

---

# LSTM XOR Project Report

Jayaram Kuchibhotla  
jayaramkuchibhotla@gmail.com

## Abstract

This project report describes the implementation and results of the LSTM model which outputs the parity of the input sequence of binary bits. It has the following sections 1.Introduction 2.Description of the datasets 3.Explanation of model architecture and implementation 4. Results 5.Improvements and Conclusion

## 1 Introduction

RNN(Recurrent Neural Networks) models are useful for training sequential data in the cases where the previous information has to be remembered to predict the output .The basic problem that is solved by RNN is prediction of the next word given few words in a meaningful sentence. Information of the immediate past word is just not sufficient to predict the current word. The context is understood only when the model has knowledge of past few words in the sentence. When it comes to practical implementation, RNN suffers from the problem of vanishing or exploding gradients for data inputs having long sequences. This problem is solved to an extent by LSTM (Long Short Term Memory) models and hence these are preferred over RNN for modeling sequential data.

In the current project, parity has to be predicted (if the count of ones is odd ,output is 1 .It is 0 if the count of ones is even). Mathematically, the bitwise operation XOR between all the bits in data helps to output the parity. The LSTM model has to be trained to implement this XOR functionality

## 2 Description of the datasets

Two types of datasets are used as input to the model in 2 different instances of execution.

Type1 : 100000 binary strings with consistent length value 50 in both input and output.  
Type2 : 100000 binary strings where each data sample length can have any value between 1 to 50 in input and the corresponding output. In both these data sets , the output data bit at index i is the parity of the data bits in the input sample from index 0 to i-1. They are split into 80000 training samples and 20000 test samples.

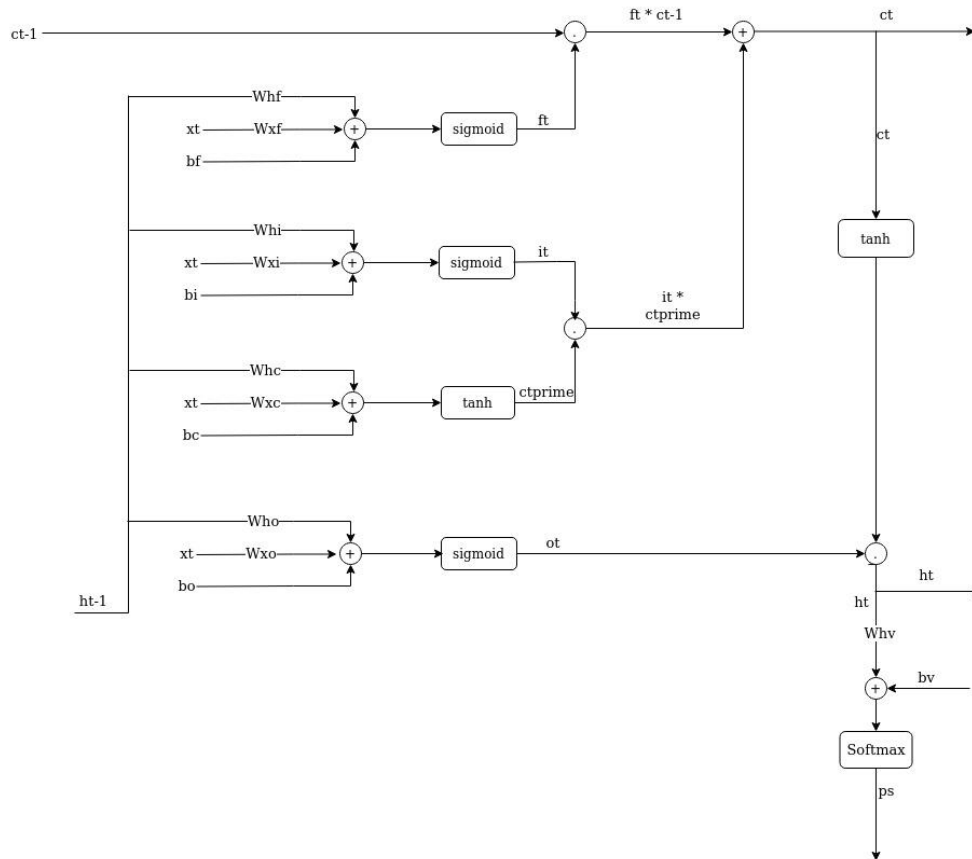
Train dataset split into batches : The 80000 data samples in type1 are divided into 2500 batches each having size of 32. In type2 , 80000 samples are divided into 2521 batches each having variable size that is  $\leq 32$ . The reason for this variability is that the data in type 2 has variable lengths and the batches are divided in such a way that each batch has equal length data samples. This helped in easier implementation of Mini batch gradient descent using 5 different processes. Training and Testing are done separately with Type1 and Type2 batches .Results obtained are independent of each other

When Type1 dataset is used for program execution , each process handled 500 batches of data and in execution with Type2 dataset, 501 batches of data are handled by each process. The 2521th batch having 30 samples of data was not used for training . The train data is equally split among 5 processes and distribution of left over data(30 samples in 2521th batch) is not handled . An effective way is to have right multiple of processes distributing

48 data equally or distribute the left over data equally among the fixed number of processes.  
 49 This is left for future improvement work

