

EMPIRICAL COMPARISON OF THREE ADAPTIVE MOMENT ESTIMATION METHODS

Elias Bughsn

MSc in Data Science & Artificial Intelligence
 Université Côte d'Azur
 eliasbughns@gmail.com

Quentin Le Roux

MSc in Data Science & Artificial Intelligence
 Université Côte d'Azur
 quentin.leroux@edhec.com

ABSTRACT

Stochastic gradient-based optimization led to the recent development of adaptive moment estimation methods, which are known to outperform traditional stochastic techniques. This paper presents the comparative analysis of three of such methods (Adam, AdamW, and AMSGrad) when applied to industry-standard, real-world datasets in the context of a multi-class classification problem.

1 INTRODUCTION

The idea behind stochastic approximation revolves around the minimization of an objective or risk function with the adjunctive use of noise as part of the optimization process (Robbins & Monro, 1951). As such, stochastic gradient descent (SGD) can be described as:

$$\underset{x}{\text{minimize}} \quad F(x) = \mathbb{E}(f(x; \zeta))$$

Given ζ , a random seed, and $f(\cdot)$, the composite of a loss function l and a prediction function h (Bottou et al., 2018).

Here the noise represents a random pick $(x^{(i)}, y^{(i)})$ from a dataset, over which the following gradient descent will be performed:

$$\theta = \theta - \eta \cdot \nabla_{\theta} F(\theta, x^{(i)}, y^{(i)})$$

Given n , the size of the dataset with $i \in \{1, \dots, n\}$, $x^{(i)}$ and $y^{(i)}$, a data point and its label (in the case of classification), θ , the model parameters, $\nabla_{\theta} F$, the gradient of the objective function, and η , the learning rate that sets the optimization's stepsize.

SGD makes use of empirical data in a more efficient way than batch methods (Ruder, 2017). Indeed, compared to batch gradient descent, SGD does not perform superfluous gradient computations, which are a growing concern as a dataset's size increases. Furthermore, the intrinsic randomness allows the descent to shift towards potentially better minimas (Bottou et al., 2018). The use of random picks from the dataset thus enables a more efficient gradient update than when all data is simultaneously iterated over. It has been shown that SGD is an efficient method that achieves fast initial improvement with low cost (Nemirovski et al., 2009).

However, the fixed learning rate (i.e. stepsize) and the use of noise impede standard SGD from converging to the exact, achievable minimum (Bottou & Bousquet, 2007). Dealing with this overshooting is a long-standing area of research, along with improving on its sublinear rate of convergence (Bottou et al., 2018).

2 ADAPTIVE MOMENT ESTIMATION METHODS

To answer SGD's overshooting problem, development in stepsize selection was undertaken. Early examples were found in Nesterov's Accelerated Gradient (Nesterov, 1983) or SGD with Momentum (Qian, 1999). More recently, however, adaptive learning rate methods arose such as with the Adaptive Gradient method or AdaGrad (Duchi et al., 2011), its variant AdaDelta (Zeiler, 2012), or the Root Mean Square Propagation or RMSProp (Tieleman & Hinton, 2012). While the two former accumulate gradients (all of them for AdaGrad, and a window selection for AdaDelta), the latter divides its gradients by a running average over its recent magnitude.

2.1 ADAM

An expansion over the concept of adaptive learning rate led to the development of a new kind of optimization algorithm: the Adaptive Moment Estimation method or Adam (Kingma & Ba, 2017). Where the standard SGD method keeps a fixed stepsize for all updates, the Adam method takes inspiration from AdaGrad and RMSProp. With Adam, the parameter learning rate is updated based on the gradient's first and second moments and the process stores an exponentially decaying average of past gradients such that:

$$\begin{aligned} m_t &= \beta_1 m_{t-1} + (1 - \beta_1) g_t \\ v_t &= \beta_2 v_{t-1} + (1 - \beta_2) g_t^2 \\ \hat{m}_t &= \frac{m_t}{1 - \beta_1^t} \\ \hat{v}_t &= \frac{v_t}{1 - \beta_2^t} \\ \theta_{t+1} &= \theta_t - \frac{\eta}{\sqrt{\hat{v}_t} + \epsilon} \hat{m}_t \end{aligned}$$

Given θ , the model parameters, β_1 and β_2 , two forgetting parameters, m_t , the estimate of the first moment, v_t , the estimate of the second moment, \hat{m}_t and \hat{v}_t , their biased-corrected counterparts, and ϵ , a very small value to avoid dividing by zero.

This method showed a strong, early success as it proved effective in practice when used on large datasets (Kingma & Ba, 2017). As such, many variants emerged.

2.2 ADAMW

One such variant takes advantage of adding a weight decay, resulting in a new algorithm: the Adam with decoupled weight decay or AdamW (Loshchilov & Hutter, 2019). The observation that led to the algorithm's development was the highlight that, contrary to standard SGD, L_2 regularization and weight decay were not identical in the case of Adam (The equivalence proof for standard SGD is reproduced in Appendix A). As L_2 regularization was found not to be as effective with Adam, the modification made with AdamW focused on the update step, inspired by the weight decay described by Hanson & Pratt (1988):

$$\theta_{t+1} = \theta_t - \alpha_t \left(\frac{\eta m_t}{\sqrt{\hat{v}_t} + \epsilon} + \lambda \theta_t \right)$$

Given λ , the rate of the weight decay at each step, and α_t , a *SetScheduleMultiplier(t)* allowing to reset the learning rate during optimization (Loshchilov & Hutter, 2019).

2.3 AMSGRAD

The final adaptive moment estimation method we are interested in began with the observation of a flaw in the Adam method: Adam can fail at converging towards an optimal solution in a convex setting (Reddi et al., 2019). Indeed, there is an underlying risk that positive definiteness could be violated for methods other than standard SGD.

Given T , a number of rounds, and Γ_t , the quantity measuring the change in the inverse of the learning rate for an adaptive method with regards to time, we have:

$$\begin{aligned} \Gamma_{t+1} &= \frac{\sqrt{V_{t+1}}}{\eta_{t+1}} - \frac{\sqrt{V_t}}{\eta_t} \\ V_t &= \text{diag}(v_t) \end{aligned}$$

For standard SGD, Reddi et al. (2019) highlighted that:

$$\forall t \in [T] \quad \Gamma_t \geq 0$$

However, they also proved that this property is not always true for methods such as Adam. Indeed, a convex optimization problem like Adam was shown to have a potentially non-zero average regret

(Shalev-Shwartz, 2012), i.e., the delta between the loss of a possible action and the action taken given an hypothesis class h^* would not converge to zero (Reddi et al., 2019) such that:

$$\frac{R_T}{T} \not\rightarrow_{T \rightarrow \infty} 0$$

$$R_T(h^*) = \sum_{t=1}^T l(p_t, y_t) - \sum_{t=1}^T l(h^*(x_t), y_t)$$

To resolve this issue (that the quantity Γ_t can be negative for Adam), a new adaptive moment estimation method was developed: AMSGrad (Reddi et al., 2019). The algorithm is meant to guarantee convergence while conserving the advantages of the Adam method. The modification over Adam impacts the computation of \hat{v}_t and of the update step (note the use of m_t rather than \hat{m}_t) such that:

$$\hat{v}_t = \max(\hat{v}_{t-1}, v_t)$$

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{\hat{v}_t} + \epsilon} m_t$$

The main difference with Adam happens with the preserving of the maximum of all past values v_t , resulting in a never-increasing stepsize, which was the underlying risk of the Adam method.

3 COMPARISON METHODOLOGY

In this section, we present an empirical comparison procedure for the three adaptive moment estimation methods we previously covered. Our experiments will focus on the problem of multi-class classification using industry-standard and real-world datasets with various kind of neural networks.

3.1 DATASETS

Our comparative analysis will rely on four industry-standard, fixed-size image datasets: MNIST (Lecun et al., 1998), Fashion-MNIST (Xiao et al., 2017), German Traffic Sign Recognition Benchmark GTSRB (Stallkamp et al., 2012), and CIFAR-10 (Krizhevsky, 2012).

The MNIST dataset corresponds to grayscale, normalized and centered fixed-size images (28×28 pixels) of handwritten digits split between 60,000 training and 10,000 testing examples, and 10 classes. The more recent Fashion-MNIST follows the same format: grayscale, normalized and centered fixed-size images (28×28 pixels) of fashion and clothing articles also split between 60,000 training and 10,000 testing examples, and 10 classes.

The GTSRB dataset corresponds to color, fixed-size images (32×32 pixels with 3 channels) split between 34,799 training and 4,410 testing examples (the provided validation set was not included), and 43 classes. Similarly, the CIFAR-10 dataset corresponds to color, fixed-size images (32×32 pixels with 3 channels) split between 50,000 training and 10,000 testing examples, and 10 classes.

As each dataset is organized into classes, they provide us a standard multi-class classification problem. The problem compatibility between the four datasets allows us to compare Adam, AdamW, and AMSGrad using a variety of interchangeable neural network models.

3.2 MODELS

We compare the three adaptive moment estimation methods using four distinct model types: a shallow neural network, a deep (fully-connected) neural network, a convolutional neural network, and a residual convolutional neural network. The corresponding model graphs for MNIST and Fashion-MNIST can be found in Appendix B, those for the GTSRB and CIFAR-10 in Appendix C. The corresponding Python code can be found in Appendix D.

Each model will use the same suite of hyper-parameters, notably being trained for 100 epochs with batch sizes of 32 elements with the categorical crossentropy loss function. Each optimization algorithm will be initialized with default parameters using the TensorFlow library (Abadi et al., 2015). The default parameters used for AdamW and AMSGrad are reproduced in Table 1.

Table 1: AdamW and AMSGrad default parameters

Parameters	Value
learning rate	$1e^{-2}$
β_1	0.9
β_2	0.999
weight decay	$1e^{-4}$
ϵ	$1e^{-8}$
decay	0.

3.3 COMPARISON APPROACH

Our comparison approach will be two-fold. We will first take note of the speed of convergence between the three different algorithms by comparing the changes in loss scores during training. Then we will study their comparative accuracy, or lack thereof, depending on the model and dataset.

4 EXPERIMENTS

Our first observation relates to how Adam performs with regards to AdamW and AMSGrad. Based on a first visual assessment, Adam’s loss convergence, when it converges, performs differently than the latter two, who display similar patterns. Train and test loss scores per epoch are reproduced in Appendix E for MNIST, Appendix F for Fashion-MNIST, Appendix G for GTSRB, and Appendix H for CIFAR-10.

Our second observation relates to how both AdamW and AMSGrad algorithms display a faster loss convergence than Adam. This observation applies for all the datasets used in this exercise. We can say that, visually, AdamW and AMSGrad will generally converge around a mean value in less than 40 epochs during training. This translates into a faster loss convergence on the test set – except in the case of a shallow neural network on CIFAR-10 where loss scores are too noisy to draw any conclusion (See Figure 21).

Our third observation highlights that the Adam algorithm does not always translate into any loss convergence on the test set (See Figures 13, 17, and 22). Furthermore, we also find that Adam may result in an increasing test set loss (See Figures 9, 10, 14, and 18). This overfitting behavior validates the observation made by Loshchilov & Hutter (2019) that Adam will not always converge.

A final observation relates to model performance, here represented by accuracy scores (See Appendix I for MNIST, Appendix J for Fashion-MNIST, Appendix K for GTSRB, and Appendix L for CIFAR-10 scores). Although AdamW and AMSGrad converge in a quicker fashion than Adam, we find that it does not always lead to a better test set accuracy. For instance, using Adam with a residual network yields a 0.79 test accuracy on GTSRB, against 0.63 and 0.61 with AdamW and AMSGrad respectively (See Figures 58, 59, and 60).

5 CONCLUSION & DISCUSSION

We saw that the AdamW and AMSGrad algorithms perform better in terms of convergence than Adam when applied to four industry-standard, real-world datasets, adding to the literature empirical proof of their overall soundness. A visible caveat, however, is that their faster convergence will not always translate into a better model accuracy on testing sets.

Beyond this short comparative analysis, further explorations can be done in two regards. First, as better convergence does not always imply better accuracy, a thorough exploration of initialization parameters for AdamW and AMSGrad should be undertaken. The goal would be to find the best mix of parameters to try to match Adam in terms of accuracy. Another area of exploration relates to other, newer adaptive moment estimation methods that could be covered as part of a future comparative analysis: Adabound (Luo et al., 2019), Radam (Liu et al., 2020), LAMB (You et al., 2020), MAS (Landro et al., 2020), and MADGRAD (Defazio & Jelassi, 2021).

REFERENCES

- Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dan Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. Tensorflow: Large-scale machine learning on heterogeneous distributed systems, 2015.
- Léon Bottou and Olivier Bousquet. The tradeoffs of large scale learning. volume 20, 01 2007.
- Léon Bottou, Frank E. Curtis, and Jorge Nocedal. Optimization methods for large-scale machine learning, 2018.
- Aaron Defazio and Samy Jelassi. Adaptivity without compromise: A momentumized, adaptive, dual averaged gradient method for stochastic optimization, 2021.
- John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12:2121–2159, 07 2011.
- Stephen José Hanson and Lorien Y. Pratt. Comparing biases for minimal network construction with back-propagation. In *Proceedings of the 1st International Conference on Neural Information Processing Systems*, NIPS’88, pp. 177–185, Cambridge, MA, USA, 1988. MIT Press.
- Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2017.
- Alex Krizhevsky. Learning multiple layers of features from tiny images. *University of Toronto*, 05 2012.
- Nicola Landro, Ignazio Gallo, and Riccardo La Grassa. Mixing adam and sgd: a combined optimization method, 2020.
- Yann Lecun, Leon Bottou, Y. Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86:2278 – 2324, 12 1998. doi: 10.1109/5.726791.
- Liyuan Liu, Haoming Jiang, Pengcheng He, Weizhu Chen, Xiaodong Liu, Jianfeng Gao, and Jiawei Han. On the variance of the adaptive learning rate and beyond, 2020.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization, 2019.
- Liangchen Luo, Yuanhao Xiong, Yan Liu, and Xu Sun. Adaptive gradient methods with dynamic bound of learning rate, 2019.
- Arkadi Nemirovski, Anatoli Juditsky, Guanghui Lan, and And Shapiro. Robust stochastic approximation approach to stochastic programming. *Society for Industrial and Applied Mathematics*, 19: 1574–1609, 01 2009.
- Y. Nesterov. A method for unconstrained convex minimization problem with the rate of convergence $o(1/k^2)$. 1983.
- Ning Qian. On the momentum term in gradient descent learning algorithms. *Neural Networks*, 12(1): 145–151, 1999.
- Sashank J. Reddi, Satyen Kale, and Sanjiv Kumar. On the convergence of adam and beyond, 2019.
- Herbert Robbins and Sutton Monro. A Stochastic Approximation Method. *The Annals of Mathematical Statistics*, 22(3):400 – 407, 1951.
- Sebastian Ruder. An overview of gradient descent optimization algorithms, 2017.
- Shai Shalev-Shwartz. Online learning and online convex optimization. *Found. Trends Mach. Learn.*, 4 (2):107–194, February 2012. ISSN 1935-8237.

- J. Stallkamp, M. Schlipsing, J. Salmen, and C. Igel. Man vs. computer: Benchmarking machine learning algorithms for traffic sign recognition. *Neural Networks*, (0):-, 2012. ISSN 0893-6080. doi: 10.1016/j.neunet.2012.02.016.
- Tijmen Tieleman and Geoffrey Hinton. Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude, 04 2012.
- Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms, 2017.
- Yang You, Jing Li, Sashank Reddi, Jonathan Hseu, Sanjiv Kumar, Srinadh Bhojanapalli, Xiaodan Song, James Demmel, Kurt Keutzer, and Cho-Jui Hsieh. Large batch optimization for deep learning: Training bert in 76 minutes, 2020.
- Matthew D. Zeiler. ADADELTA: an adaptive learning rate method. *CoRR*, abs/1212.5701, 2012.

