

EEE543 Mini Project 2 – RNN on Human Activity Recognition

1. Problem Definition

In this work, we are required to train a recurrent neural network (RNN) on a human activity recognition (HAR) dataset, with the backpropagation through time (BPTT) algorithm. Some ablation studies regarding the hyperparameters of the network and the training regime are also demonstrated with accompanying conclusions in the following sections.

2. Derivations and Implementation Details

The derivations for backpropagation through time (BPTT) for recurrent neural networks and the usage of truncated BPTT as batch learning are given in this section.

2.1. BPTT

This section includes a comprehensive derivation for BPTT using matrix notations.

$\bar{\mathbf{M}}_{K \times (M+1)}$ is the extended matrix where the +1 dimension is padded with extra 1's to accommodate for biases, and $\mathbf{M}_{K \times M}$ is its non-extended version.

The RNN equations are as follows:

$$\mathbf{h}_{50 \times 1}^t = \tanh(\bar{\mathbf{W}}_{ih_{50 \times (3+1)}} \mathbf{i}_{(3+1) \times 1}^t + \bar{\mathbf{W}}_{hh_{50 \times (50+1)}} \mathbf{h}_{(50+1) \times 1}^{t-1})$$

$$\mathbf{y}_{6 \times 1} = \text{sigmoid}(\bar{\mathbf{W}}_{ho_{6 \times (50+1)}} \mathbf{h}_{(50+1) \times 1}^t)$$

The error is determined as the multiclass cross-entropy loss:

$$E = - \sum_{k \in \text{classes}} d_k \log(y_k) + (1 - d_k) \log(1 - y_k) \equiv -\log(y_k)$$

Finding the derivative with respect to \mathbf{W}_{ho} is straightforward:

$$\frac{\partial E}{\partial \mathbf{W}_{ho}} = \frac{\partial E}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial \mathbf{o}} \frac{\partial \mathbf{o}}{\partial \mathbf{W}_{ho}} = - \frac{1}{\mathbf{y}_{6 \times 1}} (\mathbf{y}_{6 \times 1} \otimes (\mathbf{d}_{6 \times 1} - \mathbf{y}_{6 \times 1})) (\bar{\mathbf{h}}_{(50+1) \times 1}^{\text{last}})^T = (\mathbf{y}_{6 \times 1} - \mathbf{d}_{6 \times 1}) (\bar{\mathbf{h}}_{(50+1) \times 1}^{\text{last}})^T$$

Note that in $\frac{\partial \mathbf{y}^t}{\partial \mathbf{o}^t} = (\mathbf{y}_{6 \times 1} \otimes (\mathbf{d}_{6 \times 1} - \mathbf{y}_{6 \times 1}))$, \otimes is the Kronecker product. We either get $(\mathbf{y}_{6 \times 1})$ or $(\mathbf{1} - \mathbf{y}_{6 \times 1})$ from $\mathbf{y}_{6 \times 1} - \mathbf{d}_{6 \times 1}$.

For \mathbf{W}_{hh} , we will use backpropagation through time (BPTT):

$$\begin{aligned} \frac{\partial E}{\partial \mathbf{W}_{hh}} &= \frac{\partial E}{\partial \mathbf{h}^{\text{last}}} \left(\sum_t \frac{\partial \mathbf{h}^{\text{last}}}{\partial \mathbf{h}^{\text{last}-t}} \frac{\partial \mathbf{h}^{\text{last}-t}}{\partial \mathbf{W}_{hh}} \right) \\ &= \frac{\partial E}{\partial \mathbf{h}^{\text{last}}} \frac{\partial \mathbf{h}^{\text{last}}}{\partial \mathbf{W}_{hh}} + \frac{\partial E}{\partial \mathbf{h}^{\text{last}}} \frac{\partial \mathbf{h}^{\text{last}}}{\partial \mathbf{h}^{\text{last}-1}} \frac{\partial \mathbf{h}^{\text{last}-1}}{\partial \mathbf{W}_{hh}} + \frac{\partial E}{\partial \mathbf{h}^{\text{last}}} \frac{\partial \mathbf{h}^{\text{last}}}{\partial \mathbf{h}^{\text{last}-1}} \frac{\partial \mathbf{h}^{\text{last}-1}}{\partial \mathbf{h}^{\text{last}-2}} \frac{\partial \mathbf{h}^{\text{last}-2}}{\partial \mathbf{W}_{hh}} + \dots \end{aligned}$$

To understand better, write down several terms and observe the patterns:

$$\frac{\partial E}{\partial \mathbf{h}^{\text{last}}} \frac{\partial \mathbf{h}^{\text{last}}}{\partial \mathbf{W}_{hh}}$$

$$\begin{aligned}
&= \left(\text{diag} \left(\mathbf{1} - \mathbf{h}_{50 \times 1}^{\text{last}^2} \right) (\mathbf{W}_{\text{ho}_{6 \times 50}})^T (\mathbf{y}_{6 \times 1} - \mathbf{d}_{6 \times 1}) (\bar{\mathbf{h}}_{(50+1) \times 1}^{\text{last}-1})^T \right)_{50 \times (50+1)} \\
&= \mathbf{A}_{50 \times 1} (\bar{\mathbf{h}}_{(50+1) \times 1}^{\text{last}-1})^T
\end{aligned}$$

where $\text{diag}(\mathbf{1} - \mathbf{h}_{50 \times 1}^{\text{last}^2})$ is the derivative of tanh with the output $\mathbf{h}_{50 \times 1}^{\text{last}}$.