50

### 51 3 Explanation of model architecture and implementation



52 The LSTM model has an LSTM cell which is repeated for 'n' times where n is the number of  
 53 timesteps in the input data sample. The previous cell hidden state 'ht-1' and previous cell  
 54 state 'ct-1' (cell state is current long-term memory of the network) are fed as input to the  
 55 current cell along with the current time step data input

56 The above diagram shows the operations in one LSTM cell. It has 3 different gates  
 57 1.Forget gate( $ft$ ) 2.Input gate( $it$ ) 3.Output gate( $ot$ ) along with current cell state  
 58 ( $ctprime$ ).Each of these units have separate weights and biases associated with it . The  
 59 weights are in dot product with the hidden state inputs and the input data for that particular  
 60 time step . '+' symbol inside a circle is an addition operation while '.' symbol inside a circle  
 61 is the pairwise multiplication operation. Outputs from these 4 units are passed through non  
 62 linear operations such as  $\tanh$  or sigmoid as shown in the diagram above

63 Forget gate : If the previous hidden state value ( $ht-1$ ) and the current time step input ( $xt$ )are  
 64 given as input , it outputs values which, when passed through the sigmoid function, gives the  
 65 matrix values  $ft$  in the range  $[0,1]$  . $ft*ct-1$  operation helps in forgetting(making them zeros)  
 66 few values and in preserving remaining values in  $ct-1$ . Forget gate's weights  $W_{hf}$ ,  $W_{xf}$  and  
 67 bias  $b_f$  will be optimized so as to forget the right values of  $ct-1$  to contribute to the correct  
 68 output

69 Input gate: Current time step input ( $xt$ ) and previous hidden state value ( $ht-1$ ) are given as  
 70 inputs to this gate ,uses weights  $W_{hi}$ ,  $W_{xi}$  and bias  $b_i$ , the output value is passed through  
 71 sigmoid function to give values ' $it$ ' in the range  $[0,1]$  .Current cell state ( $ctprime$ ) is result of  
 72  $\tanh$  operation of factor which is again the combination of the current time step input( $xt$ )  
 73 value and previous hidden state value( $ht-1$ ) using weights  $W_{hc}$ ,  $W_{xc}$  and bias  $b_c$  .The use of

74 ctpime is to modify the information in  $x_t$  given the  $h_{t-1}$  values . Tanh range is  $(-1,1)$ , the  
75 impact of values close to lower bound is reduced while the impact of values closer to upper  
76 bound is remained as it is.

77 The operation  $it * ctpime$  nullifies the effect of few values in current cell state (  $ctprime$ ) and  
78 let other values pass through as they are . This is due to the usage of sigmoid non linear  
79 activation in input gate 'it'.

80  $ft * ct - 1 + it * ctpime$  operation produces the next cell state value 'ct'.  $ct$  is group of values  
81 which are obtained by retaining and forgetting certain values from current cell state and  
82 previous cell state.

83 Output gate: Current time step input ( $x_t$ ) and previous hidden state value ( $h_{t-1}$ ) are given as  
84 inputs to this gate ,uses weights  $W_{ho}$ ,  $W_{xo}$  and bias  $b_o$  , the output values are passed through  
85 sigmoid function to give values 'ot' in the range  $[0,1]$  . Tanh operation on the next cell  
86 state( $ct$ ) constrains the values to the range  $(-1,1)$  i.e. reduces the impact of few values and  
87 while other values are passed on as they are before.

88 The operation  $ot * \tanh(ct)$  filters out only the needed values from  $ct$  and gives them as next  
89 hidden state ( $h_t$ ) values. Thereafter, the  $h_t$  values passed through a linear layer with weight  
90  $W_{hv}$  and bias  $b_v$  and the resultant outcome is passed through a softmax layer to get the  
91 output probabilities.

92 The 'n' number of cells are cascaded one after another getting the values  $ct-1$  and  $h_{t-1}$  from  
93 previous cell.

94 Weights and biases initialization: Based on references[1][2][3] and[4] , the weights which  
95 are in dot product with hidden state values from previous cell, in every gate and current cell  
96 state are obtained by 'orthogonal initialization'. The number of units in hidden layer is 100  
97 and dimension of these matrices are (100,100). The remaining eights are consequence of  
98 'xavier uniform' initialization. Weights which multiply with input data have the dimension  
99 (batch size ,100) and the weight values in output linear layer have the dimension ( 100,2)

100 Except bias in output linear layer all other bias values have the shape (1,100) and the shape  
101 of bias in output layer is (1,2)

102 Forget gate bias values are set to ones and other bias values are set to zeros.

103 The input data and weights have rows which is equal to batch size but the bias values have  
104 single row . Addition of bias values was possible due to array broadcasting

105 Training: All batches of data are made available in every process . The start and end index of  
106 the group of batches for a process is calculated based on the rank of that process and in this  
107 way the group of batch data used differs from process to process. Mini batch gradient  
108 descent is used for optimization and training is done in parallel for each batch of data in the  
109 groups and the process can be split into 3 main steps a. Forward propagation b. Back  
110 propagation c. Optimization

111 a. Forward Propagation :

112 Before forward propagation , the weights and biases are made sure to be same across all the  
113 processes. MPI methods Send and Recv were used to get the latest parameters from process  
114 with rank 0 to the remaining processes. The previous cell state and previous hidden state  
115 values are initialized to zeros which would be used by forward propagation for 1<sup>st</sup> timestep.

