

---

# 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 approximately 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.

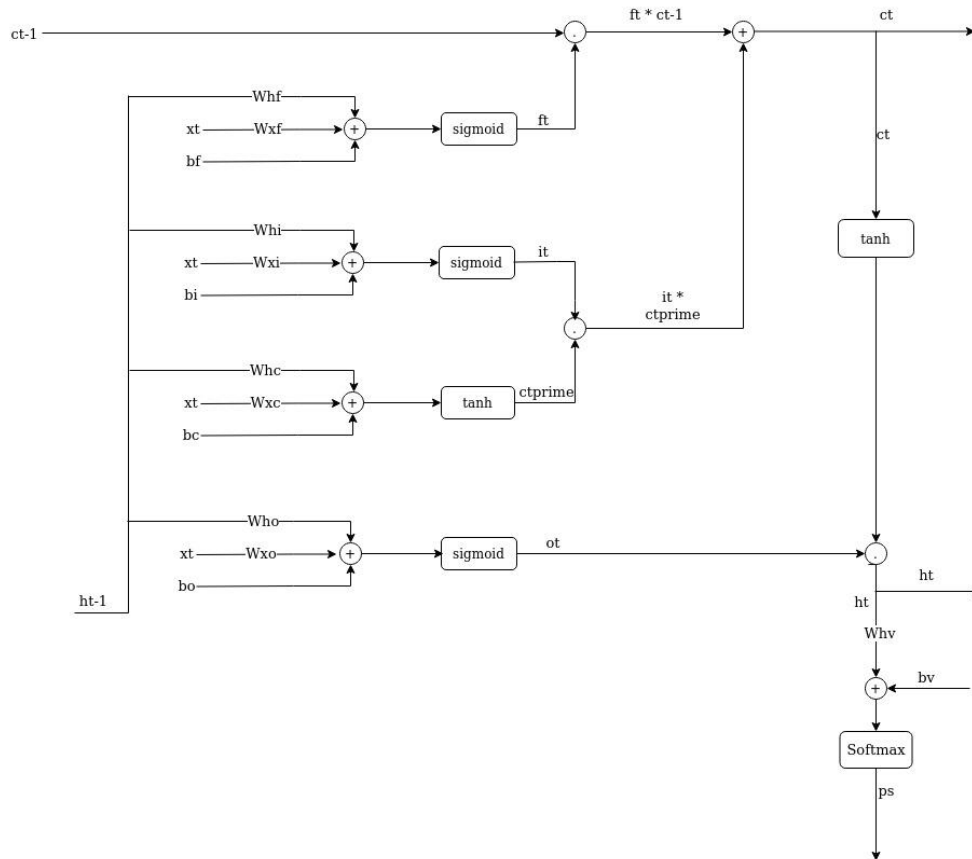
Train dataset split into batches : The 80000 data samples in type1 are divided into 2500 batches each having size of 32. In type2 , 80000 samples are divided into 2521 batches each having variable size that is  $\leq 32$ . The reason for this variability is that the data in type 2 has variable lengths 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 data

48 equally or distribute the left over data equally among the fixed number of processes. This is  
 49 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 state  
 54 'ct-1' (cell state is current long-term memory of the network) are fed as input to the current  
 55 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 weights  
 59 are in dot product with the hidden state inputs and the input data for that particular time step  
 60 . '+' symbol inside a circle is an addition operation while '.' symbol inside a circle is the  
 61 pairwise multiplication operation. Outputs from these 4 units are passed through non linear  
 62 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 Wfh,Wfx and  
 67 bias bf 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 Whi,Wxi and bias bi, 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) value  
 73 and previous hidden state value(ht-1) using weights Whc,Wxc and bias bc .The use of ctprime