A FORMAL ANALYSIS OF WEIGHT DECAY VS. L_2 REGULARIZATION FOR STANDARD SGD

Proposition (Weight decay is equivalent to L_2 regularization for standard SGD). *Standard SGD coupled with a base learning rate η executes identical steps between loss functions $f_t(\theta)$ with a weight decay λ and on regularized loss functions without weight decay $f_t^{reg}(\theta) = f_t(\theta) + \frac{\lambda'}{2} \|\theta\|_2^2$ with $\lambda' = \frac{\lambda}{\eta}$.*

Proof of Proposition

The weights update for SGD with regularization and SGD with weight decay are represented respectively by the following iterates:

$$\begin{aligned}\theta_{t+1} &\leftarrow \theta_t - \eta \nabla f_t^{reg}(\theta_t) = \theta_t - \alpha \nabla f_t(\theta_t) - \eta \lambda' \theta_t \\ \theta_{t+1} &\leftarrow (1 - \lambda) \theta_t - \alpha \nabla f_t(\theta_t)\end{aligned}$$

Resulting in the following observation:

$$\lambda' = \frac{\lambda}{\eta}$$

This proof is a modified reproduction of the one provided by Loshchilov & Hutter (2019).

B GRAPHS OF NEURAL NETWORKS USED FOR EMPIRICAL COMPARISON FOR 1-CHANNEL IMAGES (MNIST AND FASHION-MNIST)

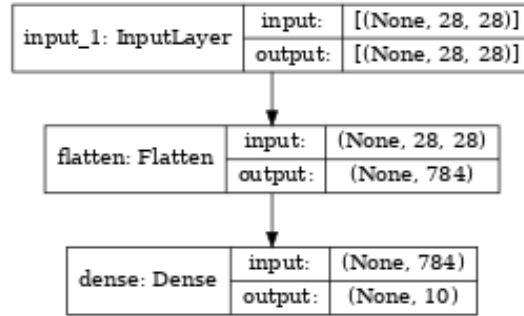


Figure 1: Shallow Neural Network

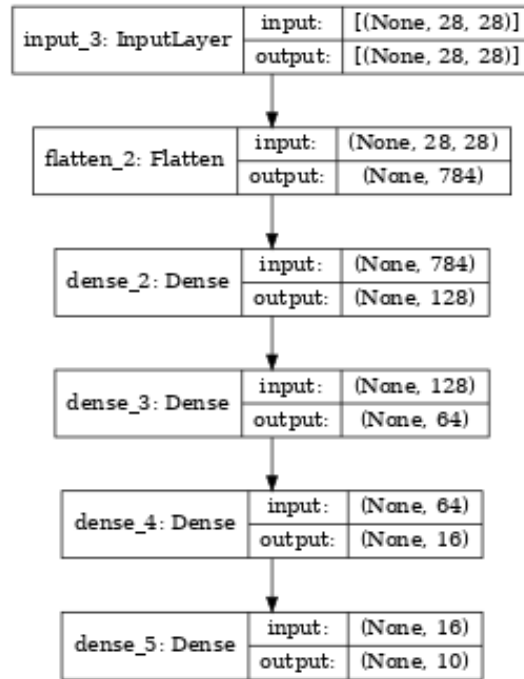


Figure 2: Deep (Fully-Connected) Neural Network

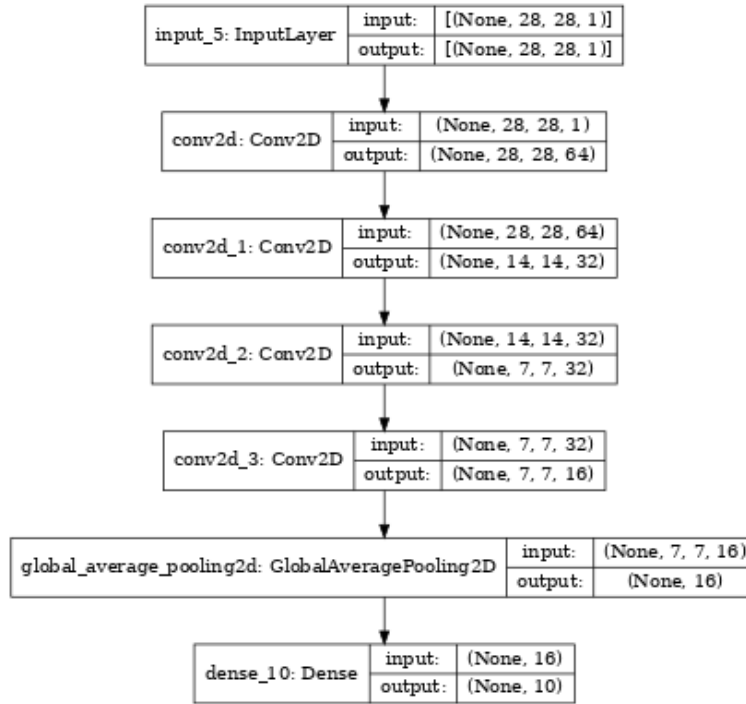


Figure 3: Convolutional Neural Network

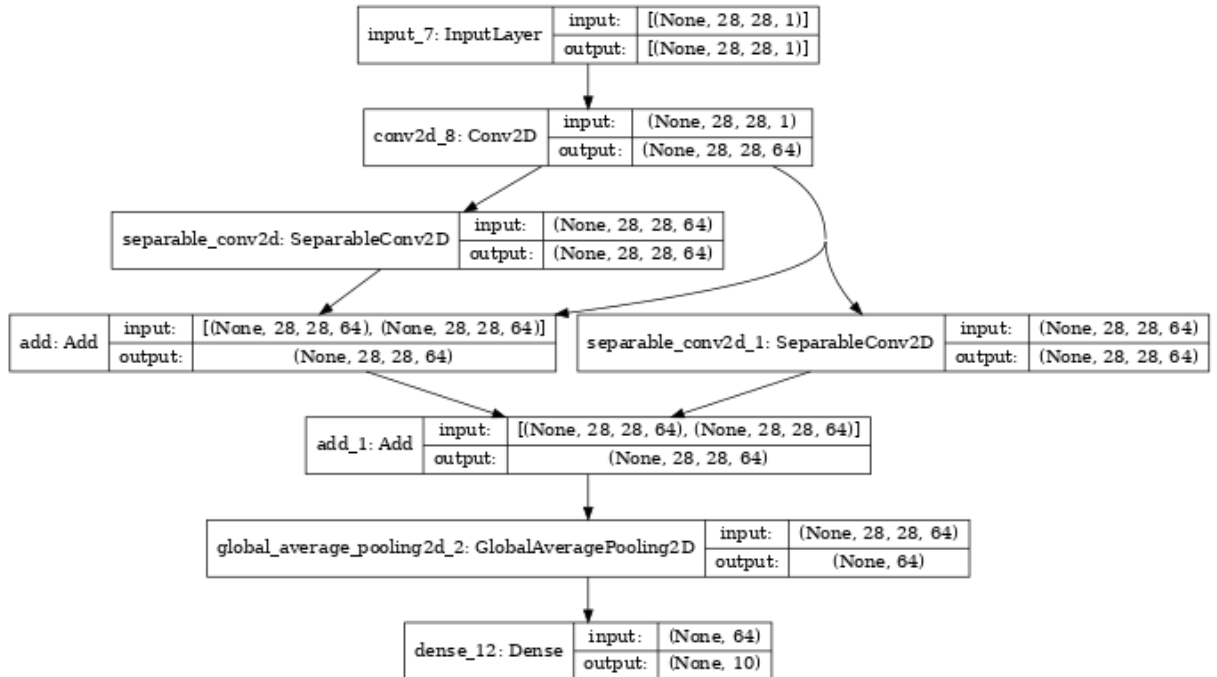


Figure 4: Residual Neural Network

C GRAPHS OF NEURAL NETWORKS USED FOR EMPIRICAL COMPARISON FOR 3-CHANNEL IMAGES (GTSRB AND CIFAR-10)

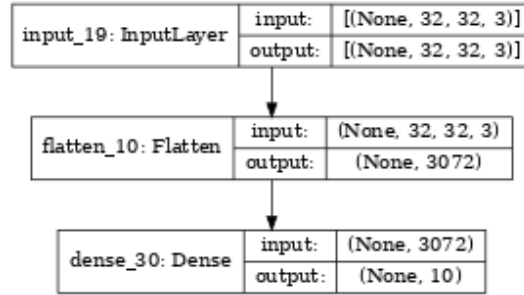


Figure 5: Shallow Neural Network

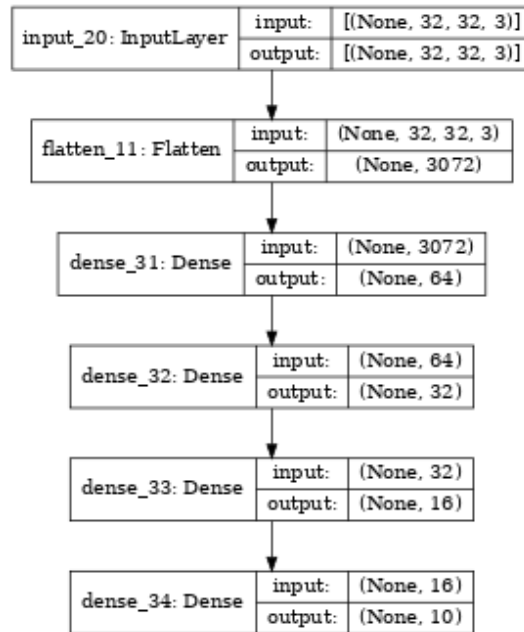


Figure 6: Deep (Fully-Connected) Neural Network

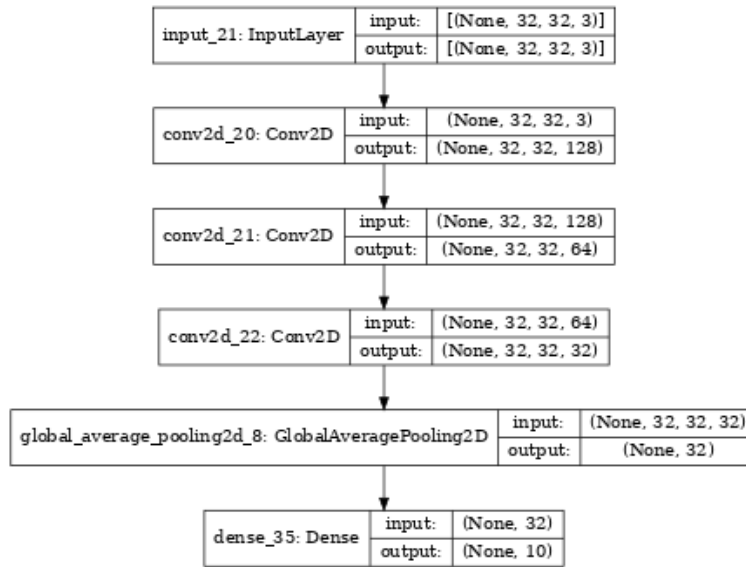


Figure 7: Convolutional Neural Network

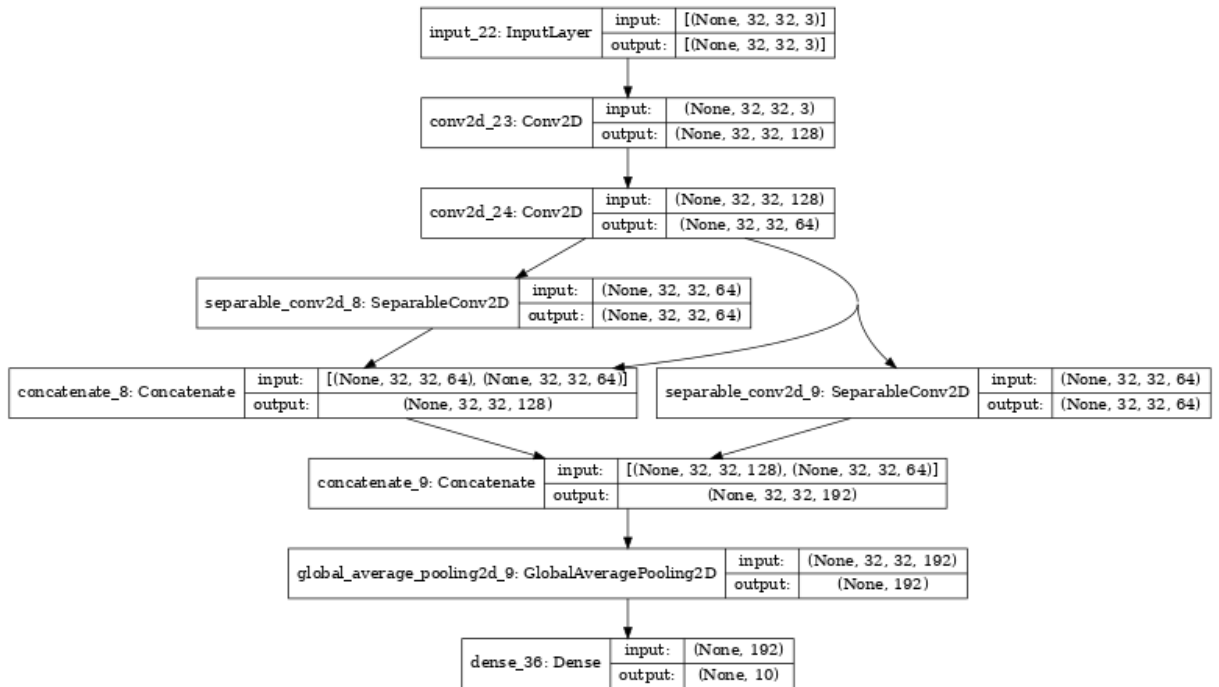


Figure 8: Residual Neural Network

D CODE IMPLEMENTATION USING PYTHON AND TENSORFLOW

Listing 1: Library Imports

```
import matplotlib.pyplot as plt
import numpy as np
import pickle
import random
import tensorflow as tf

from plot_keras_history import plot_history
from tensorflow.keras import Input, Model
from tensorflow.keras.datasets import mnist, fashion_mnist, cifar10
from tensorflow.keras.layers import Flatten, Dense, Concatenate
from tensorflow.keras.layers import Conv2D, SeparableConv2D
from tensorflow.keras.layers import BatchNormalization, Activation
from tensorflow.keras.layers import GlobalAveragePooling2D
from tensorflow.keras.optimizers import Adam
from tensorflow_addons.optimizers import AdamW
from tensorflow.keras.utils import plot_model
from keras.utils import np_utils
```