116 The input and output data consists of two types of characters i.e. '0' and '1' . For these 2  
117 classes , softmax classification can be used which necessitates one hot encoding of the input  
118 data. Input character '0' at a time step is encoded as  $[1,0]$  and character '1' is encoded as  
119  $[0,1]$ . The input and output data at every time step in the whole batch is encoded in this  
120 manner.

121 The same weight and bias values are used in the loop for n times where n is the number of  
122 time steps, to predict the output probabilities.

123 The input batch data of size (batch size, 2) is transformed into output probabilities of size  
124 (batch size,2) in the forward propagation

125 Cross entropy loss function is used to calculate the loss and the loss values of all the samples  
 126 are summed up to form batch loss based on reference[6]  
 127 The output result values are returned as dict from the forward propagation function  
 128 b. Back Propagation:  
 129 The gradients of weights and biases are declared as class members as they have same  
 130 dimension as their corresponding parameters. In the back propagation function, at the  
 131 beginning, gradients of biases are declared with row size equaling to the batch size. After the  
 132 back prop is complete for a batch ,these 2 steps are followed, 1.The gradients of biases are  
 133 summed and assigned to bias gradients which are declared before as class members. 2. The  
 134 weight and bias gradients are sent from every process to process with rank 0. Here, all the  
 135 gradients are reduced i.e. summed to form batch gradients that are used for optimization.  
 136 Reference[7] was used to arrive on this idea of implementation  
 137 References [8],[9],[10] ,[11],[12] and figure above were used to derive and implement the  
 138 back propagation  
 139 c. Optimization:  
 140 Adam optimizer is implemented from scratch using reference [13] and except the learning  
 141 rate, remaining hyper parameters for optimizer in the current implementation are same as in  
 142 the Adam optimizer paper.  
 143  
 144 Prediction: Two types of output predictions are used .  
 145 1. One method is to check if the last bit of the predicted data matches with the last bit of the  
 146 output test data simple  
 147 2. Another method is to check if every bit in the predicted data matches with every bit in the  
 148 relevant time step of the output test data sample  
 149 Accuracy Calculation:  
 150 
$$\text{Accuracy} = \frac{\text{No. of correctly predicted data samples}}{\text{No. of output data samples}} * 100$$

151

## 152 **4 Results**

153 The batch loss values from all the processes are collected in process with rank 0 and  
 154 averaged to get loss value for one iteration. The averaged batch loss from every iteration is  
 155 summed up to get the loss for the epoch  
 156

157 Training and prediction for Fix length batch data:

158 Initially, when uniform random weights were used instead of Xavier uniform and orthogonal  
 159 initialized weights , the loss value continued to increase for each epoch(Epoch:0 loss 1176  
 160 to epoch:20 loss 8843) .When the latter methods of weight initialization were used , loss  
 161 decreased steadily . This shows the exploding gradients problem has occurred when uniform  
 162 random weights are used. LSTM models do not solve the vanishing and exploding gradient  
 163 problem of RNN completely, rather they mitigate or delay this problem. There exists at least  
 164 one path in LSTM through which gradients can explode. This problem could have happened  
 165 in the current implementation and resolved after the correct weight initialization.  
 166

167 Thereafter, it took six runs to achieve the convergence. In all these runs , the learning rate  
 168 parameter and epochs are changed , other hyper parameters such as no. of units in hidden  
 169 layer =100 , batch size =32 and random seed = 0 remain the same.  
 170

171 Run1:

172 Learning rate = 0.0005, epochs =12

173 Loss does not decrease significantly even after 5 epochs

174 Training stopped  
 175

176 Run2:

177 Learning rate = 0.001 (same as in Adam optimization paper) ,epochs =12  
178 Loss increases significantly for every epoch  
179 Training stopped

180  
181 Run3:  
182 Learning rate = 0.00075, epochs =12  
183 Last epoch loss : 36  
184 min batch Loss : 29.415  
185 correct\_prediction\_count using last bit =19616  
186 Prediction Accuracy using last bit =98.080000  
187 correct\_prediction\_count for full length =7910  
188 Prediction Accuracy for full length =39.550000  
189

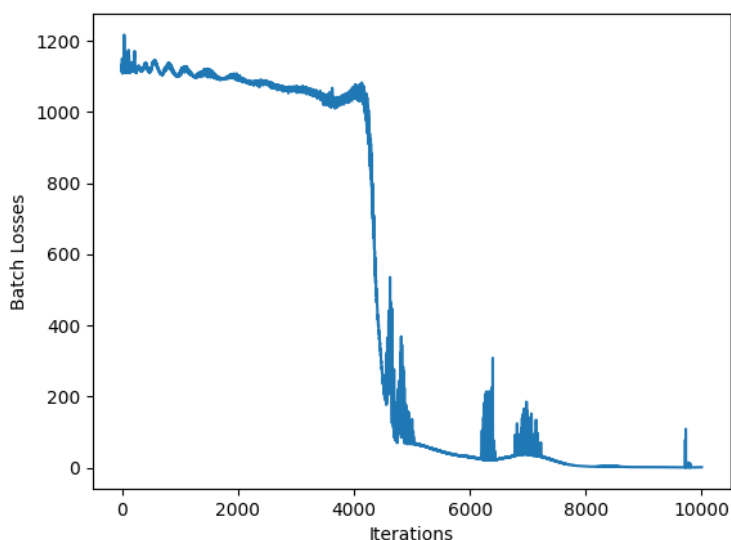
190  
191 Run4:  
192 Learning rate = 0.00080, epochs =12  
193 Last epoch loss = 1780  
194 min\_batch\_loss =1087.689  
195 correct\_prediction\_count using last bit =9922  
196 Prediction Accuracy using last bit =49.610000  
197 correct\_prediction\_count for full length =0  
198 Prediction Accuracy for full length =0.000000  
199

