
LSTM XOR Project Report

Jayaram Kuchibhotla
jayaramkuchibhotla@gmail.com

Abstract

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

1 Introduction

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

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

2 Description of the datasets

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

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

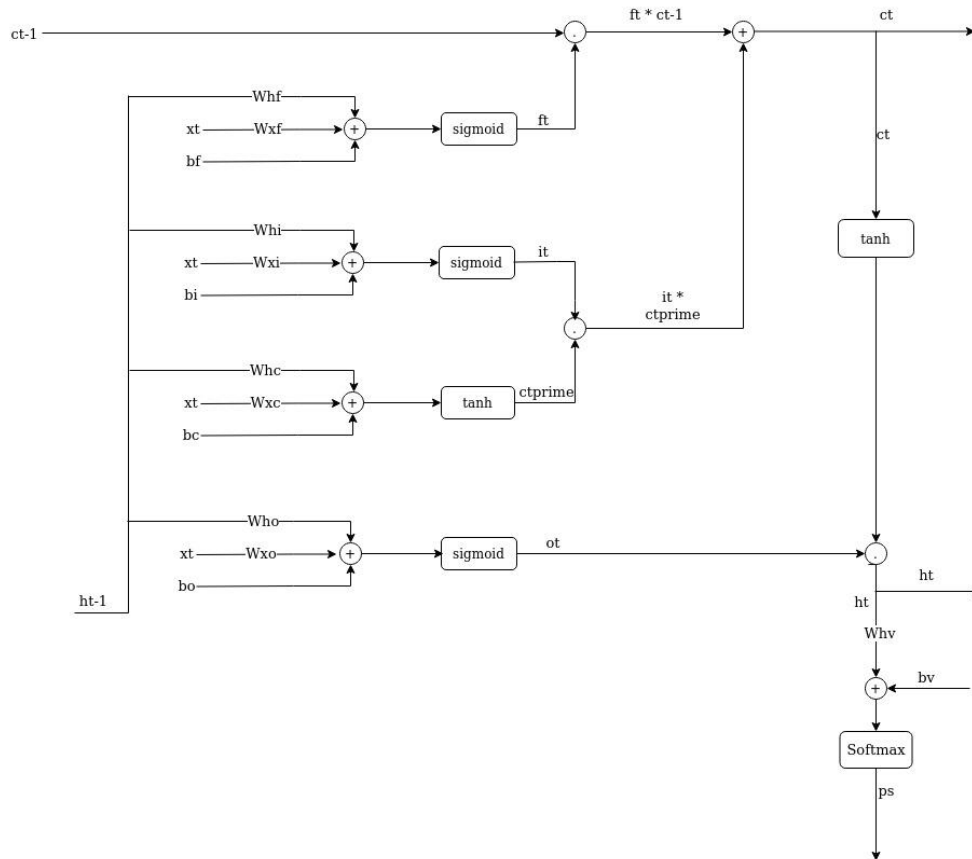
Train dataset split into batches : The 80000 data samples in type1 are divided into 2500 batches each having size of 32. In type2 , 80000 samples are divided into 2521 batches each having variable size that is ≤ 32 . The reason for this variability is that the data in type 2 has variable lengths and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

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

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

50

51 3 Explanation of model architecture and implementation



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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

111 a. Forward Propagation:

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

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

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

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

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

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

128 b. Back Propagation:

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

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

139 c. Optimization:

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

143

144 Prediction: Two types of output predictions are used .

145 1. One method is to check if the last bit of the predicted data matches with the last bit of the
146 output test data sample

147 2. Another method is to check if every bit in the predicted data matches with every bit in the
148 relevant time step of the output test data sample

149 Accuracy Calculation:

150 $\text{Accuracy} = (\text{No. of correctly predicted data samples} / \text{No. of output data samples}) * 100$

