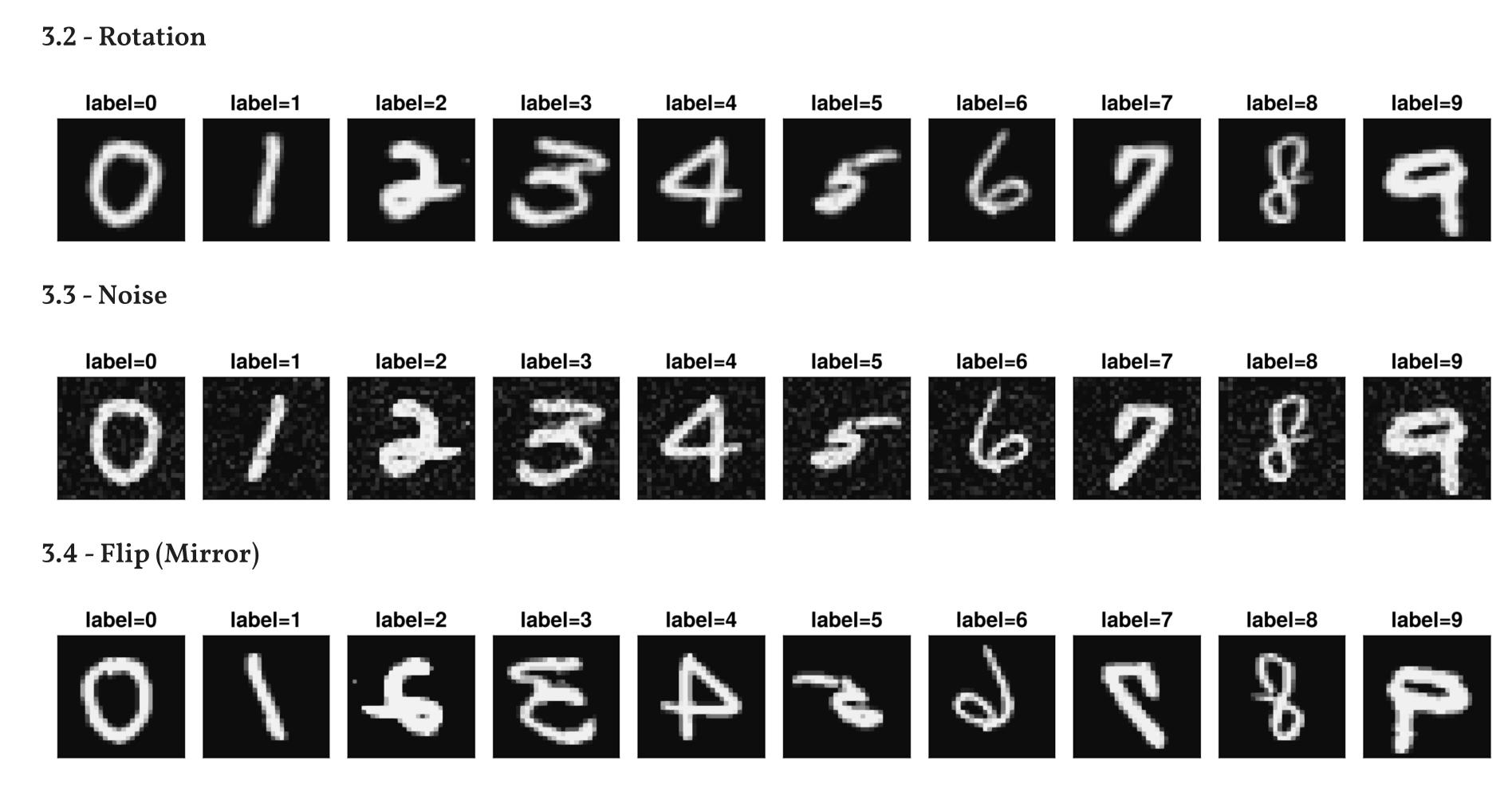
Pluto.jl -Data-Augmentation/FrontendPluto.jl **≡** Table of Contents ? project at `~/GITHUB/Scientific-Programming-Handwritten-Digit-Recognition-Data-Augmentation` **Digit Recognition - Data Augmentatio...** 1 - Introduction 2 - Motivation **Digit Recognition - Data Augmentation** 3 - Data Augmentation Methods 3.1 - No Augmentation 3.2 - Rotation 3.3 - Noise 3.4 - Flip (Mirror) 1 - Introduction 3.5 - Full Augmentation w/o Flip 4 - Trained Models 5 - Evaluation Digit recognition enables machines to interpret and process handwritten or printed numerical data, bridging human input with digital systems. 5.1 - Tables 5.2 - Fully Trained Model 5.3 - Non-Augmented Model Data augmentation helps in training for digit recognition by artificially expanding the training dataset through transformations like rotation or noise addition, 5.4 - Fully-Augmented Model allowing the model to generalize better for real-world usage. 5.5 - Only Rotation Model 5.6 - Only Noise Model 5.7 - Only Flip (Mirror) Model 2 - Motivation 6 - Conclusion 7 - Further Information 7.1 - Data Handling In this project, we train and test models implemented with a convolutional neural network (CNN) architecture based on 'LeNet-5'. 7.2 - LeNet5 7.3 - Training We first train a model on most¹ of the full MNIST dataset (60.000 images) and test it on the MNIST test set (10.000 images). This represents our 'base-case', and its values serve as guide-values. In real-world applications, usable data is often limited, and research time and computation power are constrained. Therefore our goal is to apply various augmentation methods to achieve comparable results using only fractions of the training data in a fraction of the time and processing power. ¹More on this in section 4 label=1 label=2 label=3 label=4 label=5 label=6 label=7 label=8 label=9