200  
201 Run5:  
202 Learning rate = 0.00076 , epochs =12  
203 Last epoch loss = 1960  
204 min\_batch\_loss =110.529  
205 correct\_prediction\_count using last bit =10195  
206 Prediction Accuracy using last bit =50.975000  
207 correct\_prediction\_count for full length =0  
208 Prediction Accuracy for full length =0.000000  
209

210 Run6 and Final result for fixed length data:  
211 Learning rate = 0.00075 , epochs =20  
212 Last epoch loss = 1  
213 min\_batch\_loss = 1.075  
214 correct\_prediction\_count using last bit = 19978  
215 Prediction Accuracy using last bit = 99.890000  
216 correct\_prediction\_count for full length = 18932  
217 Prediction Accuracy for full length = 94.660000  
218

219

Below graph is loss plot for Run 6 using number of iterations on x-axis  
Number of iterations =10000  
(500 values of averaged batch losses in 20 epochs)



233 Training and prediction for Variable length batch data :  
234 Varying learning rates for the ranges of epochs instead of single learning rate helped in  
235 convergence. Loss curve and decrement seem to be complicated when compared to those  
236 of the fixed length data .  
237 The number of epochs required are more in number compared to what is used for the fixed  
238 length data.  
239 It took 27 runs to figure out the correct learning rates and no. of epochs to get prediction  
240 accuracy above 92% for variable length data.  
241 Batch size =32 ,random seed =0 was used in all the runs below  
242 Except for variation in handling data set, number of epochs and learning rates , no other  
243 change is made in the model for training variable length data set.  
244  
245 Run1:  
246 Learning rate = 0.00075 , epochs =20  
247 Training loss for Epoch:11 is 1909  
248 Training stopped after 12 epochs due to high fluctuations in loss  
249  
250 Run2:  
251 Learning rate = 0.00050 ,epochs =20  
252 Training loss for Epoch:8 is 2359  
253 Training stopped after 9 epochs due to high fluctuations in loss  
254  
255 Run 3:  
256 Learning rate = 0.00040 , epochs =20  
257 Training loss for Epoch:11 is 599  
258 Training loss for Epoch:12 is 2280  
259 Training stopped  
260  
261 Run 4:  
262 Learning rate = 0.00035 , epochs = 20  
263 Last epoch loss is 44  
264 min\_batch\_loss :=4.391  
265 correct\_prediction\_count using last bit =18878  
266 Prediction Accuracy using last bit =94.390000  
267 correct\_prediction\_count for full length =11162  
268 Prediction Accuracy for full length =55.810000  
269  
270 Run 5:  
271 Learning rate = 0.00035 ,epochs = 28  
272 Last epoch loss is 464  
273 min\_batch\_loss :=3.319  
274 correct\_prediction\_count using last bit =16397  
275 Prediction Accuracy using last bit =81.985000  
276 correct\_prediction\_count for full length =9081  
277 Prediction Accuracy for full length =45.405000  
278  
279 Run 6:  
280 Learning rate = 0.00030 ,epochs = 20  
281 Training loss for Epoch:16 is 301  
282 Training loss for Epoch:17 is 434  
283 Training stopped  
284  
285 Run 7:  
286 Learning rate = 0.00025 ,epochs = 20  
287 Last epoch loss is 479  
288 min\_batch\_loss :=63.931

289 correct\_prediction\_count using last bit =10828  
290 Prediction Accuracy using last bit =54.140000  
291 correct\_prediction\_count for full length =2140  
292 Prediction Accuracy for full length =10.700000  
293  
294  
295 Run 8:  
296 Learning rate =0.00035 ,epochs=30  
297 Training loss for Epoch:29 is 2119  
298 min\_batch\_loss :=3.319  
299 correct\_prediction\_count using last bit =15151  
300 Prediction Accuracy using last bit =75.755000  
301 correct\_prediction\_count for full length =11532  
302 Prediction Accuracy for full length =57.660000  
303  
304 Run 9:  
305 Learning rate = 0.00034 ,epochs = 20  
306 Training loss for Epoch:15 is 2000  
307 Canceled after 16 epochs  
308  
309 Run 10:  
310 Learning rate = 0.000349, epochs =20  
311 Training loss for last epoch is 1138  
312  
313 Run 11:  
314 Learning rate = 0.00035 for first 20 epochs and 0.00025 for next 10 epochs ,epochs =30  
315 Training loss for Epoch:19 is 44  
316 Training loss for Epoch:20 is 170  
317 Last epoch loss is 213  
318 min\_batch\_loss :=3.289  
319 correct\_prediction\_count using last bit =18832  
320 Prediction Accuracy using last bit =94.160000  
321 correct\_prediction\_count for full length =13775  
322 Prediction Accuracy for full length =68.875000  
323  
324 Run 12:  
325 Learning rate = 0.00035 for first 20 epochs and 0.00012 for next 10 epochs , epochs =30  
326 Training loss for Epoch:19 is 44  
327 Training loss for Epoch:20 is 134  
328 last epoch loss is 48  
329 min\_batch\_loss :=3.199  
330 correct\_prediction\_count using last bit =19350  
331 Prediction Accuracy using last bit =96.750000  
332 correct\_prediction\_count for full length =11238  
333 Prediction Accuracy for full length =56.190000  
334  
335 Run 13:  
336 Learning rate = 0.00035 for first 20 epochs and 0.00006 for next 10 epochs ,epochs =30  
337 Training loss for Epoch:19 is 44  
338 Training loss for Epoch:20 is 131  
339 Last epoch loss is 16  
340 min\_batch\_loss :=3.221  
341 correct\_prediction\_count using last bit =19671  
342 Prediction Accuracy using last bit =98.355000  
343 correct\_prediction\_count for full length =13565  
344 Prediction Accuracy for full length =67.825000

