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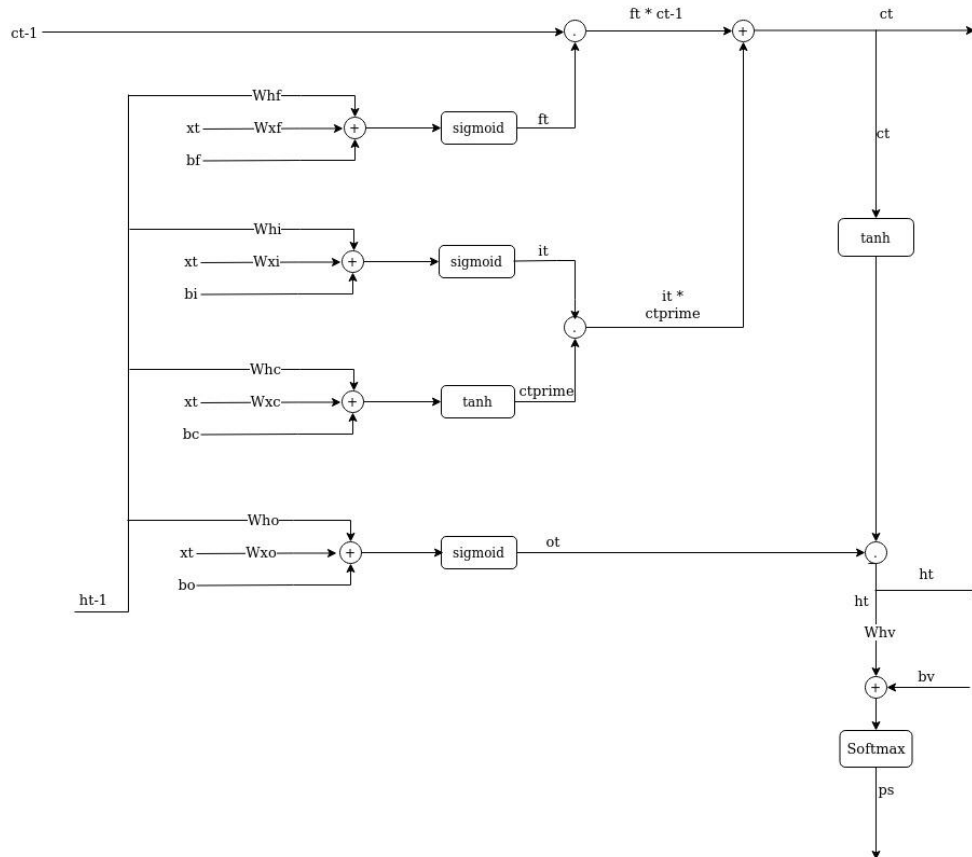
Train dataset split into batches  : The 80000 data samples in type1 are divided into 2500 batches each having size of 32. In type2 , 80000 samples are divided into 2521 batches each having  variable size that is <=32. The reason for this variability is that the data in type 2 has variable lengths and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

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 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)
73  value and previous hidden state value(ht-1) using  weights Whc,Wxc and bias bc .The use of
74  ctprime is to modify the information in xt given the ht-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*ctprime  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*ctprime 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 (xt) and previous hidden state value (ht-1) are given as
84  inputs to this gate ,uses weights Who,Wxo and bias bo , 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 (ht)  values. Thereafter, the ht values passed through a linear layer with weight
90  Whv and bias bv 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 ht-1 from
93  previous cell.

94  Weights and bias initialization: Based on references[1][2][3] and[4] , the weights which are
95  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 andtraining 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 $1^{st}$ 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    Back Propagation : The gradients of weights and biases are declared in as class members
129    as they have same dimension as their corresponding  parameters. In the back propagation
130    function ,at the beginning,gradients of biases are declared with row size equaling to the
131    batch size. After the back prop is complete for a batch ,these 2 steps are followed, a.the
132    gradients of biases are summed and assigned to bias gradients which are declared before as
133    class members. b. The weight and bias gradients are sent from every process to process with
134    rank 0. Here, all the  gradients are reduced i.e. summed to form batch gradients that are used
135    for optimization. Reference[7] was used to arrive on this idea of implementation

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

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

141

142    Prediction: Two types of output predictions are used .
143    a. One method is to check if the last bit of the predicted data matches with the last bit of the
144    output test data simple
145    b. Another method is to check if every bit in the predicted data matches with every bit in the
146    relevant time step of the output test data sample

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

149

## 150    4      Results

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

154

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

164

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

168

169    Run1:
170    Learning rate = 0.0005 , epochs =12
171    Loss does not decrease significantly even after 5 epochs
172    Training stopped

173

174    Run2:
175    Learning rate = 0.001 (same as in Adam optimization paper) ,epochs =12
176    Loss increases significantly for every epoch
177    Training stopped