151

152 **4 Results**

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

156

157 Training and prediction for Fix length batch data:

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

166

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

170

171 Run1:

172 Learning rate = 0.0005, epochs =12

173 Loss does not decrease significantly even after 5 epochs

174 Training stopped

175

176

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

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

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

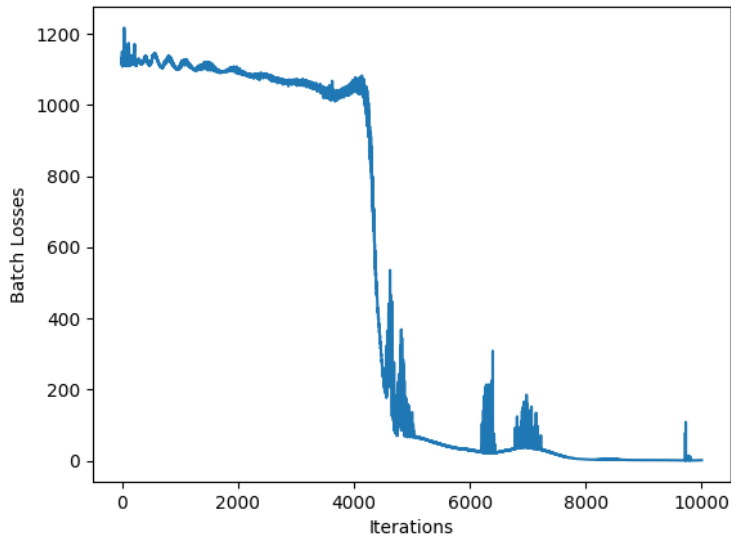
200
201
202 Run5:
203 Learning rate = 0.00076 , epochs =12
204 Last epoch loss = 1960
205 min_batch_loss =110.529
206 correct_prediction_count using last bit =10195
207 Prediction Accuracy using last bit =50.975000
208 correct_prediction_count for full length =0
209 Prediction Accuracy for full length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

210
211
212
213
214
215
216 Run6 and Final result for fixed length data:
217 Learning rate = 0.00075 , epochs =20
218 Last epoch loss = 1
219 min_batch_loss = 1.075
220 correct_prediction_count using last bit = 19978
221 Prediction Accuracy using last bit = 99.890000
222 correct_prediction_count for full length = 18932
223 Prediction Accuracy for full length = 94.660000

224
225
226
227
228
229
230
231
232

233
234
235
236

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



257

258 Training and prediction for Variable length batch data :
259 Varying learning rates for the ranges of epochs instead of single learning rate helped in
260 convergence. Loss curve and decrement seem to be complicated when compared to those
261 of the fixed length data .
262 The number of epochs required are more in number compared to what is used for the fixed
263 length data.
264 It took 27 runs to figure out the correct learning rates and no. of epochs to get prediction
265 accuracy above 92% for variable length data.
266 Batch size =32 ,random seed =0 was used in all the runs below
267 Except for variation in handling data set, number of epochs and learning rates , no other
268 change is made in the model for training variable length data set.
269
270 Run1:
271 Learning rate = 0.00075 , epochs =20
272 Training loss for Epoch:11 is 1909
273 Training stopped after 12 epochs due to high fluctuations in loss
274
275 Run2:
276 Learning rate = 0.00050 ,epochs =20
277 Training loss for Epoch:8 is 2359
278 Training stopped after 9 epochs due to high fluctuations in loss
279
280 Run 3:
281 Learning rate = 0.00040 , epochs =20
282 Training loss for Epoch:11 is 599
283 Training loss for Epoch:12 is 2280
284 Training stopped
285
286 Run 4:
287 Learning rate = 0.00035 , epochs = 20
288 Last epoch loss is 44

```
289 min_batch_loss :=4.391
290 correct_prediction_count using last bit =18878
291 Prediction Accuracy using last bit =94.390000
292 correct_prediction_count for full length =11162
293 Prediction Accuracy for full length =55.810000
294
295 Run 5:
296 Learning rate = 0.00035 ,epochs = 28
297 Last epoch loss is 464
298 min_batch_loss :=3.319
299 correct_prediction_count using last bit =16397
300 Prediction Accuracy using last bit =81.985000
301 correct_prediction_count for full length =9081
302 Prediction Accuracy for full length =45.405000
303
304 Run 6:
305 Learning rate = 0.00030 ,epochs = 20
306 Training loss for Epoch:16 is 301
307 Training loss for Epoch:17 is 434
308 Training stopped
309
310 Run 7:
311 Learning rate = 0.00025 ,epochs = 20
312 Last epoch loss is 479
313 min_batch_loss :=63.931
314 correct_prediction_count using last bit =10828
315 Prediction Accuracy using last bit =54.140000
316 correct_prediction_count for full length =2140
317 Prediction Accuracy for full length =10.700000
318
319
320 Run 8:
321 Learning rate =0.00035 ,epochs=30
322 Training loss for Epoch:29 is 2119
323 min_batch_loss :=3.319
324 correct_prediction_count using last bit =15151
325 Prediction Accuracy using last bit =75.755000
326 correct_prediction_count for full length =11532
327 Prediction Accuracy for full length =57.660000
328
329 Run 9:
330 Learning rate = 0.00034 ,epochs = 20
331 Training loss for Epoch:15 is 2000
332 Canceled after 16 epochs
333
334 Run 10:
335 Learning rate = 0.000349, epochs =20
336 Training loss for last epoch is 1138
337
338 Run 11:
339 Learning rate = 0.00035 for first 20 epochs and 0.00025 for next 10 epochs ,epochs =30
340 Training loss for Epoch:19 is 44
341 Training loss for Epoch:20 is 170
342 Last epoch loss is 213
343 min_batch_loss :=3.289
344 correct_prediction_count using last bit =18832
```

