
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.

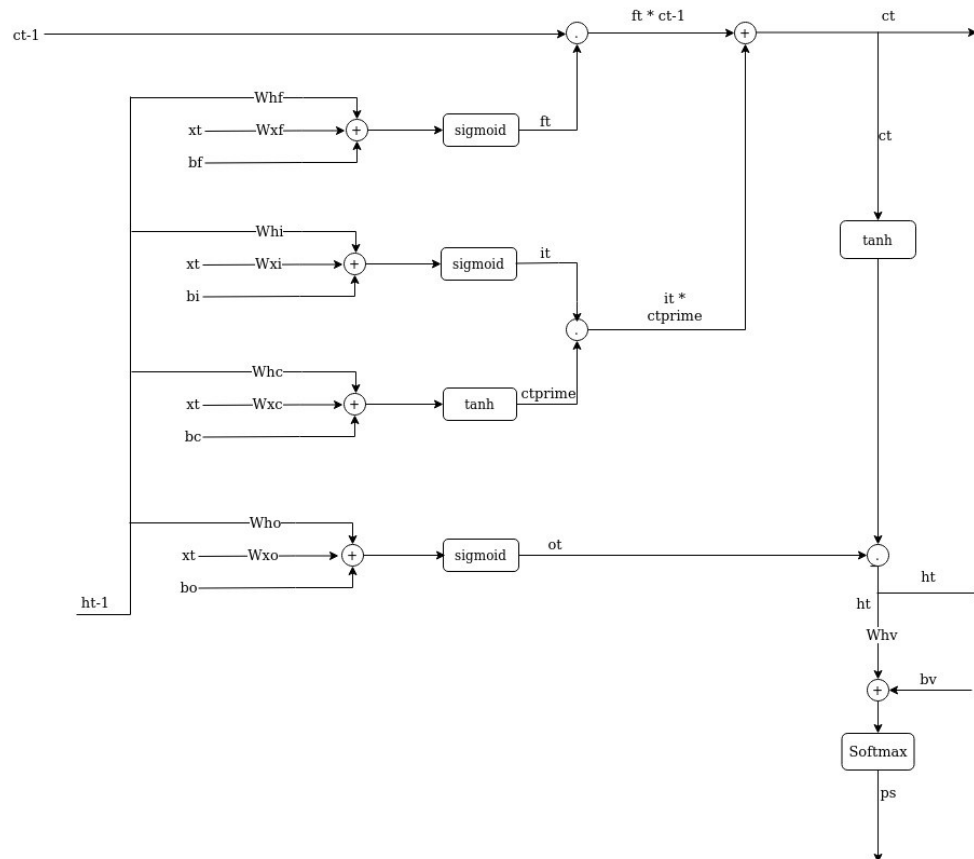
Train dataset split into batches : The 80000 data samples in type1 are divided into 2500 batches each having size of 32. In type2 , 80000 samples are divided into 2521 batches each having variable size that is ≤ 32 . The reason for this variability is that the data in type 2 has variable lengths and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

513 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_{fh} , 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 $ctprime$ is to modify the information in xt given the $ht-1$ values. Tanh range is $(-1,1)$, the

75impact of values close to lower bound is reduced while the impact of values closer to upper
76bound is remained as it is.

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

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

83Output gate: Current time step input (x_t) and previous hidden state value ($ht-1$) are given as
84inputs to this gate, uses weights W_{ho} , W_{xo} and bias b_o , the output values are passed through
85sigmoid function to give values 'ot' in the range $[0,1]$. Tanh operation on the next cell
86state(ct) constrains the values to the range $(-1,1)$ i.e. reduces the impact of few values and
87while other values are passed on as they are before.

88The operation $ot * \tanh(ct)$ filters out only the needed values from ct and gives them as next
89hidden state (ht) values. Thereafter, the ht 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
91output probabilities.

92The 'n' number of cells are cascaded one after another getting the values $ct-1$ and $ht-1$ from
93previous cell.