Listing 2: Custom Function Declarations

```
def normalize(dataset):
    """
    Normalizes the pixel values of an image stored as an array
    """
    return dataset/255.

def one_hot_labels(labels):
    """
    One-hot encodes the imported labels
    """
    return tf.one_hot(labels, depth=10)

def shallow_neural_network(input_shape, num_classes):
    """
    Creates a shallow neural network
    """
    inputs = Input(shape=input_shape)
    flatten = Flatten()(inputs)
    outputs = Dense(num_classes, activation="softmax")(flatten)
    return Model(inputs, outputs)

def deep_neural_network(input_shape, num_classes):
    """
    Creates a Deep (Fully-Connected) Neural Network
    """
    inputs = Input(shape=input_shape)
    flatten = Flatten()(inputs)
    layer_1 = Dense(64, activation="relu")(flatten)
    layer_2 = Dense(32, activation="relu")(layer_1)
    layer_3 = Dense(16, activation="relu")(layer_2)
    outputs = Dense(num_classes, activation="softmax")(layer_3)
    return Model(inputs, outputs)

def convolutional_neural_network(input_shape, num_classes):
    """
    Creates a convolutional neural network
    """
    inputs = Input(shape=input_shape)
    x = Conv2D(128, 3, padding="same", activation="relu")(inputs)
    x = Conv2D(64, 3, padding="same", activation="relu")(x)
    x = Conv2D(32, 3, padding="same", activation="relu")(x)
```

```
x = GlobalAveragePooling2D()(x)
outputs = Dense(num_classes, activation="softmax")(x)
return Model(inputs, outputs)

def residual_neural_network(input_shape, num_classes):
    """
    Creates a residual neural network
    """
    inputs = Input(shape=input_shape)
    x = Conv2D(128, 3, padding="same", activation="relu")(inputs)
    x = Conv2D(64, 3, padding="same", activation="relu")(x)
    previous_block_activation_1 = x
    previous_block_activation_2 = x
    residual = SeparableConv2D(64, 3,
                                padding="same",
                                activation="relu")(previous_block_activation_2)
    x = Concatenate()([x, residual])
    previous_block_activation = x
    residual = SeparableConv2D(64, 3,
                                padding="same",
                                activation="relu")(previous_block_activation_1)
    x = Concatenate()([x, residual])
    x = GlobalAveragePooling2D()(x)
    outputs = Dense(num_classes, activation="softmax")(x)
    return Model(inputs, outputs)

def compile_fit(model, X, y, X_val, y_val,
                optimizer, loss, metrics,
                epochs, batch_size):
    """
    Compiles and fits a model given a set of hyperparameters
    """
    model.compile(
        optimizer=optimizer,
        loss=loss,
        metrics=metrics
    )
    history = model.fit(
        x=X, y=y,
        epochs=epochs, batch_size=batch_size,
        validation_data=(X_val, y_val),
        verbose=0
    )
    return history

def plot_loss(model, dataset):
    """
    Plots a comparative representation of obtained losses for
    a specific model for all three selected optimizers
    """
    if dataset == "MNIST":
        res = results_mnist
    elif dataset == "Fashion-MNIST":
        res = results_fmnist
    elif dataset == "Road":
    else:
        res = results_cifar
    a = res[f"{model}_adam"].history["loss"]
    b = res[f"{model}_adam"].history["val_loss"]
    c = res[f"{model}_adamw"].history["loss"]
    d = res[f"{model}_adamw"].history["val_loss"]
    e = res[f"{model}_amsgrad"].history["loss"]
    f = res[f"{model}_amsgrad"].history["val_loss"]

    plt.figure(figsize=(12,10))
```

```
plt.plot(a, '--', linewidth=1, color="maroon")
plt.plot(b, linewidth=3, alpha=0.5, color="firebrick")
plt.plot(c, '--', linewidth=1, color="lightseagreen")
plt.plot(d, linewidth=3, alpha=0.5, color="teal")
plt.plot(e, '--', linewidth=1, color="purple")
plt.plot(f, linewidth=3, alpha=0.5, color="mediumvioletred")
plt.xlabel("EPOCHS", fontsize=13)
plt.ylabel("LOSS-VALUE", fontsize=13)
plt.title("Observed_loss_values_per_epoch_given_ + \
         f"a_specific_optimizer_on_the_{dataset}_dataset",
         fontsize=15)
plt.legend(["Adam_Training_Loss", "Adam_Testing_Loss",
           "AdamW_Training_Loss", "Adamw_Testing_Loss",
           "AMSGrad_Training_Loss", "AMSGrad_Testing_Loss"],
           fontsize=13)
plt.show()
```

Listing 3: Dataset Imports

```
# 1. Import dataset
# 2. Normalize X arrays (features)
# 3. One-hot encode Y arrays (labels)

# MNIST import
(mx_train, my_train), (mx_test, my_test) = mnist.load_data()
mx_train, mx_test = normalize(mx_train), normalize(mx_test)
my_train, my_test = one_hot_labels(my_train), one_hot_labels(my_test)

# Fashion-MNIST import
(fmx_train, fmy_train), (fmx_test, fmy_test) = fashion_mnist.load_data()
fmx_train, fmx_test = normalize(fmx_train), normalize(fmx_test)
fmy_train, fmy_test = one_hot_labels(fmy_train), one_hot_labels(fmy_test)

# GTSRB import
!git clone https://bitbucket.org/jadslim/german-traffic-signs
!ls german-traffic-signs
with open('german-traffic-signs/train.p', 'rb') as f:
    train_data = pickle.load(f)
with open('german-traffic-signs/valid.p', 'rb') as f:
    val_data = pickle.load(f)
sign_x_train, sign_y_train = train_data['features'], train_data['labels']
sign_x_test, sign_y_test = val_data['features'], val_data['labels']
sign_x_train, sign_x_test = normalize(sign_x_train), normalize(sign_x_test)
sign_y_train = one_hot_labels(sign_y_train, 43)
sign_y_test = one_hot_labels(sign_y_test, 43)

# CIFAR-10 import
(cifar_x_train, cifar_y_train), (cifar_x_test, cifar_y_test) = cifar10.load_data()
cifar_x_train = normalize(cifar_x_train)
cifar_x_test = normalize(cifar_x_test)
cifar_y_train = one_hot_labels(cifar_y_train)
cifar_y_test = one_hot_labels(cifar_y_test)
cifar_y_train = tf.reshape(cifar_y_train, [50000, 10])
cifar_y_test = tf.reshape(cifar_y_test, [10000, 10])
```

Listing 4: Model Declarations

```
# Global variable declarations
input_shape = (28, 28)
input_shape_cnn = (28, 28, 1)
input_shape_color = (32, 32, 3)
num_classes = 10 #or 43 for GTSRB

# Declaring Shallow Models
```

```
shallow = shallow_neural_network(input_shape=input_shape ,
                                num_classes=num_classes)

# Declaring Deep (Fully-Connected) Models
dnn = deep_neural_network(input_shape=input_shape ,
                           num_classes=num_classes)

# Declaring Convolutional Models
cnn = convolutional_neural_network(input_shape=input_shape_cnn ,
                                   num_classes=num_classes)

# Declaring ResNet Models
resnet = residual_neural_network(input_shape=input_shape_cnn ,
                                  num_classes=num_classes)
```

Listing 5: Experiments/Computing Loss Results

```
# Global variable declarations
epochs = 100
batch_size = 32
loss = "categorical_crossentropy"
metrics = ["accuracy"]

# Generic learning rate and weight decay for AdamW and AMSgrad
lr=0.001
beta_1=0.9
beta_2=0.999
weight_decay=1e-4
epsilon=1e-8
decay=0.

optimizers = [("adam", Adam()),
               ("adamw", AdamW(lr=lr, beta_1=beta_1, beta_2=beta_2,
                               weight_decay=weight_decay, epsilon=epsilon,
                               decay=decay)),
               ("amsgrad", AdamW(lr=lr, beta_1=beta_1, beta_2=beta_2,
                                  weight_decay=weight_decay, epsilon=epsilon,
                                  decay=decay,
                                  amsgrad=True)))]

results = {}

# Computes loss results for a kind of model
for name, optimizer in optimizers:
    res = compile_fit(model,
                      x_train, y_train,
                      x_test, y_test,
                      optimizer, loss, metrics,
                      epochs, batch_size)
    results[f"{name}"]=res
```

Listing 6: Plotting Results/Data Visualization example with the Shallow model and MNIST dataset

```
plot_history(results[f"{model}_adam"].history)
plot_history(results[f"{model}_adamw"].history)
plot_history(results[f"{model}_amsgrad"].history)
plot_loss(f"{model}", f"{dataset}")
```