Expanding the second and the third term yields:

$$\begin{aligned}
&\frac{\partial \mathbf{E}}{\partial \mathbf{h}^{\text{last}}} \frac{\partial \mathbf{h}^{\text{last}}}{\partial \mathbf{h}^{\text{last}-1}} \frac{\partial \mathbf{h}^{\text{last}-1}}{\partial \mathbf{W}_{\text{hh}}} \\
&= \left(\text{diag} \left(\mathbf{1} - \mathbf{h}_{50 \times 1}^{\text{last}-1^2} \right) \mathbf{W}_{\text{hh}_{50 \times 50}} \text{diag} \left(\mathbf{I} - \mathbf{h}_{50 \times 1}^{\text{last}^2} \right) (\mathbf{W}_{\text{ho}_{6 \times 50}})^T (\mathbf{y}_{6 \times 1} - \mathbf{d}_{6 \times 1}) (\bar{\mathbf{h}}_{(50+1) \times 1}^{\text{last}-2})^T \right)_{50 \times (50+1)} \\
&= \text{diag} \left(\mathbf{1} - \mathbf{h}_{50 \times 1}^{\text{last}-1^2} \right) \mathbf{W}_{\text{hh}_{50 \times 50}} \mathbf{A}_{50 \times 1} (\bar{\mathbf{h}}_{(50+1) \times 1}^{\text{last}-2})^T = \mathbf{B}_{50 \times 1} (\bar{\mathbf{h}}_{(50+1) \times 1}^{\text{last}-2})^T \\
&\frac{\partial \mathbf{E}}{\partial \mathbf{h}^{\text{last}}} \frac{\partial \mathbf{h}^{\text{last}}}{\partial \mathbf{h}^{\text{last}-1}} \frac{\partial \mathbf{h}^{\text{last}-1}}{\partial \mathbf{h}^{\text{last}-2}} \frac{\partial \mathbf{h}^{\text{last}-2}}{\partial \mathbf{W}_{\text{hh}}} \\
&= \text{diag} \left(\mathbf{1} - \mathbf{h}_{50 \times 1}^{\text{last}-2^2} \right) \mathbf{W}_{\text{hh}_{50 \times 50}} \mathbf{B}_{50 \times 1} (\bar{\mathbf{h}}_{(50+1) \times 1}^{\text{last}-3})^T = \mathbf{C}_{50 \times 1} (\bar{\mathbf{h}}_{(50+1) \times 1}^{\text{last}-3})^T
\end{aligned}$$

While propagating through time, we will utilize previously found matrices ($\mathbf{A}_{50 \times 1}, \mathbf{B}_{50 \times 1}, \mathbf{C}_{50 \times 1} \dots$). This will become handy while coding.

For \mathbf{W}_{ih} , we will again use BPTT:

$$\begin{aligned}
\frac{\partial \mathbf{E}}{\partial \mathbf{W}_{\text{ih}}} &= \frac{\partial \mathbf{E}}{\partial \mathbf{h}^{\text{last}}} \left(\sum_t \frac{\partial \mathbf{h}^{\text{last}}}{\partial \mathbf{h}^{\text{last}-t}} \frac{\partial \mathbf{h}^{\text{last}-t}}{\partial \mathbf{W}_{\text{ih}}} \right) \\
&= \frac{\partial \mathbf{E}}{\partial \mathbf{h}^{\text{last}}} \frac{\partial \mathbf{h}^{\text{last}}}{\partial \mathbf{W}_{\text{ih}}} + \frac{\partial \mathbf{E}}{\partial \mathbf{h}^{\text{last}}} \frac{\partial \mathbf{h}^{\text{last}}}{\partial \mathbf{h}^{\text{last}-1}} \frac{\partial \mathbf{h}^{\text{last}-1}}{\partial \mathbf{W}_{\text{ih}}} + \frac{\partial \mathbf{E}}{\partial \mathbf{h}^{\text{last}}} \frac{\partial \mathbf{h}^{\text{last}}}{\partial \mathbf{h}^{\text{last}-1}} \frac{\partial \mathbf{h}^{\text{last}-1}}{\partial \mathbf{h}^{\text{last}-2}} \frac{\partial \mathbf{h}^{\text{last}-2}}{\partial \mathbf{W}_{\text{ih}}} + \dots \\
&\frac{\partial \mathbf{E}}{\partial \mathbf{h}^{\text{last}}} \frac{\partial \mathbf{h}^{\text{last}}}{\partial \mathbf{W}_{\text{ih}}} \\
&= \left(\text{diag} \left(\mathbf{1} - \mathbf{h}_{50 \times 1}^{\text{last}^2} \right) (\mathbf{W}_{\text{ho}_{6 \times 50}})^T (\mathbf{y}_{6 \times 1} - \mathbf{d}_{6 \times 1}) (\bar{\mathbf{i}}_{(3+1) \times 1}^{\text{last}-1})^T \right)_{50 \times (3+1)} = \mathbf{A}_{50 \times 1} (\bar{\mathbf{i}}_{(3+1) \times 1}^{\text{last}-1})^T \\
&\frac{\partial \mathbf{E}}{\partial \mathbf{h}^{\text{last}}} \frac{\partial \mathbf{h}^{\text{last}}}{\partial \mathbf{h}^{\text{last}-1}} \frac{\partial \mathbf{h}^{\text{last}-1}}{\partial \mathbf{W}_{\text{ih}}} \\
&= \left(\text{diag} \left(\mathbf{1} - \mathbf{h}_{50 \times 1}^{\text{last}-1^2} \right) \mathbf{W}_{\text{hh}_{50 \times 50}} \text{diag} \left(\mathbf{I} - \mathbf{h}_{50 \times 1}^{\text{last}^2} \right) (\mathbf{W}_{\text{ho}_{6 \times 50}})^T (\mathbf{y}_{6 \times 1} - \mathbf{d}_{6 \times 1}) (\bar{\mathbf{i}}_{(3+1) \times 1}^{\text{last}-2})^T \right)_{50 \times (3+1)}
\end{aligned}$$

$$= \text{diag}\left(\mathbf{1} - \mathbf{h}_{50 \times 1}^{\text{last}-1^2}\right) \mathbf{W}_{\text{hh}_{50 \times 50}} \mathbf{A}_{50 \times 1} \left(\mathbf{i}_{(3+1) \times 1}^{\text{last}-2}\right)^T = \mathbf{B}_{50 \times 1} \left(\mathbf{i}_{(3+1) \times 1}^{\text{last}-2}\right)^T$$

Several iterations reveal that the only difference with $\frac{\partial E}{\partial \mathbf{W}_{\text{hh}}}$ is that the rightmost matrix is $\mathbf{i}_{(3+1) \times 1}^k$ rather than $\mathbf{h}_{(50+1) \times 1}^k$.

