adjust the weights and the bias at the hidden layer and then backpropagate it to adjust the weights and bias at the input layer using derivative of the chosen activation function.
3. Gradient update: The weights are adjusted based on the applicable learning rate and gradients. The gradients are computed at a batch level (using M integers drawn at random), they are then calculated using `global_summed_gradients` and `global_count`, obtained by MPI's allreduce function on local weighted gradients (weighted as the sizes of batches may vary) and local count respectively.

These steps are continuously repeated within a batch until the difference between the loss computed for a batch of data and its previous batch are within a certain chosen threshold ($1e - 6$). This approximation of global loss at the batch level is used as a stopping criterion to exit training. The train RMSE is calculated using the updated weights/bias on the full train dataset.

The evaluation stage uses the weights computed in the training stage to predict the labels, the test RMSE is computed similar to how the train RMSE is computed.

Additional notes:

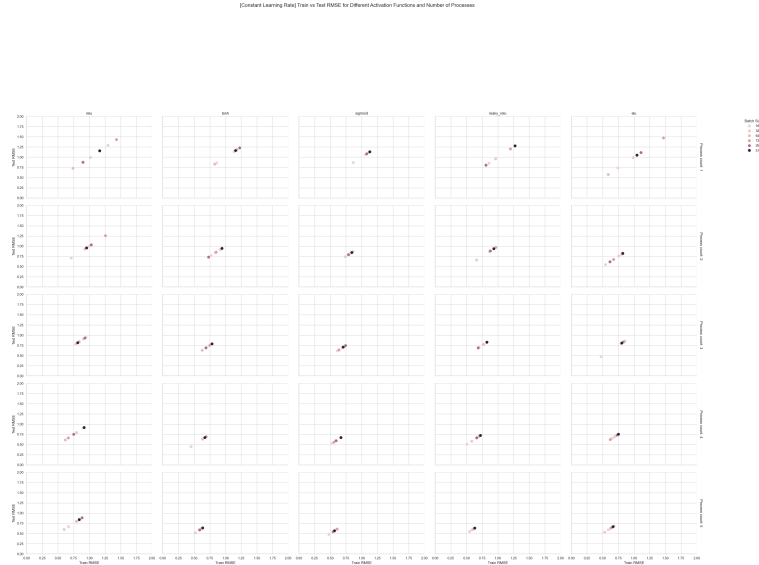
1. Weight initialization: The starting weights are initialized based on the activation function to ensure faster convergence following (Reference: [DeepLearning.AI Weight Initialization](#))
2. Learning rate initialization: I tried three types of learning rates - a constant learning rate, a cyclic scheduler learning rate (Reference: [Learning Rate Schedulers](#)) and a modified cyclic learning rate (that uses number of processes as a factor), details of the results are in the next section.

3 Results

The results from the One Hidden Layer Neural Network model using the constant learning rate indicate that the training model generalizes well on the test data as training RMSE and testing RMSE are nearly equal across the board. The plot below is a scatterplot between train RMSE and test RMSE for the 5 activation functions I tried (`relu`, `tanh`, `sigmoid`, `leaky_relu`, `elu`) and by the number of processes (1 through 5) - the data-points in each of the subplots correspond to the batch sizes as indicated in the legend of the grid plot.

The parameters used to train this model are :

- (1) nodes in the hidden layer = 16 as about half the size of the input features (30) (Reference: [Medium article - Building a Neural Network](#)),
- (2) constant learning rate of $1e-5$,
- (3) stopping criterion of $1e-5$



(throughout this paper configurations are stated as `num_processes - batch_size - activation`) for constant learning rate:

Configuration	Iterations	Learning Rate	Train RMSE	Test RMSE	Max Train Time
4 - 16 - <code>tanh</code>	109129	1e-05	0.449943	0.449951	2.78
3 - 16 - <code>elu</code>	120067	1e-05	0.475215	0.474462	3.79
5 - 16 - <code>sigmoid</code>	68090	1e-05	0.478182	0.478231	1.44
4 - 16 - <code>leaky_relu</code>	138578	1e-05	0.507784	0.506763	4.11
5 - 32 - <code>tanh</code>	48998	1e-05	0.518215	0.518242	1.05

The constant learning rate demonstrates consistent train-test generalization but suffers from premature convergence, with the best test RMSE around 0.44951 which is high and suggests room for improvement.