3 - Data Augmentation Methods 3.1 - No Augmentation label=0



4 - Trained Models

label=0

3.5 - Full Augmentation w/o Flip

label=1

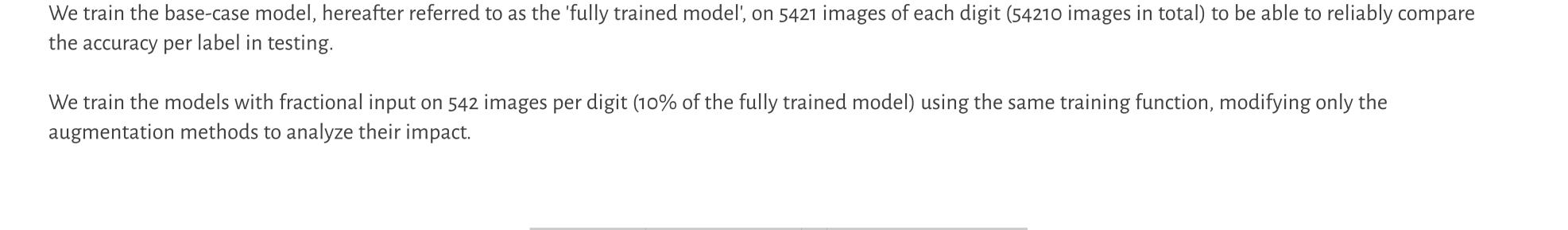
label=2

label=3

Fully Trained

Fractional

label=4



32 0 0.0003

32 0 0.0003

Epochs Batchsize λ

label=5

label=6

η Training Set Size

54210

5420

label=7

∆ Max

96.52% 96.45% 96.79% 94.25% 98.02% 96.98%

5.64% 83.06% 28.64% 33.26%

96.69%

~90s

~18s

98.96%

97.87% 96.66%

Time per Epoch Total Time

~4.5s

~0.45s

2.8% 8.66%

2.79

5.7

3.74

4.37

3.37

81.91

label=8

label=9

trained for 0.501s | reached 0.026 loss with 99.613% accuracy reached 0.017 loss with 31.089% accuracy reached 0.014 loss with 100.0% accuracy reached 0.015 loss with 99.963% accuracy