94Weights and bias initialization: Based on references[1][2][3] and[4], the weights which are
95in dot product with hidden state values from previous cell, in every gate and current cell
96state are obtained by 'orthogonal initialization'. The number of units in hidden layer is 100
97and 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)

100Except bias in output linear layer all other bias values have the shape (1,100) and the shape
101of 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
104single row. Addition of bias values was possible due to array broadcasting

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

111 a. Forward Propagation :

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

116The input and output data consists of two types of characters i.e. '0' and '1'. For these 2
117classes, softmax classification can be used which necessitates one hot encoding of the input
118data. 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
120manner.

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

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

125Cross entropy loss function is used to calculate the loss and the loss values of all the samples
126are summed up to form batch loss based on reference[6]

127The output result values are returned as dict from the forward propagation function

128Back Propagation : The gradients of weights and biases are declared in as class members
129as they have same dimension as their corresponding parameters. In the back propagation
130function ,at the beginning,gradients of biases are declared with row size equaling to the
131batch size. After the back prop is complete for a batch ,these 2 steps are followed, a.the
132gradients of biases are summed and assigned to bias gradients which are declared before as
133class members. b. The weight and bias gradients are sent from every process to process with
134rank 0. Here, all the gradients are reduced i.e. summed to form batch gradients that are used
135for optimization. Reference[7] was used to arrive on this idea of implementation

136References [8],[9],[10] ,[11],[12] and figure above were used to derive and implement the
137back propagation

138Optimization: Adam optimizer is implemented from scratch using reference [13] and except
139the learning rate, remaining hyper parameters for optimizer in the current implementation
140are same as in the Adam optimizer paper.

141

142Prediction: Two types of output predictions are used .

143a. One method is to check if the last bit of the predicted data matches with the last bit of the
144output test data sample

145b. Another method is to check if every bit in the predicted data matches with every bit in the
146relevant time step of the output test data sample

147Accuracy Calculation :

148Accuracy = No. of correctly predicted data samples /No. of output data samples *100

149

150 **4 Results**

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

154

155Training and prediction for Fix length batch data :

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

164

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

168

169Run1:

170Learning rate = 0.0005 , epochs =12

171Loss does not decrease significantly even after 5 epochs

172Training stopped

173

174Run2:

175Learning rate = 0.001 (same as in Adam optimization paper) ,epochs =12

176Loss increases significantly for every epoch

177Training stopped

178

179Run3:

180Learning rate = 0.00075 ,epochs =12

181 Last epoch loss : 36
182 min batch Loss : 29.415
183 correct_prediction_count using last bit =19616
184 Prediction Accuracy using last bit =98.080000
185 correct_prediction_count for full length =7910
186 Prediction Accuracy for full length =39.550000

