

Report for MIDAS Internship Task 2

Animesh Gupta

May 11, 2021

1 Problem Statement

Train an image classification model for classifying digits and alphabets under 3 given subtasks.

1. For the first subtask, given a dataset containing digits and alphabets (capital and small), a total of 62 classes (Fig. 1), train a neural network from scratch with random weight initialization without including any external data or pretrained weights during the training process.
2. For the second subtask, select images from the given dataset containing only digits (0-9) as labels. Compare the performance of 2 neural networks (mentioned below in detail) on the dataset containing only digits (0-9) as labels.
 - 2.1. Train on the MNIST train set a network from scratch with random weight initialization and evaluate on the MNIST test set.
 - 2.2. Train on the dataset containing only digits (0-9) as labels then save the best network. Load the best-saved network weights to then train on the MNIST train set and evaluate on the MNIST test set.
3. For the third subtask, given a new dataset containing only digits as labels but each label contains images having all the digits from 0-9. For eg, a folder with the label 0 contains images with all the digits from 0-9. (Fig. 12) Compare the performance of 2 neural networks (mentioned below in detail) on the given new dataset.
 - 3.1. Train using the given new dataset from scratch with random weight initialization and evaluate on MNIST test set.
 - 3.2. Load the best-saved network weights from subtask 1 to then train on the MNIST train set and evaluate on the MNIST test.

A well-documented source code is available on https://github.com/animesh-007/midas_internship_task2

2 Subtask 1

2.1 Implementation details

Given a dataset with images of digits and alphabets (capital and small).

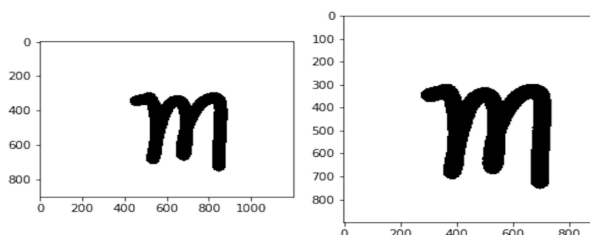
Figure 1: An overview of the dataset.



The following transformations were applied to the dataset before using it as an input to train the network.

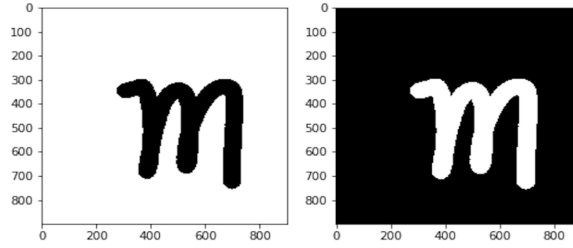
1. **Center Crop:** In the given dataset, the images are of size H=900 and W=1200, and as the images contain objects in the center, object detail is not lost. Centercrop (900,900) was applied to scale the image uniformly thus maintaining the aspect ratio. (Fig. 2)

Figure 2: A comparison between the original image and centercrop image.



2. **Image Inversion:** To maintain the same data distribution as MNIST, Invert was applied to invert the colours of the images. (Fig. 3)

Figure 3: A comparison between the original image and inversion image.



3. **Grayscale Conversion:** To maintain the same data distribution as MNIST, Grayscale was applied to convert 3-channel images to 1-channel images.
4. **Resize Images:** To train the model efficiently. The images were resized to H=28 and W=28.
5. **Normalize Images:** The training images are then converted to tensors and normalized by the mean (0.1307) and standard deviation (0.3081) to help the network get data within a range and reduce the skewness which helps learn faster and better.

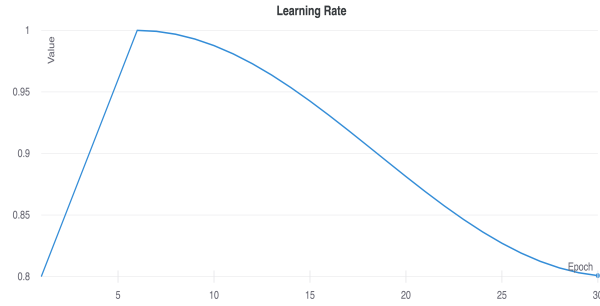
2.2 Network and Training Details

- We divide the given dataset into training and validation set in 80:20 split to cross-validate the network.
- The network consists of 2 convolutions layers with input channels [1,32] and output channels [32,64] respectively. Each followed by a relu activation function. The maxpool pooling layer was applied after the output of the second convolution layer. Dropout was applied with a value of 0.25. The output is then flattened and a fully connected layer was applied with input features 9216 and output features 128. Dropout was applied with a value of 0.5. The fully connected layer with input features 128 and output features 62 as the final layer of the model.
- We trained the network for 30 epochs with a learning rate of 1 and optimizer adadelata.
- We choose this network as it's faster+efficient for datasets like MNIST because training with deeper networks will saturate the performance.