5 - Evaluation

rained for 0.461s | reached 0.014 loss with 99.815% accuracy

trained for 0.498s | reached 0.025 loss with 99.779% accuracy

trained for 0.524s | reached 0.012 loss with 99.723% accuracy trained for 0.516s | reached 0.013 loss with 99.963% accuracy

trained for 0.498s | reached 0.025 loss with 99.797% accuracy trained for 0.48s | reached 0.016 loss with 30.664% accuracy trained for 0.518s | reached 0.012 loss with 98.579% accuracy rained for 0.514s | reached 0.013 loss with 100.0% accuracy rained for 0.461s | reached 0.011 loss with 99.852% accuracy

trained for 0.497s | reached 0.022 loss with 99.557% accuracy trained for 0.479s | reached 0.016 loss with 32.362% accuracy trained for 0.522s | reached 0.012 loss with 100.0% accuracy trained for 0.523s | reached 0.012 loss with 99.871% accuracy rained for 0.465s | reached 0.011 loss with 99.41% accuracy

Fully Trained

No Augmentation

Full Augmentation

Only Rotation

Only Noise

Only Flip

99.69%

39.59% 63.79% 11.43%

reached 0.013 loss with 99.963% accuracy

rained for 0.523s | reached 0.014 loss with 98.838% accuracy trained for 0.498s | reached 0.021 loss with 99.373% accuracy trained for 0.484s | reached 0.014 loss with 31.808% accuracy

For evaluation, we measure the model's accuracy (total and per digit) on the test dataset and compare performance across different augmentation methods. This allows us to assess which methods most effectively compensate for the reduced training data and contribute to model accuracy.

After testing with 0.1%, 1%, 10%, 25%, and 50% of the training set, we chose to focus on the 10% subset. It provides a good balance between efficiency and performance, achieving more stable and reliable results (~95% accuracy) compared to 1% (around 60–70%). At the same time, it remains distinguishable from the fully trained model, which reaches about 99% accuracy. Prediction Table: Shows the total accuracy, accuracy per digit and maximum difference of the 'accuracy per digit' Confusion Matrix: Shows the counts of true and predicted classifications for each digit, highlighting where the model gets predictions right or wrong **5.1 - Tables** 5.1.1 - Predictions

98.12% 93.18%

97.33%

1.39%

Fully Trained

Fractional (each)

