

Assignment - Fractal 3

Submitted by

Debonil Ghosh (M21AIE225)

Problem 1: Perceptron [40 points]

Following training samples are given:

x_1	x_2	Class
1	1	+1
-1	-1	-1
0	0.5	-1
0.1	0.5	-1
0.2	0.2	+1
0.9	0.5	+1

Table 1: Sample data

Assuming weight vector of initial decision boundary $w^T x = 0$ as $w=[1, 1]$, solve the following:

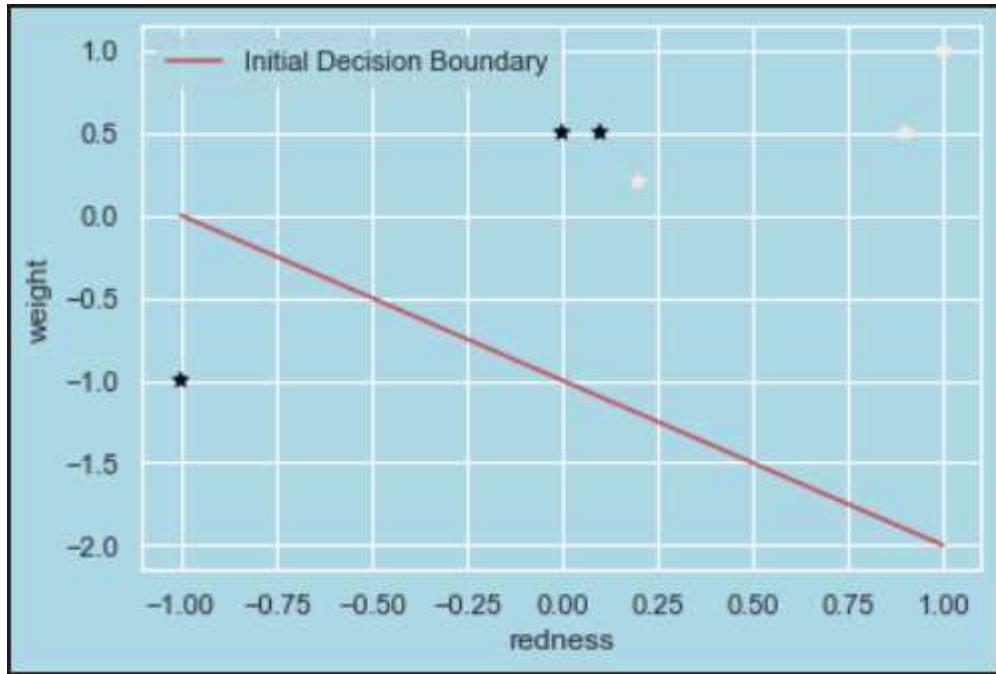
1. In how many steps perception learning algorithm will converge. (15 points)
2. What will be the final decision boundary? Show step-wise-step update of weight vector using computation as well as hand-drawn plot. (15 points)
3. Prove that Perceptron Learning Algorithm converges in a finite number of steps. (10 points)

Solution:

Program Output Part:

(Hand Written part Added later)

Initial Decision Boundary: $w=[1,1,1]$



Iteration is : 1

Sample is : [1. 1. 1.]

Class is : 1

Dot Product of W & X is : 3.0

Positive Sample is correctly classified

Sample is : [-1. -1. 1.]

Class is : -1

Dot Product of W & X is : -1.0

Negative Sample is Correctly classified

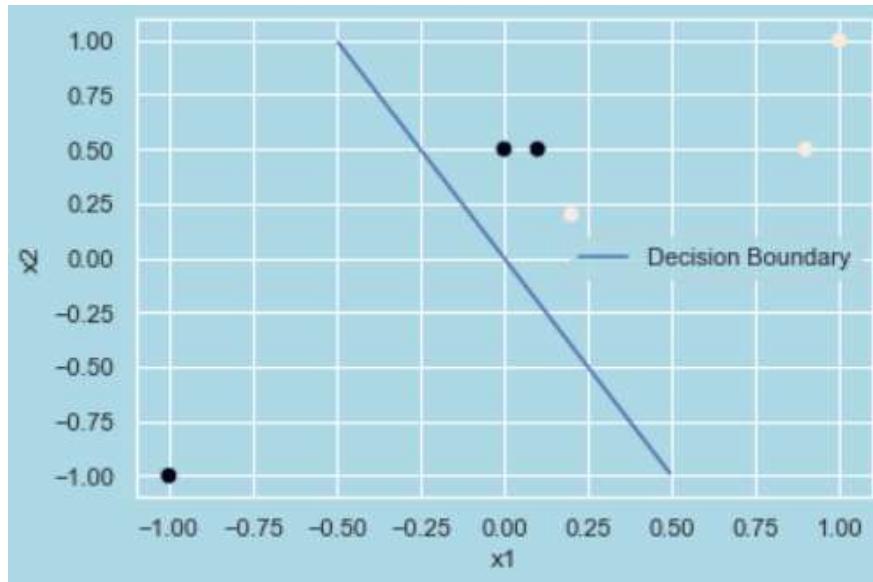
Sample is : [0. 0.5 1.]

Class is : -1

Dot Product of W & X is : 1.5

Negative Sample is classified positive

Updated W is : [1. 0.5 0.]



Sample is : [0.1 0.5 1.]

Class is : -1

Dot Product of W & X is : 0.35

Negative Sample is classified positive

Updated W is : [0.9 0. -1.]

