

MIDDLE EAST TECHNICAL UNIVERSITY

MMI 727 HOMEWORK I

MURAT SEZGİN BALOĞLU

2469518



Derivatives of Activation Functions

$$\rightarrow \text{Sigmoid} \circ \frac{1}{1+e^{-x}}$$

$$\frac{d}{dx} \rightarrow \frac{(-1+e^{-x}) \cdot 0 + (e^{-x})}{(1+e^{-x})^2} = \boxed{\frac{e^{-x}}{(1+e^{-x})^2}}$$

$$\rightarrow \text{Tanh} \circ \frac{e^{2x}-1}{e^{2x}+1}$$

$$\frac{d}{dx} \rightarrow \frac{(2e^{2x})(e^{2x}+1) - (2e^{2x}(e^{2x}-1))}{(e^{2x}+1)^2} = \boxed{\frac{4e^{4x}}{e^{4x}+2e^{2x}+1}}$$

$$\rightarrow \text{ReLU} \circ \max(0, x)$$

$$\frac{d}{dx} \rightarrow \begin{cases} 0, & x \leq 0 \\ 1, & x > 0 \end{cases}$$

→ Softmax 8

$$(x_{out})_i = \frac{e^{-\beta(x_{in})_i}}{\sum_k e^{-\beta(x_{in})_k}}$$

$$\frac{d(x_{out})_i}{d(x_{in})_j} = \text{for } i \neq j \text{ input pairs there are two general cases: } i=j \text{ and } i \neq j$$

$i=j$ case 8

$$\begin{aligned} & \frac{-\beta(x_{in})_i}{e} \sum_k e^{-\beta(x_{in})_k} - \frac{-\beta(x_{in})_j}{e} + \frac{-\beta(x_{in})_i}{e} \\ & \left(\sum_k e^{-\beta(x_{in})_k} \right)^2 \\ & = \frac{-\beta(x_{in})_i}{e} \frac{\sum_k e^{-\beta(x_{in})_k} - e^{-\beta(x_{in})_j}}{\sum_k e^{-\beta(x_{in})_k}} \\ & = \frac{-\beta(x_{in})_i}{\sum_k e^{-\beta(x_{in})_k}} \left(1 - \frac{e^{-\beta(x_{in})_j}}{\sum_k e^{-\beta(x_{in})_k}} \right) \quad \text{for } i=j \end{aligned}$$

17 J case 8

$$= \frac{\beta(x_m) - \beta(x_n)}{e^{\beta(x_m)} - e^{\beta(x_n)}}$$

$$= \frac{\sum_{k=1}^n e^{-\beta(x_m)_k} - \sum_{k=1}^m e^{-\beta(x_n)_k}}{\left(\sum_{k=1}^n e^{-\beta(x_m)_k}\right)^2}$$

⇒ The term $\beta(x_m)$ is not understood clearly. Is β a constant to $(x_m)_i$ or is $\beta(x_m)$ is a output vector of function β and $\beta(x_m)_i$ is i^{th} element of this vector.
Hence β ignored.

Q 1) Part B :