E OBSERVED LOSS VALUES FOR EACH OPTIMIZERS AND MODELS TRAINED ON THE MNIST DATASET

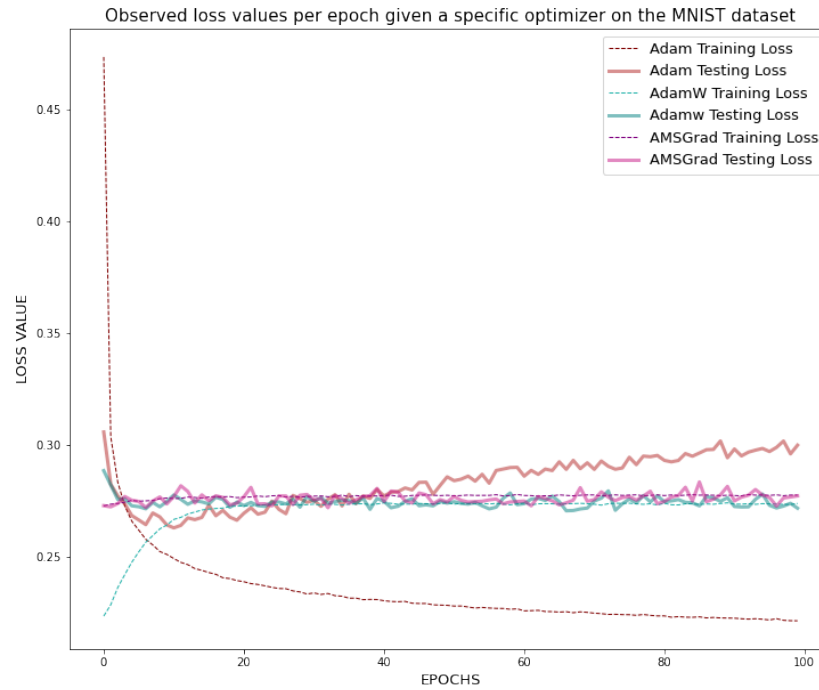


Figure 9: Results for the shallow neural network trained on MNIST

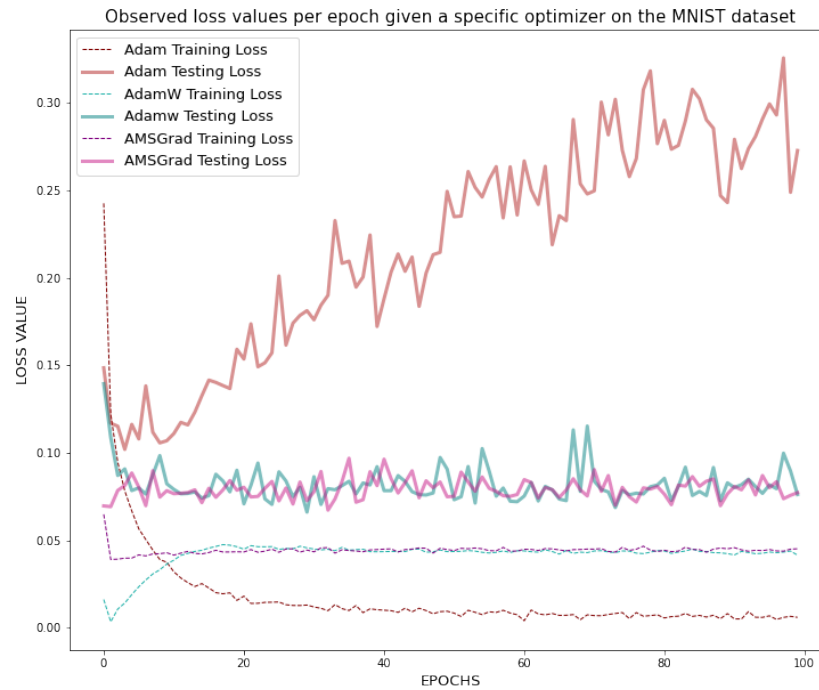


Figure 10: Results for the deep neural network trained on MNIST

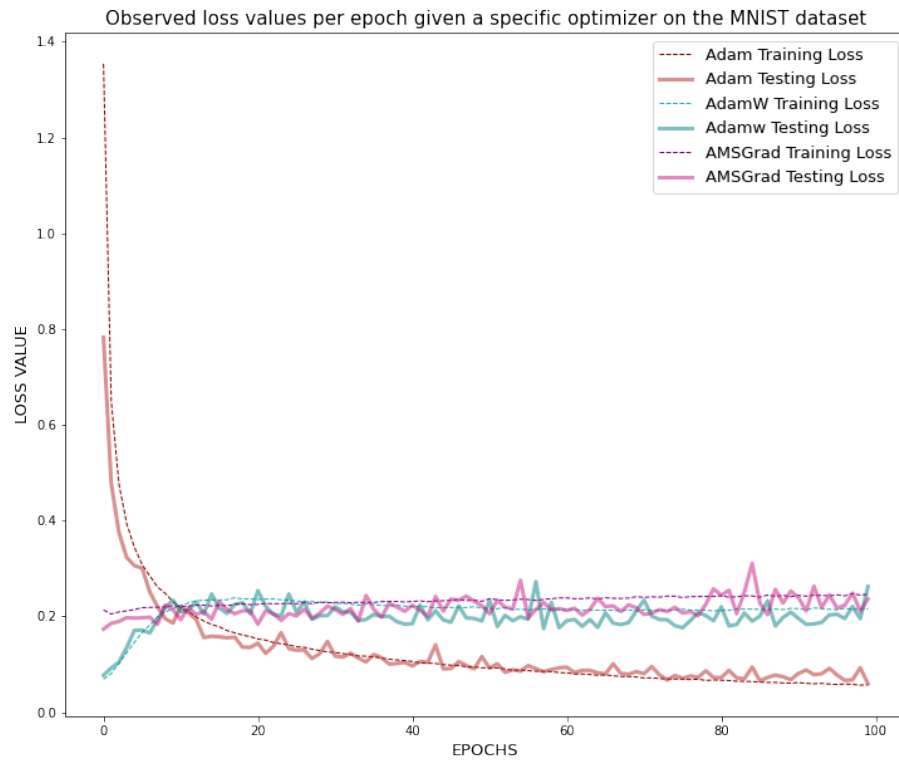


Figure 11: Results for the convolutional neural network trained on MNIST

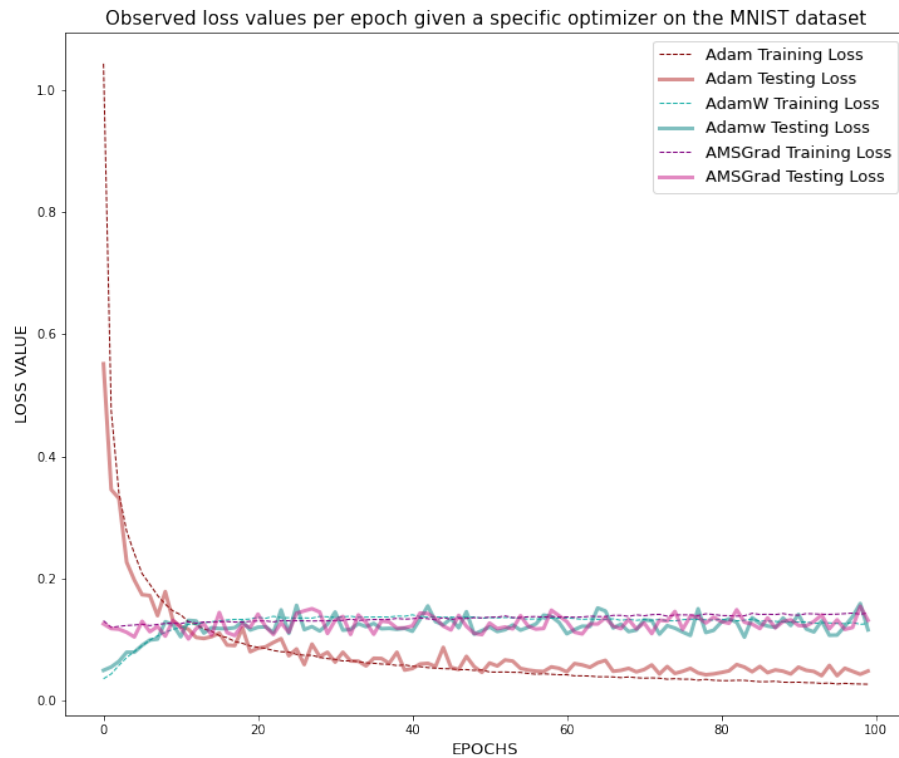


Figure 12: Results for the residual neural network trained on MNIST

F OBSERVED LOSS VALUES FOR EACH OPTIMIZERS AND MODELS TRAINED ON THE FASHION-MNIST DATASET

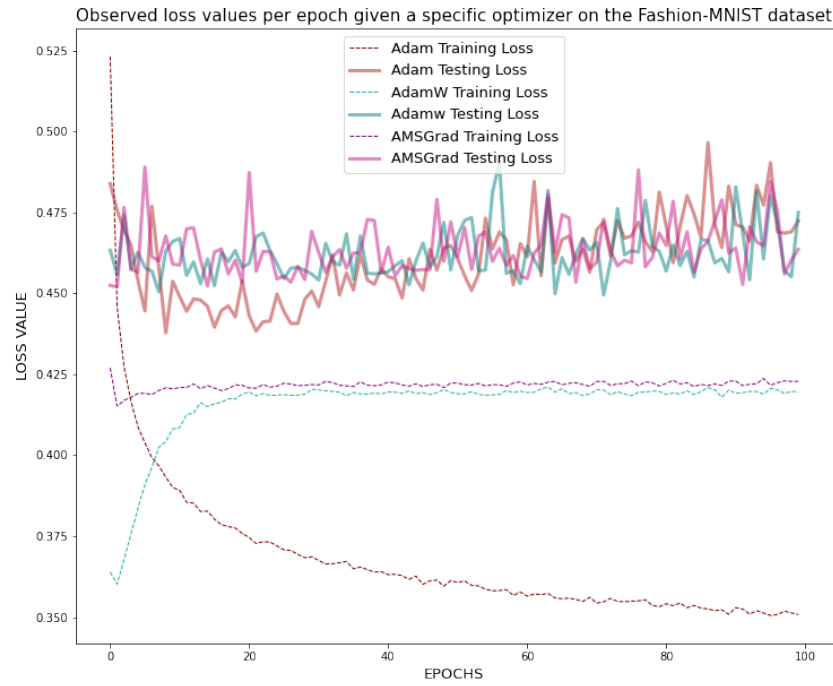


Figure 13: Results for the shallow neural network trained on Fashion-MNIST

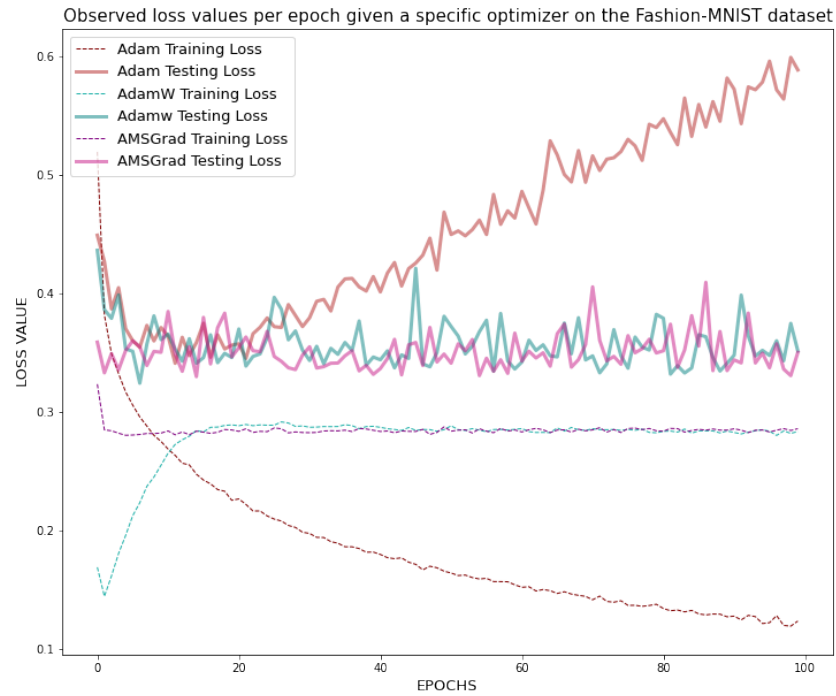


Figure 14: Results for the deep neural network trained on Fashion-MNIST

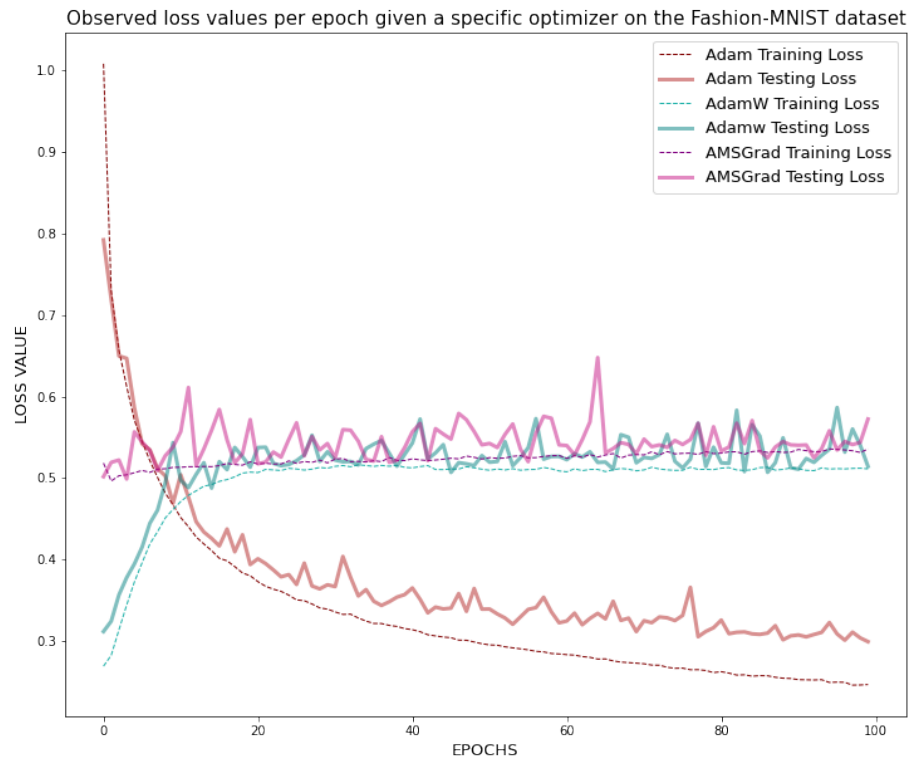


Figure 15: Results for the convolutional neural network trained on Fashion-MNIST

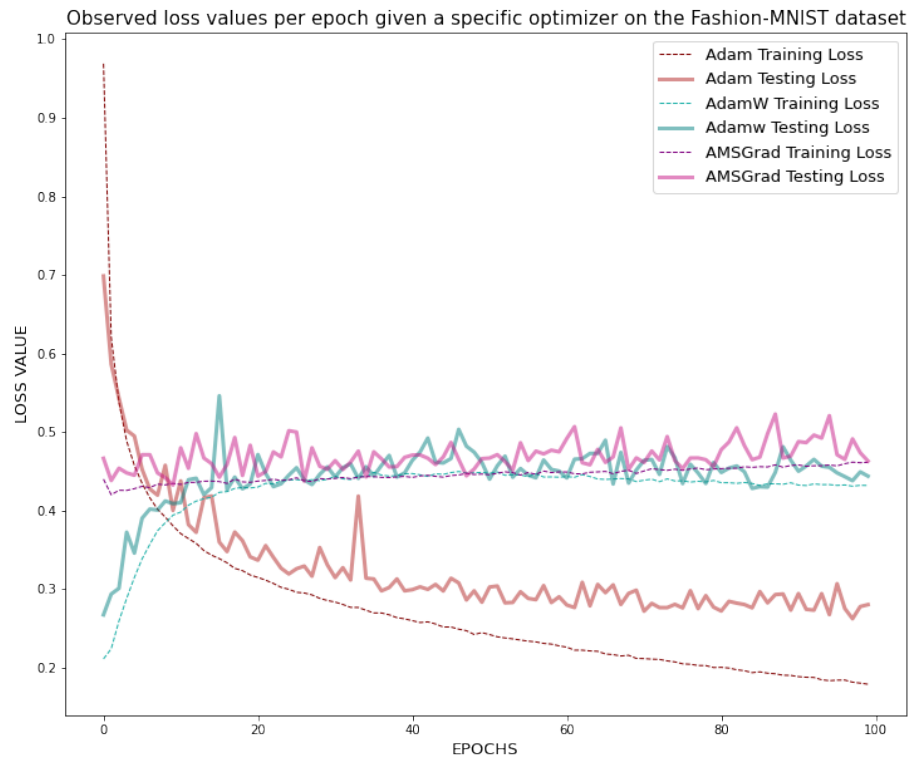


Figure 16: Results for the residual neural network trained on Fashion-MNIST

G OBSERVED LOSS VALUES FOR EACH OPTIMIZERS AND MODELS TRAINED ON THE GERMAN TRAFFIC SIGN RECOGNITION BENCHMARK DATASET

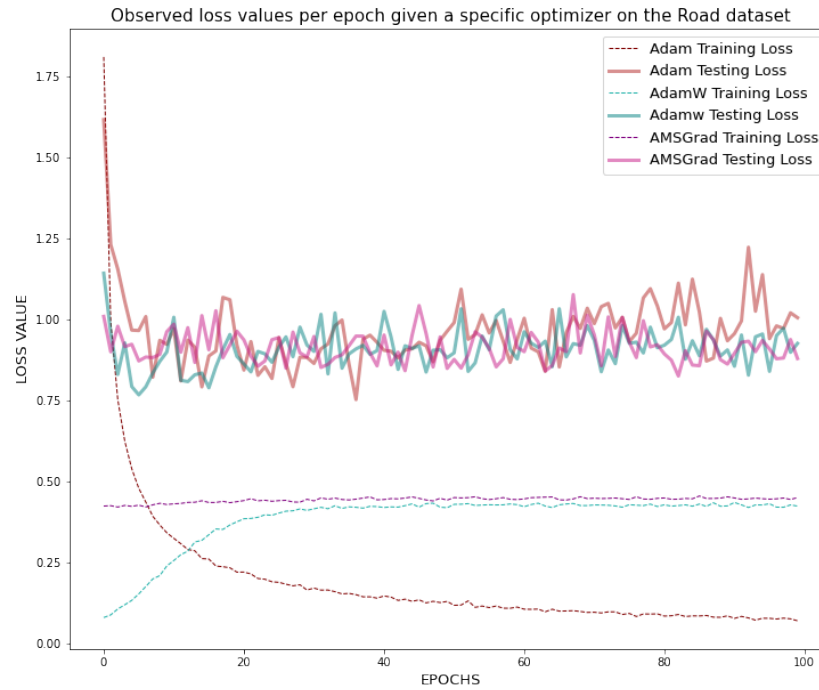


Figure 17: Results for the shallow neural network trained on GTSRB

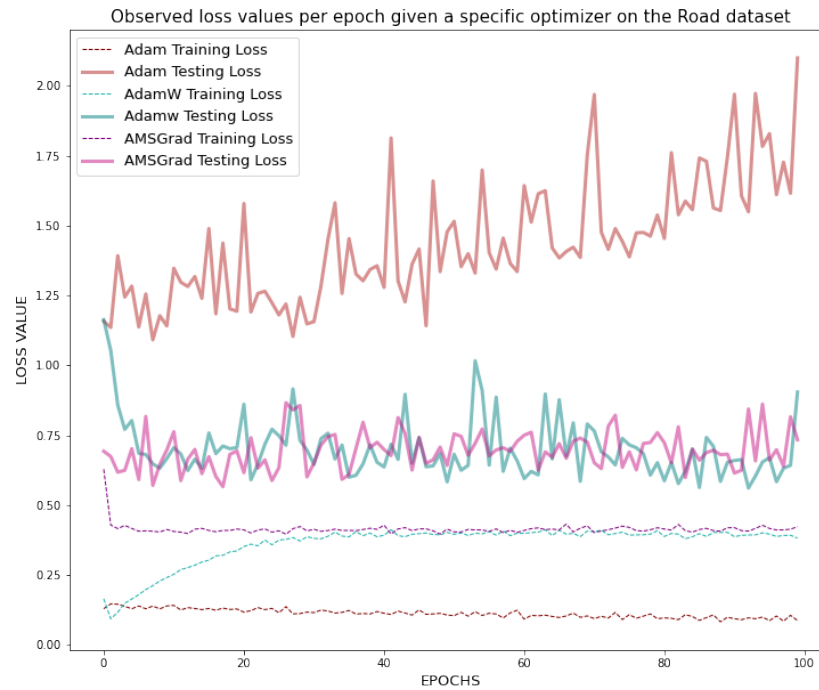


Figure 18: Results for the deep neural network trained on GTSRB

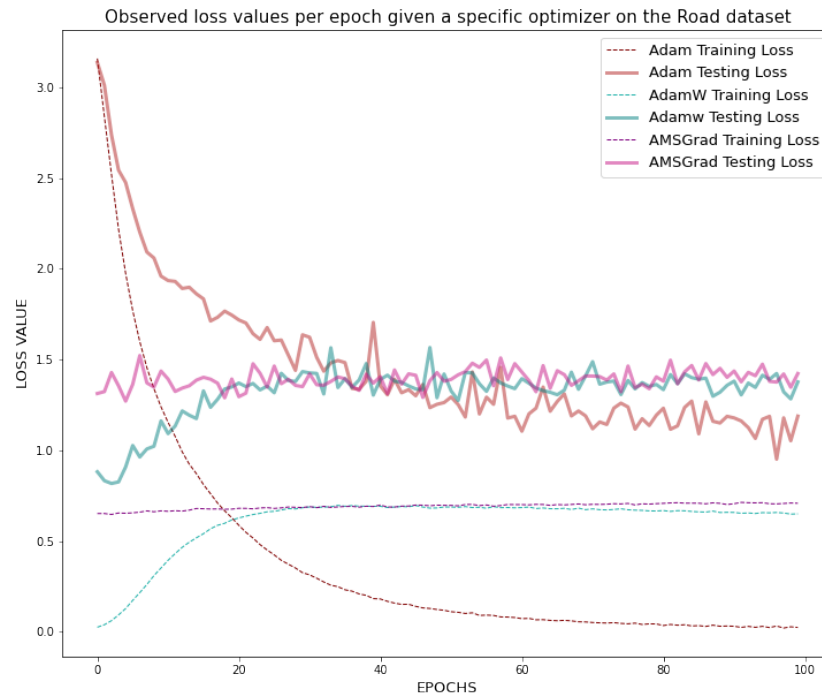


Figure 19: Results for the convolutional neural network trained on GTSRB

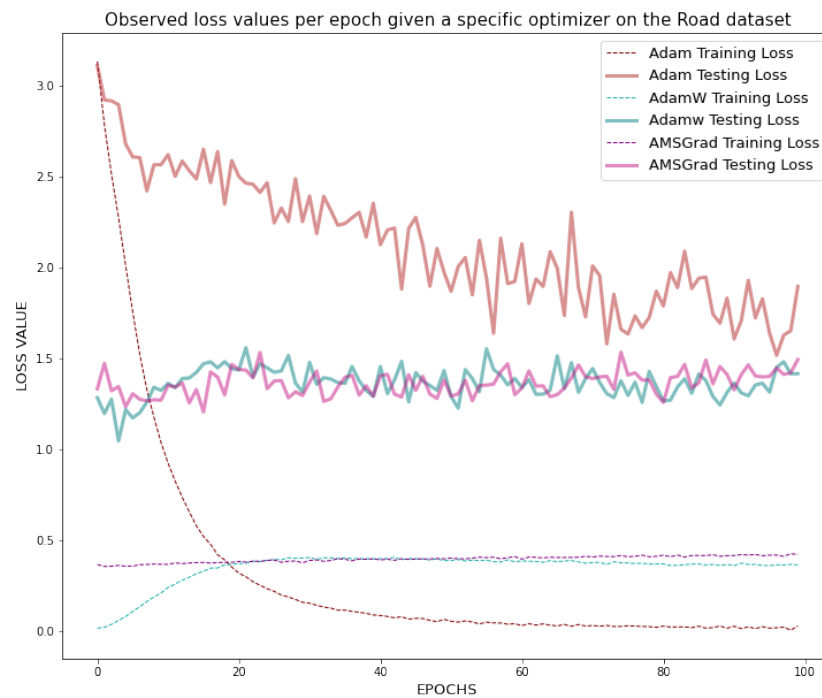


Figure 20: Results for the residual neural network trained on GTSRB