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

77 The operation  $i * ct_{prime}$  nullifies the effect of few values in current cell state ( $ct_{prime}$ ) 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  $ft * ct - 1 + i * ct_{prime}$  operation produces the next cell state value 'ct'. ct is group of values which  
81 are obtained by retaining and forgetting certain values from current cell state and previous cell  
82 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 state(ct)  
86 constrains the values to the range  $(-1,1)$  i.e. reduces the impact of few values and while other  
87 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 output  
91 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 bias initialization: Based on references[1][2][3] and[4], the weights which are in  
95 dot product with hidden state values from previous cell, in every gate and current cell state  
96 are obtained by 'orthogonal initialization'. The number of units in hidden layer is 100 and  
97 dimension of these matrices are (100,100). The remaining eight are consequence of 'xavier  
98 uniform' initialization. Weights which multiply with input data have the dimension (batch size  
99 ,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 of  
101 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 descent  
108 is used for optimization and training is done in parallel for each batch of data in the groups and  
109 the process can be split into 3 main steps a. Forward propagation b. Back propagation c.  
110 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 values  
115 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 classes  
117, softmax classification can be used which necessitates one hot encoding of the input data.  
118 Input character '0' at a time step is encoded as [1,0] and character '1' is encoded as [0,1]. The  
119 input and output data at every time step in the whole batch is encoded in this manner.

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

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

124 Cross entropy loss function is used to calculate the loss and the loss values of all the samples

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

## 149 **4 Results**

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

154 Training and prediction for Fix length batch data :

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

164 Thereafter , it took 6 runs to achieve the convergence. In all these runs , the learning rate  
165 parameter and epochs are changed , other hyper parameters such as no. of units in hidden layer  
166 =100 , batch size =32 and random seed = 0 remain the same.  
167

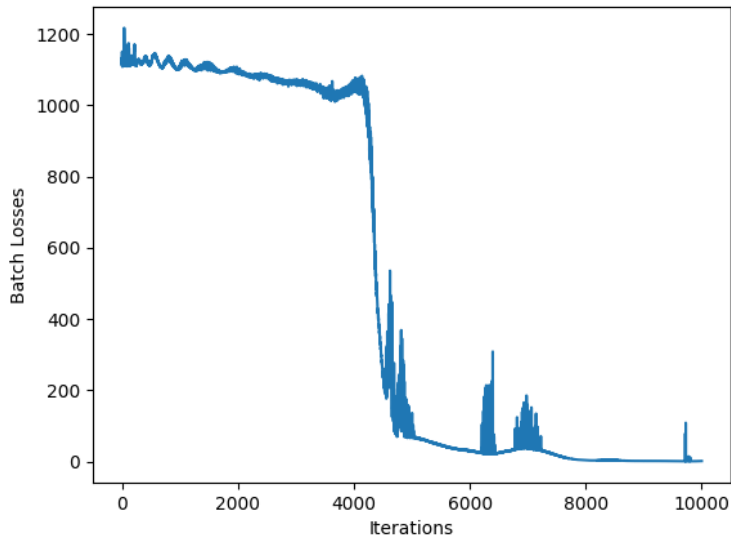
168 Run1:

169 Learning rate = 0.0005 , epochs =12  
170 Loss does not decrease significantly even after 5 epochs  
171 Training stopped  
172

173 Run2:

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

```
178 Run3:
179 Learning rate = 0.00075 ,epochs =12
180 Last epoch loss : 36
181 min batch Loss : 29.415
182 correct_prediction_count using last bit =19616
183 Prediction Accuracy using last bit =98.080000
184 correct_prediction_count for full length =7910
185 Prediction Accuracy for full length =39.550000
186
187
188 Run4:
189 Learning rate = 0.00080 ,epochs =12
190 Last epoch loss = 1780
191 min_batch_loss =1087.689
192 correct_prediction_count using last bit =9922
193 Prediction Accuracy using last bit =49.610000
194 correct_prediction_count for full length =0
195 Prediction Accuracy for full length =0.000000
196
197
198 Run5:
199 Learning rate = 0.00076 , epochs =12
200 Last epoch loss = 1960
201 min_batch_loss =110.529
202 correct_prediction_count using last bit =10195
203 Prediction Accuracy using last bit =50.975000
204 correct_prediction_count for full length =0
205 Prediction Accuracy for full length =0.000000
206
207 Run6 and Final result for fixed length data:
208 Learning rate = 0.00075 , epochs =20
209 Last epoch loss = 1
210 min_batch_loss = 1.075
211 correct_prediction_count using last bit = 19978
212 Prediction Accuracy using last bit = 99.890000
213 correct_prediction_count for full length = 18932
214 Prediction Accuracy for full length = 94.660000
215
216
217 Below graph is loss plot for Run 6 using number of iterations on x-axis
218 Number of iterations =10000 (500 values of averaged batch losses in 20 epochs)
219
```



240

241

242 Training and prediction for Variable length batch data :

243 Varying learning rates for the ranges of epochs instead of single learning rate helped in  
 244 convergence . Loss curve and decrement seem to be complicated when compared to those  
 245 of the fixed length data .