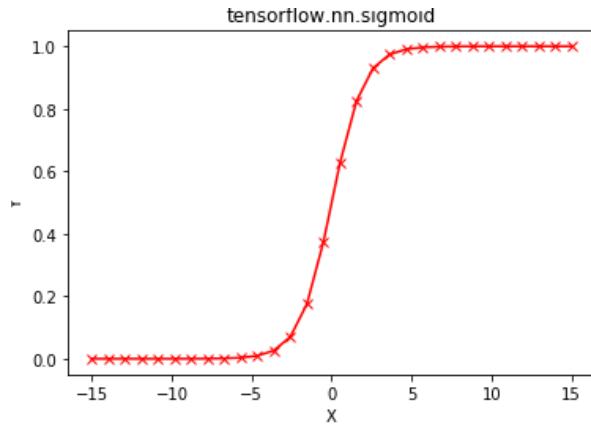


Figure 1 Sigmoid Activation Function

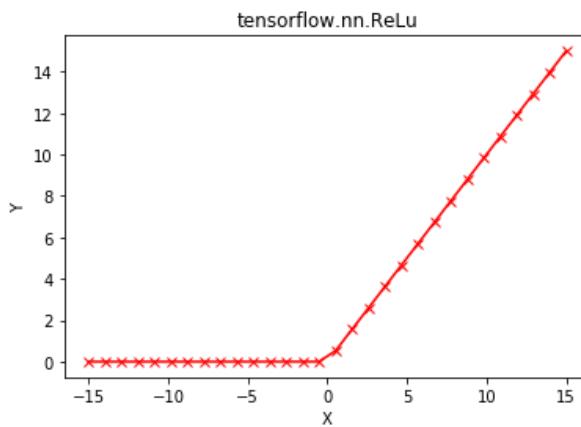


Figure 2 ReLu Activation Function

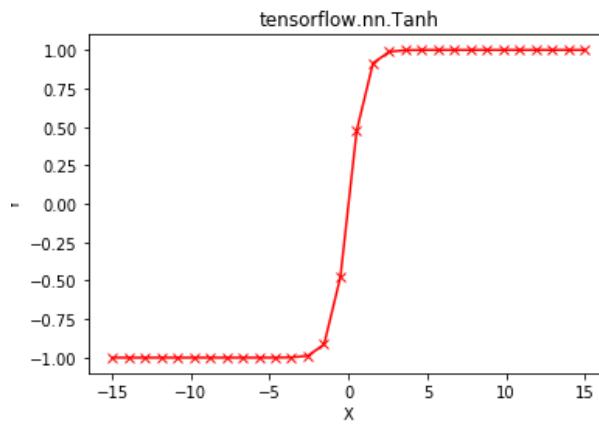


Figure 3Tanh Activation Function

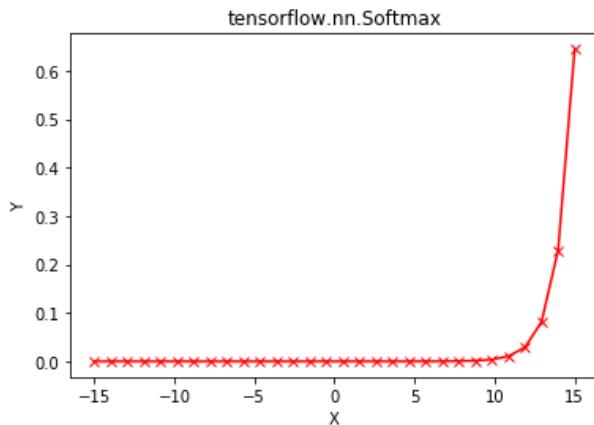


Figure 4 Softmax Activation Function

Q 1) Part C :

Sigmoid :

Advantages :

- + Easy to understand.
- + Output is bounded.
- + Differentiable everywhere

Disadvantages :

- Vanishing gradient problem
- Slow convergence
- Output is not zero centered

Where to use?

- > Popular with classification problems
- > Usually used in the output layer of a binary classification problem where result is either 0 or 1.

Tanh:

Advantages :

- + Output is zero centered
- + Output is bounded
- + Gradients are steeper with respect to sigmoid

Disadvantages :

-Vanishing gradient problem

Where to use?

Usually used in hidden layers of a NN as its values lies between -1 to +1 hence mean of hidden layer approximates to 0. This makes learning for the next layer much easier.

ReLU:

Advantages:

+Avoids vanishing gradients problem.

+Computationally less expensive compared to sigmoid and tanh.

Disadvantages :

-Gradient may die during training.

-Output is unbounded.

-It could result in dead neurons, substantial part of network may go passive.

Where to use?

It should be used in hidden layers.

Widely used in many cases especially in CNNs.

Softmax:

Advantages:

+Able to handle multiple classes.

Where to use?

Usefull for output neurons.

Widely used in multiple classification logistic regression model.

2. PARAMETER TUNING

A. LEARNING RATE EFFECT

Bigger learning rate increases the speed of convergence but may result in unsteady system. As it can be seen in figures 7-8. In figure 7, with learning rate 50, accuracy oscillated below 0.6 moreover, as learning rate increase further to 500, accuracy oscillated below 0.5 meaning larger oscillations.

On the other hand, small learning rate such as 0.005 resulted in slow convergence and it is clear in figure 6 that NN is not converged yet as of 1000th iterations.