Motivation for using Scheduler: Learning Rate (LR) Scheduler helps to adjust the LR during training. Models often benefit from this technique once learning stagnates, and helps to get better results.

- We used CosineAnnealingLR Scheduler: Min LR 0.8, Max LR 1, and warm-up 5 epochs. (Fig. 4)
- CosineAnnealingLR Scheduler increased the performance of the network as reported in Table 1.
- The trained network is cross-validated on the validation set (20% of the full dataset) to select the best model and hyperparameters.

Figure 4: Learning Rate using CosineAnnealingLR



2.3 Results

As seen in table 1, the performance is increased by 1.49% with the use of CosineAnnealingLR scheduler.

Method	Epochs	Accuracy
CNN without a scheduler.	30	67.74
CNN with a CosineAnnealingLR scheduler.	30	68.75

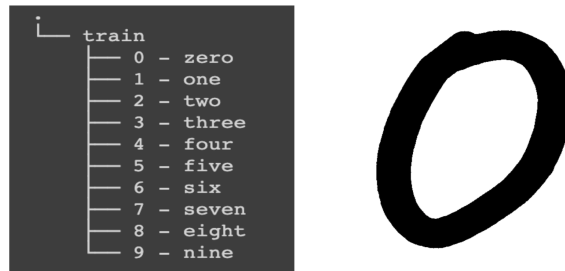
Table 1: Quantitive result comparing the scheduler, the performance is increased by **1.49%**

3 Subtask 2

3.1 Implementation details

Make a subset of the above dataset containing only digits labels.

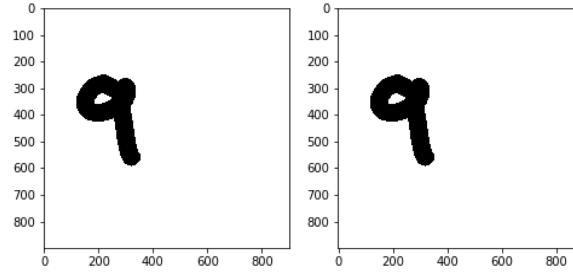
Figure 5: Overview of the directory and image inside the dataset.



The following transformations were applied to the dataset before using it as an input to train the network.

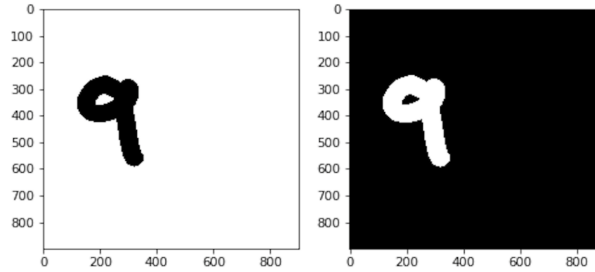
1. **Center Crop:** In the given dataset, the images are of size H=900 and W=1200, and most of the images contain objects in the center, Centercrop (900,900) was applied to scale the image uniformly thus maintaining the aspect ratio. (Fig. 6)

Figure 6: A comparison between the original image and centercrop image.



2. **Image Inversion:** To maintain the same data distribution as MNIST, Invert was applied to invert the colours of the images. (Fig. 7)

Figure 7: A comparison between the original image and inversion image.