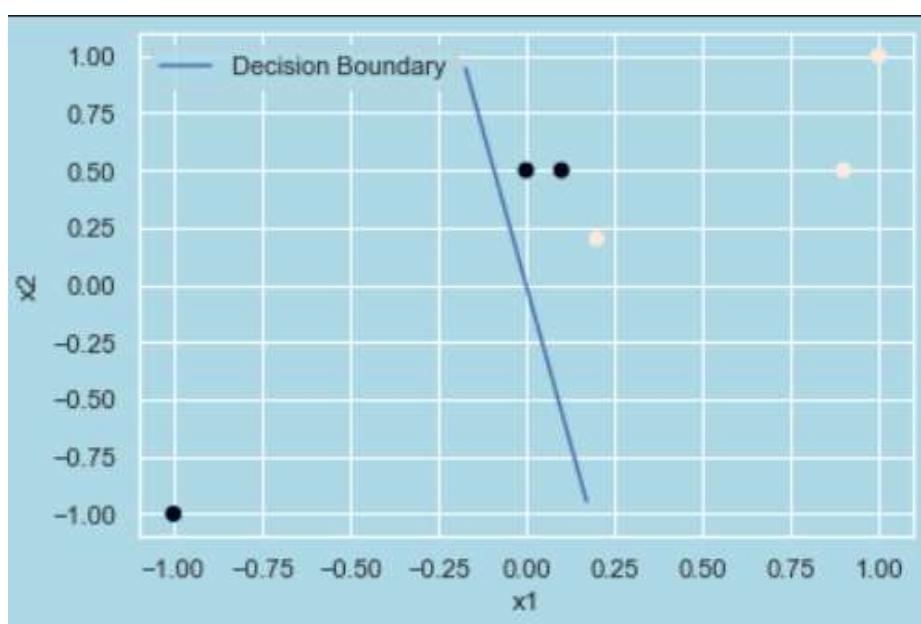
Sample is : [0.2 0.2 1.]

Class is : 1

Dot Product of W & X is : -0.82

Positive Sample is classified negative

Updated W is : [1.1 0.2 0.]



Sample is : [0.9 0.5 1.]

Class is : 1

Dot Product of W & X is : 1.09

Positive Sample is correctly classified

Iteration is : 2

Sample is : [1. 1. 1.]

Class is : 1

Dot Product of W & X is : 1.3

Positive Sample is correctly classified

Sample is : [-1. -1. 1.]

Class is : -1

Dot Product of W & X is : -1.3

Negative Sample is Correctly classified

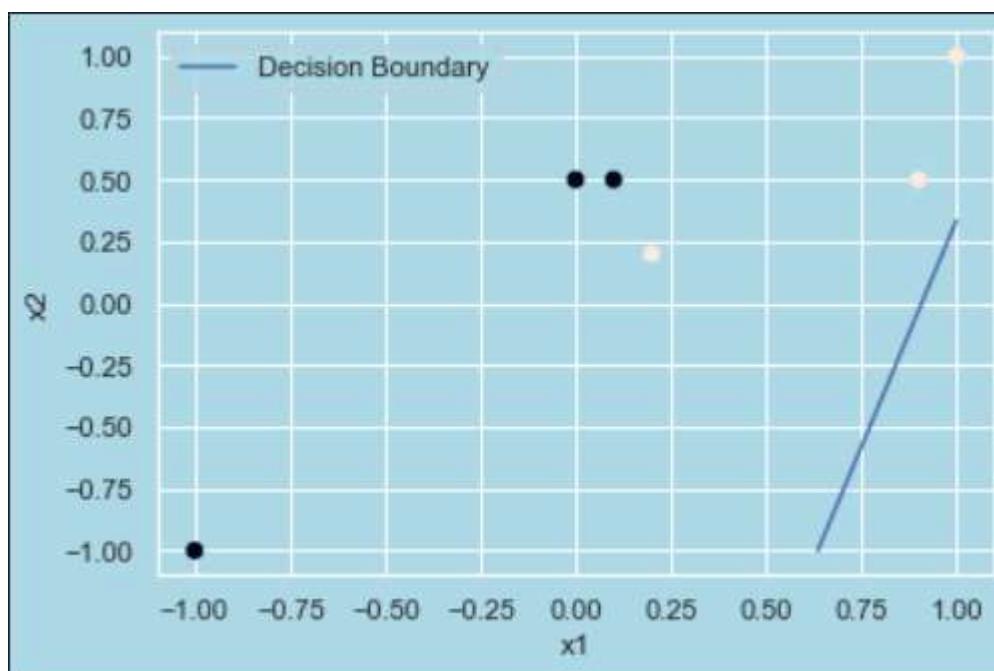
Sample is : [0. 0.5 1.]

Class is : -1

Dot Product of W & X is : 0.1

Negative Sample is classified positive

Updated W is : [1.1 -0.3 -1.]



Sample is : [0.1 0.5 1.]