187

188

189 Run4:

190 Learning rate = 0.00080 ,epochs =12

191 Last epoch loss = 1780

192 min_batch_loss =1087.689

193 correct_prediction_count using last bit =9922

194 Prediction Accuracy using last bit =49.610000

195 correct_prediction_count for full length =0

196 Prediction Accuracy for full length =0.000000

197

198

199 Run5:

200 Learning rate = 0.00076 , epochs =12

201 Last epoch loss = 1960

202 min_batch_loss =110.529

203 correct_prediction_count using last bit =10195

204 Prediction Accuracy using last bit =50.975000

205 correct_prediction_count for full length =0

206 Prediction Accuracy for full length =0.000000

207

208 Run6 and Final result for fixed length data:

209 Learning rate = 0.00075 , epochs =20

210 Last epoch loss = 1

211 min_batch_loss = 1.075

212 correct_prediction_count using last bit = 19978

213 Prediction Accuracy using last bit = 99.890000

214 correct_prediction_count for full length = 18932

215 Prediction Accuracy for full length = 94.660000

216

217

218 Below graph is loss plot for Run 6 using number of iterations on x-axis

219 Number of iterations =10000 (500 values of averaged batch losses in 20 epochs)

220

221

222

223

224

225

226

227

228

229

230

231

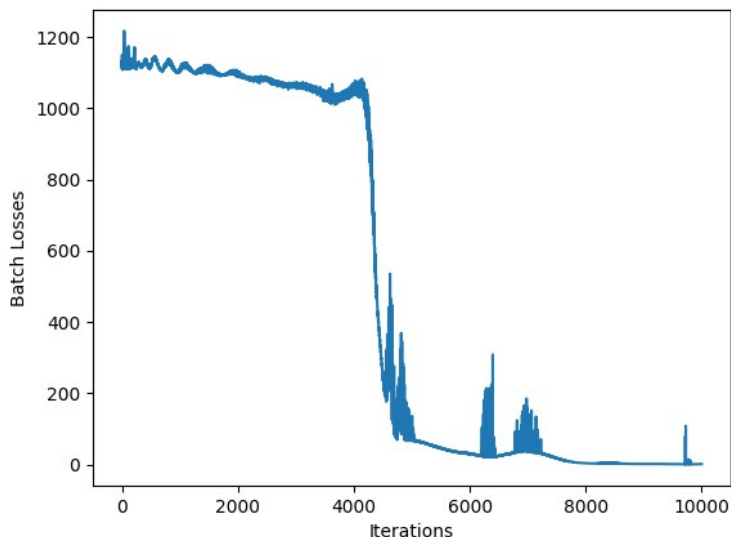
232

233

234

235

236



237
238
239
240
241
242
243Training and prediction for Variable length batch data :
244Varying learning rates for the ranges of epochs instead of single learning rate helped in
245convergence . Loss curve and decrement seem to be complicated when compared to those
246of the fixed length data .
247The number of epochs required are more in number compared to what is used for the fixed
248length data.
249It took 27 runs to figure out the correct learning rates and no. of epochs to get prediction
250accuracy above 92% for variable length data.
251Batch size =32 ,random seed =0 was used in all the runs below
252Except for variation in handling data set, number of epochs and learning rates , no other
253change is made in the model for training variable length data set.
254
255Run1:
256Learning rate = 0.00075 , epochs =20
257Training loss for Epoch:11 is 1909
258Training stopped after 12 epochs due to high fluctuations in loss
259
260Run2:
261Learning rate = 0.00050 ,epochs =20
262Training loss for Epoch:8 is 2359
263Training stopped after 9 epochs due to high fluctuations in loss
264
265Run 3:
266Learning rate = 0.00040 , epochs =20
267Training loss for Epoch:11 is 599
268Training loss for Epoch:12 is 2280
269Training stopped
270
271Run 4:
272Learning rate = 0.00035 , epochs = 20
273Last epoch loss is 44
274min_batch_loss :=4.391
275correct_prediction_count using last bit =18878
276Prediction Accuracy using last bit =94.390000
277correct_prediction_count for full length =11162
278Prediction Accuracy for full length =55.810000
279
280Run 5:
281Learning rate = 0.00035 ,epochs = 28
282Last epoch loss is 464
283min_batch_loss :=3.319
284correct_prediction_count using last bit =16397
285Prediction Accuracy using last bit =81.985000
286correct_prediction_count for full length =9081
287Prediction Accuracy for full length =45.405000
288
289Run 6:
290Learning rate = 0.00030 ,epochs = 20
291Training loss for Epoch:16 is 301
292Training loss for Epoch:17 is 434