178

179    Run3:
180    Learning 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
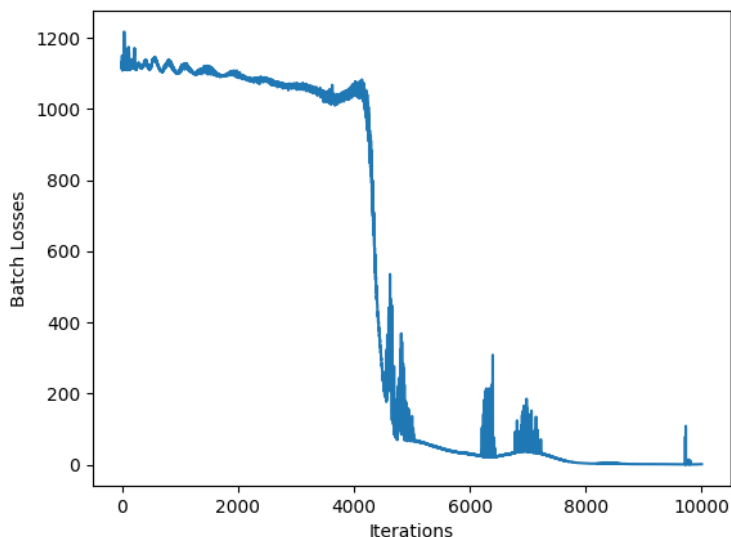207
208    <u>Run6 and Final result for fixed length data</u>:
209     Learning rate = 0.00075  , epochs =20
210     Last epoch loss  = 1
211     min_batch_loss = 1.075
212    correct_prediction_count using last bit = 19978
213    Prediction Accuracy using last bit = 99.890000
214    correct_prediction_count 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

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

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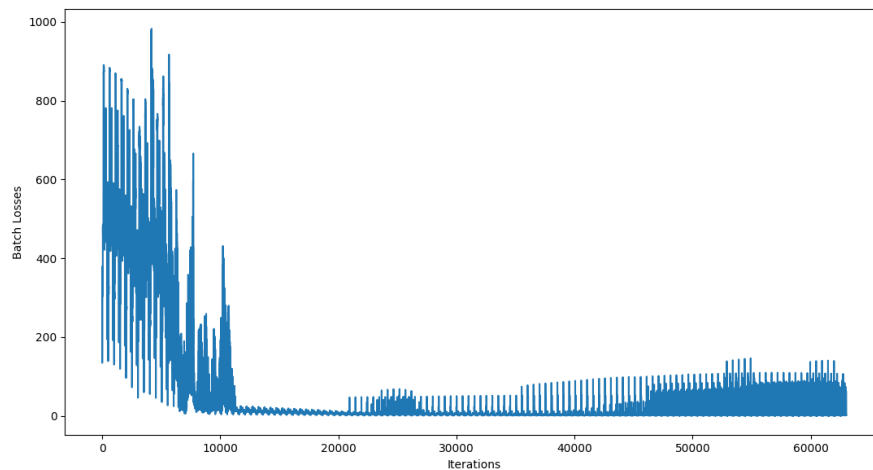348    Training loss for Epoch:20 is 131

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

```
405    Training loss for Epoch:29 is 11
406    min_batch_loss :=2.628
407    correct_prediction_count using last bit =19760
408    Prediction Accuracy using last bit =98.800000
409    correct_prediction_count for full length =14664
410    Prediction Accuracy for full length =73.320000
411
412
413    Run 19:
414    Learning rate = 0.00035 for first 20 epochs and 0.00004 for next 31 epochs ,epochs =51
415    Training loss for Epoch:19 is 44
416    Training loss for Epoch:20 is 132
417    Last epoch loss is 4
418    min_batch_loss :=0.836
419    correct_prediction_count using last bit =19894
420    Prediction Accuracy using last bit =99.470000
421    correct_prediction_count for full length =17535
422    Prediction Accuracy for full length =87.675000
423
424
425    Run 20:
426    Learning rate = 0.00035 for first 20 epochs and 0.00004 for next 45 epochs  , epochs =65
427    Training loss for Epoch:19 is 44
428    Training loss for Epoch:20 is 132
429    Last epoch loss is 50
430    min_batch_loss :=0.836
431    correct_prediction_count using last bit =19679
432    Prediction Accuracy using last bit =98.395000
433    correct_prediction_count for full length =16167
434    Prediction Accuracy for full length =80.835000
435
436
437    Run 21:
438    Learning rate = 0.00035 for first 20 epochs , 0.00004 from epoch no:20 to  epoch no:45  and
439    0.000001 from epoch no:46 to epoch no: 54 , epochs = 55
440    Training loss for Epoch:19 is 44
441    Training loss for Epoch:20 is 132
442    Training loss for Epoch:45 is 4
443    Training loss for Epoch:46 is 4
444    Last epoch loss is 4
445    min_batch_loss :=0.907
446    correct_prediction_count using last bit =19895
447    Prediction Accuracy using last bit =99.475000
448    correct_prediction_count for full length =17481
449    Prediction Accuracy for full length =87.405000
450
451    Run 22:
452    Learning rate = 0.00035 for first 20 epochs ,0.00004 for epoch no:21 to  epoch no:45 epochs
453    0.00001 from epoch no:46 to epoch no: 59 ,epochs =60
454    Training loss for Epoch:19 is 44
455    Training loss for Epoch:20 is 132
456    Training loss for Epoch:45 is 4
457    Training loss for Epoch:46 is 5
458    Last epoch loss is 4
459    min_batch_loss :=0.808
460    correct_prediction_count using last bit =19907
```