Class is : -1

Dot Product of W & X is : -1.04

Negative Sample is Correctly classified

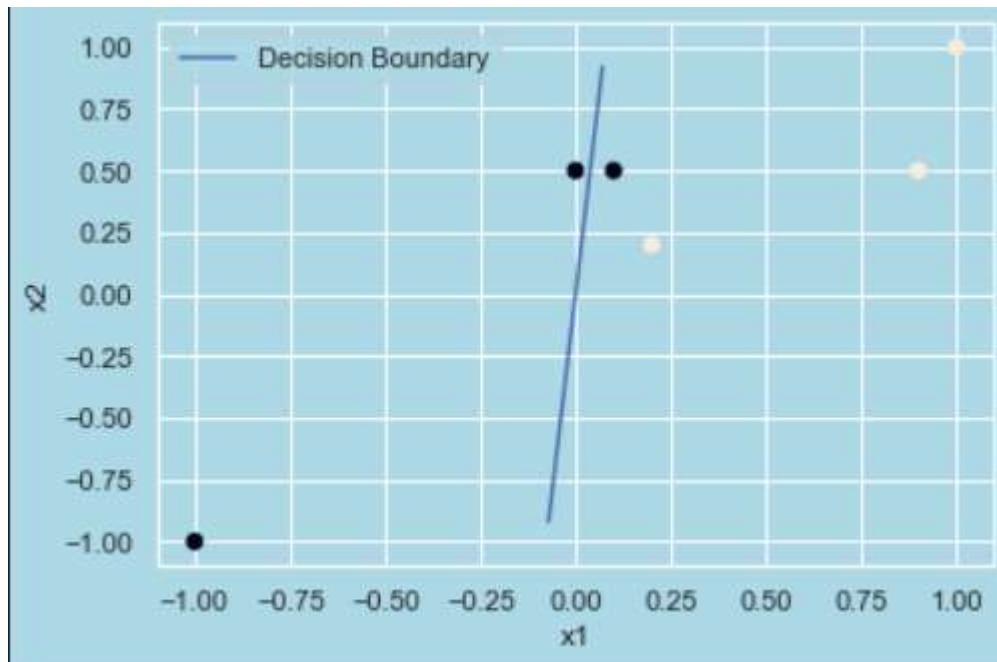
Sample is : [0.2 0.2 1.]

Class is : 1

Dot Product of W & X is : -0.84

Positive Sample is classified negative

Updated W is : [1.3 -0.1 0.]



Sample is : [0.9 0.5 1.]

Class is : 1

Dot Product of W & X is : 1.12

Positive Sample is correctly clasified

Iteration is : 3

Sample is : [1. 1. 1.]

Class is : 1

Dot Product of W & X is : 1.2000000000000002

Positive Sample is correctly clasified

Sample is : [-1. -1. 1.]

Class is : -1

Dot Product of W & X is : -1.2000000000000002

Negative Sample is Correctly classified

Sample is : [0. 0.5 1.]

Class is : -1

Dot Product of W & X is : -0.04999999999999999

Negative Sample is Correctly classified

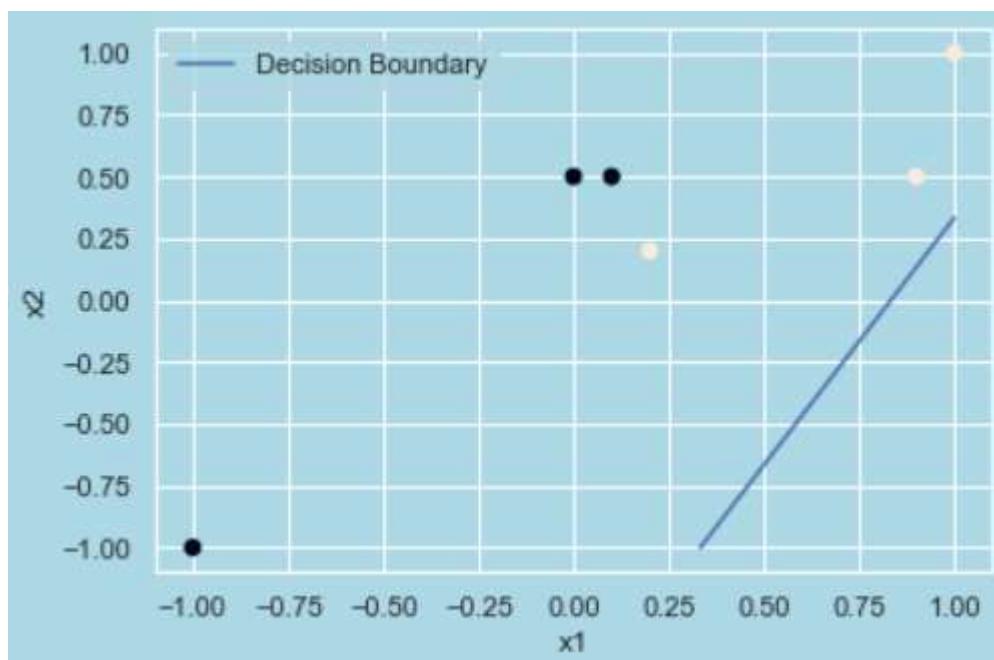
Sample is : [0.1 0.5 1.]

Class is : -1

Dot Product of W & X is : 0.08000000000000002

Negative Sample is classified positive

Updated W is : [1.2 -0.6 -1.]