345  
346  
347 Run 14:  
348 Learning rate = 0.00035 for first 20 epochs and 0.00010 for next 10 epochs , epochs =30  
349 Training loss for Epoch:19 is 44  
350 Training loss for Epoch:20 is 123  
351 Last epoch loss is 23  
352 min\_batch\_loss :=3.219  
353 correct\_prediction\_count using last bit =19597  
354 Prediction Accuracy using last bit =97.985000  
355 correct\_prediction\_count for full length =12789  
356 Prediction Accuracy for full length =63.945000  
357  
358 Run 15:  
359 Learning rate = 0.00035 for first 20 epochs and 0.00009 for next 10 epochs , epochs =30  
360 Training loss for Epoch:19 is 44  
361 Training loss for Epoch:20 is 134  
362 Last epoch loss is 22  
363 min\_batch\_loss :=3.209  
364 correct\_prediction\_count using last bit =19614  
365 Prediction Accuracy using last bit =98.070000  
366 correct\_prediction\_count for full length =12828  
367 Prediction Accuracy for full length =64.140000  
368  
369 Run 16:  
370 Learning rate = 0.00035 for first 20 epochs and 0.00001 for next 10 epochs , epochs =30  
371 Training loss for Epoch:19 is 44  
372 Training loss for Epoch:20 is 129  
373 Training loss for Epoch:29 is 14  
374 min\_batch\_loss :=3.182  
375 correct\_prediction\_count using last bit =19702  
376 Prediction Accuracy using last bit =98.510000  
377 correct\_prediction\_count for full length =13592  
378 Prediction Accuracy for full length =67.960000  
379  
380 Run 17:  
381 Learning rate = 0.00035 for first 20 epochs and 0.000001 for next 10 epochs ,epochs =30  
382 Training loss for Epoch:19 is 44  
383 Training loss for Epoch:20 is 127  
384 Training loss for Epoch:29 is 76  
385 min\_batch\_loss :=3.930  
386 correct\_prediction\_count using last bit =19238  
387 Prediction Accuracy using last bit =96.190000  
388 correct\_prediction\_count for full length =11661  
389 Prediction Accuracy for full length =58.305000  
390  
391 Run 18:  
392 Learning rate = 0.00035 for first 20 epochs and 0.00004 for next 10 epochs , epochs=30  
393 Training loss for Epoch:19 is 44  
394 Training loss for Epoch:20 is 132  
395 Training loss for Epoch:29 is 11  
396 min\_batch\_loss :=2.628  
397 correct\_prediction\_count using last bit =19760  
398 Prediction Accuracy using last bit =98.800000  
399 correct\_prediction\_count for full length =14664  
400 Prediction Accuracy for full length =73.320000