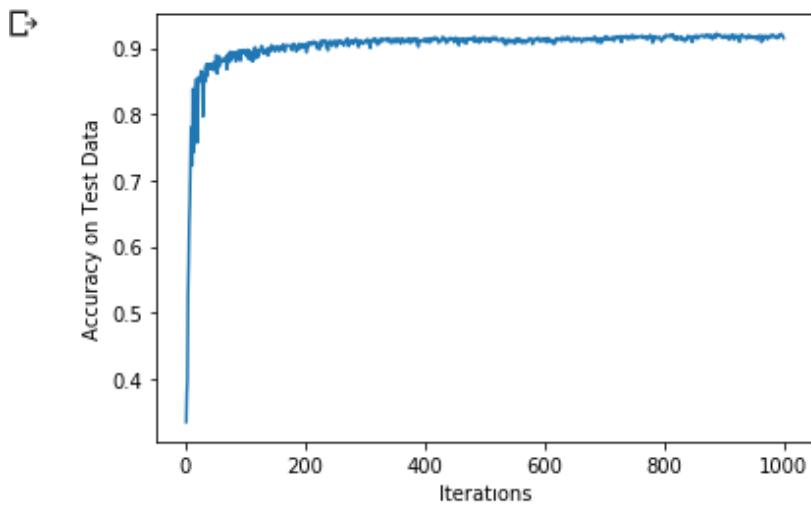


Figure 5 Learning rate 0.5

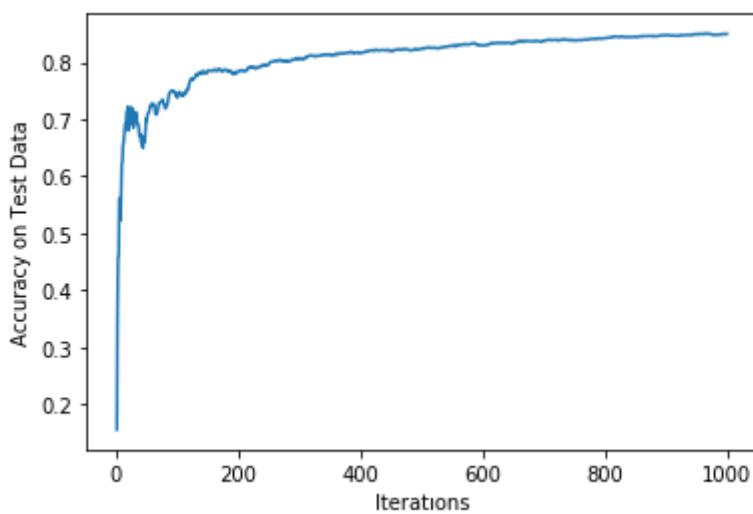


Figure 6 Learning rate 0.005

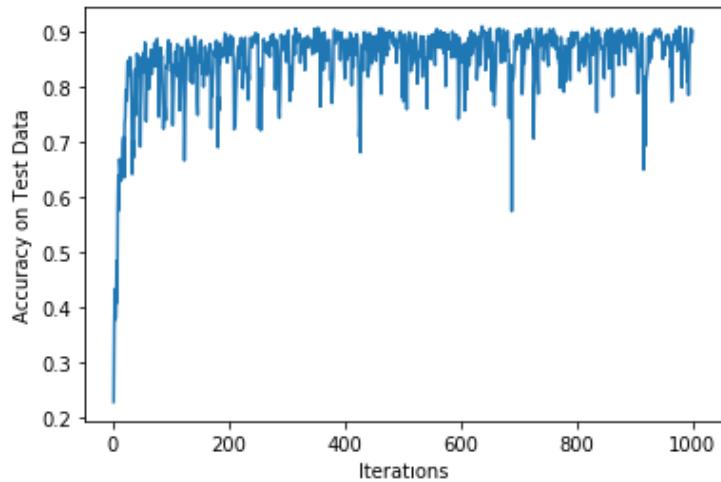


Figure 7 Learning Rate 50

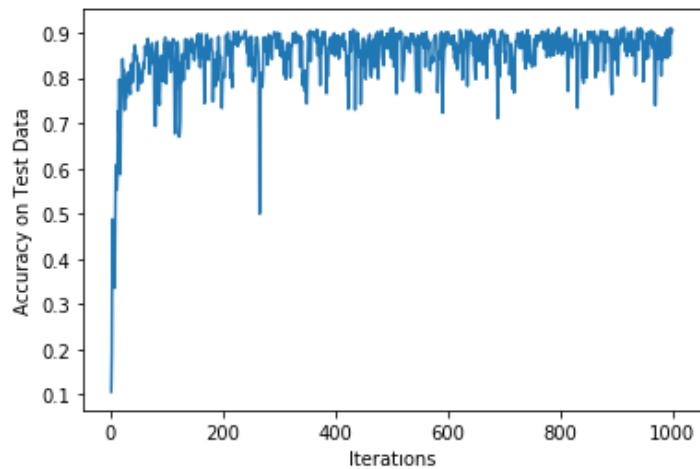


Figure 8 Learning Rate 500

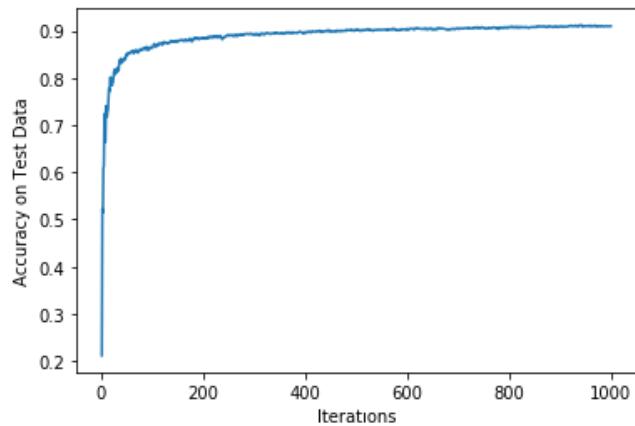


Figure 9 Learning rate 0.1

B) DIFFERENT OPTIMIZERS

Optimizer change did not result in change in computation time, convergence rate is also not drastically changed. However, Gradient Descent is more stable than compared to Adagrad and Adam. Adam optimizer's performance is weaker compared to Gradient Descent and Adagrad Optimizers.