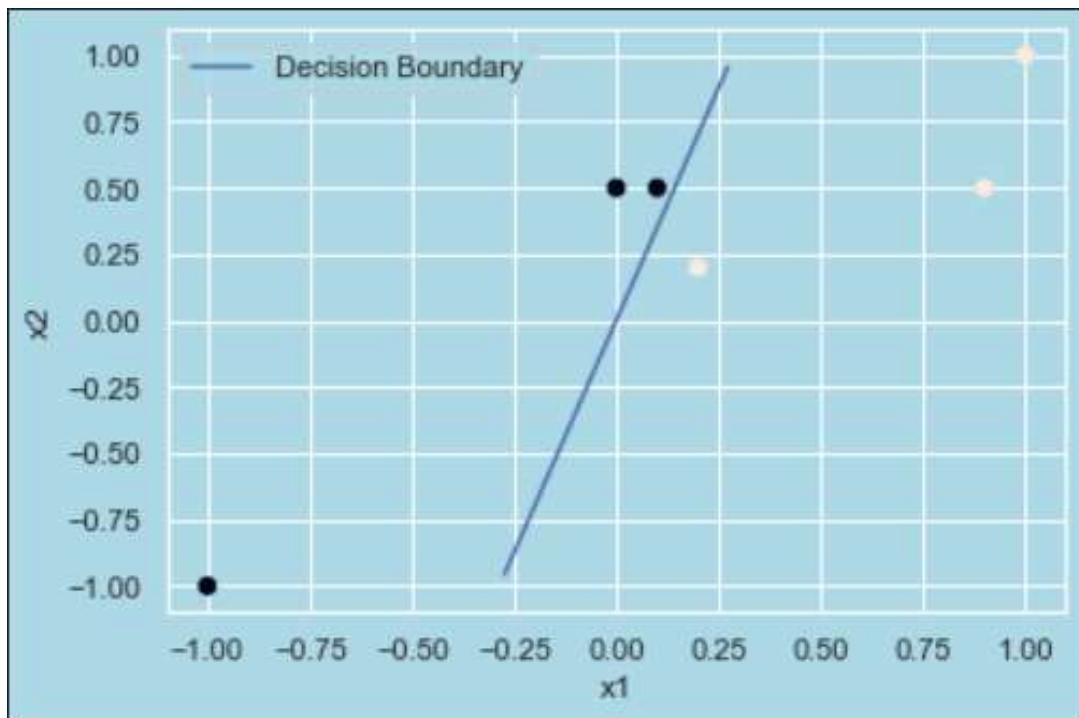
Sample is : [0.2 0.2 1.]

Class is : 1

Dot Product of W & X is : -0.88

Positive Sample is classified negative

Updated W is : [1.4 -0.4 0.]



Sample is : [0.9 0.5 1.]

Class is : 1

Dot Product of W & X is : 1.06

Positive Sample is correctly classified

Iteration is : 4

Sample is : [1. 1. 1.]

Class is : 1

Dot Product of W & X is : 1.0

Positive Sample is correctly classified

Sample is : [-1. -1. 1.]

Class is : -1

Dot Product of W & X is : -1.0

Negative Sample is Correctly classified

Sample is : [0. 0.5 1.]

Class is : -1

Dot Product of W & X is : -0.1999999999999998

Negative Sample is Correctly classified

Sample is : [0.1 0.5 1.]

Class is : -1

Dot Product of W & X is : -0.06

Negative Sample is Correctly classified

Sample is : [0.2 0.2 1.]

Class is : 1

Dot Product of W & X is : 0.1999999999999998

Positive Sample is correctly classified

Sample is : [0.9 0.5 1.]

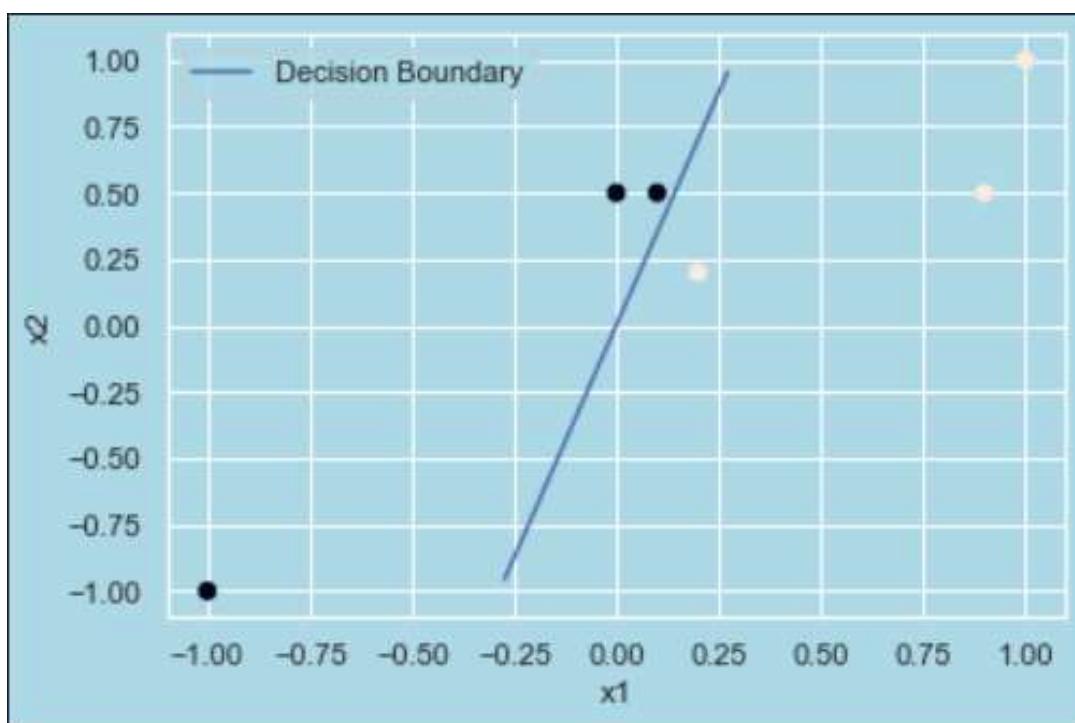
Class is : 1

Dot Product of W & X is : 1.06

Positive Sample is correctly classified

[1.4 -0.4 0.]