H OBSERVED LOSS VALUES FOR EACH OPTIMIZERS AND MODELS TRAINED ON THE CIFAR-10 DATASET

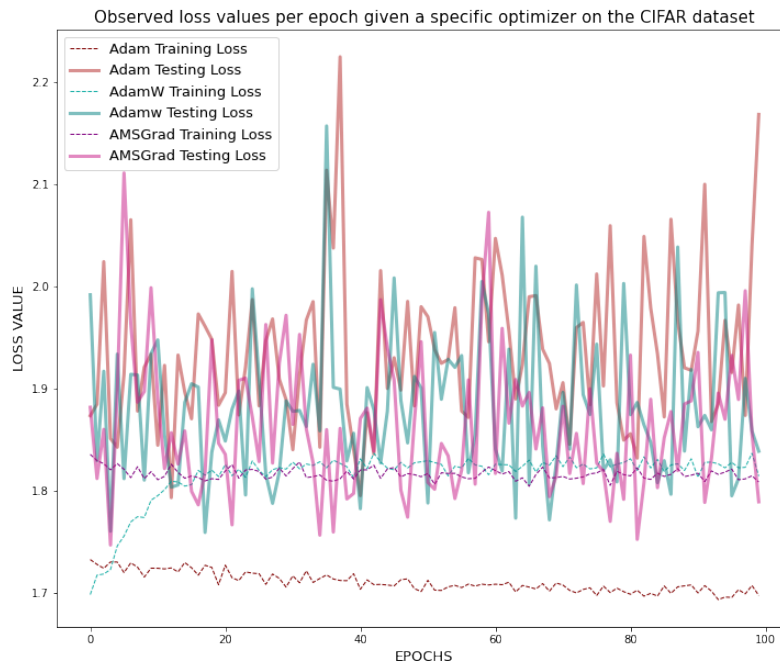


Figure 21: Results for the shallow neural network trained on CIFAR-10

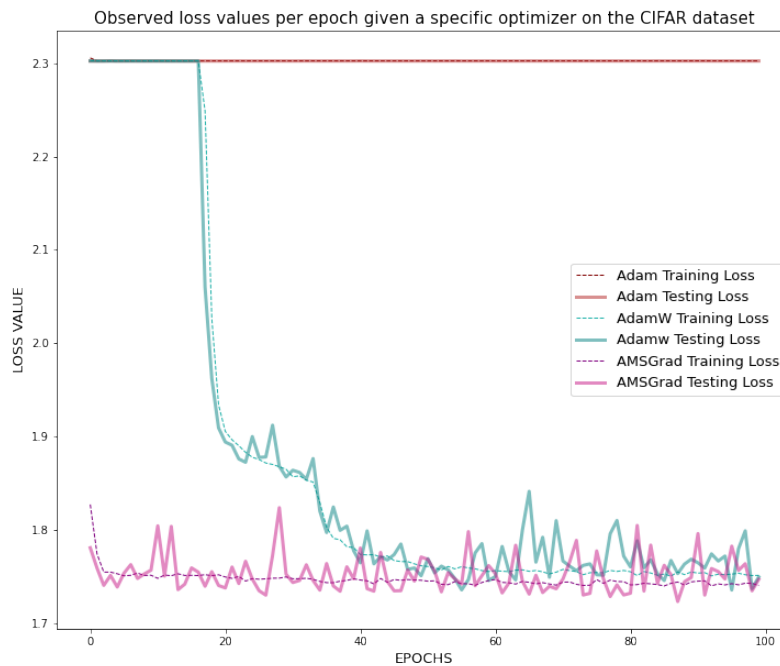


Figure 22: Results for the deep neural network trained on CIFAR-10

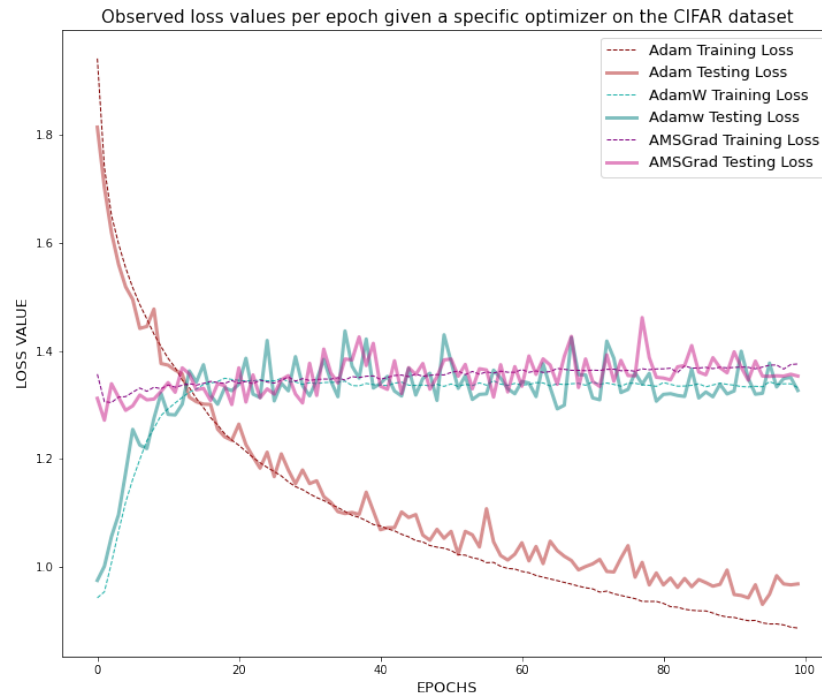


Figure 23: Results for the convolutional neural network trained on CIFAR-10

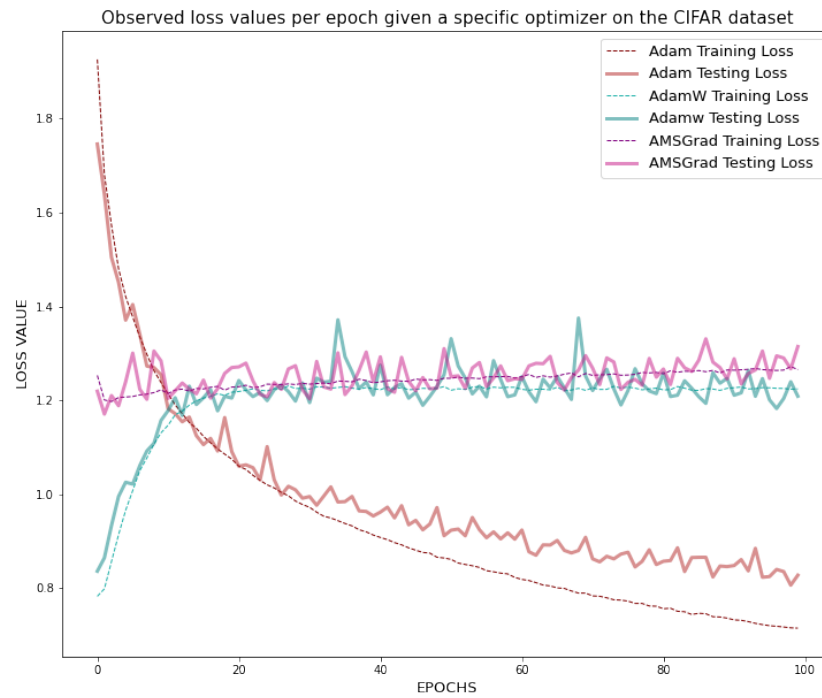


Figure 24: Results for the residual neural network trained on CIFAR-10

I OBSERVATIONS OF LOSS AND ACCURACY PER MODEL AND OPTIMIZER FOR THE MNIST DATASET

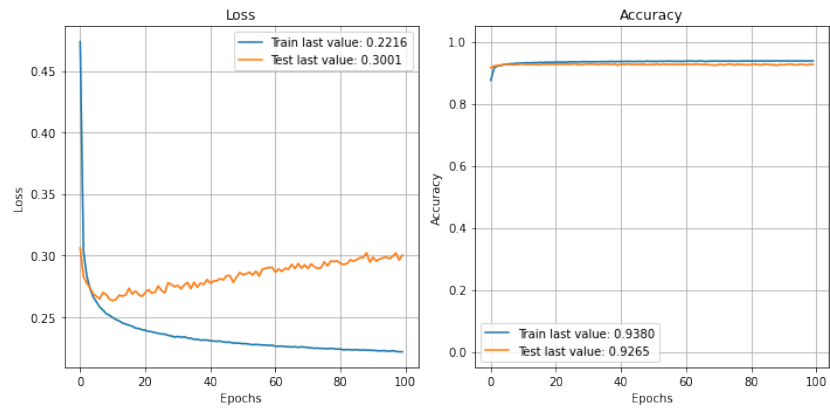


Figure 25: Loss and accuracy for the shallow neural network with Adam optimizer on MNIST

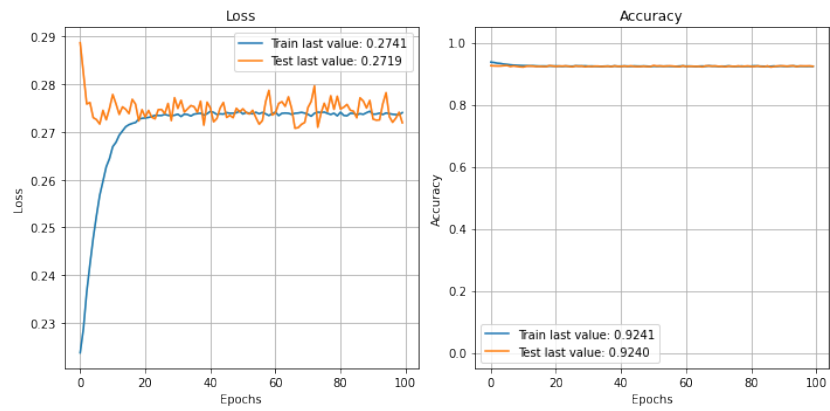


Figure 26: Loss and accuracy for the shallow neural network with AdamW optimizer on MNIST

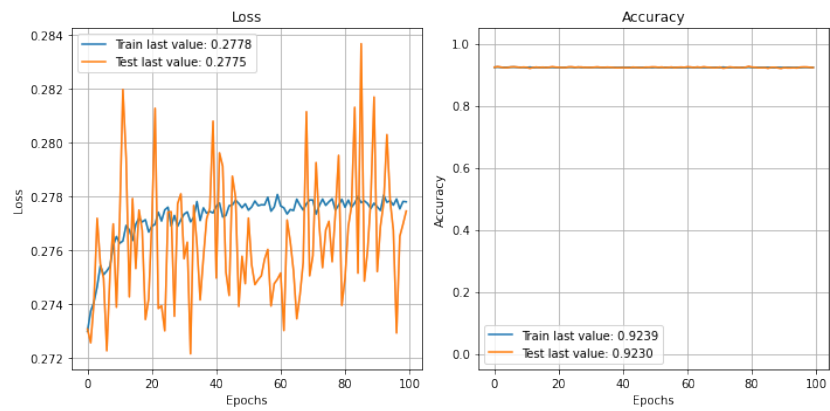


Figure 27: Loss and accuracy for the shallow neural network with AMSGrad optimizer on MNIST

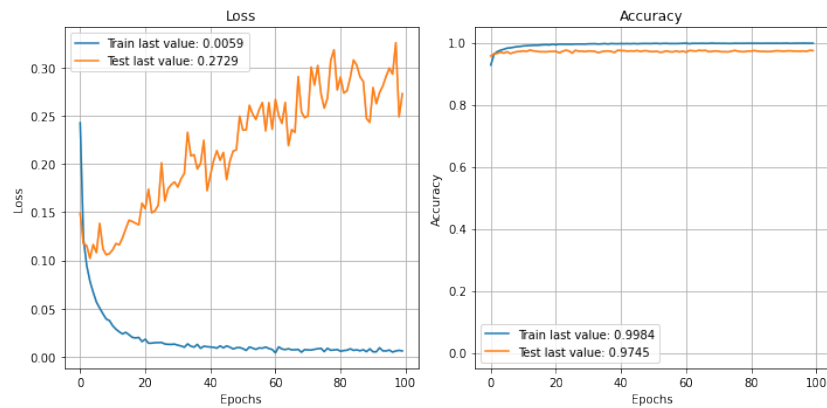


Figure 28: Loss and accuracy for the deep neural network Adam optimizer on MNIST

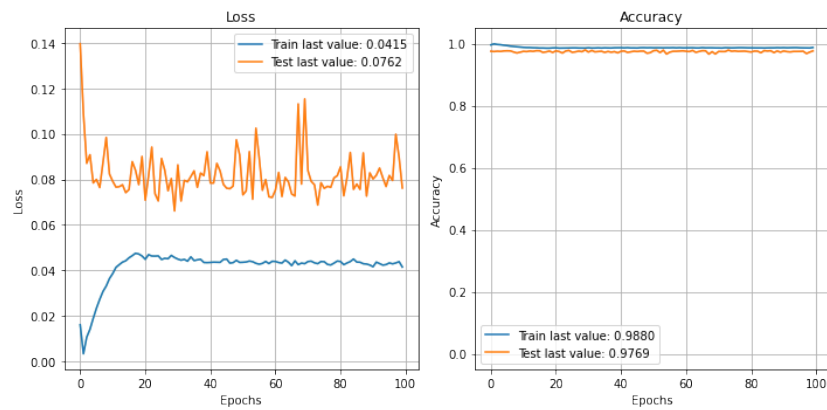


Figure 29: Loss and accuracy for the deep neural network with AdamW optimizer on MNIST



Figure 30: Loss and accuracy for the deep neural network with AMSGrad optimizer on MNIST

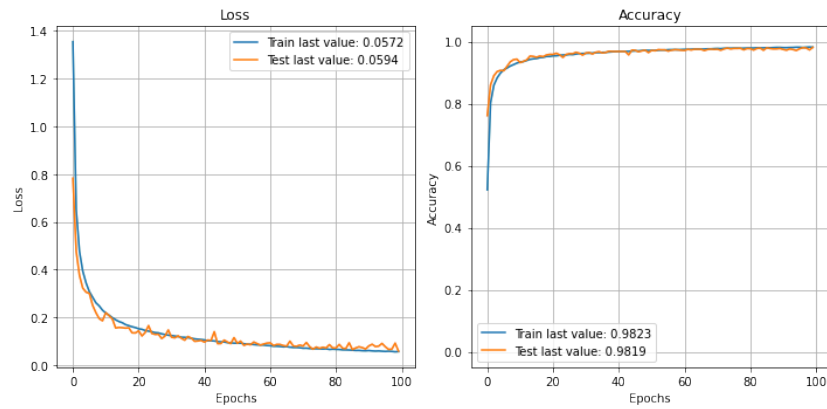


Figure 31: Loss and accuracy for the convolutional neural network with Adam optimizer on MNIST

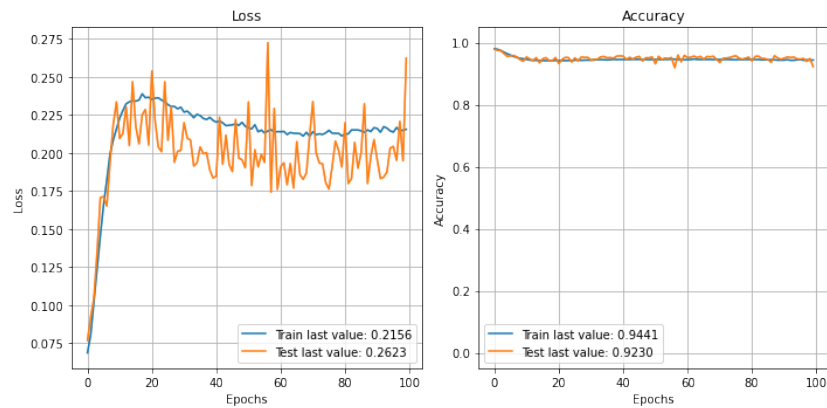


Figure 32: Loss and accuracy for the convolutional neural network with AdamW optimizer on MNIST