As Momentum optimizer's momentum increase, NN behaves more unstable, this can be compensated by decreasing the learning rate but this time convergence rate would also decreased.

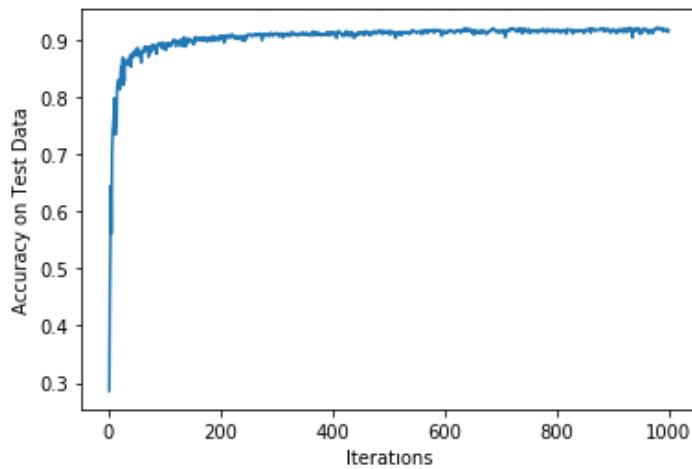


Figure 10 Default Gradient Descent Optimizer Duration 16.036s

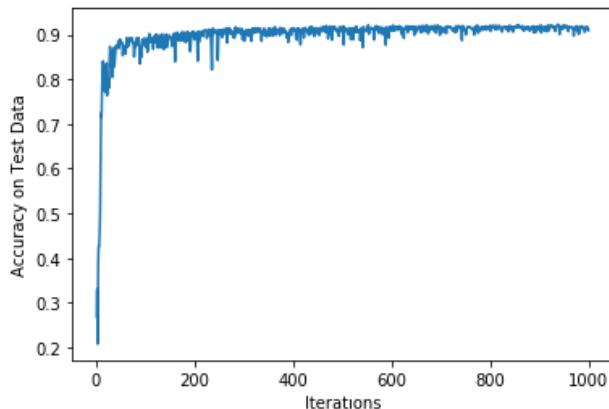


Figure 11 Adagrad Opt. Duration 16.236 s

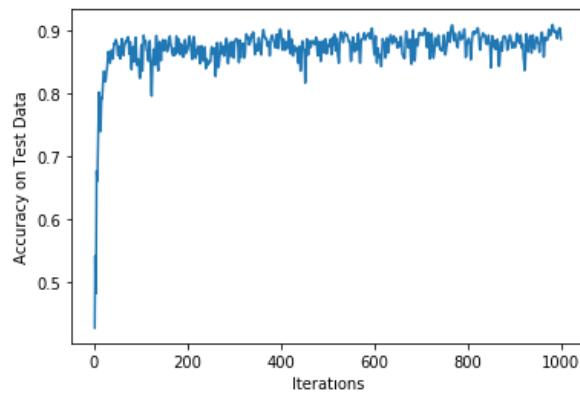


Figure 12 Adam opt. Duration 16.468s

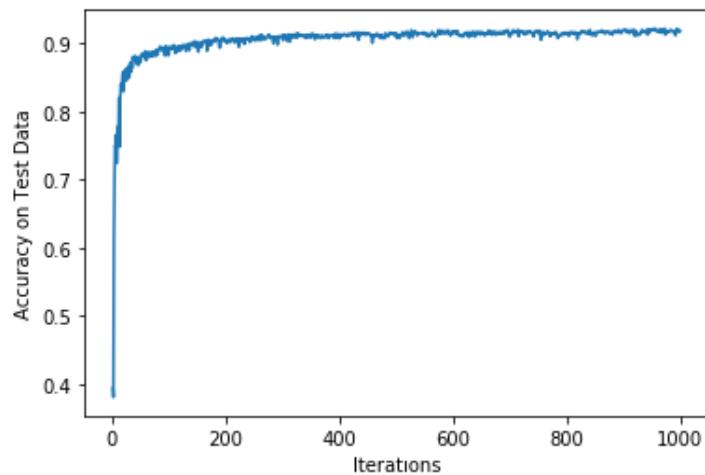


Figure 13 Momentum opt. lr 0.5 momentum 0.1 Duration 16.345 s

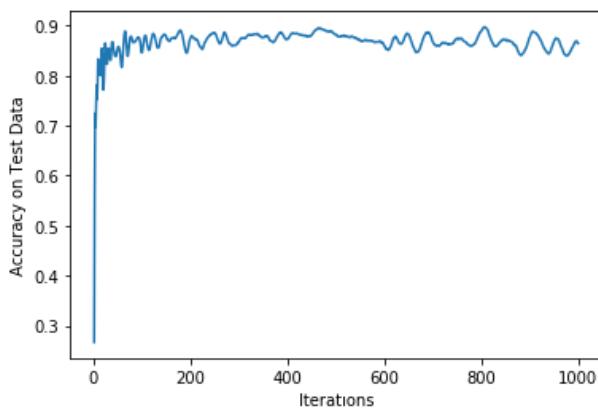


Figure 14 Momentum opt. lr 0.5 momentum 1Duration 16.257 s

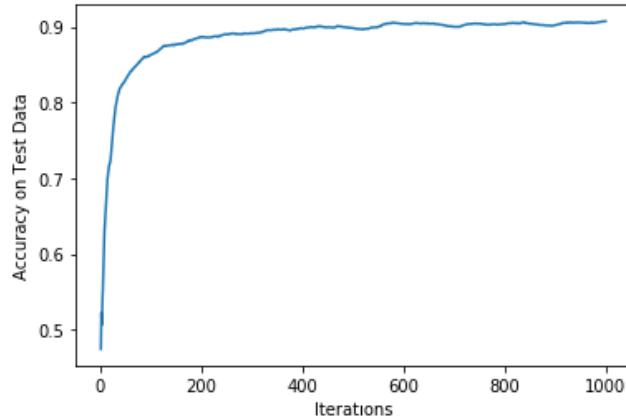


Figure 15 momentum opt. lr 0.005 momentum 1 Duration 16.778s

C) BATCH SIZE

- The plots show that small batch results generally in rapid learning but a unstable learning process with higher variance in the classification accuracy. Larger batch sizes slow down the learning process but the final stages result in a convergence to a more stable model exemplified by lower variance in classification accuracy.

Moreover, as batch size increased, accuracy is slightly increased.

Smaller batch sizes resulted in unstable networks, learning rates are decreased in order to compensate this but this led to larger convergence rates as can be seen in figures 17 and 18.