2.2. Truncated BPTT

In RNNs, batch learning is usually done with truncated BPTT. The difference between BPTT and truncated BPTT is that the forward and backward passes are “truncated”, so that the weight updates are done more frequently, avoiding vanishing/exploding gradients and encouraging short-term memory.

The important thing in truncated BPTT is that the last hidden states are saved at the end of each batch, later to be used in the next batch as the initial condition (see line 5 of **Algorithm 1**).



Figure 1. BPTT vs Truncated BPTT. Right and left arrows denote forward and backward passes, respectively. Figure adapted from the lecture notes.

2.3. Implementation Details

The implementation details are briefly given below:

- The implementation is made from scratch; hence the matrix notations are compliant with the course content. I utilized **torch** library to migrate the data to GPU and perform the training sessions faster with **CUDA**, but I did not utilize any of its automatic gradient functionalities. In other words, using **torch** is identical to using **numpy** but with GPU in this scenario; and in fact, most of the matrix functions are the same in **torch** & **numpy** in terms of valid input arguments and operation behavior.
- For stability purposes and to avoid NaNs, I **clipped the weight gradients between [-1,1]**.

The overall procedure using the derivations given in **Section 2.1** is given in **Algorithm 1**.

Algorithm 1. Pseudocode for training an RNN with truncated BPTT

```

1 Initialize all extended weights  $\bar{\mathbf{W}}$  with uniform distribution between  $[-0.01, 0.01]$ 
2 for epoch in all epochs do
3   for sample  $(\mathbf{i}, \mathbf{d})$  in training dataset do # data and labels
4     Divide  $\mathbf{i}$  into batches for truncated BPTT
5     Initialize  $\bar{\mathbf{h}}_{(50+1) \times 1}^{\text{prev}}$  with  $\mathbf{0}$  if first batch, else  $\bar{\mathbf{h}}_{(50+1) \times 1}^{\text{prev}} \leftarrow \bar{\mathbf{h}}_{(50+1) \times 1}^{\text{seq\_length}}$ 
6     for  $t=0:\text{seq\_length}$  do # Forward pass
7        $\mathbf{h}_{50 \times 1}^{\text{curr}} \leftarrow \tanh(\bar{\mathbf{W}}_{\text{ih}_{50 \times (3+1)}} \mathbf{i}_{(3+1) \times 1}^t + \bar{\mathbf{W}}_{\text{hh}_{50 \times (50+1)}} \bar{\mathbf{h}}_{(50+1) \times 1}^{\text{prev}})$ 
8        $\bar{\mathbf{h}}_{(50+1) \times 1}^{\text{prev}} \leftarrow \bar{\mathbf{h}}_{(50+1) \times 1}^{\text{curr}}$ 
9       Save  $\bar{\mathbf{h}}_{(50+1) \times 1}^{\text{curr}}$  at  $t$  and  $\mathbf{i}_{(3+1) \times 1}^t$  for backprop
10    end for
11     $\mathbf{y}_{6 \times 1} \leftarrow \text{sigmoid}(\bar{\mathbf{W}}_{\text{ho}_{6 \times (50+1)}} \bar{\mathbf{h}}_{(50+1) \times 1}^{\text{curr}})$ 
12     $\mathbf{E} = -\frac{1}{\text{seq\_length}} \sum \log(\mathbf{y})$  # log-likelihood error
13     $\frac{\partial \mathbf{E}}{\partial \mathbf{h}} \leftarrow (\mathbf{W}_{\text{ho}_{6 \times 50}})^T (\mathbf{y}_{6 \times 1} - \mathbf{d}_{6 \times 1})$ 
14     $\frac{\partial \mathbf{E}}{\partial \mathbf{W}_{\text{ho}}} \leftarrow \frac{\partial \mathbf{E}}{\partial \mathbf{W}_{\text{ho}}} + (\mathbf{y}_{6 \times 1} - \mathbf{d}_{6 \times 1}) (\bar{\mathbf{h}}_{(50+1) \times 1}^{\text{seq\_length}})^T$  # BPTT for Who
15    for  $t=\text{seq\_length}:0$  do # Backward pass
16       $\frac{\partial \mathbf{E}}{\partial \mathbf{h}} \leftarrow \text{diag}(\mathbf{1} - \mathbf{h}_{50 \times 1}^t)^2 \frac{\partial \mathbf{E}}{\partial \mathbf{h}}$  # Formation of the intermediate matrices A,B,C...
17       $\frac{\partial \mathbf{E}}{\partial \mathbf{W}_{\text{hh}}} \leftarrow \frac{\partial \mathbf{E}}{\partial \mathbf{W}_{\text{hh}}} + \frac{\partial \mathbf{E}}{\partial \mathbf{h}} (\bar{\mathbf{h}}_{(50+1) \times 1}^{t-1})^T$  # BPTT for Whh
18       $\frac{\partial \mathbf{E}}{\partial \mathbf{W}_{\text{ih}}} \leftarrow \frac{\partial \mathbf{E}}{\partial \mathbf{W}_{\text{ih}}} + \frac{\partial \mathbf{E}}{\partial \mathbf{h}} (\mathbf{i}_{(3+1) \times 1}^{t-1})^T$  # BPTT for Wih
19       $\frac{\partial \mathbf{E}}{\partial \mathbf{h}} \leftarrow \mathbf{W}_{\text{hh}_{50 \times 50}} \frac{\partial \mathbf{E}}{\partial \mathbf{h}}$  # Formation of the intermediate matrices A,B,C...
20      Save  $\frac{\partial \mathbf{E}}{\partial \mathbf{W}_{\text{ho}}}, \frac{\partial \mathbf{E}}{\partial \mathbf{W}_{\text{hh}}}, \frac{\partial \mathbf{E}}{\partial \mathbf{W}_{\text{ih}}}$  for batch update
21    Clip gradients to avoid explosion
22    end for
23     $\bar{\mathbf{W}}_{\text{ih}_{(50) \times (3+1)}} \leftarrow \bar{\mathbf{W}}_{\text{ih}_{(50) \times (3+1)}} - \eta \left( \frac{\partial \mathbf{E}}{\partial \mathbf{W}_{\text{ih}}} \right)$  # Gradient descent
24     $\bar{\mathbf{W}}_{\text{ho}_{6 \times (50+1)}} \leftarrow \bar{\mathbf{W}}_{\text{ho}_{6 \times (50+1)}} - \eta \left( \frac{\partial \mathbf{E}}{\partial \mathbf{W}_{\text{ho}}} \right)$ 
25     $\bar{\mathbf{W}}_{\text{hh}_{(50) \times (50+1)}} \leftarrow \bar{\mathbf{W}}_{\text{hh}_{(50) \times (50+1)}} - \eta \left( \frac{\partial \mathbf{E}}{\partial \mathbf{W}_{\text{hh}}} \right)$ 
26    end if
27    Track  $\mathbf{E}$  until convergence
28  end for
29  shuffle training dataset
30 end for