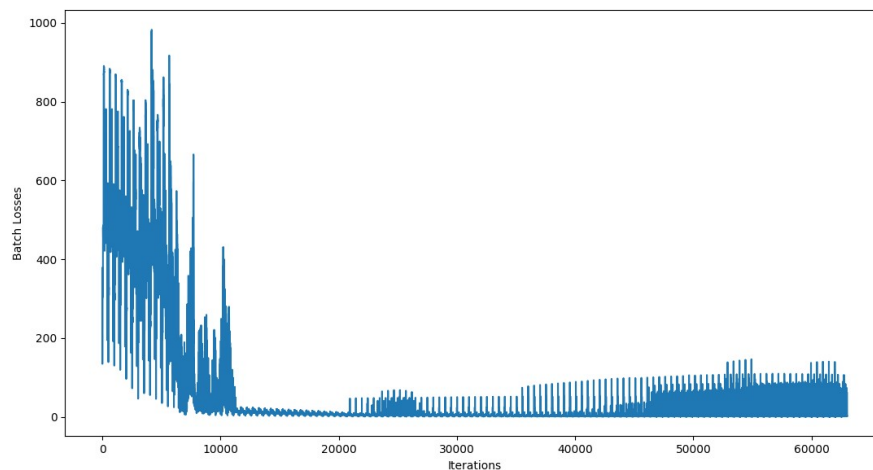
293Training stopped
294
295Run 7:
296Learning rate = 0.00025 ,epochs = 20
297Last epoch loss is 479
298min_batch_loss :=63.931
299correct_prediction_count using last bit =10828
300Prediction Accuracy using last bit =54.140000
301correct_prediction_count for full length =2140
302Prediction Accuracy for full length =10.700000
303
304
305Run 8:
306Learning rate =0.00035 ,epochs=30
307Training loss for Epoch:29 is 2119
308min_batch_loss :=3.319
309correct_prediction_count using last bit =15151
310Prediction Accuracy using last bit =75.755000
311correct_prediction_count for full length =11532
312Prediction Accuracy for full length =57.660000
313
314Run 9:
315Learning rate = 0.00034 ,epochs = 20
316Training loss for Epoch:15 is 2000
317Canceled after 16 epochs
318
319Run 10:
320Learning rate = 0.000349, epochs =20
321Training loss for last epoch is 1138
322
323Run 11:
324Learning rate = 0.00035 for first 20 epochs and 0.00025 for next 10 epochs ,epochs =30
325Training loss for Epoch:19 is 44
326Training loss for Epoch:20 is 170
327Last epoch loss is 213
328min_batch_loss :=3.289
329correct_prediction_count using last bit =18832
330Prediction Accuracy using last bit =94.160000
331correct_prediction_count for full length =13775
332Prediction Accuracy for full length =68.875000
333
334Run 12:
335Learning rate = 0.00035 for first 20 epochs and 0.00012 for next 10 epochs , epochs =30
336Training loss for Epoch:19 is 44
337Training loss for Epoch:20 is 134
338last epoch loss is 48
339min_batch_loss :=3.199
340correct_prediction_count using last bit =19350
341Prediction Accuracy using last bit =96.750000
342correct_prediction_count for full length =11238
343Prediction Accuracy for full length =56.190000
344
345Run 13:
346Learning rate = 0.00035 for first 20 epochs and 0.00006 for next 10 epochs ,epochs =30
347Training loss for Epoch:19 is 44
348Training loss for Epoch:20 is 131

349Last epoch loss is 16
350min_batch_loss :=3.221
351correct_prediction_count using last bit =19671
352Prediction Accuracy using last bit =98.355000
353correct_prediction_count for full length =13565
354Prediction Accuracy for full length =67.825000
355
356
357Run 14:
358Learning rate = 0.00035 for first 20 epochs and 0.00010 for next 10 epochs , epochs =30
359Training loss for Epoch:19 is 44
360Training loss for Epoch:20 is 123
361Last epoch loss is 23
362min_batch_loss :=3.219
363correct_prediction_count using last bit =19597
364Prediction Accuracy using last bit =97.985000
365correct_prediction_count for full length =12789
366Prediction Accuracy for full length =63.945000
367
368Run 15:
369Learning rate = 0.00035 for first 20 epochs and 0.00009 for next 10 epochs , epochs =30
370Training loss for Epoch:19 is 44
371Training loss for Epoch:20 is 134
372Last epoch loss is 22
373min_batch_loss :=3.209
374correct_prediction_count using last bit =19614
375Prediction Accuracy using last bit =98.070000
376correct_prediction_count for full length =12828
377Prediction Accuracy for full length =64.140000
378
379Run 16:
380Learning rate = 0.00035 for first 20 epochs and 0.00001 for next 10 epochs , epochs =30
381Training loss for Epoch:19 is 44
382Training loss for Epoch:20 is 129
383Training loss for Epoch:29 is 14
384min_batch_loss :=3.182
385correct_prediction_count using last bit =19702
386Prediction Accuracy using last bit =98.510000
387correct_prediction_count for full length =13592
388Prediction Accuracy for full length =67.960000
389
390Run 17:
391Learning rate = 0.00035 for first 20 epochs and 0.000001 for next 10 epochs ,epochs =30
392Training loss for Epoch:19 is 44
393Training loss for Epoch:20 is 127
394Training loss for Epoch:29 is 76
395min_batch_loss :=3.930
396correct_prediction_count using last bit =19238
397Prediction Accuracy using last bit =96.190000
398correct_prediction_count for full length =11661
399Prediction Accuracy for full length =58.305000
400
401Run 18:
402Learning rate = 0.00035 for first 20 epochs and 0.00004 for next 10 epochs , epochs=30
403Training loss for Epoch:19 is 44
404Training loss for Epoch:20 is 132