Final Accuracy is : 0.9182

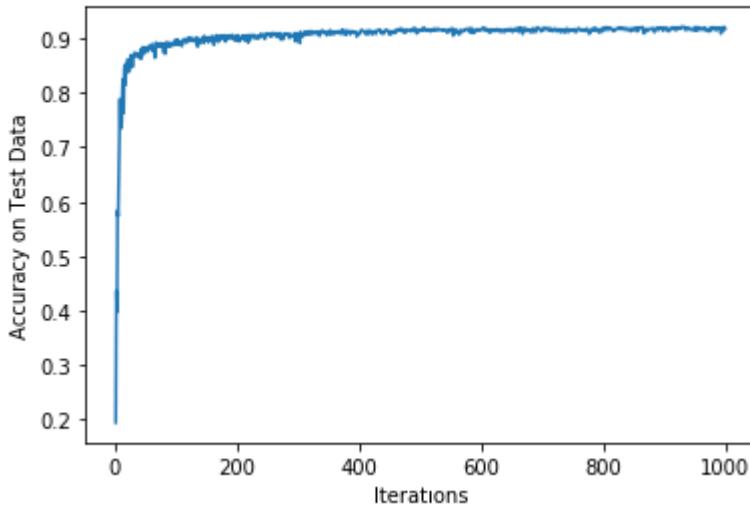


Figure 16 default one batch size = 100 Duration 16s

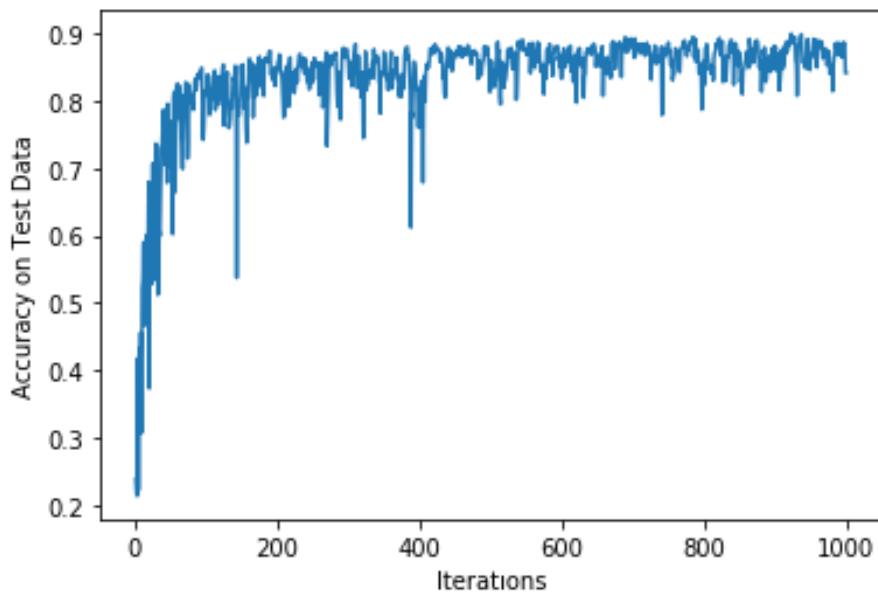


Figure 17 Batch size = 10 lr = 0.5 Duration 16s

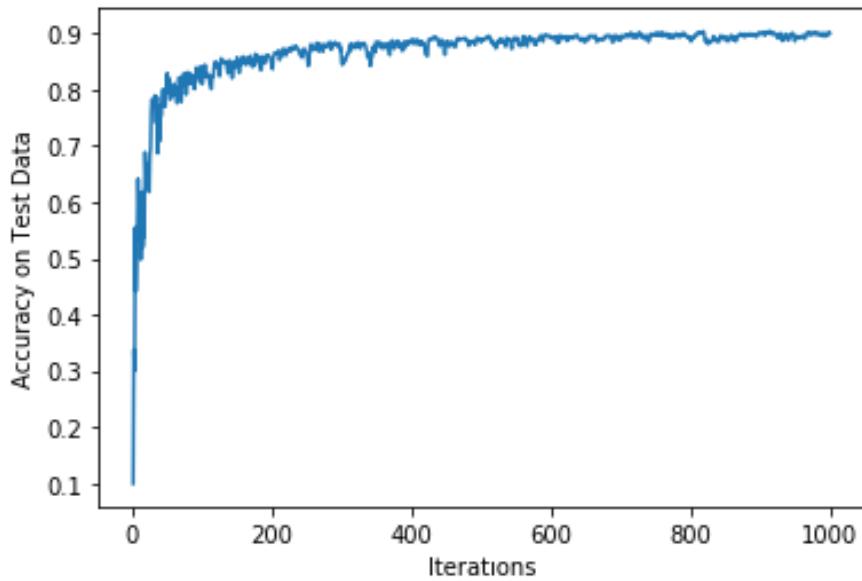


Figure 18 Batch size 10 learning rate decreased to 0.1

Final Accuracy is : 0.921

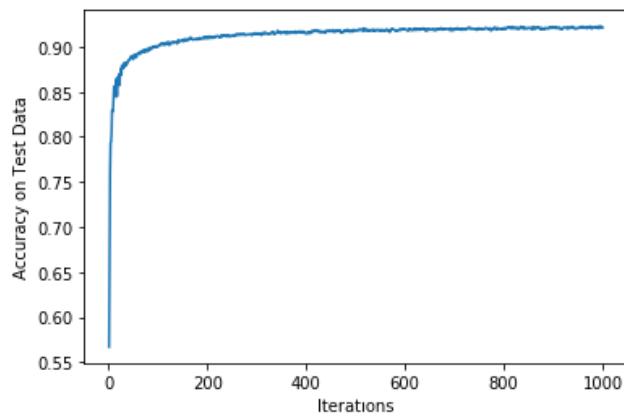


Figure 19 Batch size 1000 learning rate = 0.5 Duration 20,48s

Final Accuracy is : 0.9215