1000

- 1000

- 800

- 600

- 400

- 200

10

98.68%

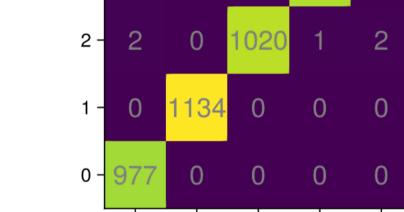
96.13%

83.3%

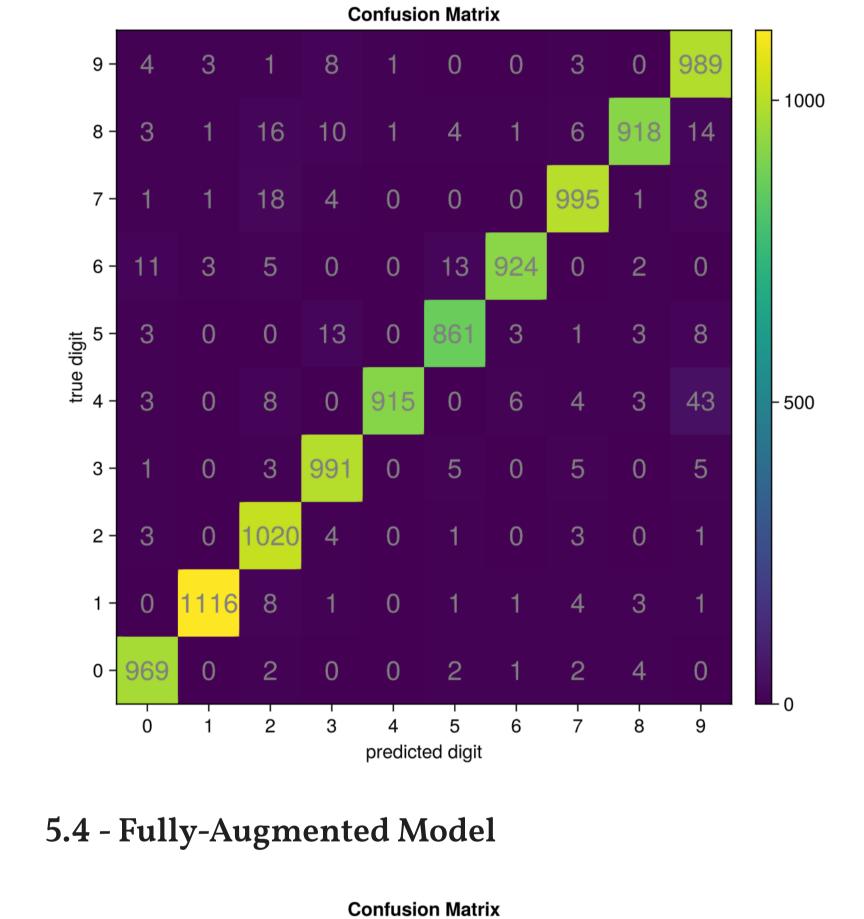
5.1.2 - Training Times

| 5.2 - Fully Trained Model | | | | | | | | | | | |
|---------------------------|------------------|---|---|---|---|---|---|------|-----|-----|--|
| | Confusion Matrix | | | | | | | | | | |
| 9 - | 1 | 1 | 1 | 2 | 8 | 3 | 2 | 5 | 6 | 980 | |
| 8 | 3 | 0 | 3 | 2 | 2 | 1 | 0 | 2 | 955 | 6 | |
| 7. | 0 | 3 | 9 | 0 | 0 | 0 | 0 | 1020 | 1 | 9 | |

Trained on M1 Pro, 16GB Ram



5.3 - Non-Augmented Model



predicted digit

predicted digit 5.5 - Only Rotation Model

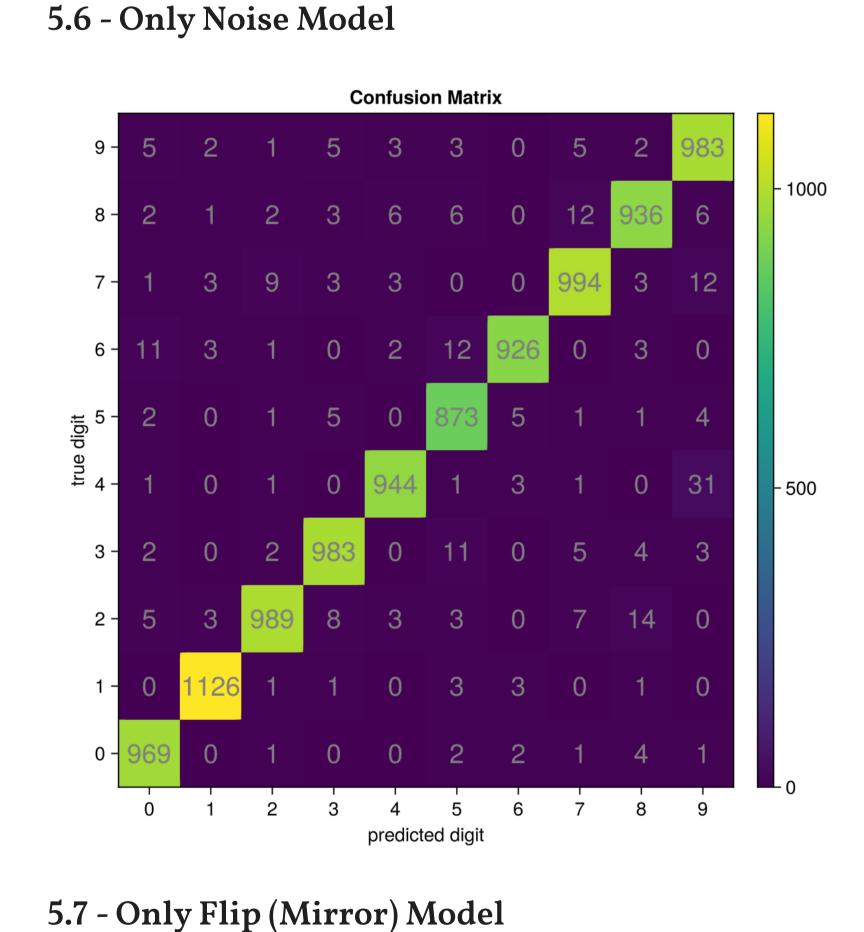
3 -

2 -

Confusion Matrix

0

14



Confusion Matrix

30

58

20

58

18

22

LeNet was originally designed to process 32×32 grayscale images; we adapt it here to handle MNIST's 28×28 images.

predicted digit

evaluation and multiple repetitions when working with limited data.