```

Test and validation steps are not given in the pseudocode. For more details, please check the attached .py file.

3. Results and Discussion

This section includes all results and conclusions.

3.1. Effects of learning rate, hidden layer size, and mini batch size without validation set (b)

The experimental results are given in Table 1 and Figure 2-4. LR, N, and B denote learning rate, hidden layer size, and mini batch size, respectively.

Table 1. Top-k scores for the models trained in Part (b). The best model is given in bold.

LR	N	B	Top-1	Top-2	Top-3	Top-4	Top-5
0.1	50	10	0.410	0.603	0.762	0.928	0.977
0.1	50	30	0.582	0.727	0.860	0.962	0.990
0.1	100	10	0.342	0.568	0.717	0.865	0.956
0.1	100	30	0.353	0.540	0.708	0.838	0.997
0.05	50	10	0.529	0.717	0.857	0.960	0.977
0.05	50	30	0.405	0.585	0.713	0.896	0.983
0.05	100	10	0.353	0.547	0.718	0.940	0.970
0.05	100	30	0.318	0.493	0.642	0.753	0.813

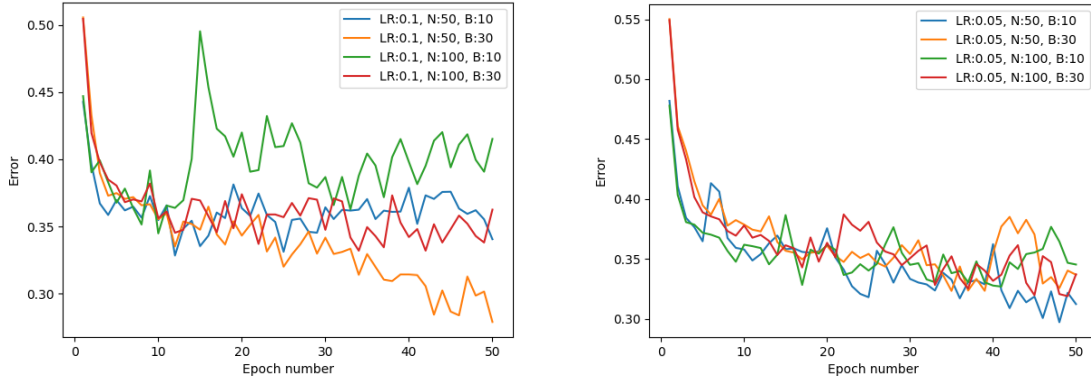


Figure 2. Comparing the models on learning rates (0.1 on left, 0.05 on right).

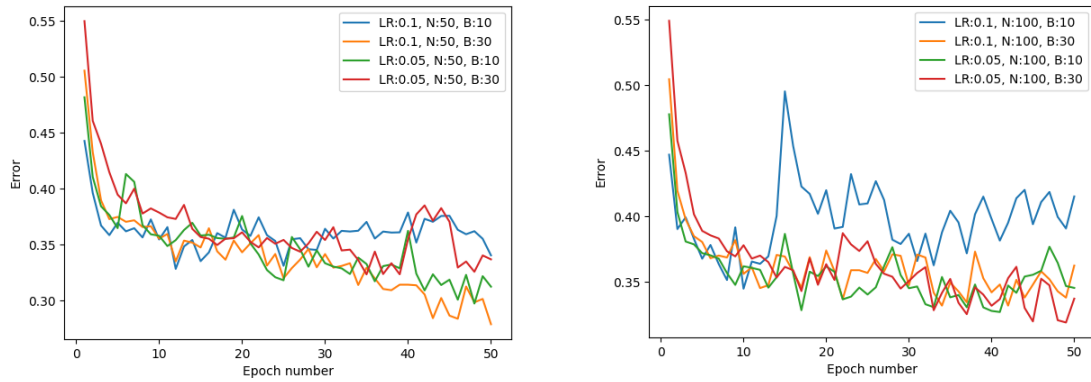


Figure 3. Comparing the models on hidden layer sizes (50 on left, 100 on right). Best viewed digitally.