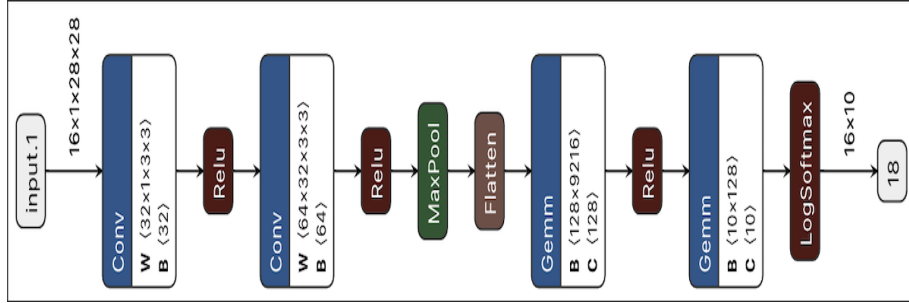
3. **Grayscale Conversion:** To maintain the same data distribution as MNIST, Grayscale was applied to convert 3-channel images to 1-channel images.
4. **Resize Images:** To train the model efficiently. The images were resized to H=28 and W=28.
5. **Normalize Images:** The training images are then converted to tensors and normalized by the mean (0.1307) and standard deviation (0.3081) to help the network get data within a range and reduce the skewness which helps learn faster and better.

3.2 Network and Training details

- The network consists of 2 convolutions layers with input channels [1,32] and output channels [32,64] respectively. Each followed by a relu activation function. The maxpool pooling layer was applied after the output of the second convolution layer. Dropout was applied with a value of 0.25. The output is then flattened and a fully connected layer was applied with input features 9216 and output features 128. Dropout was applied with a value of 0.5. The fully connected layer with input features 128 and output features 10 as the final layer of the model. (Fig. 8)
- The network is trained for 30 epochs with a learning rate of 1 and optimizer adadelata.
- We used CosineAnnealingLR Scheduler: Min LR 0.8, Max LR 1, and warm-up 5 epochs.

- CosineAnnealingLR Scheduler increased the performance of the network as reported in Table 2.
- The trained network is tested on the MNIST test set to select the best model and hyperparameters.

Figure 8: An Overview of the model architecture.



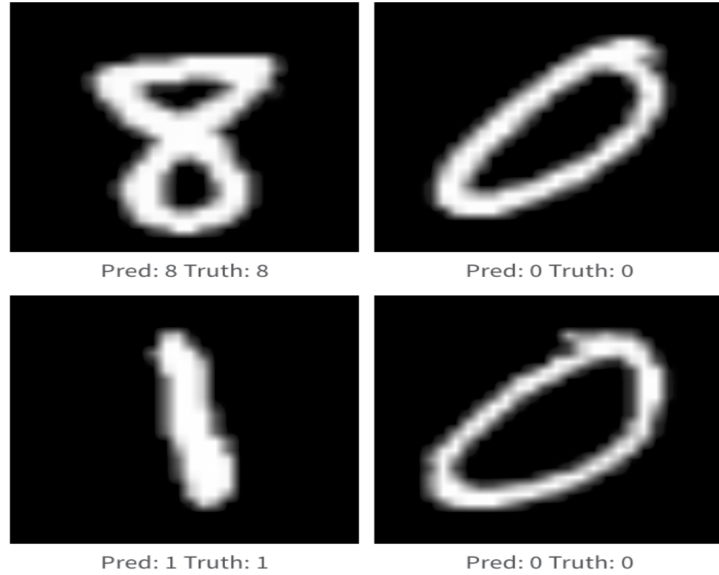
3.2.1 Training Strategy

- Making a subset of the above dataset containing only images of digits.
- Training the model on the processed dataset for 30 epochs with a learning rate of 1.

3.3 Results

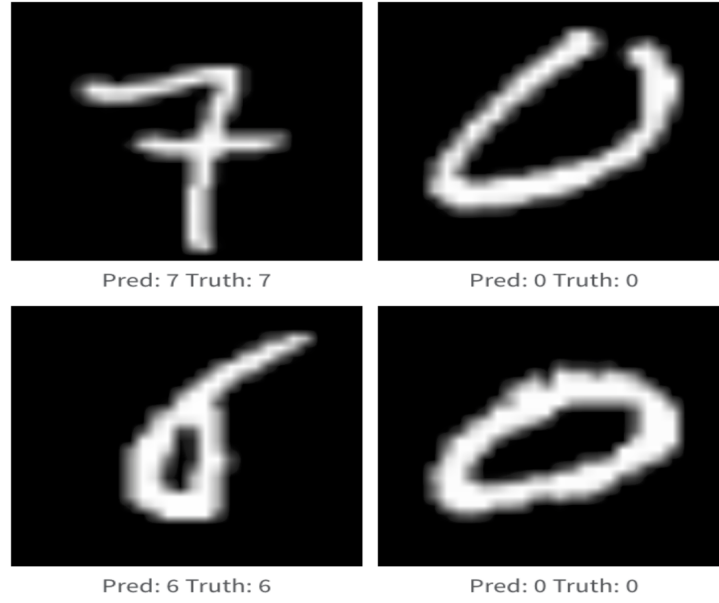
3.3.1 Training on the MNIST dataset with random weights

Figure 9: Prediction and Ground truth by training the model with random weights.



