

HW2

Single-layer neural network:

1. Implement a single-layer fully connected neural network model to classify MNIST images. It inputs a raw image as a 1D vector of length $784=28 \times 28$ and outputs a vector of length 10 (each dimension corresponds to a class). You may want to add a bias term to the input, but that is optional for this assignment. The output is connected to the first layer simply by a set of linear weights. The second layer uses Softmax function as the non-linearity function. Softmax is a simple function that converts a n -dimensional input (z) to a n -dimensional output (o) where the output sums to one and each element is a value between 0 and 1. It is defined as

$$y_i = \frac{\exp(x_i)}{\sum_{j=1}^n \exp(x_j)}$$

When we apply this function to the output of the network, o , it predicts a vector which can be seen as the probability of each category given the input x :

$$P(c_i|x) = \frac{\exp(o_i)}{\sum_{j=1}^n \exp(o_j)}$$

where n is the number of categories, 10, in our case. We want the i 'th output to mimic $P(c_i|x)$, the probability of the input x belonging to the category i . We can represent the desired probability distribution as the vector gt where $gt(i)$ is one only if the input is from the i 'th category and zero otherwise. This is called one-hot encoding. Assuming x is from y 'th category, $gt(y)$ is the only element in gt that is equal to one. Then, we want the output probability distribution to be similar to the desired one (ground-truth). Hence, we use cross-entropy loss to compare these two probability distributions, P and gt :

$$L(x, y, w) = \sum_{i=1}^n -gt(i) \log(P(c_i|x))$$

where n is the number of categories. Since gt is one hot encoding, we can remove the terms of gt that are zero, keeping only the y 'th term. Since $gt(y) = 1$, we can remove it in the multiplication to come up with the following loss which is identical to the above one:

$$L(x, y, w) = -\log(P(c_y|x))$$

This is the loss for one input only, so the total loss on a mini-batch is:

$$L = \sum_{k=1}^N -\log(P(c_{y_k}|x_k))$$

where N is the size of mini-batch, number of training data being processed at this iteration.

Please implement stochastic gradient descend (SGD) algorithm from scratch to train the model. You may use NumPy, but should not use PyTorch, TensorFlow, or any similar deep learning framework. Use mini-batch of 10 images per iteration. Then, train it on all MNIST train dataset and plot the accuracy on all test data for every n iteration. Please choose n small enough so that the graph shows the progress of learning and large enough so that testing does not take a lot of time. You may use smaller n initially and then increase it gradually as the learning progresses. Choose a learning rate so that the loss goes down.

Closely adapted from [2]

→ Derivative of softmax and cross entropy loss

$$\text{softmax function} \rightarrow S_i = \frac{e^{z_i}}{\sum_{l=1}^n e^{z_l}}, \quad \forall i = 1, \dots, n$$

As softmax outputs are always positive:

$$\frac{\partial}{\partial z_j} \log(S_i) = \frac{1}{S_i} \cdot \frac{\partial S_i}{\partial z_j}$$

$$\frac{\partial S_i}{\partial z_j} = S_i \cdot \frac{\partial}{\partial z_j} (\log(S_i))$$

$$\log(S_i) = \log\left(\frac{e^{z_i}}{\sum_{l=1}^n e^{z_l}}\right) = z_i - \log\left(\sum_{l=1}^n e^{z_l}\right)$$

$$\frac{\partial}{\partial z_j} \log S_i = \frac{\partial z_i}{\partial z_j} - \frac{\partial}{\partial z_j} \log\left(\sum_{l=1}^n e^{z_l}\right)$$

$$\frac{\partial z_i}{\partial z_j} = \begin{cases} 1, & \text{if } i=j \\ 0, & \text{otherwise} \end{cases}$$

$$\frac{\partial}{\partial z_j} \log S_i = \mathbb{1}\{i=j\} - \frac{1}{\sum_{l=1}^n e^{z_l}} \cdot \left(\frac{\partial}{\partial z_j} \sum_{l=1}^n e^{z_l}\right)$$

$$\frac{\partial}{\partial z_j} \log S_i = \mathbb{1}\{i=j\} - S_j$$

$$\frac{\partial S_i}{\partial z_j} = S_i \cdot (\mathbb{1}\{i=j\} - S_j) \text{ ————— (i)}$$

$$\text{cross entropy loss function} \rightarrow L(\gamma, s) = - \sum_{i=1}^C \gamma_i \cdot \log(s_i)$$

$$\frac{\partial L}{\partial z_j} = - \frac{\partial}{\partial z_j} \sum_{i=1}^C \gamma_i \cdot \log(s_i) = - \sum_{i=1}^C \frac{\gamma_i}{s_i} \cdot \frac{\partial s_i}{\partial z_j} \quad \text{--- (ii)}$$