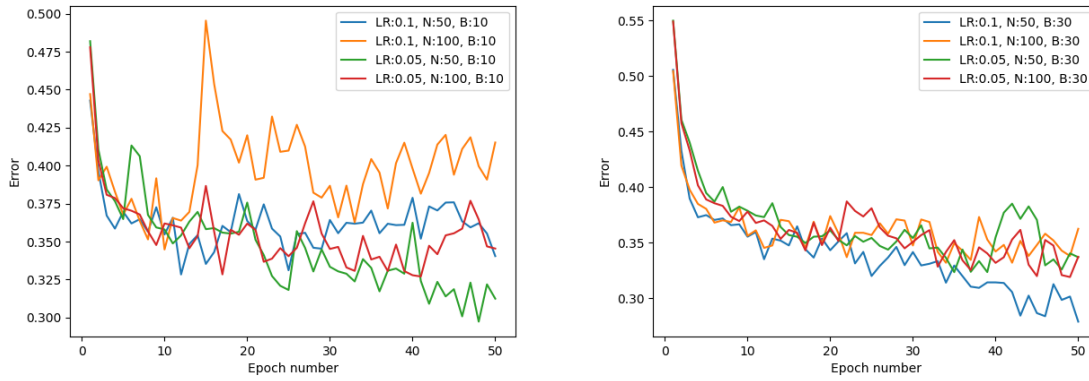
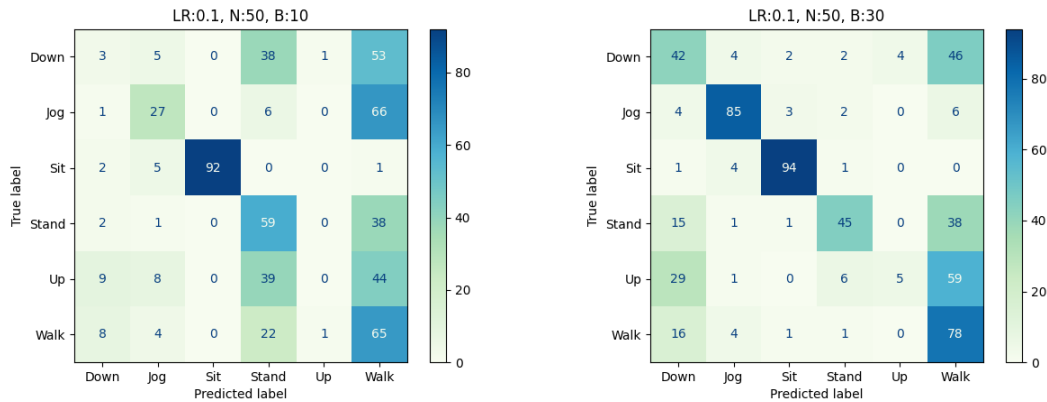


Figure 4. Comparing the models on mini batch sizes (10 on left, 30 on right).

- **Figure 2** reveals that higher batch size with a smaller number of hidden layers work the best among the tried configurations for learning rate = 0.1 (left plot). Higher learning rate enables the models converge to a point with less training error but is more unstable compared to the ones with lower learning rate.
- **Figure 3** suggests using a lower learning rate when model size increases (notice the blue curve on the right plot).
- **Figure 4** reveals using a larger batch size (hence enabling longer-term memory) may be more stable during training (right plot), especially if the learning rate is larger.