405Training loss for Epoch:29 is 11
406min_batch_loss :=2.628
407correct_prediction_count using last bit =19760
408Prediction Accuracy using last bit =98.800000
409correct_prediction_count for full length =14664
410Prediction Accuracy for full length =73.320000
411
412
413Run 19:
414Learning rate = 0.00035 for first 20 epochs and 0.00004 for next 31 epochs ,epochs =51
415Training loss for Epoch:19 is 44
416Training loss for Epoch:20 is 132
417Last epoch loss is 4
418min_batch_loss :=0.836
419correct_prediction_count using last bit =19894
420Prediction Accuracy using last bit =99.470000
421correct_prediction_count for full length =17535
422Prediction Accuracy for full length =87.675000
423
424
425Run 20:
426Learning rate = 0.00035 for first 20 epochs and 0.00004 for next 45 epochs , epochs =65
427Training loss for Epoch:19 is 44
428Training loss for Epoch:20 is 132
429Last epoch loss is 50
430min_batch_loss :=0.836
431correct_prediction_count using last bit =19679
432Prediction Accuracy using last bit =98.395000
433correct_prediction_count for full length =16167
434Prediction Accuracy for full length =80.835000
435
436
437Run 21:
438Learning rate = 0.00035 for first 20 epochs , 0.00004 from epoch no:20 to epoch no:45 and
4390.000001 from epoch no:46 to epoch no: 54 , epochs = 55
440Training loss for Epoch:19 is 44
441Training loss for Epoch:20 is 132
442Training loss for Epoch:45 is 4
443Training loss for Epoch:46 is 4
444Last epoch loss is 4
445min_batch_loss :=0.907
446correct_prediction_count using last bit =19895
447Prediction Accuracy using last bit =99.475000
448correct_prediction_count for full length =17481
449Prediction Accuracy for full length =87.405000
450
451Run 22:
452Learning rate = 0.00035 for first 20 epochs ,0.00004 for epoch no:21 to epoch no:45 epochs
4530.00001 from epoch no:46 to epoch no: 59 ,epochs =60
454Training loss for Epoch:19 is 44
455Training loss for Epoch:20 is 132
456Training loss for Epoch:45 is 4
457Training loss for Epoch:46 is 5
458Last epoch loss is 4
459min_batch_loss :=0.808
460correct_prediction_count using last bit =19907

461Prediction Accuracy using last bit =99.535000
462correct_prediction_count for full length =17667
463Prediction Accuracy for full length =88.335000
464
465Run 23:
466Learning rate = 0.00035 for first 20 epochs, 0.00004 from epoch no:20 to epoch no:45
467and 0.000004 from epoch no:46 to epoch no: 64 ,epochs = 65
468Training loss for Epoch:19 is 44
469Training loss for Epoch:20 is 132
470Training loss for Epoch:45 is 4
471Training loss for Epoch:46 is 4
472Last epoch loss is 4
473min_batch_loss :=0.826
474correct_prediction_count using last bit =19900
475Prediction Accuracy using last bit =99.500000
476correct_prediction_count for full length =17643
477Prediction Accuracy for full length =88.215000
478
479Run 24:
480Learning rate = 0.00035 for first 20 epochs , 0.00004 from epoch no:20 to epoch no:45,
4810.00002 from epoch no:46 to epoch no: 64 ,epochs = 65
482Training loss for Epoch:19 is 44
483Training loss for Epoch:20 is 132
484Training loss for Epoch:45 is 4
485Training loss for Epoch:46 is 5
486Last epoch loss is 6
487min_batch_loss :=0.741
488correct_prediction_count using last bit =19890
489Prediction Accuracy using last bit =99.450000
490correct_prediction_count for full length =17863
491Prediction Accuracy for full length =89.31500
492
493Run 25:
494Learning rate = 0.00035 for first 20 epochs , 0.00004 from epoch no:21 to epoch no:45 and
4950.000014 from epoch no:46 to epoch no: 64 ,epochs = 65
496Training loss for Epoch:19 is 44
497Training loss for Epoch:20 is 132
498Training loss for Epoch:45 is 4
499Training loss for Epoch:46 is 5
500Last epoch loss is 4
501min_batch_loss :=0.795
502correct_prediction_count using last bit =19903
503Prediction Accuracy using last bit =99.515000
504correct_prediction_count for full length =17753
505Prediction Accuracy for full length =88.765000
506
507Run 26:
508Learning rate = 0.00035 for first 20 epochs , 0.00004 for epoch no:20 to epoch no: 45 and
5090.00001 from epoch no:46 to epoch no: 79 ,epochs =80
510Training loss for Epoch:19 is 44
511Training loss for Epoch:20 is 132
512Training loss for Epoch:45 is 4
513Training loss for Epoch:46 is 5
514Training loss for last epoch is 3
515min_batch_loss :=0.596
516correct_prediction_count using last bit =19931