Replacing (i) in (ii)

$$\frac{\partial L}{\partial z_j} = - \sum_{i=1}^C \frac{\gamma_i}{s_i} \cdot s_i (1_{\{i=j\}} - s_j) = - \sum_{i=1}^C \gamma_i \cdot (1_{\{i=j\}} - s_j)$$

$$\Rightarrow \sum_{i=1}^C \gamma_i \cdot s_j - \sum_{i=1}^C \gamma_i \cdot 1_{\{i=j\}}$$

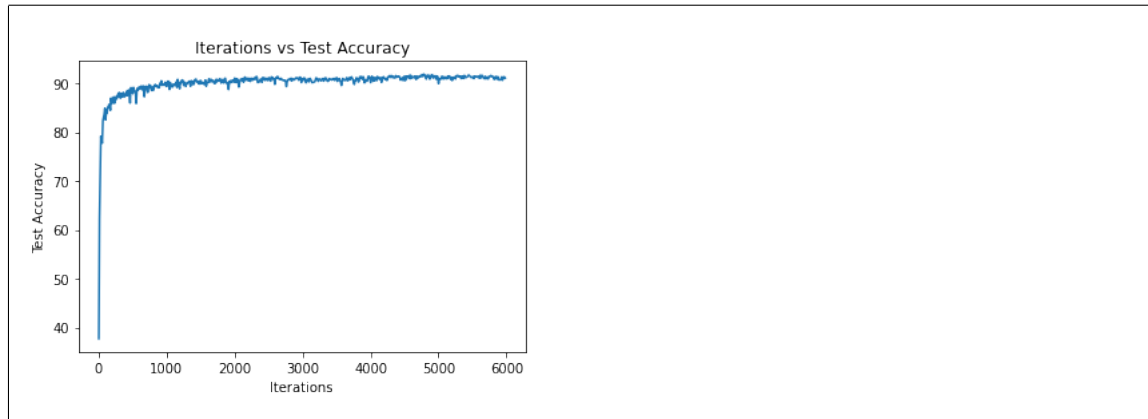
$$\Rightarrow \sum_{i=1}^C \gamma_i \cdot s_j - \gamma_j = s_j \sum_{i=1}^C \gamma_i - \gamma_j = s_j - \gamma_j$$

$$\therefore \frac{\partial L}{\partial z} = s - \gamma \quad \text{--- (iii)}$$

we need to find $\frac{\partial L}{\partial w}$, we know that $z = w^T x$

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial z} \times \frac{\partial z}{\partial w} = (s - \gamma) \times \frac{\partial}{\partial w} (w^T x)$$

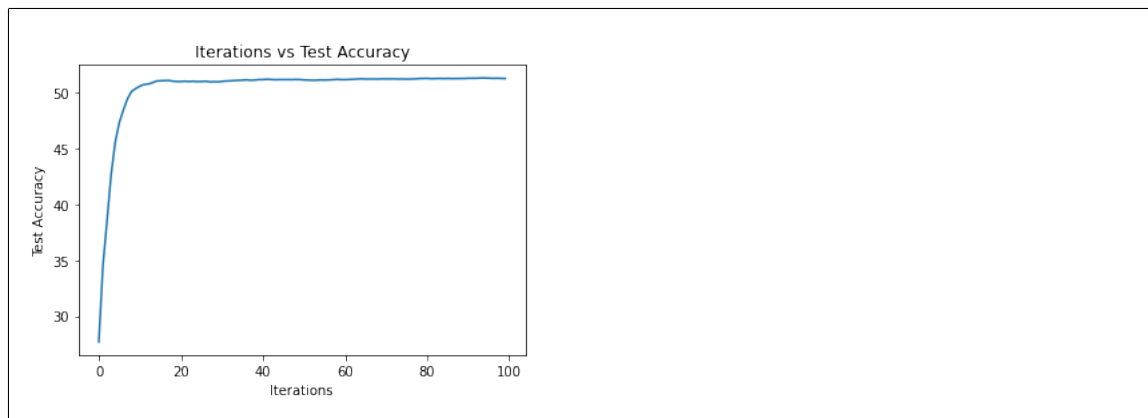
$$\boxed{\frac{\partial L}{\partial w} = (s - \gamma) \times x} \quad \text{--- (iv)}$$