246 The number of epochs required are more in number compared to what is used for the fixed  
 247 length data.

248 It took 27 runs to figure out the correct learning rates and no. of epochs to get prediction  
 249 accuracy above 92% for variable length data.

250 Batch size =32 ,random seed =0 was used in all the runs below

251 Except for variation in handling data set, number of epochs and learning rates , no other  
 252 change is made in the model for training variable length data set.

253

254 Run1:

255 Learning rate = 0.00075 , epochs =20

256 Training loss for Epoch:11 is 1909

257 Training stopped after 12 epochs due to high fluctuations in loss

258

259 Run2:

260 Learning rate = 0.00050 ,epochs =20

261 Training loss for Epoch:8 is 2359

262 Training stopped after 9 epochs due to high fluctuations in loss

263

264 Run 3:

265 Learning rate = 0.00040 , epochs =20

266 Training loss for Epoch:11 is 599

267 Training loss for Epoch:12 is 2280

268 Training stopped

269

270 Run 4:

271 Learning rate = 0.00035 , epochs = 20

272 Last epoch loss is 44

273 min\_batch\_loss :=4.391

274 correct\_prediction\_count using last bit =18878

275 Prediction Accuracy using last bit =94.390000

```
276 correct_prediction_count for full length =11162
277 Prediction Accuracy for full length =55.810000
278
279 Run 5:
280 Learning rate = 0.00035 ,epochs = 28
281 Last epoch loss is 464
282 min_batch_loss :=3.319
283 correct_prediction_count using last bit =16397
284 Prediction Accuracy using last bit =81.985000
285 correct_prediction_count for full length =9081
286 Prediction Accuracy for full length =45.405000
287
288 Run 6:
289 Learning rate = 0.00030 ,epochs = 20
290 Training loss for Epoch:16 is 301
291 Training loss for Epoch:17 is 434
292 Training stopped
293
294 Run 7:
295 Learning rate = 0.00025 ,epochs = 20
296 Last epoch loss is 479
297 min_batch_loss :=63.931
298 correct_prediction_count using last bit =10828
299 Prediction Accuracy using last bit =54.140000
300 correct_prediction_count for full length =2140
301 Prediction Accuracy for full length =10.700000
302
303
304 Run 8:
305 Learning rate =0.00035 ,epochs=30
306 Training loss for Epoch:29 is 2119
307 min_batch_loss :=3.319
308 correct_prediction_count using last bit =15151
309 Prediction Accuracy using last bit =75.755000
310 correct_prediction_count for full length =11532
311 Prediction Accuracy for full length =57.660000
312
313 Run 9:
314 Learning rate = 0.00034 ,epochs = 20
315 Training loss for Epoch:15 is 2000
316 Canceled after 16 epochs
317
318 Run 10:
319 Learning rate = 0.000349, epochs =20
320 Training loss for last epoch is 1138
321
322 Run 11:
323 Learning rate = 0.00035 for first 20 epochs and 0.00025 for next 10 epochs ,epochs =30
324 Training loss for Epoch:19 is 44
325 Training loss for Epoch:20 is 170
326 Last epoch loss is 213
327 min_batch_loss :=3.289
328 correct_prediction_count using last bit =18832
329 Prediction Accuracy using last bit =94.160000
330 correct_prediction_count for full length =13775
331 Prediction Accuracy for full length =68.875000
```