0

predicted digit

25 3

6 - Conclusion

trained on the full dataset.

7 - Further Information

7.1 - Data Handling

architecture.

the number of images.

45

13

2 -

33

14

The results of our models could improve if we limit specific augmentation methods to specific digits (e.g. to only use '0' and '8' in our flip (mirror) method), because these digits remain visually similar when mirrored, whereas others (like '2' or '5') can resemble different digits, potentially confusing the model. Additionally, even more augmentation methods, training for multiple runs to test the stability of the accuracies, and futher improving the models and training algorithms may improve these results.

We only tested on the MNIST dataset, which limits our ability to assess generalization, as it does not reflect the diversity and complexity of real-world digit

recognition tasks. Future work could invole testing on more challenging data sets with custom handwritten digits, to better evaluate model robustness.

We found that augmenting the dataset with rotation and noise, individually and in conjunction, helps stabilize the results and reduce variance between runs on

On some runs the non augmented model achieved higher accuracy and less variance than the augmented models. This highlights the importance of consistent

These augmentations should improve the model's ability to generalize from a smaller dataset, making performance more consistent and closer to that of models

average. Still, with only 10% of the training data, the model's performance is more sensitive to randomness. This leads to greater variability in accuracy across runs.

The labels are represented using one-hot encoding: a binary vector of length 10, where the position of the 1 indicates the corresponding digit class (e.g., [0, 0, 0, 0, 0, 0, 0, 0, 0, 0] represents the digit 3). This encoding is a standard format used for classification tasks in machine learning. 7.2 - LeNet5

 $f: \mathbb{R}^{28 \times 28 \times 1} \to \mathbb{R}^{10}$, with the whole function being:

Each image is stored as a tensor of shape (28, 28, 1) — the third dimension representing the grayscale feature map. The dataset is also batched, so the full training set has the shape (28, 28, 1, N), where N is

MNIST dataset: A collection of 28×28 grayscale images of handwritten digits (0–9). Using the Julia package MLDatasets.jl, the data is preprocessed to match the input format expected by the LeNet

Our implementation consists of two convolutional-pooling blocks, followed by three fully connected layers. Each hidden layer uses the ReLU (Rectified Linear Unit) activation function, which introduces nonlinearity and helps the network learn complex patterns by zeroing out negative values. The final output layer uses the softmax function to convert the network's raw outputs into a probability distribution over the 10 digit classes.

 $\mathbb{R}^{28 \times 28 \times 1} \xrightarrow{\text{Conv 1}} \mathbb{R}^{24 \times 24 \times 6} \xrightarrow{\text{Pool 1}} \mathbb{R}^{12 \times 12 \times 6} \xrightarrow{\text{Conv 2}} \mathbb{R}^{8 \times 8 \times 16} \xrightarrow{\text{Pool 2}} \mathbb{R}^{4 \times 4 \times 16} \xrightarrow{\text{Dense}} \mathbb{R}^{256} \xrightarrow{\text{Dense}} \mathbb{R}^{120} \xrightarrow{\text{Dense}} \mathbb{R}^{84} \xrightarrow{\text{Softmax}} \mathbb{R}^{10}$ where the input matrix $\mathbb{R}^{28 \times 28 \times 1}$ is the MNIST image (28 × 28), and the output vector \mathbb{R}^{10} represents the class probabilities for digits 0 through 9 calculated from the last softmax step. 7.3 - Training

ground truth label. The cross-entropy loss over a batch $\mathbb N$ of samples is defined as: $-rac{1}{N}\sum_{i=1}^N y_i \cdot \log(\hat{y}_i)$

Cross-entropy loss function: Commonly used for multi-class classification problems. Let $\hat{\mathbf{y}}_i$ denote the predicted probability distribution for the i-th sample, and \mathbf{y}_i the corresponding one-hot encoded