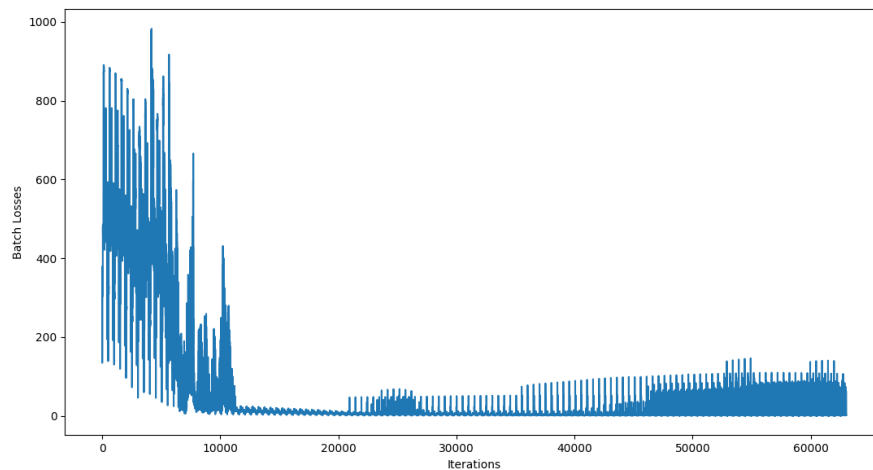
401  
402  
403 Run 19:  
404 Learning rate = 0.00035 for first 20 epochs and 0.00004 for next 31 epochs ,epochs =51  
405 Training loss for Epoch:19 is 44  
406 Training loss for Epoch:20 is 132  
407 Last epoch loss is 4  
408 min\_batch\_loss :=0.836  
409 correct\_prediction\_count using last bit =19894  
410 Prediction Accuracy using last bit =99.470000  
411 correct\_prediction\_count for full length =17535  
412 Prediction Accuracy for full length =87.675000  
413  
414  
415 Run 20:  
416 Learning rate = 0.00035 for first 20 epochs and 0.00004 for next 45 epochs , epochs =65  
417 Training loss for Epoch:19 is 44  
418 Training loss for Epoch:20 is 132  
419 Last epoch loss is 50  
420 min\_batch\_loss :=0.836  
421 correct\_prediction\_count using last bit =19679  
422 Prediction Accuracy using last bit =98.395000  
423 correct\_prediction\_count for full length =16167  
424 Prediction Accuracy for full length =80.835000  
425  
426  
427 Run 21:  
428 Learning rate = 0.00035 for first 20 epochs , 0.00004 from epoch no:20 to epoch no:45 and  
429 0.000001 from epoch no:46 to epoch no: 54 , epochs = 55  
430 Training loss for Epoch:19 is 44  
431 Training loss for Epoch:20 is 132  
432 Training loss for Epoch:45 is 4  
433 Training loss for Epoch:46 is 4  
434 Last epoch loss is 4  
435 min\_batch\_loss :=0.907  
436 correct\_prediction\_count using last bit =19895  
437 Prediction Accuracy using last bit =99.475000  
438 correct\_prediction\_count for full length =17481  
439 Prediction Accuracy for full length =87.405000  
440  
441 Run 22:  
442 Learning rate = 0.00035 for first 20 epochs ,0.00004 for epoch no:21 to epoch no:45 epochs  
443 0.00001 from epoch no:46 to epoch no: 59 ,epochs =60  
444 Training loss for Epoch:19 is 44  
445 Training loss for Epoch:20 is 132  
446 Training loss for Epoch:45 is 4  
447 Training loss for Epoch:46 is 5  
448 Last epoch loss is 4  
449 min\_batch\_loss :=0.808  
450 correct\_prediction\_count using last bit =19907  
451 Prediction Accuracy using last bit =99.535000  
452 correct\_prediction\_count for full length =17667  
453 Prediction Accuracy for full length =88.335000  
454  
455 Run 23:  
456 Learning rate = 0.00035 for first 20 epochs, 0.00004 from epoch no:20 to epoch no:45

```

457 and 0.000004 from epoch no:46 to epoch no: 64 ,epochs = 65
458 Training loss for Epoch:19 is 44
459 Training loss for Epoch:20 is 132
460 Training loss for Epoch:45 is 4
461 Training loss for Epoch:46 is 4
462 Last epoch loss is 4
463 min_batch_loss :=0.826
464 correct_prediction_count using last bit =19900
465 Prediction Accuracy using last bit =99.500000
466 correct_prediction_count for full length =17643
467 Prediction Accuracy for full length =88.215000
468
469 Run 24:
470 Learning rate = 0.00035 for first 20 epochs , 0.00004 from epoch no:20 to epoch no:45,
471 0.00002 from epoch no:46 to epoch no: 64 ,epochs = 65
472 Training loss for Epoch:19 is 44
473 Training loss for Epoch:20 is 132
474 Training loss for Epoch:45 is 4
475 Training loss for Epoch:46 is 5
476 Last epoch loss is 6
477 min_batch_loss :=0.741
478 correct_prediction_count using last bit =19890
479 Prediction Accuracy using last bit =99.450000
480 correct_prediction_count for full length =17863
481 Prediction Accuracy for full length =89.31500
482
483 Run 25:
484 Learning rate = 0.00035 for first 20 epochs , 0.00004 from epoch no:21 to epoch no:45 and
485 0.000014 from epoch no:46 to epoch no: 64 ,epochs = 65
486 Training loss for Epoch:19 is 44
487 Training loss for Epoch:20 is 132
488 Training loss for Epoch:45 is 4
489 Training loss for Epoch:46 is 5
490 Last epoch loss is 4
491 min_batch_loss :=0.795
492 correct_prediction_count using last bit =19903
493 Prediction Accuracy using last bit =99.515000
494 correct_prediction_count for full length =17753
495 Prediction Accuracy for full length =88.765000
496
497 Run 26:
498 Learning rate = 0.00035 for first 20 epochs , 0.00004 for epoch no:20 to epoch no: 45 and
499 0.00001 from epoch no:46 to epoch no: 79 ,epochs =80
500 Training loss for Epoch:19 is 44
501 Training loss for Epoch:20 is 132
502 Training loss for Epoch:45 is 4
503 Training loss for Epoch:46 is 5
504 Training loss for last epoch is 3
505 min_batch_loss :=0.596
506 correct_prediction_count using last bit =19931
507 Prediction Accuracy using last bit =99.655000
508 correct_prediction_count for full length =18317
509 Prediction Accuracy for full length =91.585000
510
511
512