345 Prediction Accuracy using last bit =94.160000
346 correct_prediction_count for full length =13775
347 Prediction Accuracy for full length =68.875000
348
349 Run 12:
350 Learning rate = 0.00035 for first 20 epochs and 0.00012 for next 10 epochs , epochs =30
351 Training loss for Epoch:19 is 44
352 Training loss for Epoch:20 is 134
353 last epoch loss is 48
354 min_batch_loss :=3.199
355 correct_prediction_count using last bit =19350
356 Prediction Accuracy using last bit =96.750000
357 correct_prediction_count for full length =11238
358 Prediction Accuracy for full length =56.190000
359
360 Run 13:
361 Learning rate = 0.00035 for first 20 epochs and 0.00006 for next 10 epochs ,epochs =30
362 Training loss for Epoch:19 is 44
363 Training loss for Epoch:20 is 131
364 Last epoch loss is 16
365 min_batch_loss :=3.221
366 correct_prediction_count using last bit =19671
367 Prediction Accuracy using last bit =98.355000
368 correct_prediction_count for full length =13565
369 Prediction Accuracy for full length =67.825000
370
371
372 Run 14:
373 Learning rate = 0.00035 for first 20 epochs and 0.00010 for next 10 epochs , epochs =30
374 Training loss for Epoch:19 is 44
375 Training loss for Epoch:20 is 123
376 Last epoch loss is 23
377 min_batch_loss :=3.219
378 correct_prediction_count using last bit =19597
379 Prediction Accuracy using last bit =97.985000
380 correct_prediction_count for full length =12789
381 Prediction Accuracy for full length =63.945000
382
383 Run 15:
384 Learning rate = 0.00035 for first 20 epochs and 0.00009 for next 10 epochs , epochs =30
385 Training loss for Epoch:19 is 44
386 Training loss for Epoch:20 is 134
387 Last epoch loss is 22
388 min_batch_loss :=3.209
389 correct_prediction_count using last bit =19614
390 Prediction Accuracy using last bit =98.070000
391 correct_prediction_count for full length =12828
392 Prediction Accuracy for full length =64.140000
393
394 Run 16:
395 Learning rate = 0.00035 for first 20 epochs and 0.00001 for next 10 epochs , epochs =30
396 Training loss for Epoch:19 is 44
397 Training loss for Epoch:20 is 129
398 Training loss for Epoch:29 is 14
399 min_batch_loss :=3.182
400 correct_prediction_count using last bit =19702