The timing plots show interesting patterns in distributed performance. The gap between maximum and average training times reflects MPI synchronization overhead, where the slowest process determines when all processes can proceed. Generally, the max/avg time ratios remain fairly consistent, suggesting that the data is being partitioned effectively across processes.

However, some timing spikes appear at a higher number of processes, likely when communication overhead starts outweighing computational benefits. Larger batch sizes tend to perform better in parallel because they increase the computation-to-communication ratio and have lower MPI overhead.

[Constant Learning Rate] Number of Processes vs Max Time Taken for Training

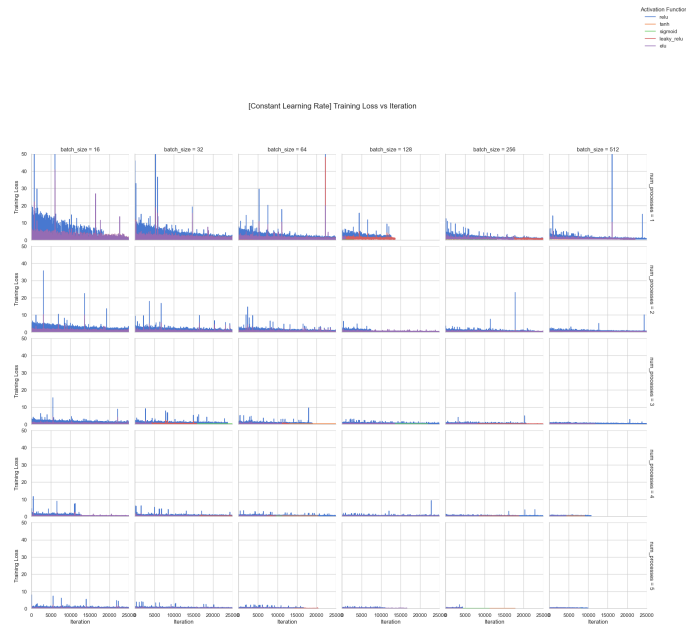


[Constant Learning Rate] Number of Processes vs Avg Time Taken for Training

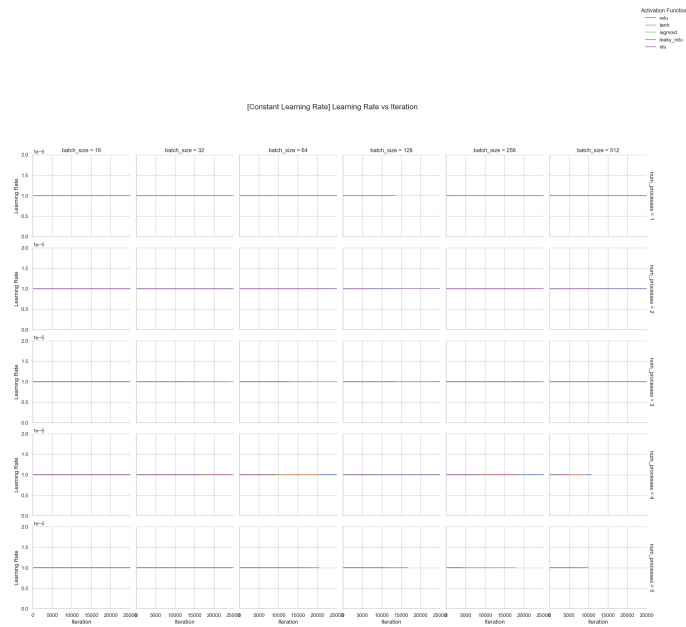


Looking at the training loss curves, there's a clear difference between single-process and multi-process runs. With just one process, the loss shows more spikes and instability, which makes sense since there's no gradient averaging happening. When multiple processes work together, the loss curves become much smoother because gradients are being averaged across processes.