```

513 Run 27 and Final epoch for variable length data:  
 514 Learning rate = 0.00035 for first 20 epochs ,0.00004 for epoch no:20 to epoch no: 45 and  
 515 0.00001 from epoch no:46 to epoch no: 125 ,epochs = 125  
 516 Training loss for Epoch:19 is 44  
 517 Training loss for Epoch:20 is 132  
 518 Training loss for Epoch:45 is 4  
 519 Training loss for Epoch:46 is 5  
 520 Training loss for last epoch is 5  
 521 min\_batch\_loss :=0.405  
 522 correct\_prediction\_count using last bit =19938  
 523 Prediction Accuracy using last bit =99.690000  
 524 correct\_prediction\_count for full length =18838  
 525 Prediction Accuracy for full length =94.190000  
 526  
 527 Below graph is loss plot for Run 27 using number of iterations on x-axis  
 528 Number of iterations = 63000 (504 values of averaged batch losses in 125 epochs)  
 529



530 The point to note is that the epochs are increased from 80 to 125 and other hyper parameters  
 531 are same when Run 26 and Run27 are compared . Final epoch loss is more in Run 27 by 1  
 532 but the accuracy is high by 2.6 %  
 533  
 534

## 535 **5 Improvements and Conclusion**

536 Following improvements can be further made which are left for future work

- 537 a. Tuning learning rate using learning rate scheduler perhaps can reduce the training time to
- 538 get good results
- 539 b. MPI parallel programming with 5 processes reduced the training time of a single epoch to
- 540 14 min compared to 20 min training time of one epoch when MPI was not used.
- 541 Using GPU based cloud infrastructure such as GCP or AWS can further reduce the
- 542 training time to a large extent . This can help to train using efficient hyper parameter
- 543 tuning methods
- 544 c. All experiments are run with seed 0 . Few other random seeds can be used in different
- 545 trials.
- 546 d. Batch size 32 was used and other batch sizes such as 64,128 ,256 could be tried.

547  
 548 This project implementation helped to understand LSTM model and application of parallel  
 549 programming using MPI

## 550 References

- 551 [1]Orthogonal initialization for hidden to hidden , Xavier uniform initialization for rest of the weights  
552 <https://www.kaggle.com/code/junkoda/pytorch-lstm-with-tensorflow-like-initialization/notebook>
- 553 [2]Weight initialization  
554 [https://www.deeplearningwizard.com/deep\\_learning/boosting\\_models\\_pytorch/weight\\_initialization\\_activation\\_functions/](https://www.deeplearningwizard.com/deep_learning/boosting_models_pytorch/weight_initialization_activation_functions/)  
555 <https://wandb.ai/sauravmaheshkar/initialization/reports/A-Gentle-Introduction-To-Weight-Initialization-for-Neural-Networks--Vmlldzo2ODExMTg>  
556 [https://www.reddit.com/r/deeplearning/comments/9vjckm/what\\_is\\_the\\_best\\_way\\_to\\_initialize\\_lstm\\_weights/](https://www.reddit.com/r/deeplearning/comments/9vjckm/what_is_the_best_way_to_initialize_lstm_weights/)  
557 [https://www.reddit.com/r/deeplearning/comments/9vjckm/what\\_is\\_the\\_best\\_way\\_to\\_initialize\\_lstm\\_weights/](https://www.reddit.com/r/deeplearning/comments/9vjckm/what_is_the_best_way_to_initialize_lstm_weights/)
- 558 [3] Orthogonal weight initialization  
559 <https://github.com/kastnerkyle/net/blob/d66d533b1ebcc25dd442bdb41431d334b617e80e/net.py#L168>  
560 <https://github.com/nyu-dl/dl4mt-tutorial/blob/master/session3/lm.py>
- 561 [4] Fan in Fan out for xavier uniform  
562 <https://stackoverflow.com/questions/42670274/how-to-calculate-fan-in-and-fan-out-in-xavier-initialization-for-neural-networks>  
563 <https://ml-cheatsheet.readthedocs.io/en/latest/forwardpropagation.html>  
564 <https://visualstudiomagazine.com/articles/2019/09/05/neural-network-glorot.aspx>  
565 <https://visualstudiomagazine.com/articles/2019/09/05/neural-network-glorot.aspx>
- 566 [5] One hot encoding :  
567 <https://datascience.stackexchange.com/questions/45803/should-we-use-only-one-hot-vector-for-lstm-input-outputs>  
568 <https://datascience.stackexchange.com/questions/45803/should-we-use-only-one-hot-vector-for-lstm-input-outputs>
- 569 [6]Batch Loss  
570 <https://stackoverflow.com/questions/54053868/how-do-i-get-a-loss-per-epoch-and-not-per-batch>  
571 <https://stackoverflow.com/questions/54053868/how-do-i-get-a-loss-per-epoch-and-not-per-batch>
- 572 [7] Distributed training and optimization  
573 Diagram : [https://github.com/vlimant/mpi\\_learn/blob/master/docs/downpour.png](https://github.com/vlimant/mpi_learn/blob/master/docs/downpour.png)  
574 [https://github.com/vlimant/mpi\\_learn/blob/master/docs/downpour.png](https://github.com/vlimant/mpi_learn/blob/master/docs/downpour.png)
- 575 MPI.AllReduce:  
576 [https://github.com/openai/baselines/blob/master/baselines/common/mpi\\_adam\\_optimizer.py](https://github.com/openai/baselines/blob/master/baselines/common/mpi_adam_optimizer.py)  
577 [https://github.com/openai/baselines/blob/master/baselines/common/mpi\\_adam\\_optimizer.py](https://github.com/openai/baselines/blob/master/baselines/common/mpi_adam_optimizer.py)
- 578 [8]Back propagation flow diagrams  
579 <https://towardsdatascience.com/backpropagation-in-rnn-explained-bdf853b4e1c2>  
580 [https://github.com/christinakouridi/scratchML/blob/master/LSTM/LSTMForward\\_diagram.png](https://github.com/christinakouridi/scratchML/blob/master/LSTM/LSTMForward_diagram.png)  
581 [https://github.com/christinakouridi/scratchML/blob/master/LSTM/LSTMForward\\_diagram.png](https://github.com/christinakouridi/scratchML/blob/master/LSTM/LSTMForward_diagram.png)
- 582 [9] Back propagation for batch  
583 <https://ai.stackexchange.com/questions/11667/is-back-propagation-applied-for-each-data-point-or-for-a-batch-of-data-points>  
584 <https://ai.stackexchange.com/questions/11667/is-back-propagation-applied-for-each-data-point-or-for-a-batch-of-data-points>
- 585 [10] Modularizing code and RNN Implementation :  
586 <https://github.com/JY-Yoon/RNN-Implementation-using-NumPy/blob/master/RNN%20Implementation%20using%20NumPy.ipynb>  
587 <https://github.com/JY-Yoon/RNN-Implementation-using-NumPy/blob/master/RNN%20Implementation%20using%20NumPy.ipynb>
- 588 [11] Back prop chain rule, matrix multiply dimension matching and Matrix -Matrix Multiply gradient  
589 : <https://cs231n.github.io/optimization-2/>  
590 <https://cs231n.github.io/optimization-2/>
- 591 [12] Adam Optimizer  
592 <https://towardsdatascience.com/how-to-implement-an-adam-optimizer-from-scratch-76e7b217f1cc>  
593 <https://towardsdatascience.com/how-to-implement-an-adam-optimizer-from-scratch-76e7b217f1cc>
- 594 <https://stackoverflow.com/questions/70225531/time-step-in-adam-optimizer>  
595 <https://arxiv.org/pdf/1412.6980.pdf>  
596 <https://arxiv.org/pdf/1412.6980.pdf>
- 597 [13]Vanishing and exploding gradients in LSTM ( This problem is not entirely solved in LSTM ,rather it is mitigated and delayed)  
598 <https://www.quora.com/How-does-LSTM-help-prevent-the-vanishing-and-exploding-gradient-problem-in-a-recurrent-neural-network>  
599 <https://www.quora.com/How-does-LSTM-help-prevent-the-vanishing-and-exploding-gradient-problem-in-a-recurrent-neural-network>
- 600 [14] Cross Entropy  
601 <https://www.quora.com/Whats-an-intuitive-way-to-think-of-cross-entropy>  
602 <https://www.quora.com/Whats-an-intuitive-way-to-think-of-cross-entropy>