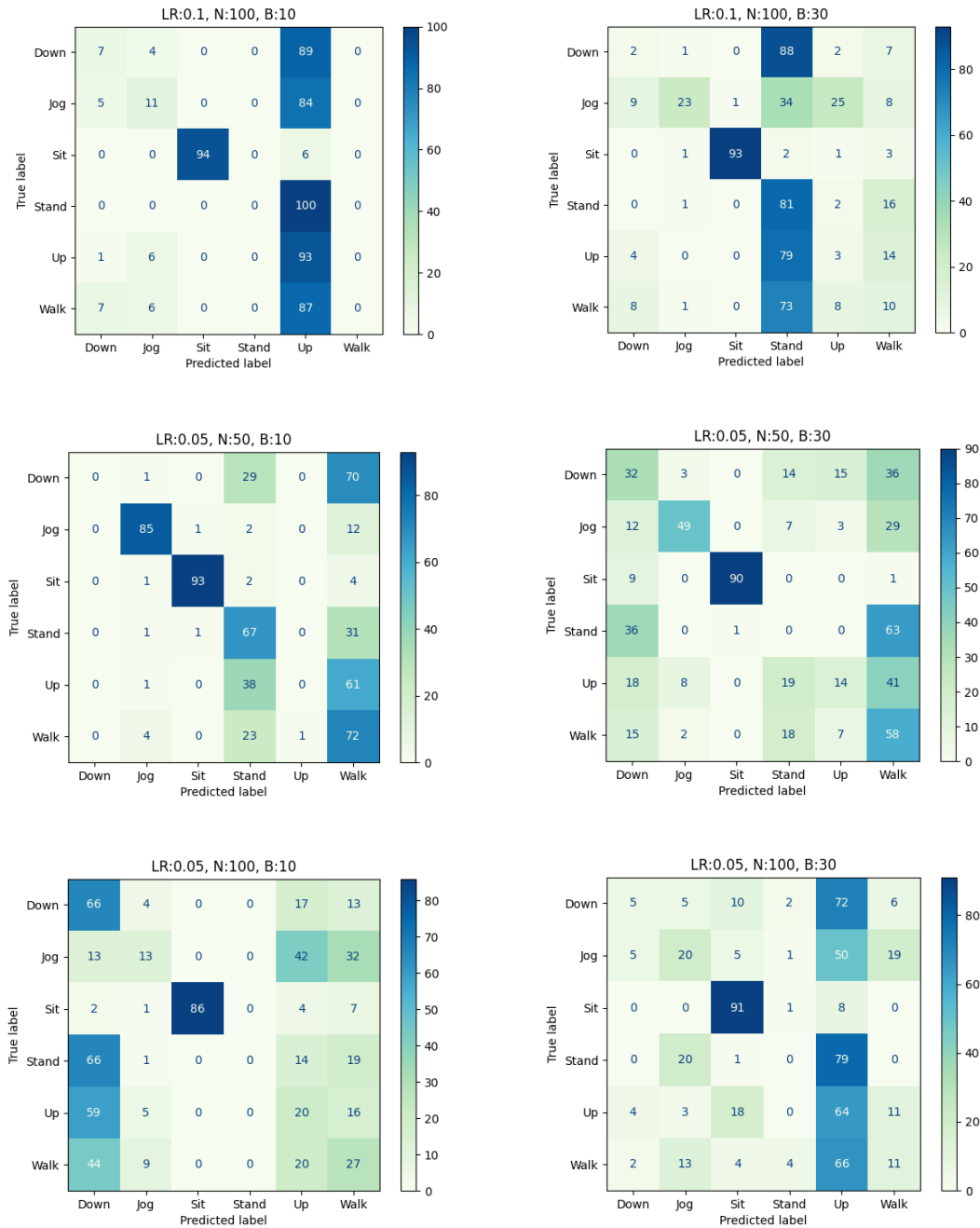


Figure 5. Confusion matrices for all trained 8 models in Part (b). The ideal case is to get a diagonal matrix.

- All models got the sitting motion correct. Notice some models are biased towards a specific motion due to bad training. The best model (LR:0.1, N:50, B:30) is the closest to a diagonal confusion matrix.
- Walking is sometimes mixed up with all other actions except sitting - which makes sense.
- The distribution of the dataset classes (i.e., inherent bias) may play a role in the apparent biases in the confusion matrices as well.

3.2. Early stopping with validation set (c)

Table 2. Top-k scores for the models trained in Part (c). The improvements with early stopping compared to without early stopping are shown in bold, and worse results are shown in italic.

Note that the Top-6 score is always 1.

LR	N	B	Without early stopping (last checkpoint)					With early stopping					Early Stop Epoch
			Top-1	Top-2	Top-3	Top-4	Top-5	Top-1	Top-2	Top-3	Top-4	Top-5	
0.1	50	10	0.358	0.525	0.706	0.872	0.975	0.483	0.657	0.843	0.965	0.987	38
0.1	50	30	0.555	0.732	0.867	0.980	0.995	0.555	0.732	0.867	0.980	0.995	50
0.1	100	10	0.418	0.613	0.730	0.855	0.938	0.418	0.625	0.770	0.915	0.972	12
0.1	100	30	0.450	0.652	0.785	0.912	0.977	0.480	0.677	0.810	0.945	<i>0.973</i>	35
0.05	50	10	0.410	0.675	0.792	0.888	0.915	0.410	0.675	0.792	0.888	0.915	50
0.05	50	30	0.522	0.645	0.783	0.955	0.973	0.522	0.645	0.783	0.955	0.973	50
0.05	100	10	0.490	0.626	0.813	0.913	0.925	<i>0.440</i>	<i>0.622</i>	<i>0.780</i>	0.923	0.978	46
0.05	100	30	0.283	0.480	0.667	0.883	0.980	0.283	0.480	0.667	0.883	0.980	50

Table 3. Top-k scores for the best models in Part (b) and Part (c). Note that the Top-6 score is always 1.

LR	N	B	Best results of Part (b)					Best results of Part (c)				
			Top-1	Top-2	Top-3	Top-4	Top-5	Top-1	Top-2	Top-3	Top-4	Top-5
0.1	50	10	0.410	0.603	0.762	0.928	0.977	0.483	0.657	0.843	0.965	0.987
0.1	50	30	0.582	0.727	0.860	0.962	0.990	0.555	0.732	0.867	0.980	0.995
0.1	100	10	0.342	0.568	0.717	0.865	0.956	0.418	0.625	0.770	0.915	0.972
0.1	100	30	0.353	0.540	0.708	0.838	0.997	0.480	0.677	0.810	0.945	0.973
0.05	50	10	0.529	0.717	0.857	0.960	0.977	0.410	0.675	0.792	0.888	0.915
0.05	50	30	0.405	0.585	0.713	0.896	0.983	0.522	0.645	0.783	0.955	0.973
0.05	100	10	0.353	0.547	0.718	0.940	0.970	0.440	0.622	0.780	0.923	0.978
0.05	100	30	0.318	0.493	0.642	0.753	0.813	0.283	0.480	0.667	0.883	0.980

- **Table 2** reveals the positive effect of early stopping. Since training RNNs are not very stable, in most of the considered hyperparameter combinations the early stopping helps.
- **Table 3** combines the best results of Part (b) and (c). The improvements are thanks to early stopping, and worse results may be due to the decrease in the training samples.