- For each class, visualize the 10 images that are misclassified with the highest score along with their predicted label and score. These are very confident wrong predictions.

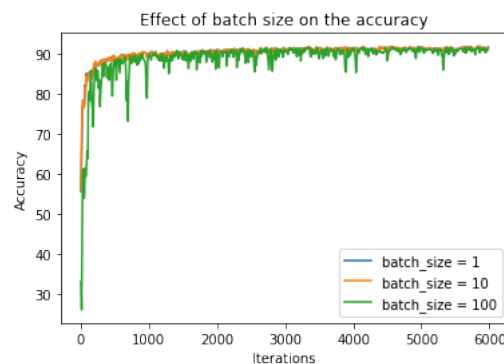
Figure present in the .ipynb file

- Please reduce the number of training data to 1 example per class (chosen randomly from training data) and plot the curve (accuracy vs, iterations). The whole training data will be 10 images only.



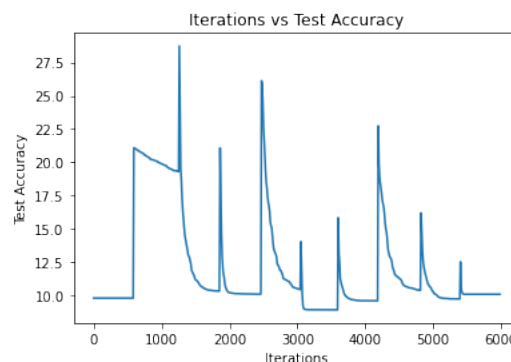
- Try different mini-batch sizes (1, 10, 100) for the original case and plot the results. Which one is better and why?

In comparison to the models with batch sizes of 10 and 100, the model with a mini-batch size of 1 had the highest final accuracy. According to [1], models with smaller batch sizes typically outperform those with bigger batch sizes. The study goes on to explain why this is possible. Big batch models have significantly lower generalization ability when compared to small batch models; this difference is exacerbated by the fact that large batch models converge to *sharp minimizers*, whereas small batch models converge to *flat minimizers*. A *flat minimizer* is one in which the function fluctuates gradually across a large region. A *sharp minimizer*, on the other hand, allows the function to rise quickly in a small region.

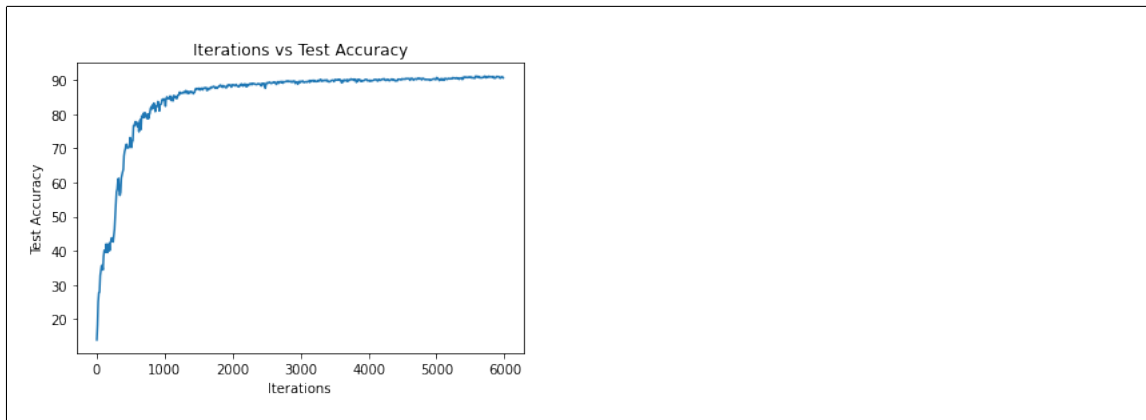


5. Instead of using random sampling, sort the data before training so that all "1"s appear before "2"s and so on. Then, sample sequentially in running SGD instead of random sampling. Does this work well, why?