461    Prediction Accuracy using last bit =99.535000
462    correct_prediction_count for full length =17667
463    Prediction Accuracy for full length =88.335000
464
465    Run 23:
466    Learning rate = 0.00035 for first 20 epochs, 0.00004 from epoch no:20 to  epoch no:45
467    and 0.000004 from epoch no:46 to epoch no: 64 ,epochs = 65
468    Training loss for Epoch:19 is 44
469    Training loss for Epoch:20 is 132
470    Training loss for Epoch:45 is 4
471    Training loss for Epoch:46 is 4
472    Last epoch loss is 4
473    min_batch_loss :=0.826
474    correct_prediction_count using last bit =19900
475    Prediction Accuracy using last bit =99.500000
476    correct_prediction_count for full length =17643
477    Prediction Accuracy for full length =88.215000
478
479    Run 24:
480    Learning rate = 0.00035 for first 20 epochs , 0.00004 from epoch no:20 to  epoch no:45,
481    0.00002 from epoch no:46 to epoch no: 64 ,epochs = 65
482    Training loss for Epoch:19 is 44
483    Training loss for Epoch:20 is 132
484    Training loss for Epoch:45 is 4
485    Training loss for Epoch:46 is 5
486    Last epoch loss is 6
487    min_batch_loss :=0.741
488    correct_prediction_count using last bit =19890
489    Prediction Accuracy using last bit =99.450000
490    correct_prediction_count for full length =17863
491    Prediction Accuracy for full length =89.31500
492
493    Run 25:
494    Learning rate = 0.00035 for first 20 epochs , 0.00004 from epoch no:21 to  epoch no:45 and
495    0.000014 from epoch no:46 to epoch no: 64 ,epochs =  65
496    Training loss for Epoch:19 is 44
497    Training loss for Epoch:20 is 132
498    Training loss for Epoch:45 is 4
499    Training loss for Epoch:46 is 5
500    Last epoch loss is 4
501    min_batch_loss :=0.795
502    correct_prediction_count using last bit =19903
503    Prediction Accuracy using last bit =99.515000
504    correct_prediction_count for full length =17753
505    Prediction Accuracy for full length =88.765000
506
507    Run 26:
508    Learning rate = 0.00035 for first 20 epochs , 0.00004 for epoch no:20 to  epoch no: 45  and
509    0.00001 from epoch no:46 to epoch no: 79 ,epochs =80
510    Training loss for Epoch:19 is 44
511    Training loss for Epoch:20 is 132
512    Training loss for Epoch:45 is 4
513    Training loss for Epoch:46 is 5
514    Training loss for last epoch is 3
515    min_batch_loss :=0.596
516    correct_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    <u>Run 27 and Final epoch for variable length data:</u>
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

## 5 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
554    c. All experiments are run with seed 0 . Few other random seeds can be used in different

555    trials.
556    d. 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

## References

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/weight_initialization_a
565    ctivation_functions/

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

568    https://www.reddit.com/r/deeplearning/comments/9vjckm/what_is_the_best_way_to_initialize_lstm_w
569    eights/

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-
575    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-
580    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

584    Diagram : https://github.com/vlimant/mpi_learn/blob/master/docs/downpour.png

585    MPI.AllReduce:
586    .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

589    https://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-
592    a-batch-of-data-points

593    [10] Modularizing code and RNN Implementation :

594    https://github.com/JY-Yoon/RNN-Implementation-using-
595    NumPy/blob/master/RNN%20Implementation%20using%20NumPy.ipynb

596    [11] Back prop chain rule, matrix multiply dimension matching  and Matrix -Matrix Multiply gradient
597    : 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

603     it is mitigated and delayed)

604     https://www.quora.com/How-does-LSTM-help-prevent-the-vanishing-and-exploding-gradient-
605     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_self/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-
615     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-
618     margani.medium.com/how-to-speed-up-the-training-of-the-sequence-model-using-bucketing-
619     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-
624     global-norm-for-rnns-and-how