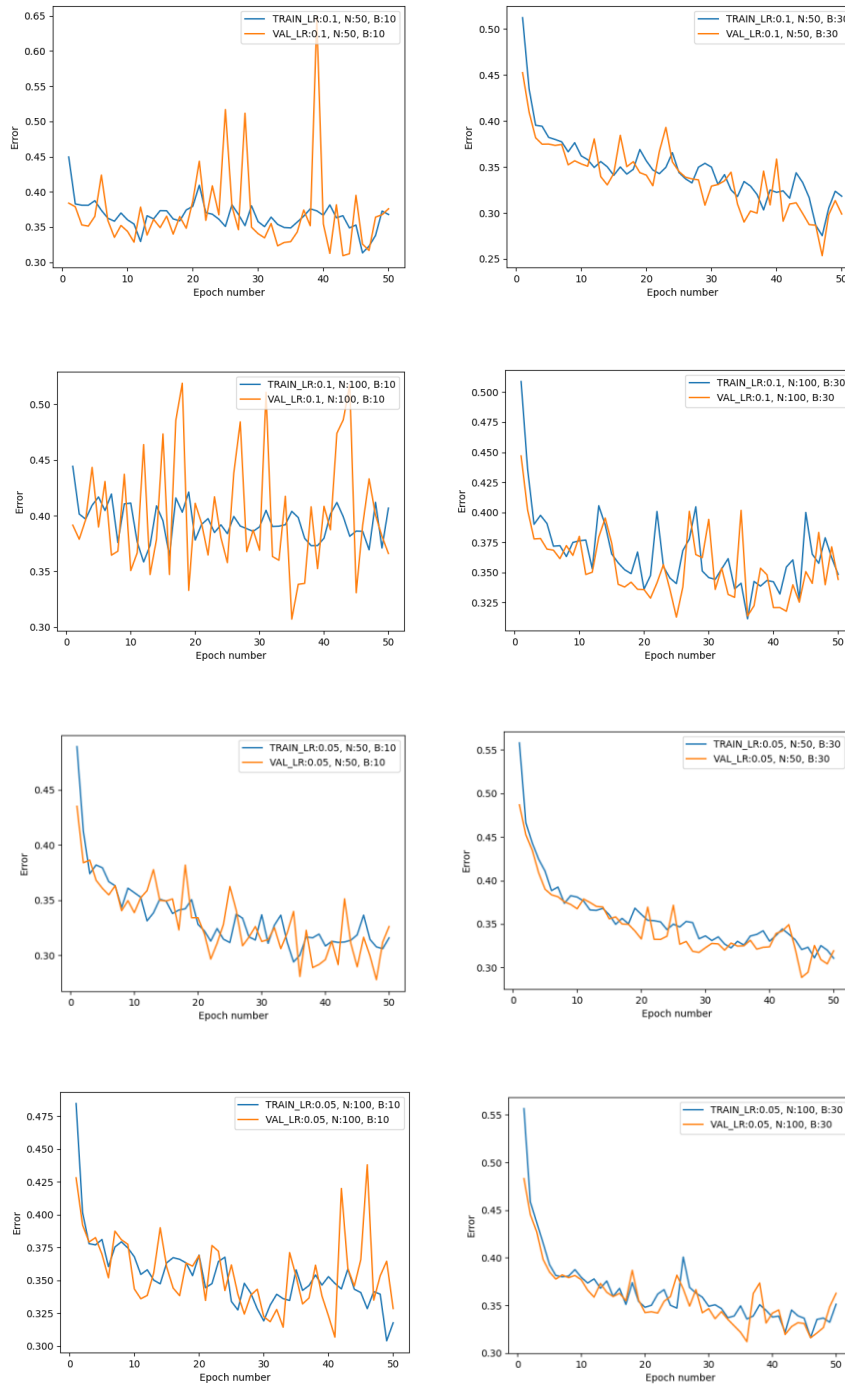
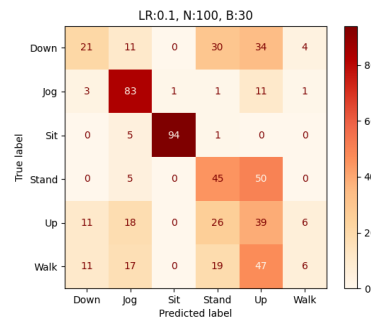
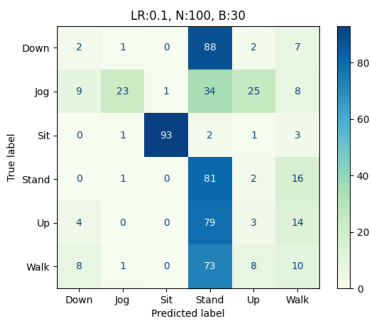
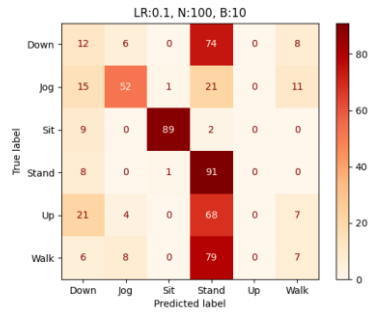
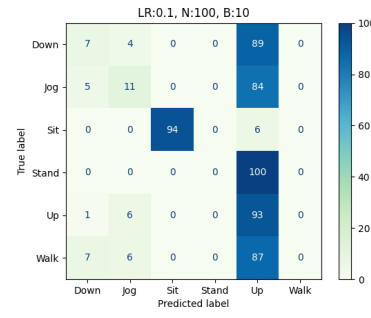
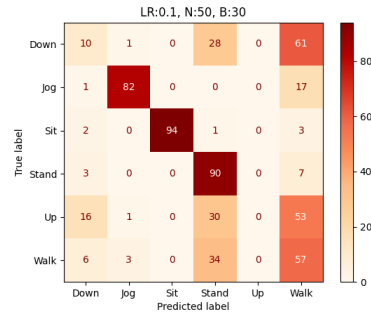
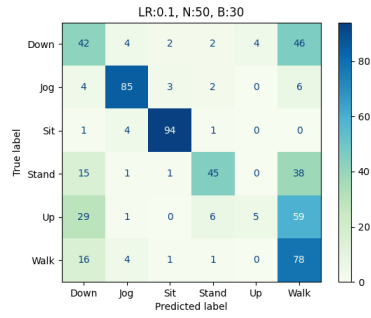
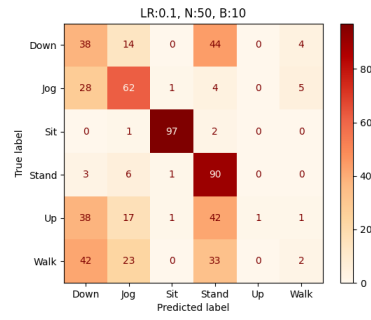
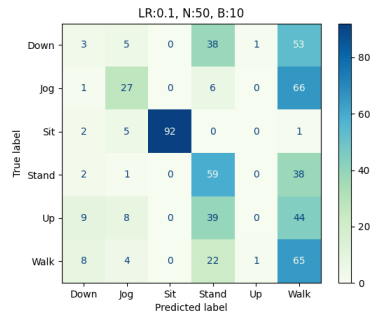