Final Decision Boundary is:



Machine Learning Assignment (Fract-3)

Debonil Ghosh (M21AIE225)

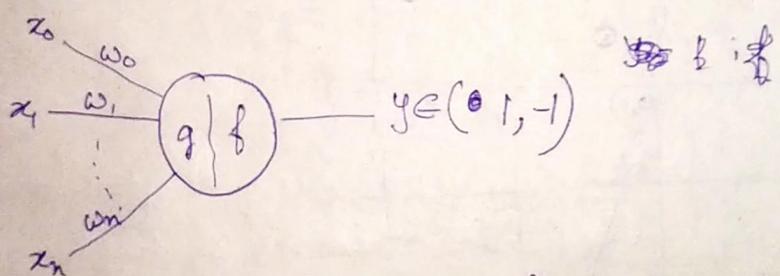
[Hand Written Part]

Q1. Perceptron

x_1	x_2	Class
1	1	+1
-1	-1	-1
0	0.5	-1
0.1	0.5	-1
0.2	0.2	+1
0.9	0.5	+1

Sample data

Initial decision boundary $w^T x = 0$ as $w = [1, 1]$

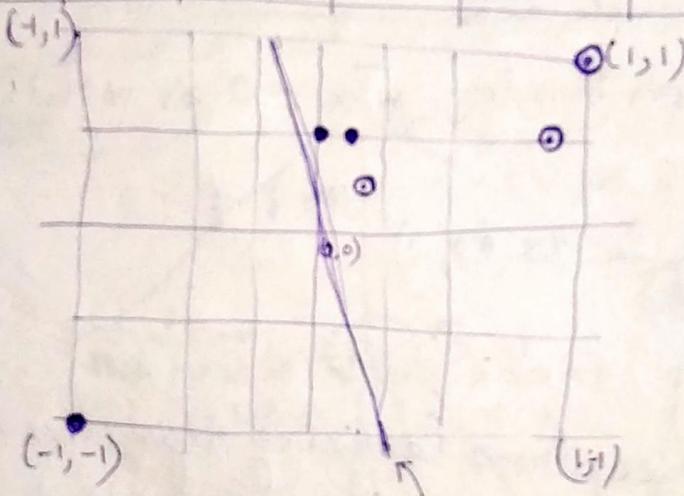


- 1.1. The perceptron learning algorithm will converge in 4 steps (iteration), which can be verified from the next answer.

1.2

Step | Iteration 1:

x_1	x_2	x_0	Class	w_1	w_2	w_0	$w^T \cdot x$	Correct	w -updated
1	1	1	1	1	1	1	3	✓	—
-1	-1	1	-1	1	1	1	-1	✓	—
0.5	1	-1	-1	1	1	1	1.5	✗	$[1, 0.5, 0]$
0.1	0.5	1	-1	1	1	0	0.35	✗	$[0.9, 0, -1]$
0.2	0.2	1	1	0.9	0	-1	-0.82	✗	$[1.1, 0.2, 0]$
0.9	0.5	1	1	1	0.2	0	1.09	✓	—

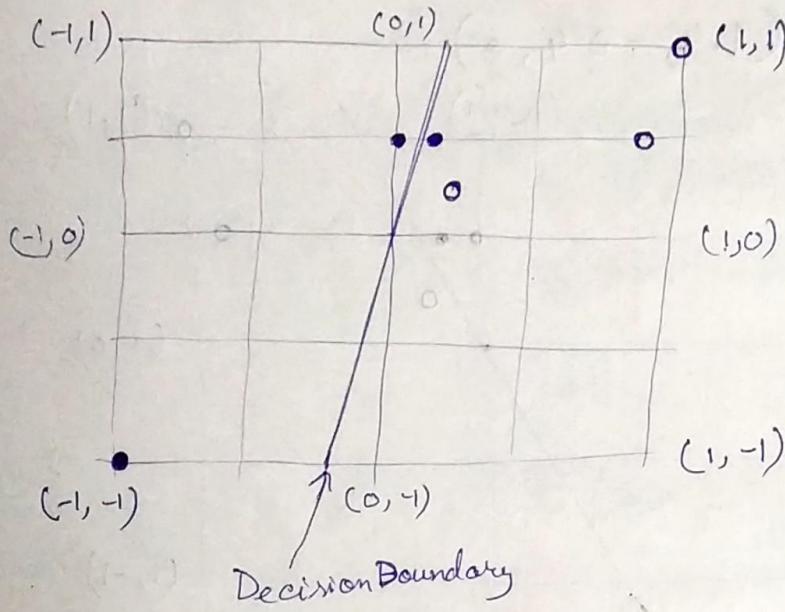


decision boundary after iteration 1
 $(w = [1.1, 0.2, 0])$

Step-2:

x_1	x_2	x_0	Class	w_1	w_2	w_0	$w^T \cdot x$	Correct	w updated
1	1	1	1	1.1	0.2	0	1.3	✓	—
-1	-1	1	-1	1.1	0.2	0	-1.3	✓	—
0.5	1	-1	-1	1.1	0.2	0	0.1	✗	$[1.1, -0.3, -1]$
0.1	0.5	1	-1	1.1	-0.3	-1	-1.04	✓	—
0.2	0.2	1	1	1.1	-0.3	-1	-0.84	✗	$[1.3, -0.1, 0]$
0.9	0.5	1	1	1.3	-0.1	0	1.12	✓	—