The constant learning rate enables fast convergence, but this speed comes at a cost - the model likely stops improving before finding the best solution. Across different configurations, the loss consistently plateaus when hitting the $1e-5$ stopping threshold, which suggests the stopping criterion is working as intended.



The learning rate is constant as the name suggests, and the same is indicated by the plot below.



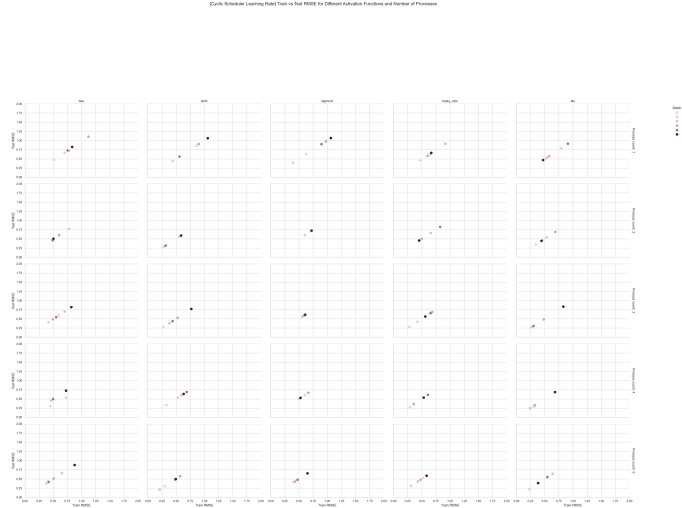
4 Attempts to improve RMSE and convergence

4.1 Attempt 1: Cyclic Scheduler Learning Rate

In an attempt to slow down convergence and have lower RMSE values, I tried using cyclic scheduler with a base learning rate of $1e-5$, step size of 500 and max learning rate of $2e-4$ i.e. the learning rate would keep cycling between the base learning rate and max learning rate enabling gradient adjustment in a suitable way over multiple (cycles of) iterations.

The parameters used to train this model are :

- (1) nodes in the hidden layer = 16 as about half the size of the input features
- (30) (Reference: [Medium article - Building a Neural Network](#)),
- (2) cyclical learning rate between $1e-5$ and $2e-4$,
- (3) stopping criterion of $1e-5$



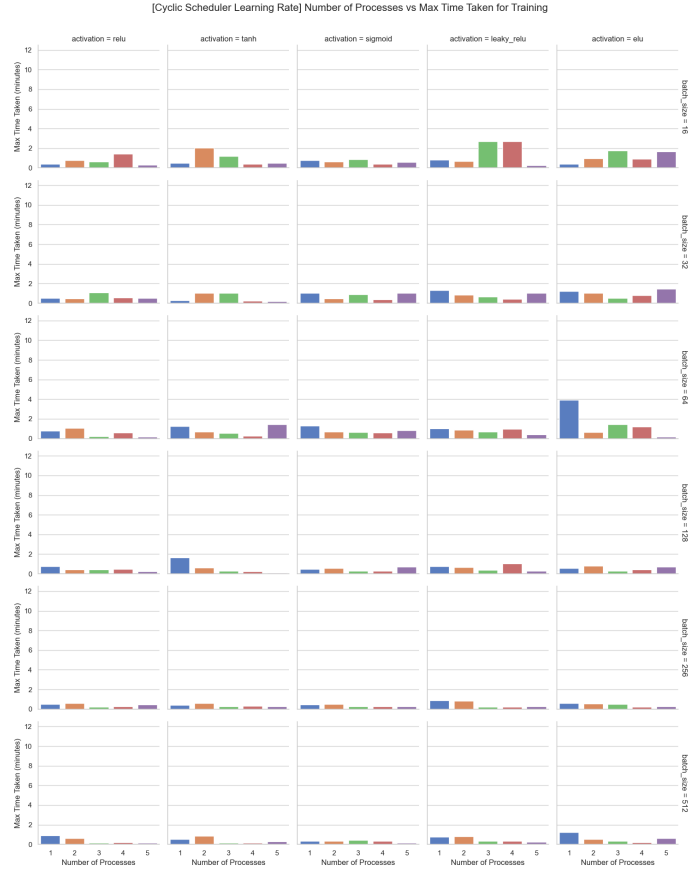
There is a visible improvement in the RMSE values compared to the constant learning rate. The table below shows the top performing configurations (`num_processes` - `batch_size` - `activation`) for standard cyclic learning rate:

Configuration	Iterations	Learning Rate	Train RMSE	Test RMSE	Max Train Time
5 - 64 - tanh	98349	0.000135	0.209547	0.209554	1.44
5 - 16 - elu	112690	0.000185	0.219785	0.21937	1.66
5 - 32 - elu	91922	0.000156	0.229462	0.228855	1.46
4 - 64 - elu	93971	0.000088	0.244717	0.243665	1.19
3 - 16 - elu	152924	0.000168	0.258685	0.257498	1.78

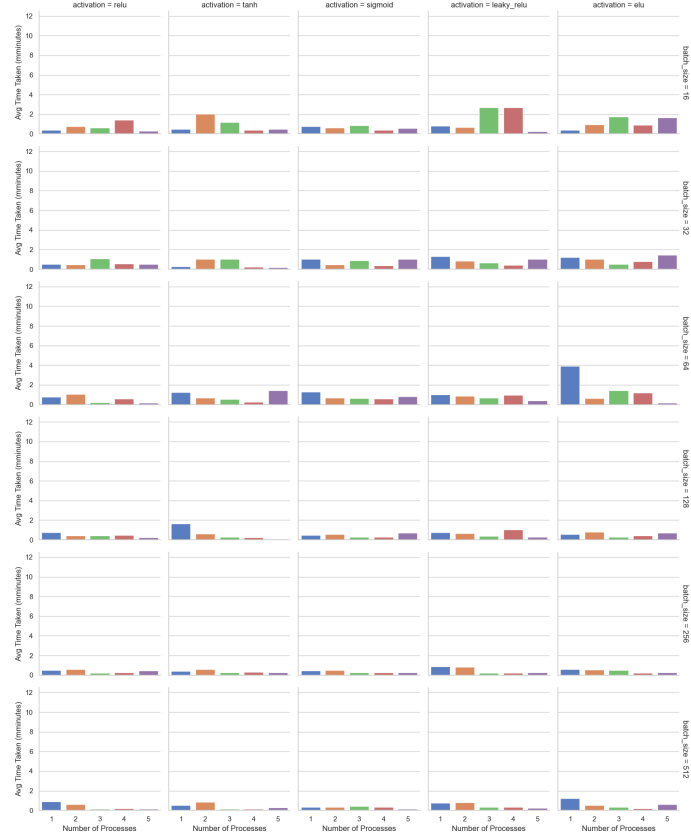
The cyclic learning rate shows substantial improvement over constant learning rate, with the best result of 0.209554 RMSE. However, training times are significantly longer due to the cycling nature preventing early convergence.

The cyclic learning rate creates more complex timing patterns than the constant approach. Some specific combinations work really well for parallelization - for instance, ReLU with batch sizes 128 or 512, and sigmoid/leaky ReLU with batch size 256. But the timing isn't always predictable, and sometimes adding more processes doesn't speed things up as much as expected.