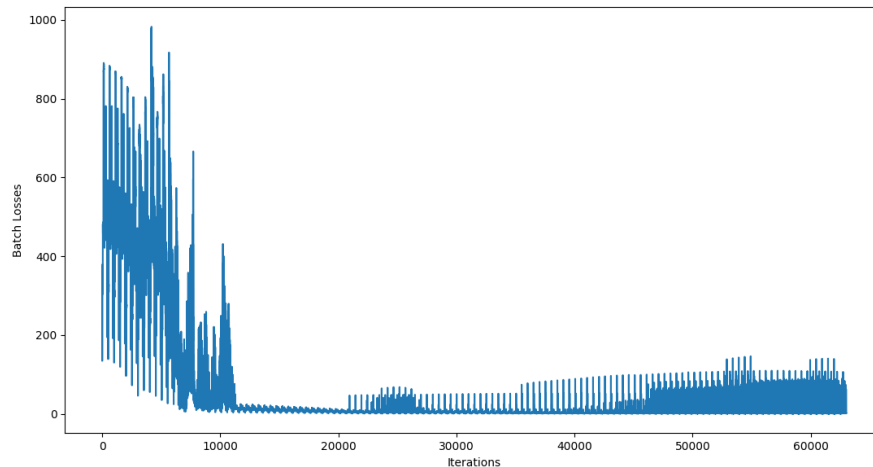
401 Prediction Accuracy using last bit =98.510000
402 correct_prediction_count for full length =13592
403 Prediction Accuracy for full length =67.960000
404
405 Run 17:
406 Learning rate = 0.00035 for first 20 epochs and 0.000001 for next 10 epochs ,epochs =30
407 Training loss for Epoch:19 is 44
408 Training loss for Epoch:20 is 127
409 Training loss for Epoch:29 is 76
410 min_batch_loss :=3.930
411 correct_prediction_count using last bit =19238
412 Prediction Accuracy using last bit =96.190000
413 correct_prediction_count for full length =11661
414 Prediction Accuracy for full length =58.305000
415
416 Run 18:
417 Learning rate = 0.00035 for first 20 epochs and 0.00004 for next 10 epochs , epochs=30
418 Training loss for Epoch:19 is 44
419 Training loss for Epoch:20 is 132
420 Training loss for Epoch:29 is 11
421 min_batch_loss :=2.628
422 correct_prediction_count using last bit =19760
423 Prediction Accuracy using last bit =98.800000
424 correct_prediction_count for full length =14664
425 Prediction Accuracy for full length =73.320000
426
427
428 Run 19:
429 Learning rate = 0.00035 for first 20 epochs and 0.00004 for next 31 epochs ,epochs =51
430 Training loss for Epoch:19 is 44
431 Training loss for Epoch:20 is 132
432 Last epoch loss is 4
433 min_batch_loss :=0.836
434 correct_prediction_count using last bit =19894
435 Prediction Accuracy using last bit =99.470000
436 correct_prediction_count for full length =17535
437 Prediction Accuracy for full length =87.675000
438
439
440 Run 20:
441 Learning rate = 0.00035 for first 20 epochs and 0.00004 for next 45 epochs , epochs =65
442 Training loss for Epoch:19 is 44
443 Training loss for Epoch:20 is 132
444 Last epoch loss is 50
445 min_batch_loss :=0.836
446 correct_prediction_count using last bit =19679
447 Prediction Accuracy using last bit =98.395000
448 correct_prediction_count for full length =16167
449 Prediction Accuracy for full length =80.835000
450
451
452 Run 21:
453 Learning rate = 0.00035 for first 20 epochs , 0.00004 from epoch no:20 to epoch no:45 and
454 0.000001 from epoch no:46 to epoch no: 54 , epochs = 55
455 Training loss for Epoch:19 is 44
456 Training loss for Epoch:20 is 132

457 Training loss for Epoch:45 is 4
458 Training loss for Epoch:46 is 4
459 Last epoch loss is 4
460 min_batch_loss :=0.907
461 correct_prediction_count using last bit =19895
462 Prediction Accuracy using last bit =99.475000
463 correct_prediction_count for full length =17481
464 Prediction Accuracy for full length =87.405000
465
466 Run 22:
467 Learning rate = 0.00035 for first 20 epochs ,0.00004 for epoch no:21 to epoch no:45 epochs
468 0.00001 from epoch no:46 to epoch no: 59 ,epochs =60
469 Training loss for Epoch:19 is 44
470 Training loss for Epoch:20 is 132
471 Training loss for Epoch:45 is 4
472 Training loss for Epoch:46 is 5
473 Last epoch loss is 4
474 min_batch_loss :=0.808
475 correct_prediction_count using last bit =19907
476 Prediction Accuracy using last bit =99.535000
477 correct_prediction_count for full length =17667
478 Prediction Accuracy for full length =88.335000
479
480 Run 23:
481 Learning rate = 0.00035 for first 20 epochs, 0.00004 from epoch no:20 to epoch no:45
482 and 0.000004 from epoch no:46 to epoch no: 64 ,epochs = 65
483 Training loss for Epoch:19 is 44
484 Training loss for Epoch:20 is 132
485 Training loss for Epoch:45 is 4
486 Training loss for Epoch:46 is 4
487 Last epoch loss is 4
488 min_batch_loss :=0.826
489 correct_prediction_count using last bit =19900
490 Prediction Accuracy using last bit =99.500000
491 correct_prediction_count for full length =17643
492 Prediction Accuracy for full length =88.215000
493
494 Run 24:
495 Learning rate = 0.00035 for first 20 epochs , 0.00004 from epoch no:20 to epoch no:45,
496 0.00002 from epoch no:46 to epoch no: 64 ,epochs = 65
497 Training loss for Epoch:19 is 44
498 Training loss for Epoch:20 is 132
499 Training loss for Epoch:45 is 4
500 Training loss for Epoch:46 is 5
501 Last epoch loss is 6
502 min_batch_loss :=0.741
503 correct_prediction_count using last bit =19890
504 Prediction Accuracy using last bit =99.450000
505 correct_prediction_count for full length =17863
506 Prediction Accuracy for full length =89.31500
507
508 Run 25:
509 Learning rate = 0.00035 for first 20 epochs , 0.00004 from epoch no:21 to epoch no:45 and
510 0.000014 from epoch no:46 to epoch no: 64 ,epochs = 65
511 Training loss for Epoch:19 is 44
512 Training loss for Epoch:20 is 132