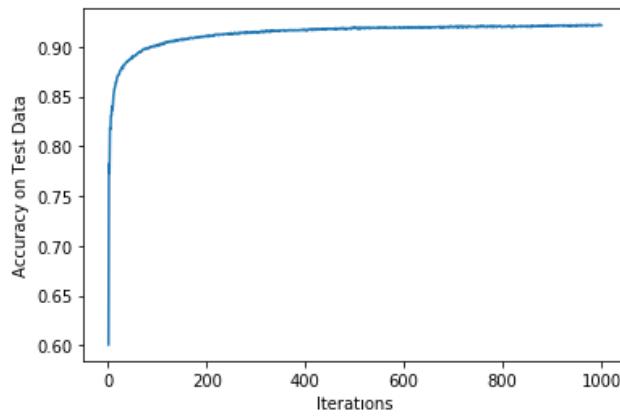


Figure 20 Batch size 10000 learning rate = 0.5 Duration 63s

Final Accuracy is : 0.9218

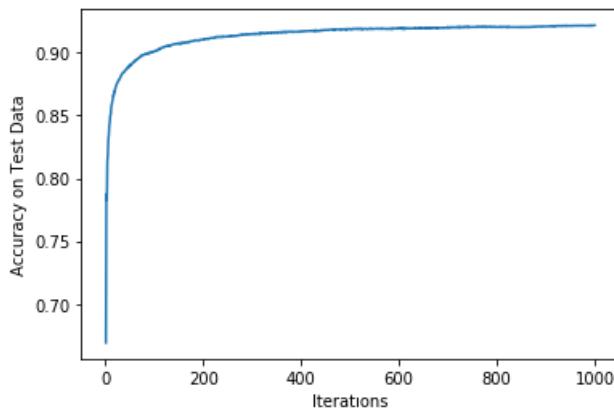


Figure 21 Batch size 50000 Learning Rate 0.5 Duration 270s

3) NETWORK MODEL

A) ORIGINAL SINGLE LAYER NETWORK

Final Accuracy is : 0.9197

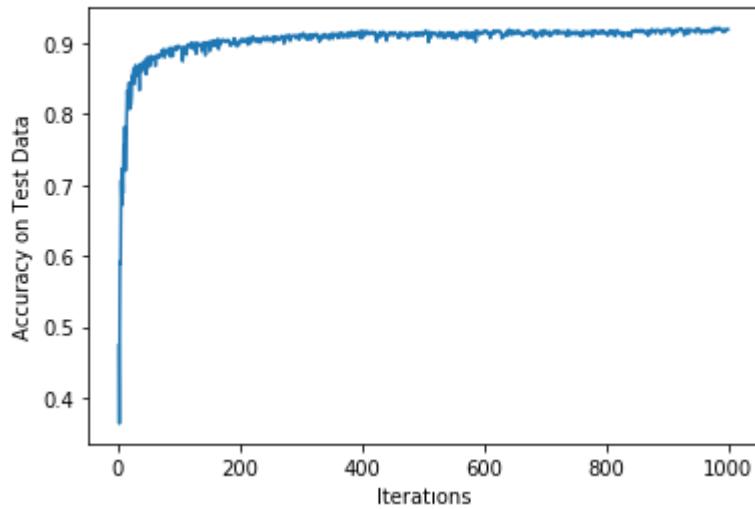


Figure 22 Original Single Layer Network

B) TWO LAYER NETWORK