Decision Boundary after step 2 is $\omega = (1.3, -0.1, 0)$

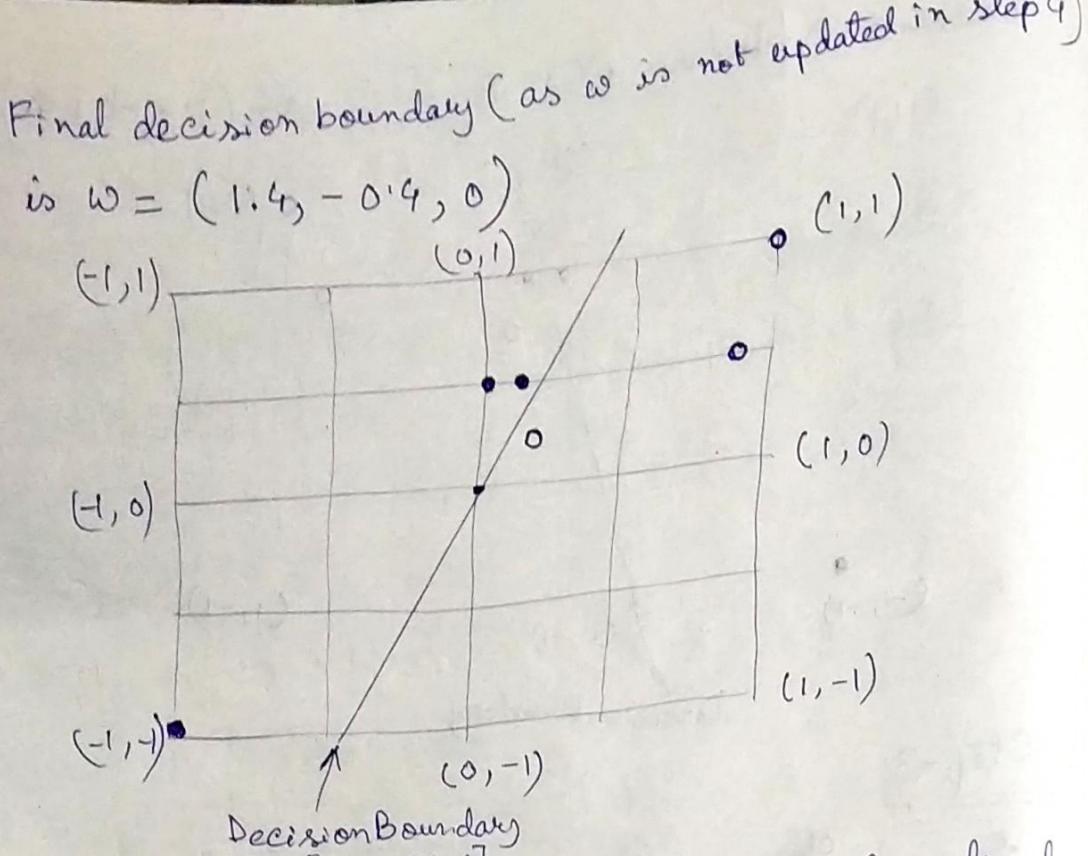


Step-3:

x	Class	ω	$w^T \cdot x$	Correct	ω updated
x_1	x_2	$w_1 \quad w_2 \quad w_0$			
1	1	1 1.3 -0.1 0	1.2	✓	-
-1	-1	1 1.3 -0.1 0	-1.2	✓	-
0	0.5	1 1.3 -0.1 0	-0.05	✓	-
0.1	0.5	1 1.3 -0.1 0	0.08	X	[1.2, -0.6, -1]
0.2	0.2	1 1.2, -0.6, -1	-0.88	X	[1.4, -0.4, 0]
0.9	0.5	1 1.4, -0.4, 0	1.06	✓	-

Step-4

x	Class	ω	$w^T \cdot x$	Correct	ω updated
x_1	x_2	$w_1 \quad w_2 \quad w_0$			
1	1	1 1.4, -0.4 0	1	✓	-
-1	-1	1 1.4 -0.4 0	-1	✓	-
0	0.5	1 1.4 -0.4 0	-0.2	✓	-
0.1	0.5	1 1.4 -0.4 0	-0.06	✓	-
0.2	0.2	1 1.4 -0.4 0	0.2	✓	-
0.9	0.5	1 1.4 -0.4 0	1.06	✓	-



After step 3, all samples are correctly classified in step 4.

1.3 Prove that Perceptron Learning Algorithm converges in a finite number of steps.

Ans: Suppose at timestamp t , we inspected the point p_i , and found that $\omega^T \cdot p_i \leq 0$

$$\text{we make a correction } \omega_{t+1} = \omega_t + p_i$$

Let β be the angle between ω^* and ω_{t+1}

$$\text{then, } \cos \beta = \frac{\omega^* \cdot \omega_{t+1}}{\|\omega_{t+1}\|}$$

$$\begin{aligned}
 \text{Numerator} &= w^* \cdot w_{t+1} = w^*(w_t + p_i) \\
 &= w^* w_t + w^* p_i \\
 &\geq w^* w_t + \delta \left(\delta = \min\{w^* \cdot r_i | v_i\} \right) \\
 &> w^*(w_{t-1} + p_i) + \delta \\
 &\geq w^* w_{t-1} + w^* p_i + \delta \\
 &\geq w^* w_{t-1} + 2\delta
 \end{aligned}$$

- Everytime we make a correction a quantity δ is added. Say if k corrections are made, a quantity of $k\delta$ is added.

$$\cos \beta = \frac{(w^* w_{t+1})}{\|w_{t+1}\|}$$

$$\text{Numerator} \geq w^* w_0 + k\delta \quad (\text{proved by induction})$$

$$\begin{aligned}
 \text{Denominator} &= \|w_{t+1}\|^2 \\
 &= (w_t + p_i)^2 \\
 &= w_t^2 + 2w_t p_i + p_i^2 \leq \|w_t\|^2 + \|p_i\|^2 \\
 &\leq \|w_{t-1}\|^2 + 1 \quad (\because w_i, p \leq 0) \\
 &\leq \|w_{t-1}\|^2 + k
 \end{aligned}$$