3.3.2 Train on the MNIST dataset with pretrained weights

Figure 10: Prediction and Ground truth by training the model with pretrained weights



Method	Epochs	Accuracy
CNN on MIDAS dataset containing only digits, with a CosineAnnealingLR scheduler.	30	66.56
CNN on MNIST dataset with random weights, with a CosineAnnealingLR scheduler.	30	99.39
CNN on MNIST dataset with pretrained weights, with a CosineAnnealingLR scheduler.	30	99.34

Table 2: Results after using the CosineAnnealingLR scheduler.

4 Subtask 3

4.1 Implementation details

Given a dataset with images of digits.

Figure 11: The directory structure of the folder.

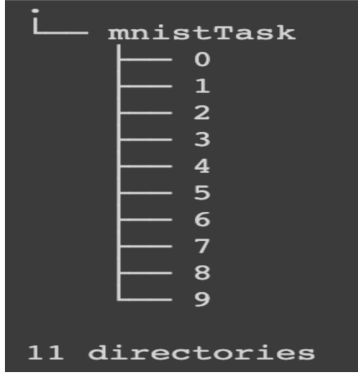


Figure 12: Images inside the folder of "0-zero".



The following transformations were applied to the dataset before using it as an input to train the network.

1. **Grayscale Conversion:** To maintain the same data distribution as MNIST, Grayscale was applied to convert 3-channel images to 1-channel images.
2. **Resize Images:** To train the model efficiently. The images were resized to H=28 and W=28.
3. **Normalize Images:** The training images are then converted to tensors and normalized by the mean (0.1307) and standard deviation (0.3081) to help the network get data within a range and reduce the skewness which helps learn faster and better.

4.1.1 Network details

- The network consists of 2 convolutions layers with input channels [1,32] and output channels [32,64] respectively. Each followed by a relu activation function. The maxpool pooling layer was applied after the output of the second convolution layer. Dropout was applied with a value of 0.25. The output is then flattened and a fully connected layer was applied with input features 9216 and output features 128. Dropout was applied with a value of 0.5. The fully connected layer with input features 128 and output features 10 as the final layer of the model.
- The network is trained for 30 epochs with a learning rate of 1 and optimizer adadelata.
- We used CosineAnnealingLR Scheduler: Min LR 0.8, Max LR 1, and warm-up 5.
- The trained network is tested on the MNIST test set to select the best model and hyperparameters.

4.1.2 Training Strategy

- Training the model on the given dataset.
- Training the model on the processed dataset for 30 epochs with a learning rate of 1

4.2 Results

Method	Epochs	Accuracy
CNN on MIDAS Dataset with random weights.	30	1.74
CNN on MIDAS Dataset with pretrained weights of MIDAS dataset containing only digits.	30	10.32

Table 3: Results of the model after training on the dataset.

4.2.1 Analysis

- After visualizing the images inside the dataset, I found that images inside the folder are present in a particular manner i.e images of "0-zero" are very less present inside the folder "0".
- This made me clear that the labels of the images are shuffled which results in poor accuracy on the MNIST test after training the model on the dataset.

4.3 Future Goals

- Load the best-saved network weights from subtask 2.2, then train on the MNIST train set and evaluate on the MNIST test.
- Train a CNN with 224 by 224 as an input image for subtask 1.

5 References

I want to give shoutouts to these works which enabled me to compile this report with various code examples and experimental results:

- [Training a Classifier](#) by Pytorch.
- [Transfer Learning For Computer Vision Tutorial](#) by Pytorch.
- [Pytorch MNIST example](#) by Pytorch.
- [Weights & Biases Documentation](#) by Weights & Biases.
- [CosineAnnealingLR](#)