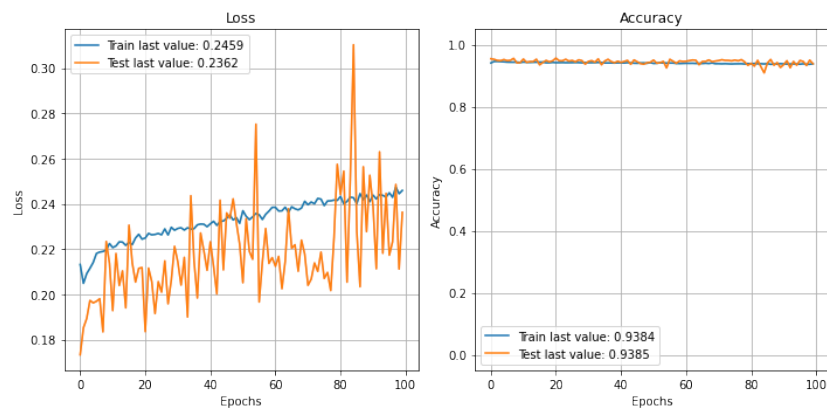


Figure 33: Loss and accuracy for the convolutional neural network with AMSGrad optimizer on MNIST

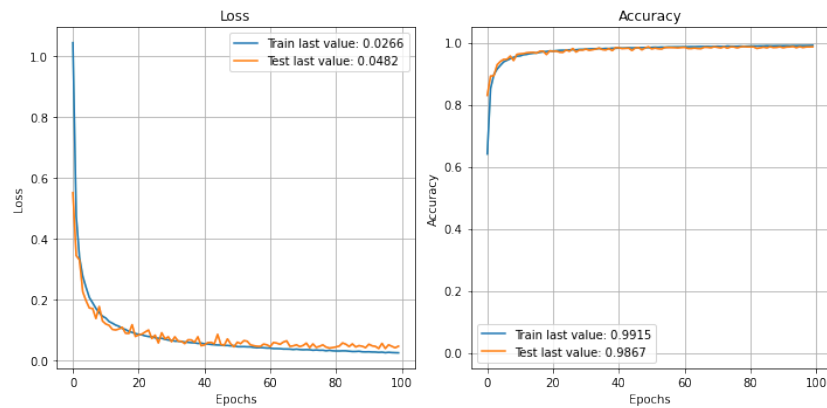


Figure 34: Loss and accuracy for the residual neural network with Adam optimizer on MNIST

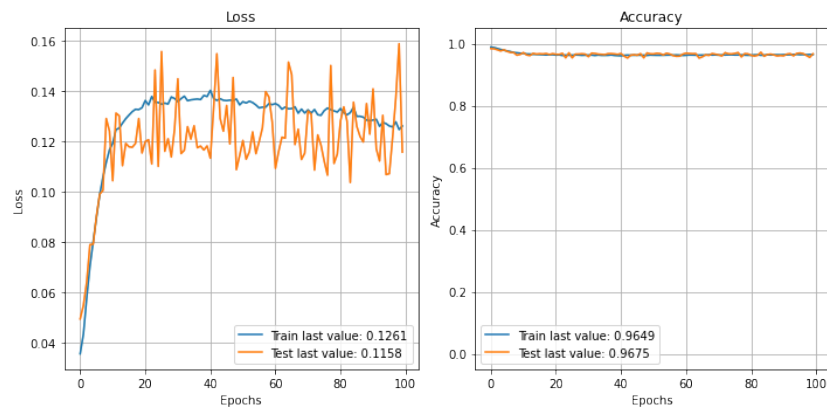


Figure 35: Loss and accuracy for the residual neural network with AdamW optimizer on MNIST

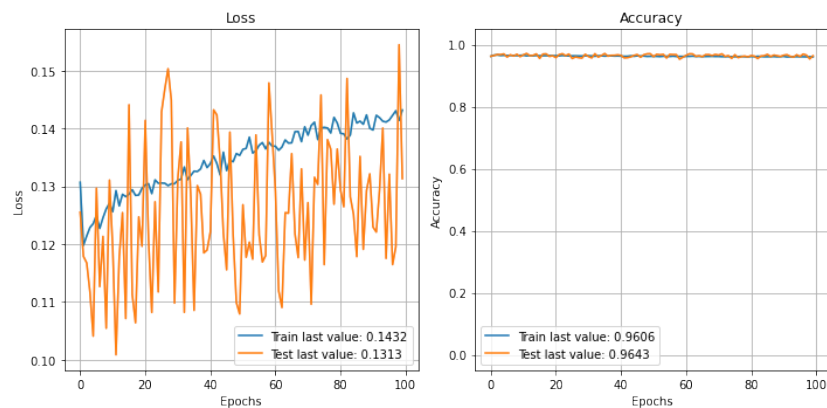


Figure 36: Loss and accuracy for the residual neural network with AMSGrad optimizer on MNIST

J OBSERVATIONS OF LOSS AND ACCURACY PER MODEL AND OPTIMIZER FOR THE FASHION-MNIST DATASET

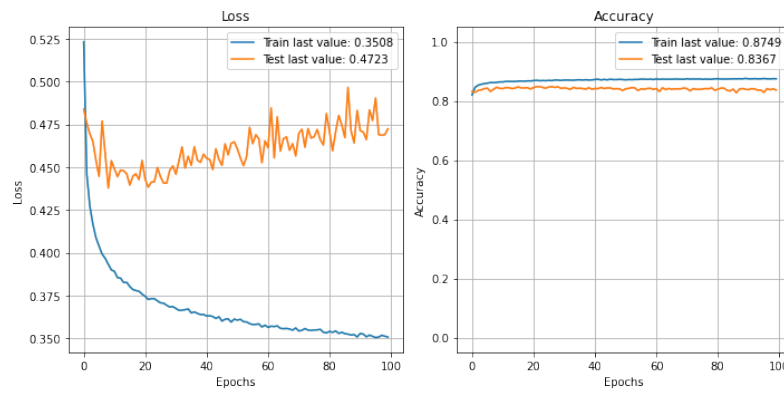


Figure 37: Loss and accuracy for the shallow neural network with Adam optimizer on Fashion-MNIST

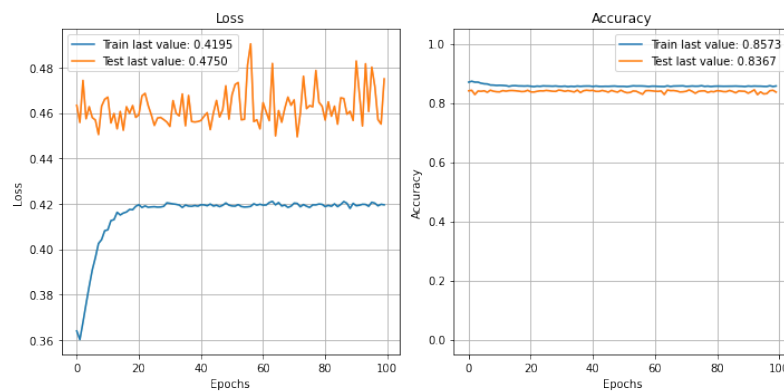


Figure 38: Loss and accuracy for the shallow neural network with AdamW optimizer on Fashion-MNIST

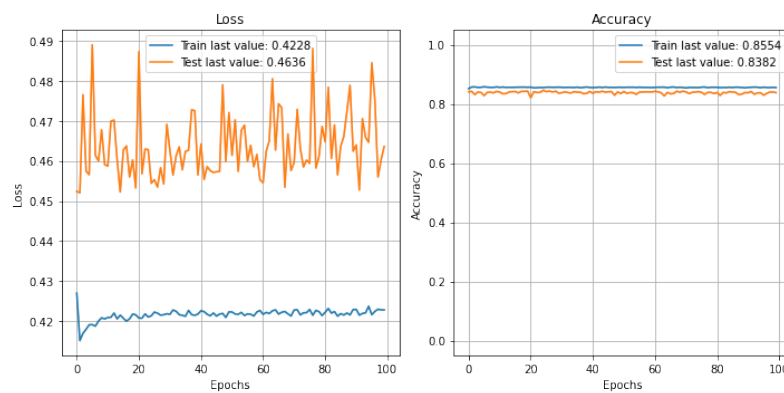


Figure 39: Loss and accuracy for the shallow neural network with AMSGrad optimizer on Fashion-MNIST

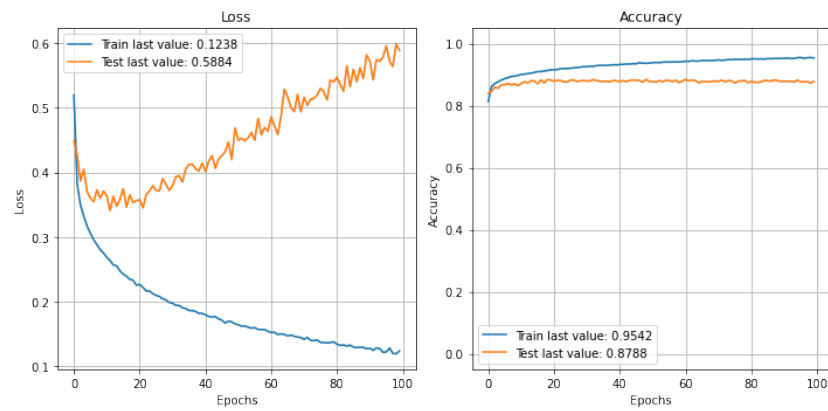


Figure 40: Loss and accuracy for the deep neural network Adam optimizer on Fashion-MNIST

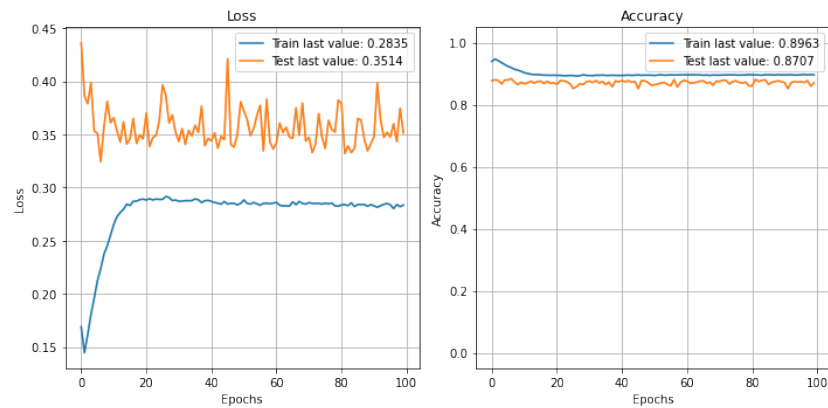


Figure 41: Loss and accuracy for the deep neural network with AdamW optimizer on Fashion-MNIST

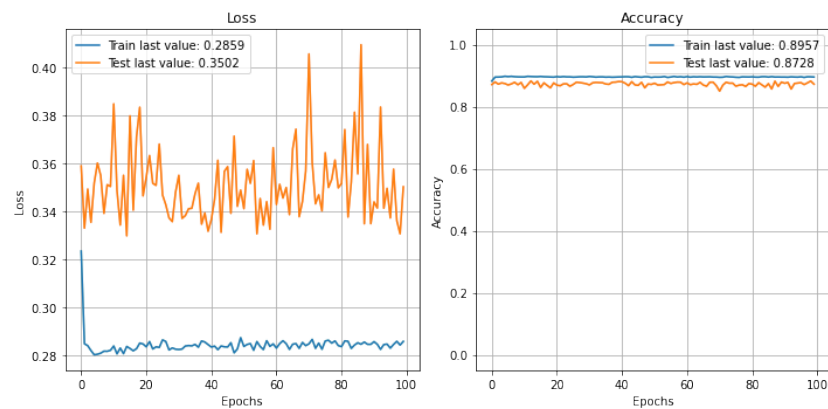


Figure 42: Loss and accuracy for the deep neural network with AMSGrad optimizer on Fashion-MNIST

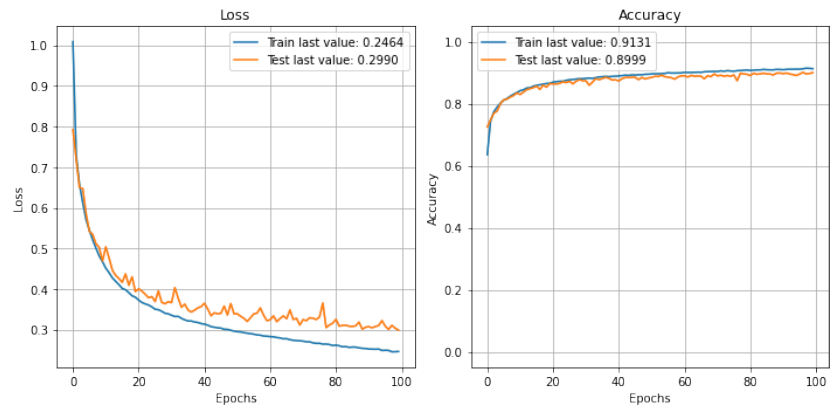


Figure 43: Loss and accuracy for the convolutional neural network with Adam optimizer on Fashion-MNIST

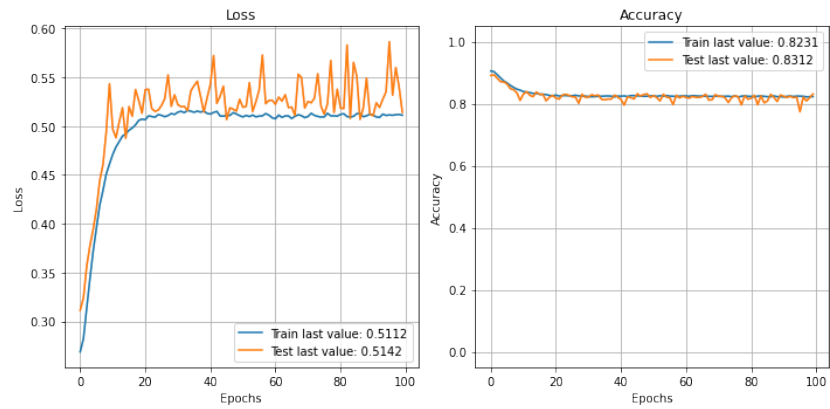


Figure 44: Loss and accuracy for the convolutional neural network with AdamW optimizer on Fashion-MNIST

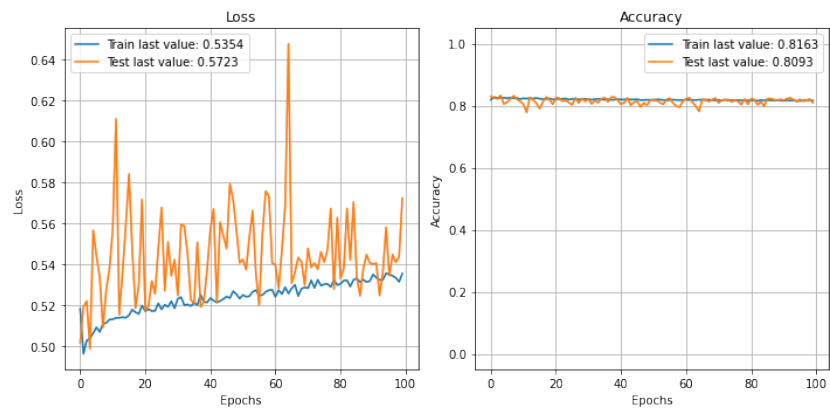


Figure 45: Loss and accuracy for the convolutional neural network with AMSGrad optimizer on Fashion-MNIST

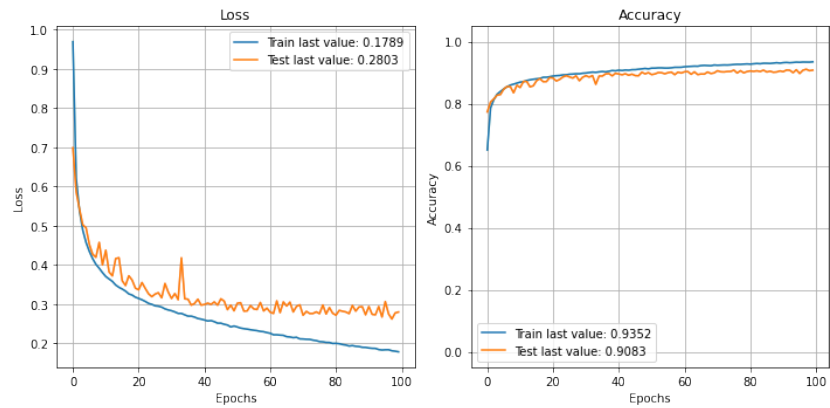


Figure 46: Loss and accuracy for the residual neural network with Adam optimizer on Fashion-MNIST