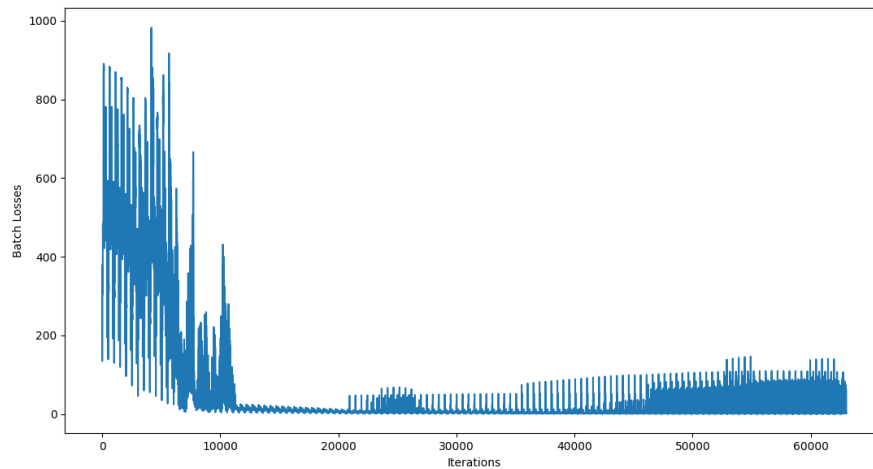
```
332
333 Run 12:
334 Learning rate = 0.00035 for first 20 epochs and 0.00012 for next 10 epochs , epochs =30
335 Training loss for Epoch:19 is 44
336 Training loss for Epoch:20 is 134
337 last epoch loss is 48
338 min_batch_loss :=3.199
339 correct_prediction_count using last bit =19350
340 Prediction Accuracy using last bit =96.750000
341 correct_prediction_count for full length =11238
342 Prediction Accuracy for full length =56.190000
343
344 Run 13:
345 Learning rate = 0.00035 for first 20 epochs and 0.00006 for next 10 epochs ,epochs =30
346 Training loss for Epoch:19 is 44
347 Training loss for Epoch:20 is 131
348 Last epoch loss is 16
349 min_batch_loss :=3.221
350 correct_prediction_count using last bit =19671
351 Prediction Accuracy using last bit =98.355000
352 correct_prediction_count for full length =13565
353 Prediction Accuracy for full length =67.825000
354
355
356 Run 14:
357 Learning rate = 0.00035 for first 20 epochs and 0.00010 for next 10 epochs , epochs =30
358 Training loss for Epoch:19 is 44
359 Training loss for Epoch:20 is 123
360 Last epoch loss is 23
361 min_batch_loss :=3.219
362 correct_prediction_count using last bit =19597
363 Prediction Accuracy using last bit =97.985000
364 correct_prediction_count for full length =12789
365 Prediction Accuracy for full length =63.945000
366
367 Run 15:
368 Learning rate = 0.00035 for first 20 epochs and 0.00009 for next 10 epochs , epochs =30
369 Training loss for Epoch:19 is 44
370 Training loss for Epoch:20 is 134
371 Last epoch loss is 22
372 min_batch_loss :=3.209
373 correct_prediction_count using last bit =19614
374 Prediction Accuracy using last bit =98.070000
375 correct_prediction_count for full length =12828
376 Prediction Accuracy for full length =64.140000
377
378 Run 16:
379 Learning rate = 0.00035 for first 20 epochs and 0.00001 for next 10 epochs , epochs =30
380 Training loss for Epoch:19 is 44
381 Training loss for Epoch:20 is 129
382 Training loss for Epoch:29 is 14
383 min_batch_loss :=3.182
384 correct_prediction_count using last bit =19702
385 Prediction Accuracy using last bit =98.510000
386 correct_prediction_count for full length =13592
387 Prediction Accuracy for full length =67.960000
```



```
388
389 Run 17:
390 Learning rate = 0.00035 for first 20 epochs and 0.000001 for next 10 epochs ,epochs =30
391 Training loss for Epoch:19 is 44
392 Training loss for Epoch:20 is 127
393 Training loss for Epoch:29 is 76
394 min_batch_loss :=3.930
395 correct_prediction_count using last bit =19238
396 Prediction Accuracy using last bit =96.190000
397 correct_prediction_count for full length =11661
398 Prediction Accuracy for full length =58.305000
399
400 Run 18:
401 Learning rate = 0.00035 for first 20 epochs and 0.00004 for next 10 epochs , epochs=30
402 Training loss for Epoch:19 is 44
403 Training loss for Epoch:20 is 132
404 Training loss for Epoch:29 is 11
405 min_batch_loss :=2.628
406 correct_prediction_count using last bit =19760
407 Prediction Accuracy using last bit =98.800000
408 correct_prediction_count for full length =14664
409 Prediction Accuracy for full length =73.320000
410
411
412 Run 19:
413 Learning rate = 0.00035 for first 20 epochs and 0.00004 for next 31 epochs ,epochs =51
414 Training loss for Epoch:19 is 44
415 Training loss for Epoch:20 is 132
416 Last epoch loss is 4
417 min_batch_loss :=0.836
418 correct_prediction_count using last bit =19894
419 Prediction Accuracy using last bit =99.470000
420 correct_prediction_count for full length =17535
421 Prediction Accuracy for full length =87.675000
422
423
424 Run 20:
425 Learning rate = 0.00035 for first 20 epochs and 0.00004 for next 45 epochs , epochs =65
426 Training loss for Epoch:19 is 44
427 Training loss for Epoch:20 is 132
428 Last epoch loss is 50
429 min_batch_loss :=0.836
430 correct_prediction_count using last bit =19679
431 Prediction Accuracy using last bit =98.395000
432 correct_prediction_count for full length =16167
433 Prediction Accuracy for full length =80.835000
434
435
436 Run 21:
437 Learning rate = 0.00035 for first 20 epochs , 0.00004 from epoch no:20 to epoch no:45 and
438 0.000001 from epoch no:46 to epoch no: 54 , epochs = 55
439 Training loss for Epoch:19 is 44
440 Training loss for Epoch:20 is 132
441 Training loss for Epoch:45 is 4
442 Training loss for Epoch:46 is 4
443 Last epoch loss is 4
```