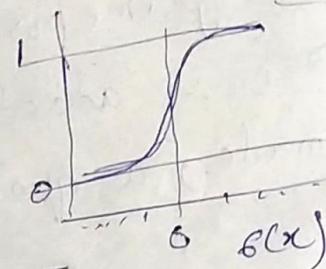
$$\cos \beta \Rightarrow \frac{w^* w_0 + k\delta}{\sqrt{\|w_0\|^2 + k}}$$

thus $\cos \beta$ is proportional to \sqrt{K}
 $\therefore K$ has to be finite so that $\cos \beta \leq 1$
 \therefore there can only be a finite number of
 corrections (K) to w and the perceptron algorithm
 will converge.

Q.3

3.1. Compute Derivatives of the following activation
 functions:

1. Sigmoid:



$$\text{Sigmoid}(x) = \sigma(x) = \frac{1}{1 + e^{-x}}$$

$$\frac{d \sigma(x)}{dx} = \frac{d}{dx} \left(\frac{1}{1 + e^{-x}} \right)$$

$$= \frac{(1 + e^{-x}) \frac{d}{dx} \cdot 1 - 1 \times \frac{d}{dx} e^{-x}}{(1 + e^{-x})^2} \quad (\text{applying quotient rule})$$

$$= \frac{e^{-x}}{(1 + e^{-x})^2} = \frac{1 + e^{-x} - 1}{(1 + e^{-x})^2}$$

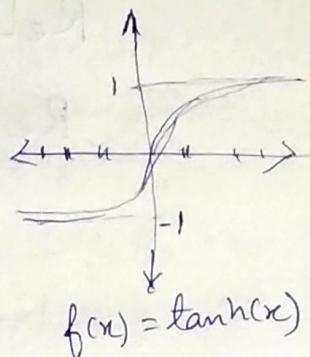
$$= \frac{1 + e^{-x}}{(1 + e^{-x})^2} - \frac{1}{(1 + e^{-x})^2} = \frac{1}{1 + e^{-x}} - \frac{1}{(1 + e^{-x})^2}$$

$$= \frac{1}{1 + e^{-x}} \left(1 - \frac{1}{1 + e^{-x}} \right) = \sigma(x) (1 - \sigma(x))$$

$$\therefore \frac{d}{dx} f(x) = f(x)(1-f(x)) \quad \underline{\text{Ans}}$$

2. tanh:

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} = T(x)$$



$$\frac{d}{dx} T(x) = \frac{d}{dx} \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

$$f'(x) = \tanh'(x)$$

applying quotient rule,

$$\begin{aligned} \frac{d}{dx} \frac{e^x - e^{-x}}{e^x + e^{-x}} &= \frac{e^x + e^{-x} \left(\frac{d}{dx}(e^x - e^{-x}) \right) - e^x - e^{-x} \left(\frac{d}{dx}(e^x + e^{-x}) \right)}{(e^x + e^{-x})^2} \\ &= \frac{(e^x + e^{-x})(e^x + e^{-x}) - \{(e^x - e^{-x})(e^x - e^{-x})\}^2}{(e^x + e^{-x})^2} \\ &= \frac{(e^x + e^{-x})^2 - (e^x - e^{-x})^2}{(e^x + e^{-x})^2} \\ &= 1 - \frac{(e^x - e^{-x})^2}{(e^x + e^{-x})^2} \end{aligned}$$

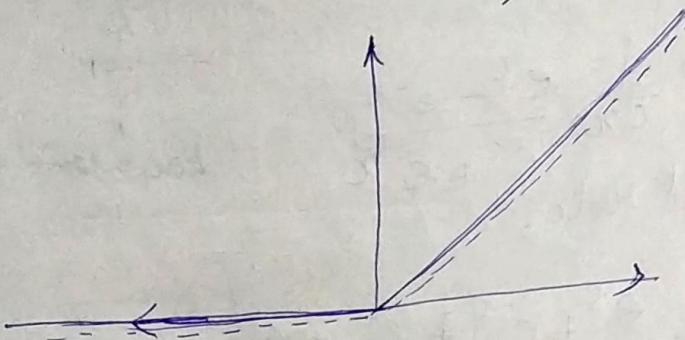
$$\frac{d}{dx} T(x) = 1 - (T(x))^2$$

$$\therefore \frac{d}{dx} \tanh = 1 - (\tanh)^2$$

3. ReLU:

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

$$\text{ReLU} = \begin{cases} x, & x > 0 \\ 0, & x \leq 0 \end{cases}$$



$$f(x) = \text{ReLU}$$

$$f'(x) = \lim_{\Delta x \rightarrow 0} \frac{\max(0, x + \Delta x) - \max(0, x)}{\Delta x}$$

At $x > 0$

$$f'(x) = \lim_{\Delta x \rightarrow 0} \frac{x + \Delta x - x}{\Delta x} = 1$$

At $x < 0$

$$f'(x) = \lim_{\Delta x \rightarrow 0} \frac{0 - 0}{\Delta x} = 0$$

At $x = 0$

$$f'(x) = \lim_{\Delta x \rightarrow 0^+} \frac{0 + \Delta x - 0}{\Delta x} = 1$$

$$f'(x) = \lim_{\Delta x \rightarrow 0^-} \frac{0 - 0}{\Delta x} = 0$$

LHL \neq RHL

$$\therefore \frac{d}{dx} \text{ReLU} = \frac{d}{dx} \max(0, x)$$

$$= \begin{cases} 1, & x > 0 \\ 0, & x < 0 \\ \text{Does not exist, } & x = 0 \end{cases}$$

3.2 What are the strategies you will follow to avoid overfitting in neural network.

Ans:

When a model of neural network learns training data very well but does not perform well on testing dataset, then the situation is called Overfitting. Several precautions may be taken to avoid overfitting. and these are following:

① Simplifying the model: By reducing the number of

Internal ~~out~~ layers and nodes in a layer, making network smaller.