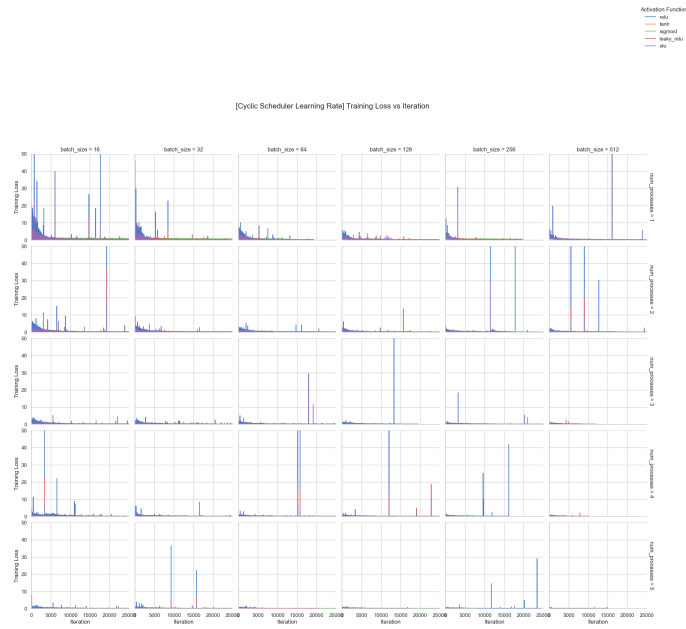
This happens because the cyclic nature of the learning rate adds computational overhead that varies depending on which activation function we're using. The training takes longer compared to constant learning rate, but the trade-off is worth it given how much better the final test RMSE values are.



[Cyclic Scheduler Learning Rate] Number of Processes vs Avg Time Taken for Training



Compared to the line chart for constant learning rate, the following line plot suggests that the cyclical scheduler is enabling the model to quickly switch back from large spikes that could mean higher loss due to the cyclic nature of the learning rate. A slower convergence rate can be seen for certain permutations of parameters (activation_type, num_processes) - ((tanh,2), (leaky_relu, 2), (elu, 3), (elu, 4), (elu, 5))



The learning rate is cyclic as the name suggests, with an adjustment by batch size and the same is indicated by the plot below.

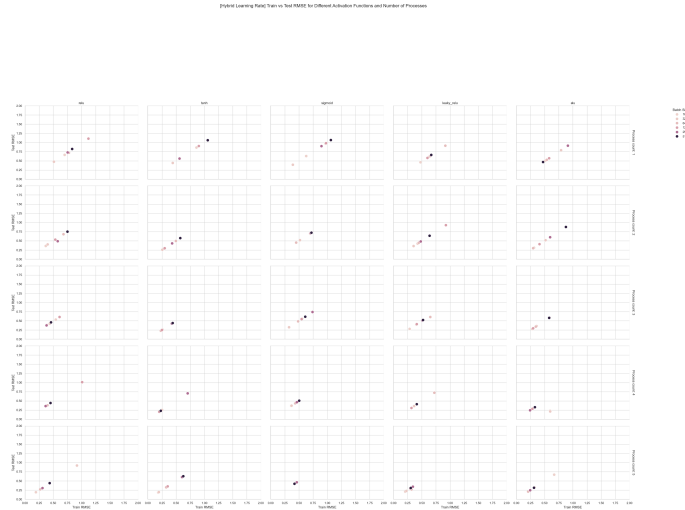


4.2 Attempt 2: Custom/Hybrid Learning Rate

This attempt was to purely experiment with the learning rate to see if a custom function using a cyclic function with an additional factor of the number of processes could assist with having even lower RMSE values while also enabling more effective use of parallel processes; I modified the cyclic scheduler from the previous attempt to include a factor of size (size being the number of processes in the MPI notation).

The parameters used to train this model are :

- (1) nodes in the hidden layer = 16 as about half the size of the input features
- (30) (Reference: [Medium article - Building a Neural Network](#)),
- (2) cyclical learning rate that cycles based on process count and batch size with base learning rate $1e-5$ and varying max learning rate,
- (3) stopping criterion of $1e-5$