Figure 6. Training loss versus validation loss for all trained models in Part (c).

- Notice the validation versus training loss curve may suggest about the stability of the training besides the overfitting problem (LR:0.1, N:50, B:10 & LR:0.1, N:100, B:10).
- Some models could have been trained for more than 50 epochs since the validation curve does not show any sign of overfitting (LR:0.05, N:50, B:30 & LR:0.05, N:100, B:30).



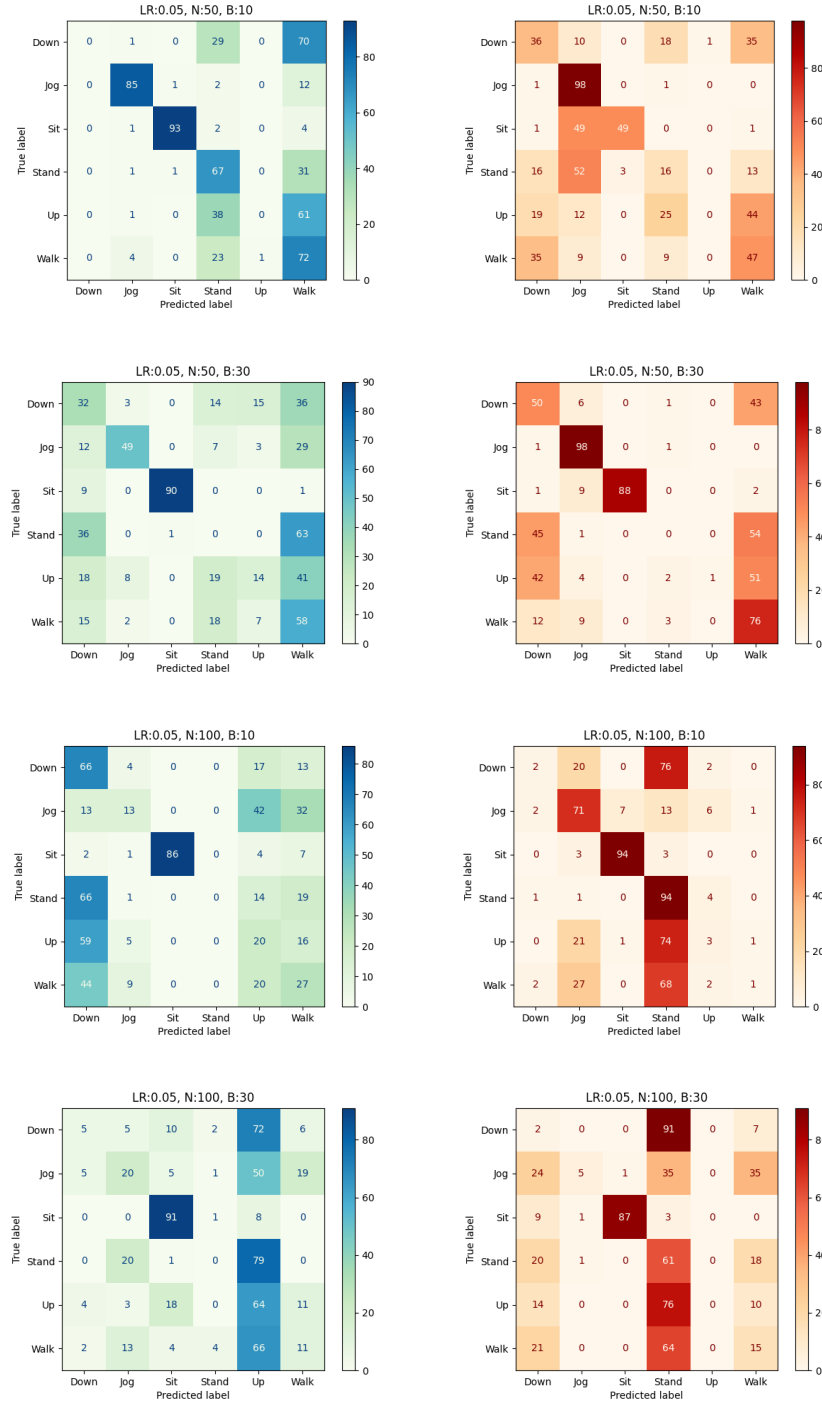


Figure 7. Confusion matrices for all trained 8 models in Part (b) (left) versus the models in Part (c) (right). The ideal case is to get a diagonal matrix.

- Notice how early stopping helps the confusion matrices to “move towards being more diagonal”. This enables the models to be more generalizable, with the probable cost of a lesser accuracy score on the test set, especially on the Top-1 score.