513 Training loss for Epoch:45 is 4
514 Training loss for Epoch:46 is 5
515 Last epoch loss is 4
516 min_batch_loss :=0.795
517 correct_prediction_count using last bit =19903
518 Prediction Accuracy using last bit =99.515000
519 correct_prediction_count for full length =17753
520 Prediction Accuracy for full length =88.765000
521
522
523
524
525 Run 26:
526 Learning rate = 0.00035 for first 20 epochs , 0.00004 for epoch no:20 to epoch no: 45 and
527 0.00001 from epoch no:46 to epoch no: 79 ,epochs =80
528 Training loss for Epoch:19 is 44
529 Training loss for Epoch:20 is 132
530 Training loss for Epoch:45 is 4
531 Training loss for Epoch:46 is 5
532 Training loss for last epoch is 3
533 min_batch_loss :=0.596
534 correct_prediction_count using last bit =19931
535 Prediction Accuracy using last bit =99.655000
536 correct_prediction_count for full length =18317
537 Prediction Accuracy for full length =91.585000
538
539
540
541
542
543
544
545
546 Run 27 and Final epoch for variable length data:
547 Learning rate = 0.00035 for first 20 epochs ,0.00004 for epoch no:20 to epoch no: 45 and
548 0.00001 from epoch no:46 to epoch no: 125 ,epochs = 125
549 Training loss for Epoch:19 is 44
550 Training loss for Epoch:20 is 132
551 Training loss for Epoch:45 is 4
552 Training loss for Epoch:46 is 5
553 Training loss for last epoch is 5
554 min_batch_loss :=0.405
555 correct_prediction_count using last bit =19938
556 Prediction Accuracy using last bit =99.690000
557 correct_prediction_count for full length =18838
558 Prediction Accuracy for full length =94.190000
559
560
561
562
563
564
565
566
567
568

569
570
571
572

Below graph is loss plot for Run 27 using number of iterations on x-axis
Number of iterations = 63000 (504 values of averaged batch losses in 125 epochs)



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

576
577

578 **5 Improvements and Conclusion**

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

- 580 a. Tuning learning rate using learning rate scheduler perhaps can reduce the training time to
- 581 get good results
- 582 b. MPI parallel programming with 5 processes reduced the training time of a single epoch to
- 583 14 min compared to 20 min training time of one epoch when MPI was not used.
- 584 Using GPU based cloud infrastructure such as GCP or AWS can further reduce the
- 585 training time to a large extent . This can help to train using efficient hyper parameter
- 586 tuning methods
- 587 c. All experiments are run with seed 0 . Few other random seeds can be used in different
- 588 trials.
- 589 d. Batch size 32 was used and other batch sizes such as 64,128 ,256 could be tried.