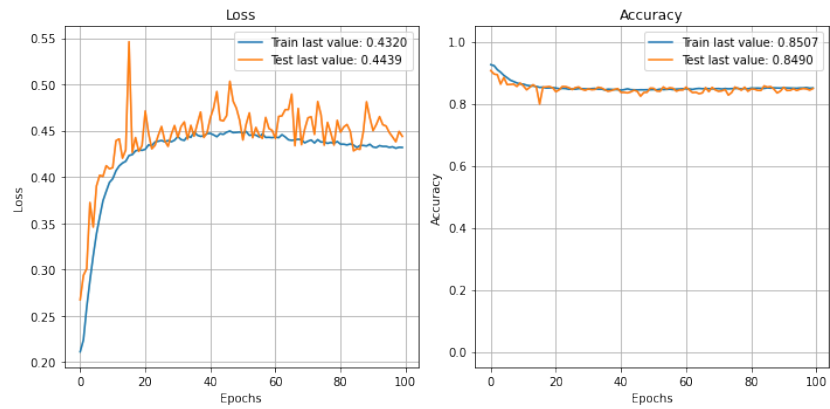


Figure 47: Loss and accuracy for the residual neural network with AdamW optimizer on Fashion-MNIST

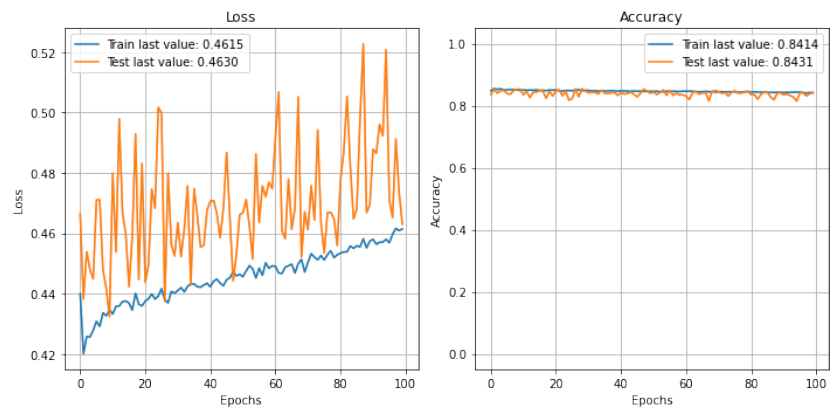


Figure 48: Loss and accuracy for the residual neural network with AMSGrad optimizer on Fashion-MNIST

K OBSERVATIONS OF LOSS AND ACCURACY PER MODEL AND OPTIMIZER FOR THE GERMAN TRAFFIC SIGN RECOGNITION BENCHMARK DATASET

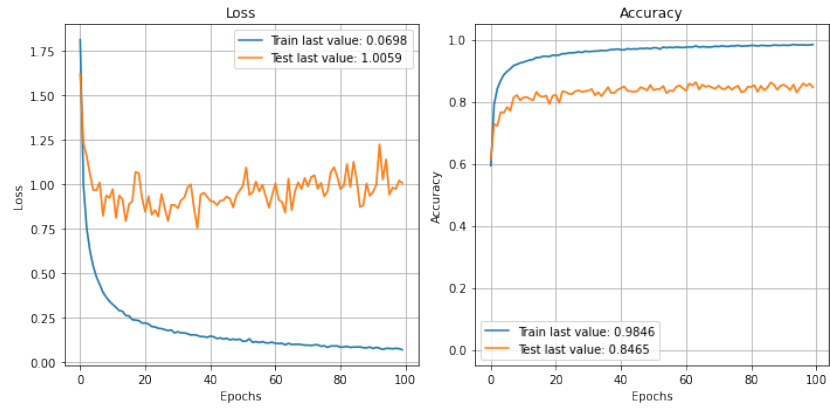


Figure 49: Loss and accuracy for the shallow neural network with Adam optimizer on GTSRB

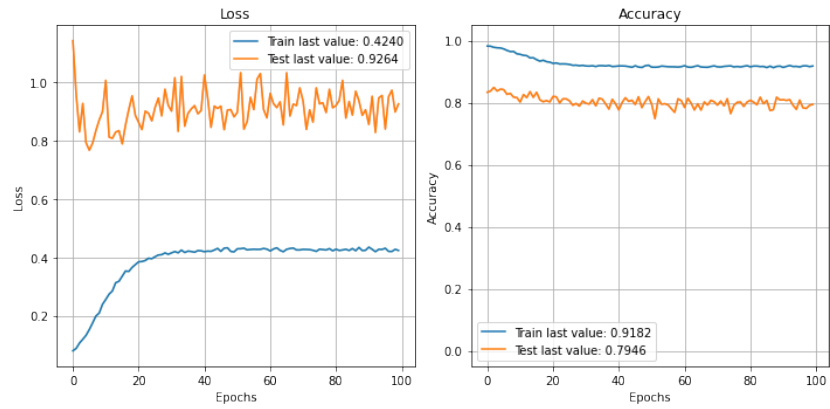


Figure 50: Loss and accuracy for the shallow neural network with AdamW optimizer on GTSRB

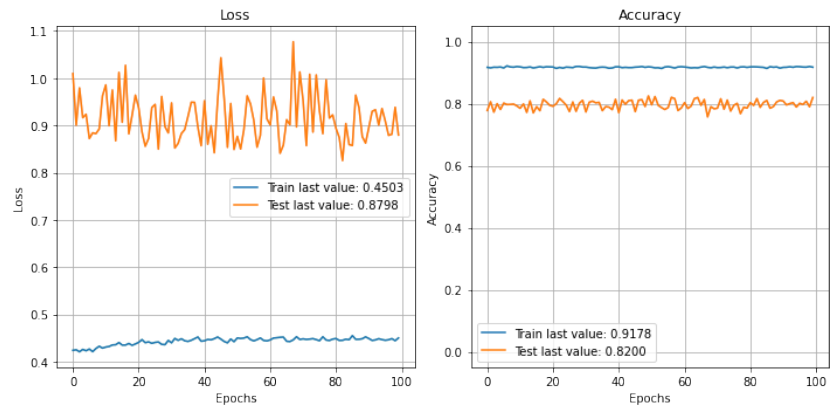


Figure 51: Loss and accuracy for the shallow neural network with AMSGrad optimizer on GTSRB

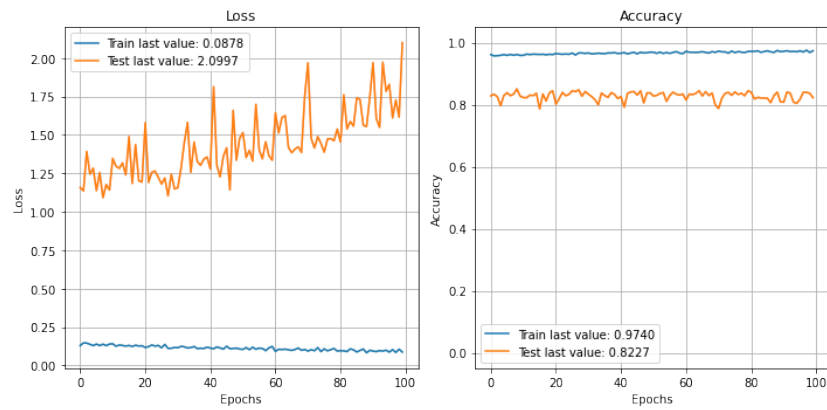


Figure 52: Loss and accuracy for the deep neural network Adam optimizer on GTSRB

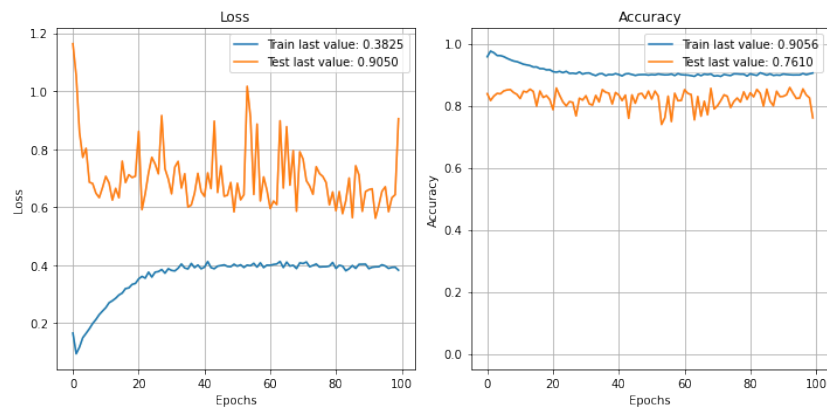


Figure 53: Loss and accuracy for the deep neural network with AdamW optimizer on GTSRB

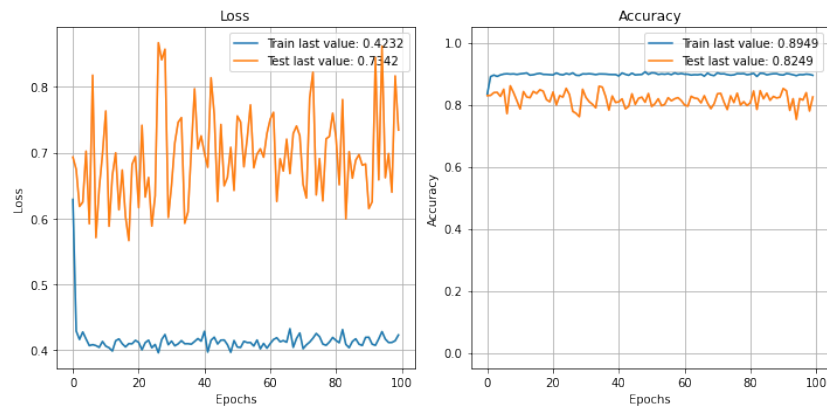


Figure 54: Loss and accuracy for the deep neural network with AMSGrad optimizer on GTSRB

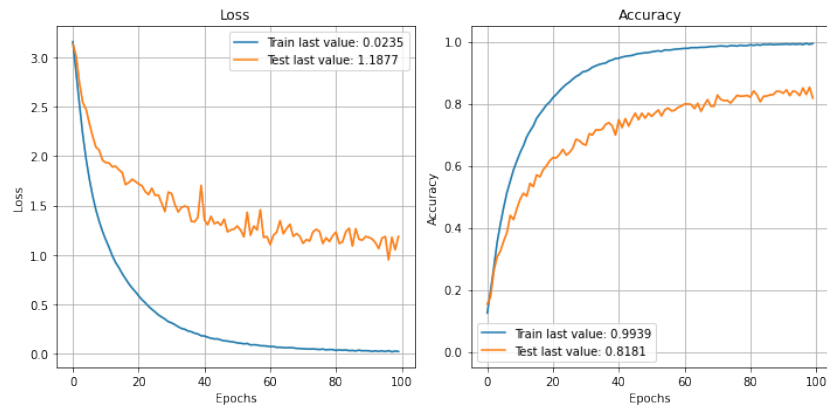


Figure 55: Loss and accuracy for the convolutional neural network with Adam optimizer on GTSRB

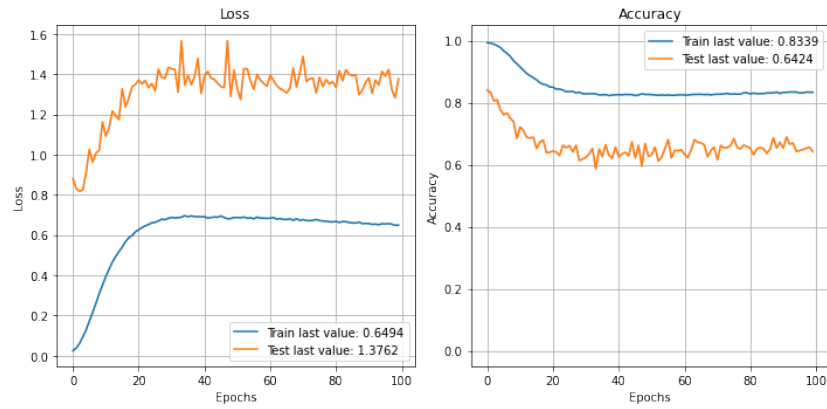


Figure 56: Loss and accuracy for the convolutional neural network with AdamW optimizer on GTSRB

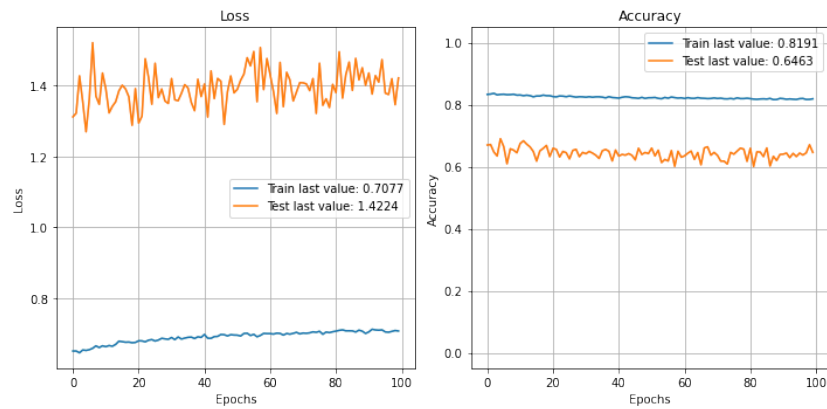


Figure 57: Loss and accuracy for the convolutional neural network with AMSGrad optimizer on GTSRB

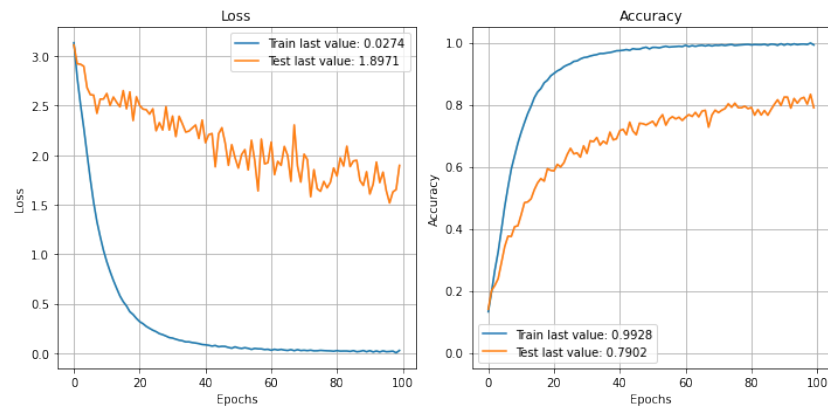


Figure 58: Loss and accuracy for the residual neural network with Adam optimizer on GTSRB

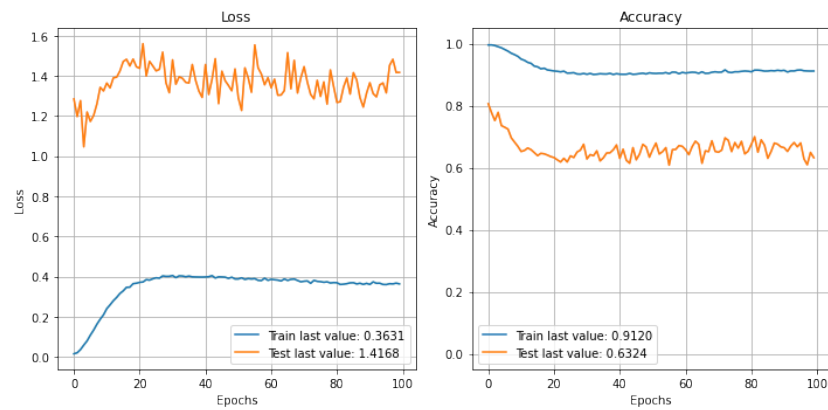


Figure 59: Loss and accuracy for the residual neural network with AdamW optimizer on GTSRB