517 Prediction Accuracy using last bit =99.655000
 518 correct_prediction_count for full length =18317
 519 Prediction Accuracy for full length =91.585000
 520
 521
 522
 523 Run 27 and Final epoch for variable length data:
 524 Learning rate = 0.00035 for first 20 epochs ,0.00004 for epoch no:20 to epoch no: 45 and
 525 0.00001 from epoch no:46 to epoch no: 125 ,epochs = 125
 526 Training loss for Epoch:19 is 44
 527 Training loss for Epoch:20 is 132
 528 Training loss for Epoch:45 is 4
 529 Training loss for Epoch:46 is 5
 530 Training loss for last epoch is 5
 531 min_batch_loss :=0.405
 532 correct_prediction_count using last bit =19938
 533 Prediction Accuracy using last bit =99.690000
 534 correct_prediction_count for full length =18838
 535 Prediction Accuracy for full length =94.190000
 536
 537 Below graph is loss plot for Run 27 using number of iterations on x-axis
 538 Number of iterations = 63000 (504 values of averaged batch losses in 125 epochs)
 539



540 The point to note is that the epochs are increased from 80 to 125 and other hyper parameters
 541 are same when Run 26 and Run27 are compared . Final epoch loss is more in Run 27 by 1
 542 but the accuracy is high by 2.6 %
 543
 544

545 Improvements and Conclusion

546 Following improvements can be further made which are left for future work
 547 a. Tuning learning rate using learning rate scheduler perhaps can reduce the training time to
 548 get good results
 549 b. MPI parallel programming with 5 processes reduced the training time of a single epoch to
 550 14 min compared to 20 min training time of one epoch when MPI was not used.
 551 Using GPU based cloud infrastructure such as GCP or AWS can further reduce the
 552 training time to a large extent . This can help to train using efficient hyper parameter
 553 tuning methods

554c. All experiments are run with seed 0 . Few other random seeds can be used in different
555 trials.
556d. Batch size 32 was used and other batch sizes such as 64,128 ,256 could be tried.
557
558 This project implementation helped to understand LSTM model and application of parallel
559 programming using MPI