444 min\_batch\_loss :=0.907  
445 correct\_prediction\_count using last bit =19895  
446 Prediction Accuracy using last bit =99.475000  
447 correct\_prediction\_count for full length =17481  
448 Prediction Accuracy for full length =87.405000  
449  
450 Run 22:  
451 Learning rate = 0.00035 for first 20 epochs ,0.00004 for epoch no:21 to epoch no:45 epochs  
452 0.00001 from epoch no:46 to epoch no: 59 ,epochs =60  
453 Training loss for Epoch:19 is 44  
454 Training loss for Epoch:20 is 132  
455 Training loss for Epoch:45 is 4  
456 Training loss for Epoch:46 is 5  
457 Last epoch loss is 4  
458 min\_batch\_loss :=0.808  
459 correct\_prediction\_count using last bit =19907  
460 Prediction Accuracy using last bit =99.535000  
461 correct\_prediction\_count for full length =17667  
462 Prediction Accuracy for full length =88.335000  
463  
464 Run 23:  
465 Learning rate = 0.00035 for first 20 epochs, 0.00004 from epoch no:20 to epoch no:45  
466 and 0.000004 from epoch no:46 to epoch no: 64 ,epochs = 65  
467 Training loss for Epoch:19 is 44  
468 Training loss for Epoch:20 is 132  
469 Training loss for Epoch:45 is 4  
470 Training loss for Epoch:46 is 4  
471 Last epoch loss is 4  
472 min\_batch\_loss :=0.826  
473 correct\_prediction\_count using last bit =19900  
474 Prediction Accuracy using last bit =99.500000  
475 correct\_prediction\_count for full length =17643  
476 Prediction Accuracy for full length =88.215000  
477  
478 Run 24:  
479 Learning rate = 0.00035 for first 20 epochs , 0.00004 from epoch no:20 to epoch no:45,  
480 0.00002 from epoch no:46 to epoch no: 64 ,epochs = 65  
481 Training loss for Epoch:19 is 44  
482 Training loss for Epoch:20 is 132  
483 Training loss for Epoch:45 is 4  
484 Training loss for Epoch:46 is 5  
485 Last epoch loss is 6  
486 min\_batch\_loss :=0.741  
487 correct\_prediction\_count using last bit =19890  
488 Prediction Accuracy using last bit =99.450000  
489 correct\_prediction\_count for full length =17863  
490 Prediction Accuracy for full length =89.31500  
491  
492 Run 25:  
493 Learning rate = 0.00035 for first 20 epochs , 0.00004 from epoch no:21 to epoch no:45 and  
494 0.000014 from epoch no:46 to epoch no: 64 ,epochs = 65  
495 Training loss for Epoch:19 is 44  
496 Training loss for Epoch:20 is 132  
497 Training loss for Epoch:45 is 4  
498 Training loss for Epoch:46 is 5  
499 Last epoch loss is 4