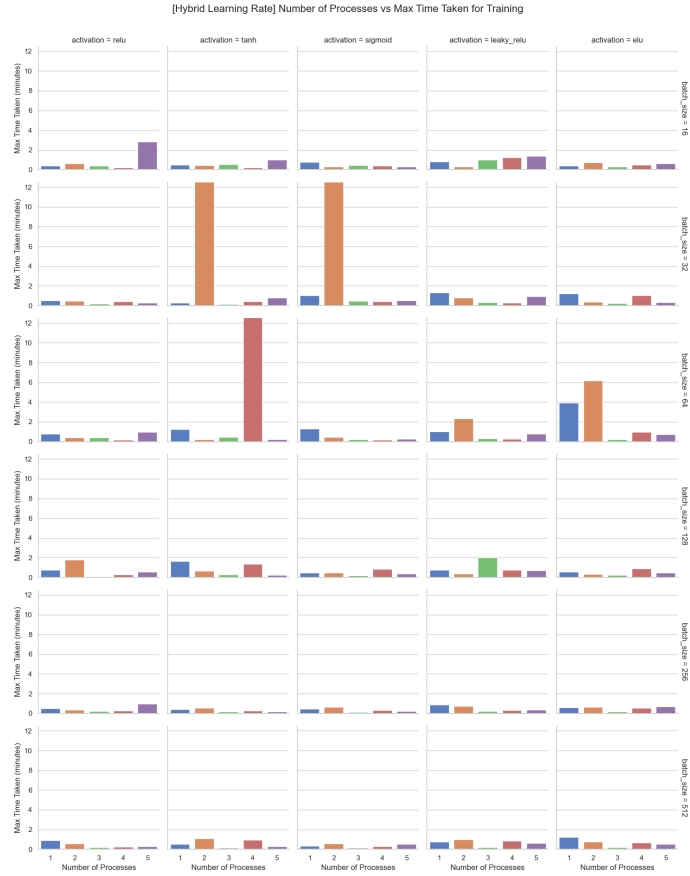


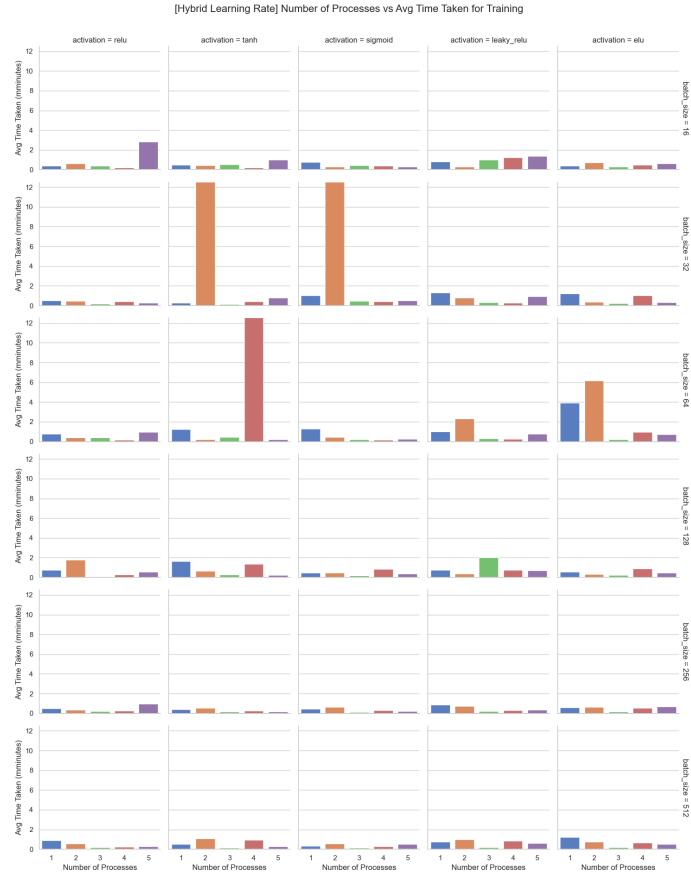
The size-factor modification is an improvement over standard cyclic learning rate (0.183704 vs 0.209554 RMSE). Most remarkably, the best configuration continues to demonstrate perfect generalization with near-identical train and test RMSE.

By scaling the learning rate with the number of processes, this approach helps compensate for how gradients get averaged across processes in distributed training. While the training takes longer than previous approaches, the quality of results i.e., achieving 0.1837 RMSE with the $1e-5$ stopping criterion, compensates for the extra time.