560References

561[1]Orthogonal initialization for hidden to hidden , Xavier uniform initialization for rest of the weights
562<https://www.kaggle.com/code/junkoda/pytorch-lstm-with-tensorflow-like-initialization/notebook>
563[2]Weight initialization
564[https://www.deeplearningwizard.com/deep_learning/boosting_models_pytorch/](https://www.deeplearningwizard.com/deep_learning/boosting_models_pytorch/weight_initialization_activation_functions/)
565[weight_initialization_activation_functions/](https://www.deeplearningwizard.com/deep_learning/boosting_models_pytorch/weight_initialization_activation_functions/)
566[https://wandb.ai/sauravmaheshkar/initialization/reports/A-Gentle-Introduction-To-Weight-](https://wandb.ai/sauravmaheshkar/initialization/reports/A-Gentle-Introduction-To-Weight-Initialization-for-Neural-Networks--Vmlldzo2ODExMTg)
567[Initialization-for-Neural-Networks--Vmlldzo2ODExMTg](https://wandb.ai/sauravmaheshkar/initialization/reports/A-Gentle-Introduction-To-Weight-Initialization-for-Neural-Networks--Vmlldzo2ODExMTg)
568[https://www.reddit.com/r/deeplearning/comments/9vjckm/](https://www.reddit.com/r/deeplearning/comments/9vjckm/what_is_the_best_way_to_initialize_lstm_weights/)
569[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/)
570[3] Orthogonal weight initialization
571<https://github.com/kastnerkyle/net/blob/d66d533b1ebcc25dd442bdb41431d334b617e80e/net.py#L168>
572<https://github.com/nyu-dl/dl4mt-tutorial/blob/master/session3/lm.py>
573[4] Fan in Fan out for xavier uniform
574[https://stackoverflow.com/questions/42670274/how-to-calculate-fan-in-and-fan-out-in-xavier-](https://stackoverflow.com/questions/42670274/how-to-calculate-fan-in-and-fan-out-in-xavier-initialization-for-neural-networks)
575[initialization-for-neural-networks](https://stackoverflow.com/questions/42670274/how-to-calculate-fan-in-and-fan-out-in-xavier-initialization-for-neural-networks)
576<https://ml-cheatsheet.readthedocs.io/en/latest/forwardpropagation.html>
577<https://visualstudiomagazine.com/articles/2019/09/05/neural-network-glorot.aspx>
578[5] One hot encoding :
579[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)
580[input-outputs](https://datascience.stackexchange.com/questions/45803/should-we-use-only-one-hot-vector-for-lstm-input-outputs)
581[6]Batch Loss
582<https://stackoverflow.com/questions/54053868/how-do-i-get-a-loss-per-epoch-and-not-per-batch>
583[7] Distibuted training and optimization
584Diagram : https://github.com/vlimant/mpi_learn/blob/master/docs/downpour.png
585MPI.AllReduce: [https://github.com/openai/baselines/blob/master/baselines/common/](https://github.com/openai/baselines/blob/master/baselines/common/mpi_adam_optimizer.py)
586[mpi_adam_optimizer.py](https://github.com/openai/baselines/blob/master/baselines/common/mpi_adam_optimizer.py)
587[8]Back propagation flow diagrams
588<https://towardsdatascience.com/backpropagation-in-rnn-explained-bdf853b4e1c2>
589https://github.com/christinakouridi/scratchML/blob/master/LSTM/LSTMForward_diagram.png
590[9] Back propagation for batch
591[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)
592[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)
593[10] Modularizing code and RNN Implementation :
594[https://github.com/JY-Yoon/RNN-Implementation-using-NumPy/blob/master/RNN](https://github.com/JY-Yoon/RNN-Implementation-using-NumPy/blob/master/RNN%20Implementation%20using%20NumPy.ipynb)
595[%20Implementation%20using%20NumPy.ipynb](https://github.com/JY-Yoon/RNN-Implementation-using-NumPy/blob/master/RNN%20Implementation%20using%20NumPy.ipynb)
596[11] Back prop chain rule, matrix multiply dimension matching and Matrix -Matrix Multiply
597gradient : <https://cs231n.github.io/optimization-2/>
598[12] Adam Optimizer
599<https://towardsdatascience.com/how-to-implement-an-adam-optimizer-from-scratch-76e7b217f1cc>
600<https://stackoverflow.com/questions/70225531/time-step-in-adam-optimizer>
601<https://arxiv.org/pdf/1412.6980.pdf>

602[13] Vanishing and exploding gradients in LSTM (This problem is not entirely solved in LSTM ,rather
603it is mitigated and delayed)

604[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)
605[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)

606 [14] Cross Entropy

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

608[15] Matrix and batch dimensions:

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

610[16] Bias gradients sum along axis =0 :

611<https://github.com/Erlemar/cs231n/blob/master/assignment2/cs231n/layers.py#L116>

612[17] RNN Theory : <http://karpathy.github.io/2015/05/21/rnn-effectiveness/>

613[18] LSTM Theory : <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

614 [19] LSTM theory explanation : [https://towardsdatascience.com/lstm-networks-a-detailed-](https://towardsdatascience.com/lstm-networks-a-detailed-explanation-8fae6aefc7f9)

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

616[20] Batch size : <https://ai.stackexchange.com/questions/8560/how-do-i-choose-the-optimal-batch-size>

617[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)

618[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)

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

620[22] Python binary string generation

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

622[23] Gradient clipping

623[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)

624[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)