500 min\_batch\_loss :=0.795  
501 correct\_prediction\_count using last bit =19903  
502 Prediction Accuracy using last bit =99.515000  
503 correct\_prediction\_count for full length =17753  
504 Prediction Accuracy for full length =88.765000  
505  
506 Run 26:  
507 Learning rate = 0.00035 for first 20 epochs , 0.00004 for epoch no:20 to epoch no: 45 and  
508 0.00001 from epoch no:46 to epoch no: 79 ,epochs =80  
509 Training loss for Epoch:19 is 44  
510 Training loss for Epoch:20 is 132  
511 Training loss for Epoch:45 is 4  
512 Training loss for Epoch:46 is 5  
513 Training loss for last epoch is 3  
514 min\_batch\_loss :=0.596  
515 correct\_prediction\_count using last bit =19931  
516 Prediction Accuracy using last bit =99.655000  
517 correct\_prediction\_count for full length =18317  
518 Prediction Accuracy for full length =91.585000  
519  
520  
521  
522 Run 27 and Final epoch for variable length data:  
523 Learning rate = 0.00035 for first 20 epochs ,0.00004 for epoch no:20 to epoch no: 45 and  
524 0.00001 from epoch no:46 to epoch no: 125 ,epochs = 125  
525 Training loss for Epoch:19 is 44  
526 Training loss for Epoch:20 is 132  
527 Training loss for Epoch:45 is 4  
528 Training loss for Epoch:46 is 5  
529 Training loss for last epoch is 5  
530 min\_batch\_loss :=0.405  
531 correct\_prediction\_count using last bit =19938  
532 Prediction Accuracy using last bit =99.690000  
533 correct\_prediction\_count for full length =18838  
534 Prediction Accuracy for full length =94.190000  
535  
536 Below graph is loss plot for Run 27 using number of iterations on x-axis  
537 Number of iterations = 63000 (504 values of averaged batch losses in 125 epochs)  
538



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

542

543

544

## 5 Improvements and Conclusion

545

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

546

a. Tuning learning rate using learning rate scheduler perhaps can reduce the training time to  
 547 get good results

548

b. MPI parallel programming with 5 processes reduced the training time of a single epoch to  
 549 14 min compared to 20 min training time of one epoch when MPI was not used.

550

Using GPU based cloud infrastructure such as GCP or AWS can further reduce the  
 551 training time to a large extent . This can help to train using efficient hyper parameter  
 552 tuning methods

553

c. All experiments are run with seed 0 . Few other random seeds can be used in different  
 554 trials.

555

d. Batch size 32 was used and other batch sizes such as 64,128 ,256 could be tried.

556

557

This project implementation helped to understand LSTM model and application of parallel  
 558 programming using MPI

559

## References

560

[1]Orthogonal initialization for hidden to hidden , Xavier uniform initialization for rest of the weights

561

<https://www.kaggle.com/code/junkoda/pytorch-lstm-with-tensorflow-like-initialization/notebook>

562

[2]Weight initialization

563

[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/)

564

<https://wandb.ai/sauravmaheshkar/initialization/reports/A-Gentle-Introduction-To-Weight-Initialization-for-Neural-Networks--Vmlldzo2ODExMTg>

565

[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/)

566

[3] Orthogonal weight initialization  
<https://github.com/kastnerkyle/net/blob/d66d533b1ebcc25dd442bdb41431d334b617e80e/net.py#L168>

567

<https://github.com/nyu-dl/dl4mt-tutorial/blob/master/session3/lm.py>

568

[4] Fan in Fan out for xavier uniform

569

<https://stackoverflow.com/questions/42670274/how-to-calculate-fan-in-and-fan-out-in-xavier->

570

<https://stackoverflow.com/questions/42670274/how-to-calculate-fan-in-and-fan-out-in-xavier->

571

<https://stackoverflow.com/questions/42670274/how-to-calculate-fan-in-and-fan-out-in-xavier->

572

<https://stackoverflow.com/questions/42670274/how-to-calculate-fan-in-and-fan-out-in-xavier->

573

<https://stackoverflow.com/questions/42670274/how-to-calculate-fan-in-and-fan-out-in-xavier->

574 [initialization-for-neural-networks](#)

575 <https://ml-cheatsheet.readthedocs.io/en/latest/forwardpropagation.html>

576 <https://visualstudiomagazine.com/articles/2019/09/05/neural-network-glorot.aspx>

577 [5] One hot encoding :

578 [https://datascience.stackexchange.com/questions/45803/should-we-use-only-one-hot-vector-for-lstm-](https://datascience.stackexchange.com/questions/45803/should-we-use-only-one-hot-vector-for-lstm-input-outputs)

579 [input-outputs](https://datascience.stackexchange.com/questions/45803/should-we-use-only-one-hot-vector-for-lstm-input-outputs)