Two layer network didn't work with its previous hyper parameters because weights in the first layer change very little and error propagates. Implemented solution is to initialize weights with random uniform distribution.

0.1135

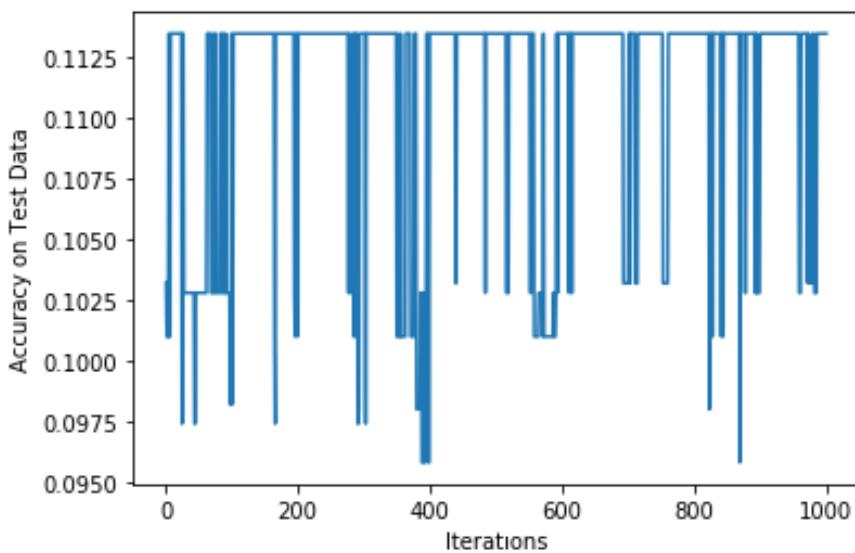


Figure 23 First version of 2-layer network, network did not learn, Duration 140s

Final Accuracy is : 0.9112

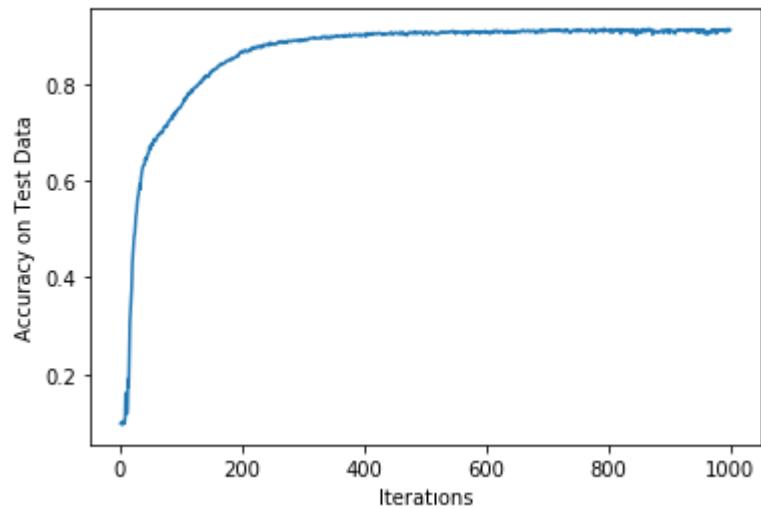


Figure 24 Weights are randomly initialized, batch size increased to 1000 Duration 244s, Learning rate 0.1, Gradient Descent Opt.