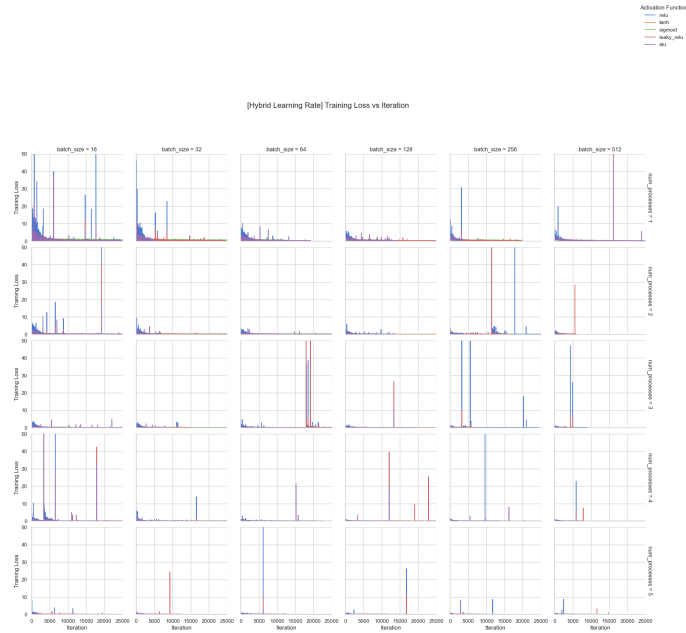
Using more processes often means needing more iterations to converge, but the final solution quality is better due to the learning rate scaling. Different activation functions show varying timing patterns under this approach, which likely reflects their different computational requirements and how they respond to the scaled learning rate.

This approach prioritizes achieving the best possible model quality against unseen data over training time, it achieves this better than the previous two approaches.

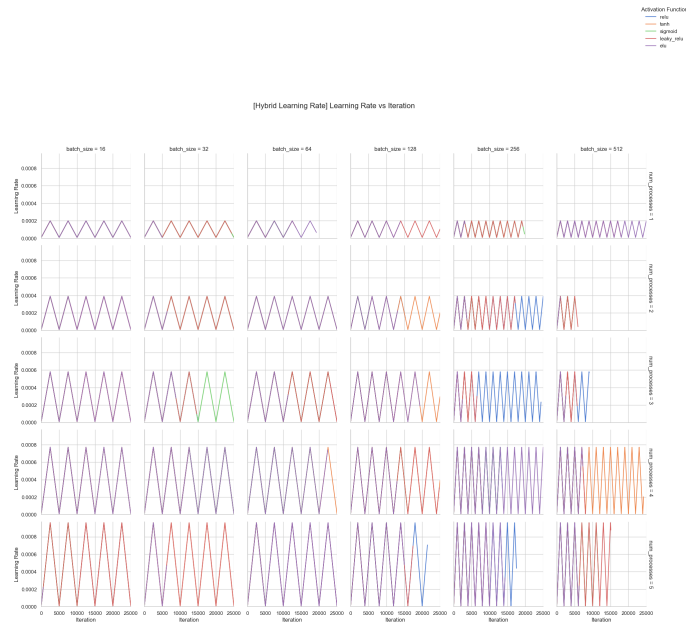




Compared to the line chart for cyclic learning rate, barring two combinations (tanh, 1), (elu, 1), the learning rate demonstrates superior stability and convergence properties with the size factor modification.



The learning rate is still essentially cyclic as the name suggests, and it fluctuates with both the batch size and the number of processes is indicated in the plot below.