598 [15] Matrix and batch dimensions:  
599 <https://www.quora.com/In-LSTM-how-do-you-figure-out-what-size-the-weights-are-supposed-to-be>  
600 [16] Bias gradients sum along axis =0 :  
601 [https://github.com/Erleamar/cs231n\\_self/blob/master/assignment2/cs231n/layers.py#L116](https://github.com/Erleamar/cs231n_self/blob/master/assignment2/cs231n/layers.py#L116)  
602 [17] RNN Theory : <http://karpathy.github.io/2015/05/21/rnn-effectiveness/>  
603 [18] LSTM Theory : <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>  
604 [19] LSTM theory explanation : [https://towardsdatascience.com/lstm-networks-a-detailed-](https://towardsdatascience.com/lstm-networks-a-detailed-explanation-8fae6aefc7f9)  
605 [explanation-8fae6aefc7f9](https://towardsdatascience.com/lstm-networks-a-detailed-explanation-8fae6aefc7f9)  
606 [20] Batch size : <https://ai.stackexchange.com/questions/8560/how-do-i-choose-the-optimal-batch-size>  
607 [21] Bucketing into batches of data for irregular length sequences: [https://rashmi-](https://rashmi-margani.medium.com/how-to-speed-up-the-training-of-the-sequence-model-using-bucketing-techniques-9e302b0fd976)  
608 [margani.medium.com/how-to-speed-up-the-training-of-the-sequence-model-using-bucketing-](https://rashmi-margani.medium.com/how-to-speed-up-the-training-of-the-sequence-model-using-bucketing-techniques-9e302b0fd976)  
609 [techniques-9e302b0fd976](https://rashmi-margani.medium.com/how-to-speed-up-the-training-of-the-sequence-model-using-bucketing-techniques-9e302b0fd976)  
610 [22] Python binary string generation  
611 <https://www.geeksforgeeks.org/python-program-to-generate-random-binary-string/>  
612 [23] Gradient clipping  
613 [https://stackoverflow.com/questions/44796793/difference-between-tf-clip-by-value-and-tf-clip-by-](https://stackoverflow.com/questions/44796793/difference-between-tf-clip-by-value-and-tf-clip-by-global-norm-for-rnns-and-how)  
614 [global-norm-for-rnns-and-how](https://stackoverflow.com/questions/44796793/difference-between-tf-clip-by-value-and-tf-clip-by-global-norm-for-rnns-and-how)