590
591
592
593
594
595

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

596 **References**

- 597 [1]Orthogonal initialization for hidden to hidden , Xavier uniform initialization for rest of the weights
- 598 <https://www.kaggle.com/code/junkoda/pytorch-lstm-with-tensorflow-like-initialization/notebook>
- 599 [2]Weight initialization
- 600 [https://www.deeplearningwizard.com/deep_learning/boosting_models_pytorch/weight_initialization_a](https://www.deeplearningwizard.com/deep_learning/boosting_models_pytorch/weight_initialization_activation_functions/)
- 601 [ctivation_functions/](https://www.deeplearningwizard.com/deep_learning/boosting_models_pytorch/weight_initialization_activation_functions/)
- 602 [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)
- 603 [Initialization-for-Neural-Networks--Vmlldzo2ODExMTg](https://wandb.ai/sauravmaheshkar/initialization/reports/A-Gentle-Introduction-To-Weight-Initialization-for-Neural-Networks--Vmlldzo2ODExMTg)
- 604 https://www.reddit.com/r/deeplearning/comments/9vjckm/what_is_the_best_way_to_initialize_lstm_w

605 [eights/](#)

606 [3] Orthogonal weight initialization

607 <https://github.com/kastnerkyle/net/blob/d66d533b1ebcc25dd442bdb41431d334b617e80e/net.py#L168>

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

609 [4] Fan in Fan out for xavier uniform

610 <https://stackoverflow.com/questions/42670274/how-to-calculate-fan-in-and-fan-out-in-xavier-initialization-for-neural-networks>

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

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

613 [5] One hot encoding :

614 <https://datascience.stackexchange.com/questions/45803/should-we-use-only-one-hot-vector-for-lstm-input-outputs>

615 <https://datascience.stackexchange.com/questions/45803/should-we-use-only-one-hot-vector-for-lstm-input-outputs>

616 [6] Batch Loss

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

618 [7] Distributed training and optimization

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

620 MPI.AllReduce:

621 https://github.com/openai/baselines/blob/master/baselines/common/mpi_adam_optimizer.py

622 [8] Back propagation flow diagrams

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

624 https://github.com/christinakouridi/scratchML/blob/master/LSTM/LSTMForward_diagram.png

625 [9] Back propagation for batch

626 <https://ai.stackexchange.com/questions/11667/is-back-propagation-applied-for-each-data-point-or-for-a-batch-of-data-points>

627 <https://ai.stackexchange.com/questions/11667/is-back-propagation-applied-for-each-data-point-or-for-a-batch-of-data-points>

628 [10] Modularizing code and RNN Implementation :

629 <https://github.com/JY-Yoon/RNN-Implementation-using-NumPy/blob/master/RNN%20Implementation%20using%20NumPy.ipynb>

630 <https://github.com/JY-Yoon/RNN-Implementation-using-NumPy/blob/master/RNN%20Implementation%20using%20NumPy.ipynb>

631 [11] Back prop chain rule, matrix multiply dimension matching and Matrix -Matrix Multiply gradient

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

633 [12] Adam Optimizer

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

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

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

637 [13] Vanishing and exploding gradients in LSTM (This problem is not entirely solved in LSTM ,rather it is mitigated and delayed)

638 <https://www.quora.com/How-does-LSTM-help-prevent-the-vanishing-and-exploding-gradient-problem-in-a-recurrent-neural-network>

639 <https://www.quora.com/How-does-LSTM-help-prevent-the-vanishing-and-exploding-gradient-problem-in-a-recurrent-neural-network>

640 [14] Cross Entropy

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

642 [15] Matrix and batch dimensions:

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

644 [16] Bias gradients sum along axis =0 :

645 https://github.com/Erleamar/cs231n_self/blob/master/assignment2/cs231n/layers.py#L116

646 [17] RNN Theory : <http://karpathy.github.io/2015/05/21/rnn-effectiveness/>

647 [18] LSTM Theory : <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

648 [19] LSTM theory explanation : <https://towardsdatascience.com/lstm-networks-a-detailed-explanation-8fae6aefc7f9>

649 <https://towardsdatascience.com/lstm-networks-a-detailed-explanation-8fae6aefc7f9>

650 <https://towardsdatascience.com/lstm-networks-a-detailed-explanation-8fae6aefc7f9>

651 <https://towardsdatascience.com/lstm-networks-a-detailed-explanation-8fae6aefc7f9>

652 [20] Batch size : <https://ai.stackexchange.com/questions/8560/how-do-i-choose-the-optimal-batch-size>
653 [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)
654 [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)
655 [techniques-9e302b0fd976](https://rashmi-margani.medium.com/how-to-speed-up-the-training-of-the-sequence-model-using-bucketing-techniques-9e302b0fd976)
656 [22] Python binary string generation
657 <https://www.geeksforgeeks.org/python-program-to-generate-random-binary-string/>
658 [23] Gradient clipping
659 [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)
660 [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)