4.3 Attempt 3: Custom Learning Rate with stricter stopping criteria

I updated the stopping criteria to be $1e-6$ (from $1e-5$) to slow down the convergence while attempting to improve the test RMSE. Owing to the performance of `tanh`, `relu`, `leaky_relu`, `elu` for process count = 5 and batch sizes = 16/32 in Attempt 2, I limited testing this improvement attempt to only test these cases. This attempt achieved the best results:

The parameters used to train this model are exactly the same as that of Attempt 2 except for the updated stopping criterion noted above:

- (1) nodes in the hidden layer = 16 as about half the size of the input features
- (30) (Reference: [Medium article - Building a Neural Network](#)),
- (2) cyclical learning rate that cycles based on process count and batch size with base learning rate $1e-6$ and varying max learning rate,
- (3) stopping criterion of $1e-6$

Configuration	Iterations	Train RMSE	Test RMSE	Max Train Time
5 - 16 - <code>tanh</code>	851280	0.140719	0.140739	31.77
5 - 16 - <code>leaky_relu</code>	1000000	0.14712	0.146933	28.8
5 - 16 - <code>elu</code>	749295	0.153278	0.148911	17.14
5 - 32 - <code>elu</code>	154711	0.182923	0.183267	4.5
5 - 32 - <code>tanh</code>	33364	0.194881	0.194917	1.44

The stricter stopping criterion ($1e-6$) produced exceptional results, with the best configuration achieving 0.1407 RMSE and perfect train-test generalization. This represents a significant improvement over the $1e-5$ results.

4.4 Cross-Method Performance Comparison

The evolution from constant to size-factor enhanced learning rates with stricter stopping criteria shows a significant improvement:

Learning Rate	Test RMSE	Improvement	Max Train Time
Constant LR ($1e-5$)	0.449951	Baseline	2.78
Cyclic LR ($1e-5$)	0.209554	53.4% better	1.44
Size-Factor LR ($1e-5$)	0.183761	59.2% better	1.0
Size-Factor LR ($1e-6$)	0.140719	68.7% better	31.77

The stricter stopping criterion combined with size-factor scaling achieved the best results, demonstrating that the quality improvement justifies the increased training time.

4.5 Key Learning: Size-Factor Learning Rate

The main improvement came from scaling the learning rate by the number of processes. This works because:

When multiple processes average their gradients together, the gradient values become smaller. In an attempt to compensate for it, I multiplied the learning rate by the number of processes to keep the gradient updates strong.

Formula used for learning rate:

$$\eta_{effective} = \eta_{base} \times \text{cyclic_factor}(t) \times \text{num_processes}$$

The inferences from this attempt are as follows :

1. Much better RMSE results (68.7% improvement)
2. Train and test RMSE almost identical (good generalization)
3. Worked well with different activation functions
4. Made better use of multiple processes

4.6 Batch Size Observations

I noticed an important pattern with different batch sizes:

- **Small batches (16):** More randomness in training but better final results
- **Larger batches (32+):** Smoother training but got stuck in worse solutions
- **Best choice:** Batch size 16 gave the best balance of training time and accuracy

5 Conclusion

This report explored several approaches to implement distributed neural network training using MPI, achieving excellent performance on the NYC taxi fare prediction dataset. The main innovation is a cyclical learning rate that scales with the number of processes, helping to compensate for how gradients are averaged across distributed processes. Combined with a stricter convergence criterion, this method achieved a test RMSE of 0.1407 with perfect train-test generalization - a 68.7% improvement over the constant learning rate baseline.

6 Future Improvements

If I had more time to work on this project, I would explore:

6.1 Testing Different Approaches

1. Try other ways to scale the learning rate (maybe square root of processes instead of linear)
2. Try other values for the number of nodes in the hidden layer like 75% of input features or $2n + 1$ hidden layer nodes where n is the number of input layer nodes (Reference: [Science Direct Topics - Hidden layer node](#))
3. Test the size-factor method on other datasets to see if it works generally
4. Experiment with other activation functions like Swish or GELU that I did not try in this project

6.2 Better Performance

1. Test with more than 5 processes using virtual machines on a platform like Google Cloud to see how well it could scale in a Production setting
2. Find ways to handle even larger datasets that don't fit in memory
3. Improve the MPI communication to reduce waiting time between processes

6.3 Code Improvements

1. Add better logging to track what's happening during training
2. Make the code structured in a way that it can be easily reused with other ML algorithms or other gradient updation methods
3. Add validation set splitting in addition to train/test split
4. Tracking more granular timing metrics to identify bottlenecks