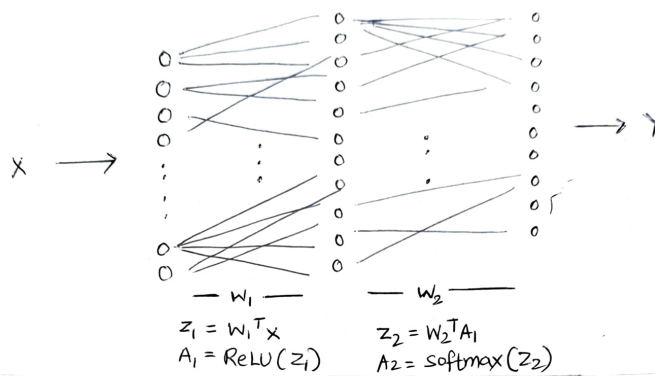
As seen in the graph, sorting the training data before training the model does not work well. The input to the model is intended to provide the model with an overall representation of the present data. By sorting the data, the model assumes that the only type of input is the class it has seen so far. During testing, the model is more inclined to predict the image it has seen the most from input rather than what the image actually is. Another issue is that because each image processed by the model is of the same type until a specific class is exhausted, the model will be unable to generalize the data and will also be unable to discover patterns within the data. Random sampling overcomes these problems since the model now understands what sorts of inputs to expect and the nuances with which it must adjust its weights in order to categorize a specific image as another class.



6. (Bonus point) Add a hidden layer with 11 hidden neurons and ReLU activation function. Then, plot the accuracy curve to see how the accuracy changes.



→ Derivation of $\frac{\partial L}{\partial w_1}, \frac{\partial L}{\partial w_2}$



$$w_2 = w_2 - \alpha \cdot \frac{\partial L}{\partial w_2}$$

— From (iv) —

$$\frac{\partial L}{\partial w} = (s - y) \times X$$

— Replacing the terms —

$$\frac{\partial L}{\partial w_2} = (A_2 - y) \times A_1$$

$$\Rightarrow w_2 = w_2 - \alpha [(A_2 - y) \times A_1]$$

$$w_1 = w_1 - \alpha \cdot \frac{\partial L}{\partial w_1}$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial z_2} \times \frac{\partial z_2}{\partial A_1} \times \frac{\partial A_1}{\partial z_1} \times \frac{\partial z_1}{\partial w_1}$$

— From (iii) —

$$\frac{\partial L}{\partial z} = (s - y)$$

— Replacing the terms —

$$\frac{\partial L}{\partial z_2} = A_2 - y \quad \frac{\partial z_2}{\partial A_1} = \frac{\partial w_2^T A_1}{\partial A_1} = w_2^T$$

$$\frac{\partial A_1}{\partial z_1} = \frac{\partial}{\partial z_1} (\text{ReLU}(z_1)) = \begin{cases} 1, & \text{if } z_1 > 0 \\ 0, & \text{otherwise} \end{cases}$$

$$\frac{\partial z_1}{\partial w_1} = \frac{\partial}{\partial w_1} (w_1^T X) = X$$

$$\therefore \frac{\partial L}{\partial w_1} = (A_2 - y) \times (w_2^T) \times (z_1 > 0) \times X$$

$$\Rightarrow w_1 = w_1 - \alpha (A_2 - y) \times w_2^T \times (z_1 > 0) \times X$$

References

- [1] Nitish Shirish Keskar et al. “On Large-Batch Training for Deep Learning: Generalization Gap and Sharp Minima”. In: *CoRR* abs/1609.04836 (2016). arXiv: 1609.04836. URL: <http://arxiv.org/abs/1609.04836>.
- [2] Thomas Kurbiel. *Derivative of the Softmax function and the categorical cross-entropy loss*. URL: <https://towardsdatascience.com/derivative-of-the-softmax-function-and-the-categorical-cross-entropy-loss-ffceefc081d1>.