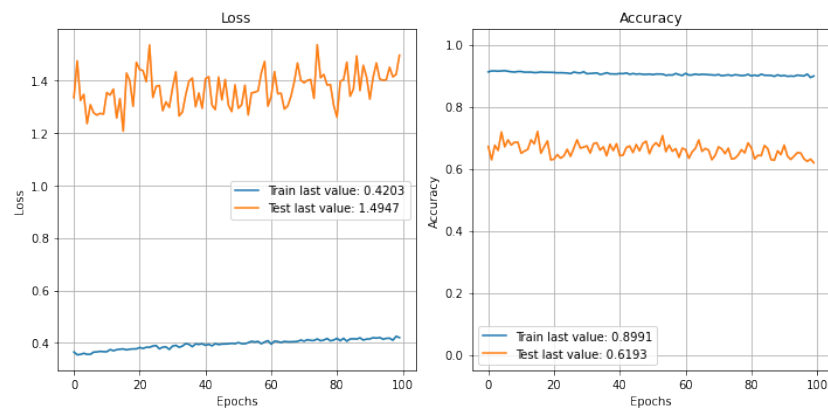


Figure 60: Loss and accuracy for the residual neural network with AMSGrad optimizer on GTSRB

L OBSERVATIONS OF LOSS AND ACCURACY PER MODEL AND OPTIMIZER FOR THE CIFAR-10 DATASET

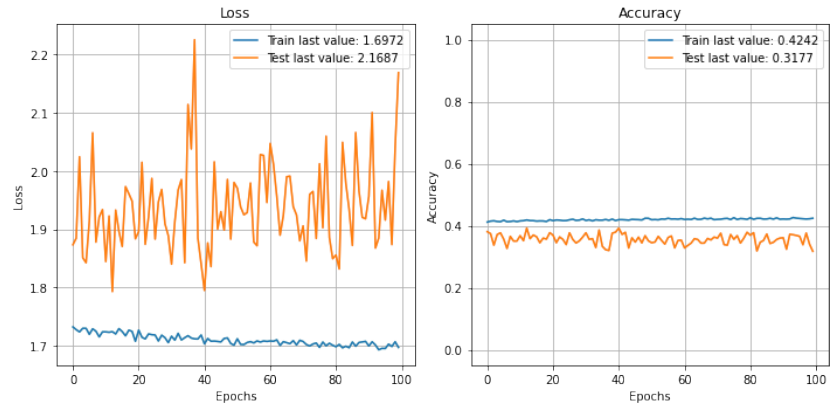


Figure 61: Loss and accuracy for the shallow neural network with Adam optimizer on CIFAR-10

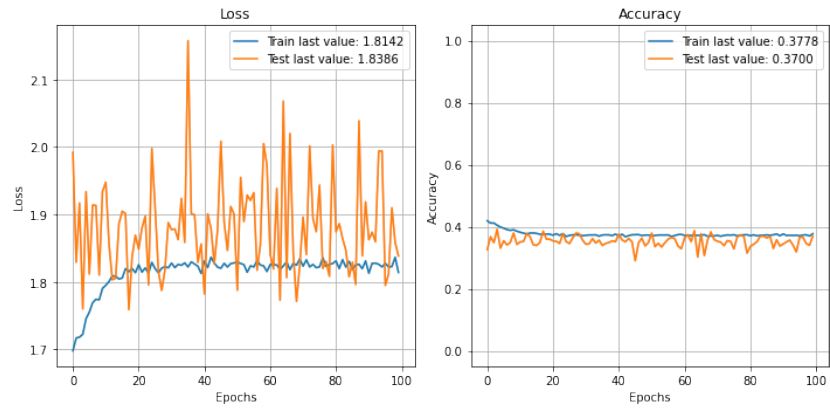


Figure 62: Loss and accuracy for the shallow neural network with AdamW optimizer on CIFAR-10

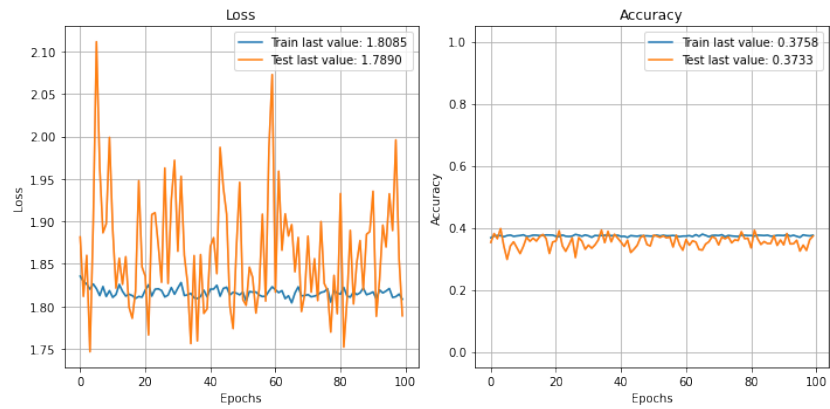


Figure 63: Loss and accuracy for the shallow neural network with AMSGrad optimizer on CIFAR-10

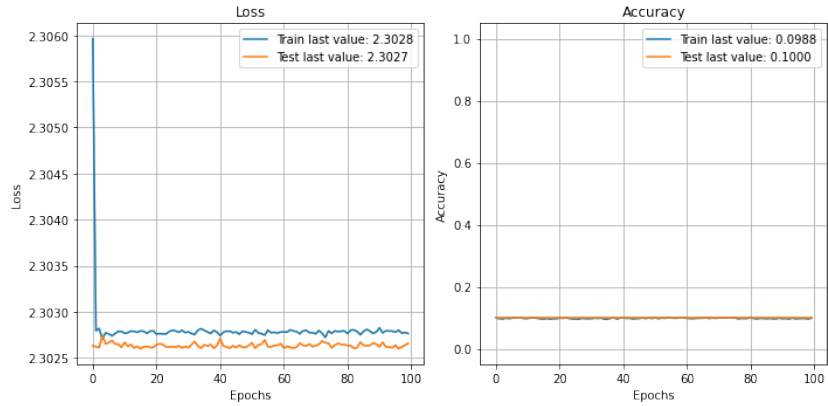


Figure 64: Loss and accuracy for the deep neural network Adam optimizer on CIFAR-10

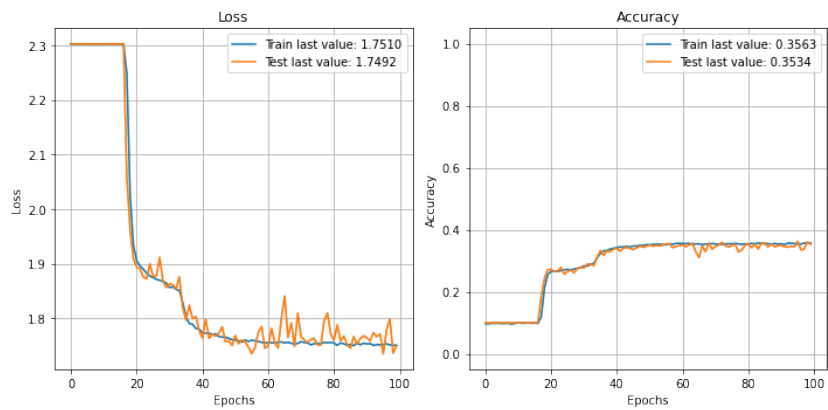


Figure 65: Loss and accuracy for the deep neural network with AdamW optimizer on CIFAR-10

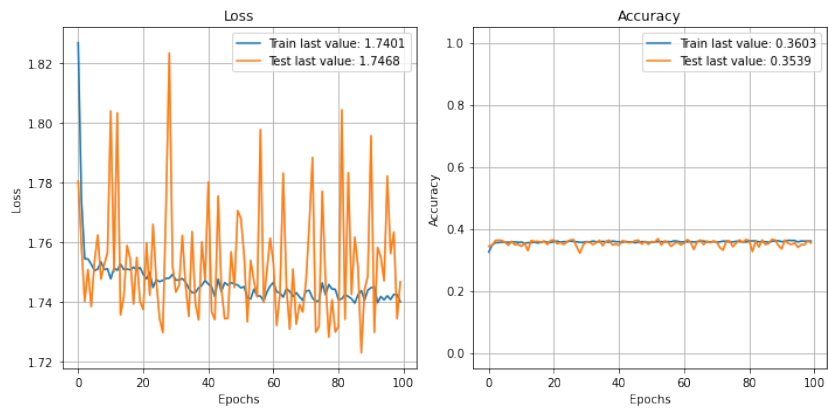


Figure 66: Loss and accuracy for the deep neural network with AMSGrad optimizer on CIFAR-10

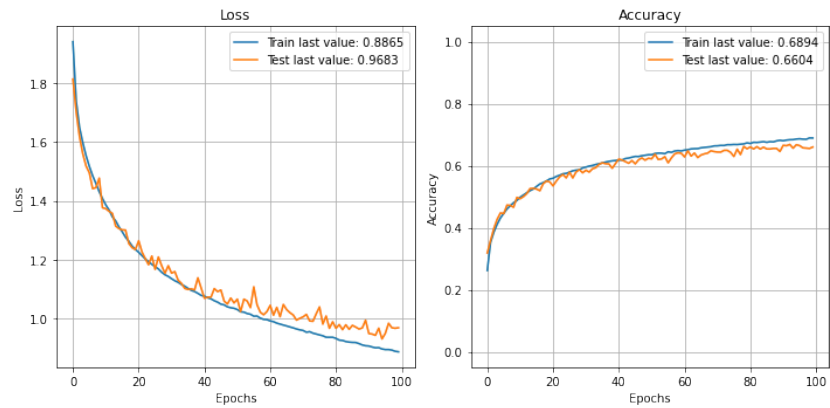


Figure 67: Loss and accuracy for the convolutional neural network with Adam optimizer on CIFAR-10

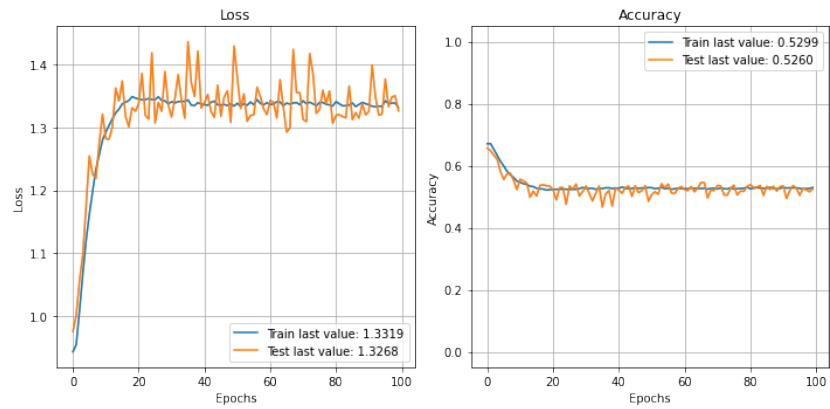


Figure 68: Loss and accuracy for the convolutional neural network with AdamW optimizer on CIFAR-10

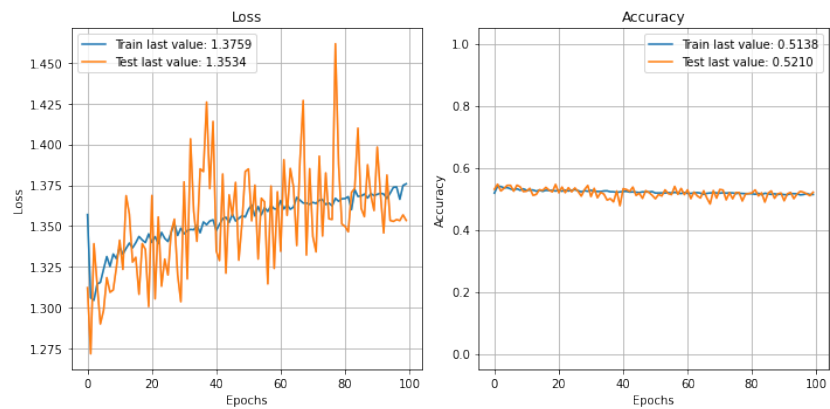


Figure 69: Loss and accuracy for the convolutional neural network with AMSGrad optimizer on CIFAR-10

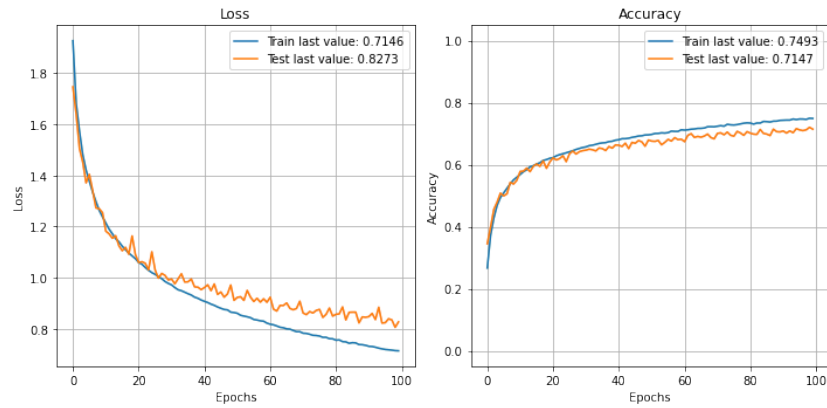


Figure 70: Loss and accuracy for the residual neural network with Adam optimizer on CIFAR-10

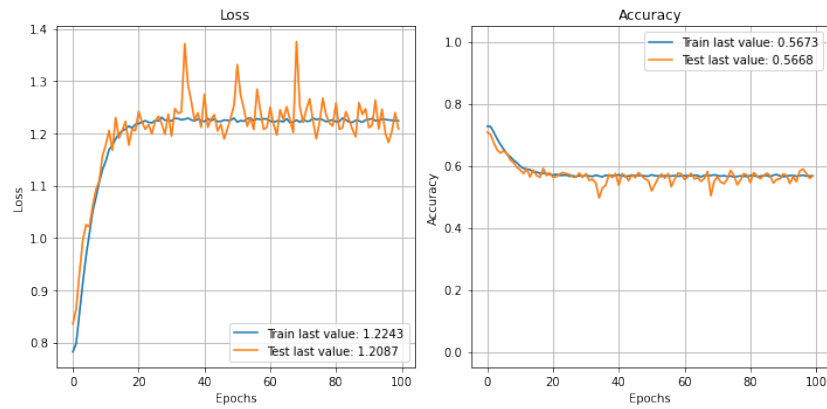


Figure 71: Loss and accuracy for the residual neural network with AdamW optimizer on CIFAR-10

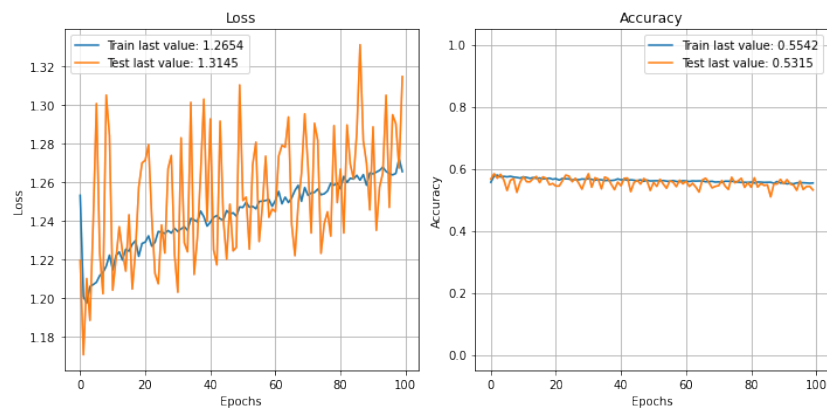


Figure 72: Loss and accuracy for the residual neural network with AMSGrad optimizer on CIFAR-10