580 [6] Batch Loss

581 <https://stackoverflow.com/questions/54053868/how-do-i-get-a-loss-per-epoch-and-not-per-batch>

582 [7] Distributed training and optimization

583 Diagram : [https://github.com/vlimant/mpi\\_learn/blob/master/docs/downpour.png](https://github.com/vlimant/mpi_learn/blob/master/docs/downpour.png)

584 MPI.AllReduce:

585 [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)

586 [8] Back propagation flow diagrams

587 <https://towardsdatascience.com/backpropagation-in-rnn-explained-bdf853b4e1c2>

588 [https://github.com/christinakouridi/scratchML/blob/master/LSTM/LSTMForward\\_diagram.png](https://github.com/christinakouridi/scratchML/blob/master/LSTM/LSTMForward_diagram.png)

589 [9] Back propagation for batch

590 [https://ai.stackexchange.com/questions/11667/is-back-propagation-applied-for-each-data-point-or-for-](https://ai.stackexchange.com/questions/11667/is-back-propagation-applied-for-each-data-point-or-for-a-batch-of-data-points)

591 [a-batch-of-data-points](https://ai.stackexchange.com/questions/11667/is-back-propagation-applied-for-each-data-point-or-for-a-batch-of-data-points)

592 [10] Modularizing code and RNN Implementation :

593 [https://github.com/JY-Yoon/RNN-Implementation-using-](https://github.com/JY-Yoon/RNN-Implementation-using-NumPy/blob/master/RNN%20Implementation%20using%20NumPy.ipynb)

594 [NumPy/blob/master/RNN%20Implementation%20using%20NumPy.ipynb](https://github.com/JY-Yoon/RNN-Implementation-using-NumPy/blob/master/RNN%20Implementation%20using%20NumPy.ipynb)

595 [11] Back prop chain rule, matrix multiply dimension matching and Matrix -Matrix Multiply gradient

596 : <https://cs231n.github.io/optimization-2/>

597 [12] Adam Optimizer

598 <https://towardsdatascience.com/how-to-implement-an-adam-optimizer-from-scratch-76e7b217f1cc>

599 <https://stackoverflow.com/questions/70225531/time-step-in-adam-optimizer>

600 <https://arxiv.org/pdf/1412.6980.pdf>

601 [13] Vanishing and exploding gradients in LSTM ( This problem is not entirely solved in LSTM ,rather

602 it is mitigated and delayed)

603 [https://www.quora.com/How-does-LSTM-help-prevent-the-vanishing-and-exploding-gradient-](https://www.quora.com/How-does-LSTM-help-prevent-the-vanishing-and-exploding-gradient-problem-in-a-recurrent-neural-network)

604 [problem-in-a-recurrent-neural-network](https://www.quora.com/How-does-LSTM-help-prevent-the-vanishing-and-exploding-gradient-problem-in-a-recurrent-neural-network)

605 [14] Cross Entropy

606 <https://www.quora.com/Whats-an-intuitive-way-to-think-of-cross-entropy>

607 [15] Matrix and batch dimensions:

608 <https://www.quora.com/In-LSTM-how-do-you-figure-out-what-size-the-weights-are-supposed-to-be>

609 [16] Bias gradients sum along axis =0 :

610 [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)

611 [17] RNN Theory : <http://karpathy.github.io/2015/05/21/rnn-effectiveness/>

612 [18] LSTM Theory : <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

613 [19] LSTM theory explanation : [https://towardsdatascience.com/lstm-networks-a-detailed-](https://towardsdatascience.com/lstm-networks-a-detailed-explanation-8fae6aefc7f9)

614 [explanation-8fae6aefc7f9](https://towardsdatascience.com/lstm-networks-a-detailed-explanation-8fae6aefc7f9)

615 [20] Batch size : <https://ai.stackexchange.com/questions/8560/how-do-i-choose-the-optimal-batch-size>

616 [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)

617 [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)

618 [techniques-9e302b0fd976](https://rashmi-margani.medium.com/how-to-speed-up-the-training-of-the-sequence-model-using-bucketing-techniques-9e302b0fd976)

619 [22] Python binary string generation

620 <https://www.geeksforgeeks.org/python-program-to-generate-random-binary-string/>

621 [23] Gradient clipping

622 <https://stackoverflow.com/questions/44796793/difference-between-tf-clip-by-value-and-tf-clip-by->  
623 [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)