C) THREE LAYER NETWORK

3 layers NN did not work at first, then optimizer changed to Adagrad Optimizer then, other hyperparameters tuned.

Final Accuracy is : 0.9163

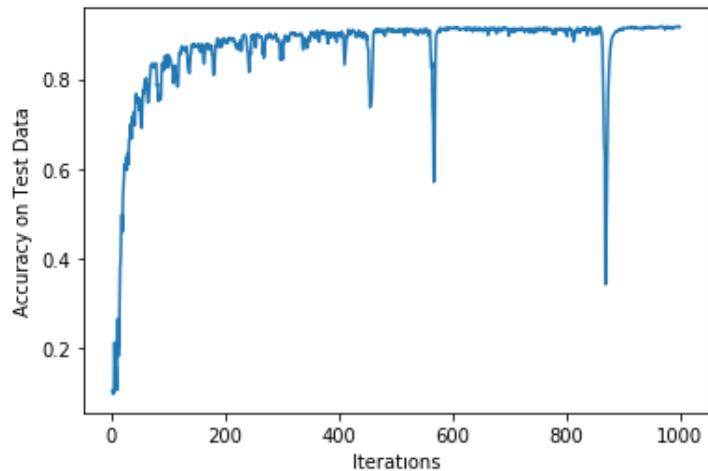


Figure 25 Adagrad opt. Lr 0.5 batch size 1000

D) THREE LAYER IMPLEMENTATION USING TF.LAYER

Final Accuracy is : 0.895

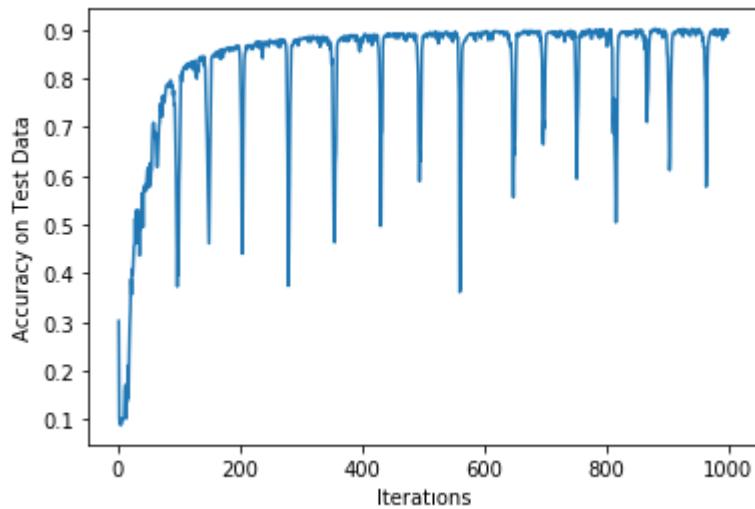


Figure 2 Same network with same hyper parameters in part c implemented with tf.layers module

Final Accuracy is : 0.9066

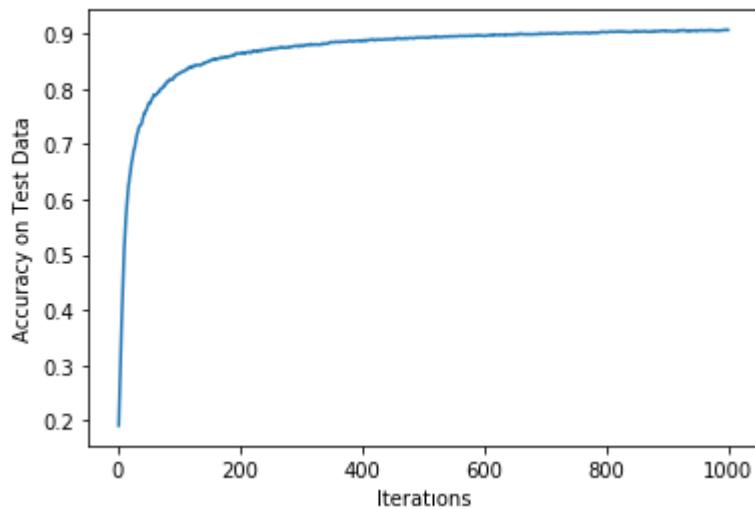


Figure 3 3 layer network implemented with tf.layers module with parameters : GradientDescentOptimizer, lr = 0.01, batch_size = 500