② Early Stopping: Stopping the learning or training process of the model before it crosses overfitting.

③ Data Augmentation: By making small changes into the existing dataset, we make the model consider each augmented model as a distinct dataset. This will reduce overfitting of a ~~bad~~ neural network.

④ Regularization: It can be used to reduce the complexity of a model. The best way is to add penalty to the loss function in proportion to the weight of the model.

Q. 3.3

$$\cancel{x} = [1, 1]^T, y = [1, 1]^T \in \mathbb{R}^2$$

$f: \mathbb{R}^2 \rightarrow \mathbb{R}^2$ with $f(z) = z_1 x + z_2 y$ for any

$$z = [z_1, z_2]^T \in \mathbb{R}^2$$

$$z = g(r) = [r^2, r^3]^T \text{ where } r \in \mathbb{R}.$$

$$\frac{\partial f}{\partial r} = \frac{\partial f}{\partial z} \times \frac{\partial z}{\partial r} = \begin{bmatrix} x \\ y \end{bmatrix} \begin{bmatrix} 2r & 3r^2 \end{bmatrix}$$

$$= \begin{bmatrix} 2rx_1 & 3ry_1 \\ 2rx_2 & 3ry_2 \end{bmatrix} \underset{x=[1,1]}{=} \begin{bmatrix} 4 & 12 \\ 4 & 12 \end{bmatrix}$$

$y = [1, 1]$

$r = 2$

Q. 3.4

① Mean Square Error (MSE): Mean Square Error calculates the mean

of distance of set data points from a regression line. It calculates the difference between the truth and prediction and square it. The reason behind the squaring the result is to make the negative difference not relevant.

$$MSE = (Prediction - Truth)^2$$

for n number of data points,

$$MSE = \frac{(Pred_1 - Truth_1)^2 + (Pred_2 - Truth_2)^2 + \dots + (Pred_n - Truth_n)^2}{n}$$

② Binary Cross-Entropy (BCE): Binary Cross Entropy is the

negative average of the log of corrected predicted probabilities. It compares each of the predicted probabilities to actual class output which can be either 0, or 1. It then calculates the score that penalizes the probabilities based on the distance from the expected value. That means how close of far from the actual value.

$$\text{Binary Cross Entropy} = -\frac{1}{N} \sum_{i=1}^N (\log(p_i))$$

③ Categorical Cross Entropy: It is a loss function that is used in multi-class classification task. These are tasks where an example can only belong to one out of many possible categories and the model must decide which one.

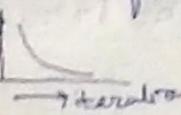
$$\text{Categorical Cross Entropy} = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^M y_{ij} \log(p_{ij})$$

where, N is the number of rows
 M is the number of classes

Q. 3.5: Explain the following variants of Gradient Descent in brief:

(i) Batch Gradient Descent: Batch Gradient Descent computes the gradient of the cost function w.r.t. to the parameter θ for the entire training dataset

$$\theta = \theta - \eta \cdot \nabla_{\theta} J(\theta)$$

With Batch Gradient Descent function cost function reduces smoothly. 

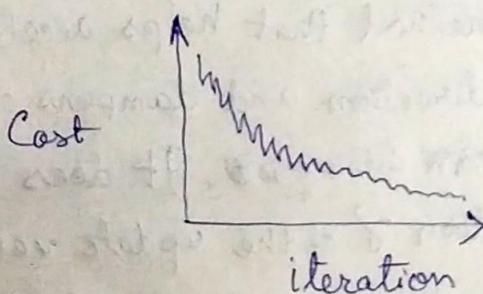
(ii) Stochastic Gradient Descent: Stochastic Gradient Descent

performs a parameter update for each training example x^i and label y^i :

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

Batch GD performs redundant computations for large datasets, as it recomputes gradients for similar examples before each parameter update.

SGD does away with this redundancy by performing one update at a time it is therefore usually much faster and can also be used to learn online.



Parameter updation in case of SGD is not that smooth, as it is possible to learn noise as well.

(iii) Mini Batch Gradient Descent: Mini-batch Gradient Descent

takes the best of both worlds (BGD and SGD) and performs an update for every mini-batch of n training examples:

$$\theta = \theta - \eta \cdot \nabla_{\theta} J(\theta; x^{(i:i+n)}, y^{(i:i+n)})$$

this way it: (a) reduces the variance of the parameter update, which can lead to more stable convergence; (b) can make use of highly optimized matrix optimizations common to deep learning libraries.

(iv) Momentum Based Gradient Descent:

Momentum is a method that helps accelerate SGD in the relevant direction and dampens oscillations. ~~can be seen as~~. It does this by adding a fraction γ of the update vector of the past time step to the current update vector:

$$v_t = \gamma v_{t-1} + \eta \nabla_{\theta} J(\theta)$$

$$\theta = \theta - v_t$$

(v) Adam Solver: Adaptive Moment Estimation
(ADAM) is another method that computes adaptive learning rates for each parameter.
In addition to storing an exponentially decaying average of past squared gradients v_t like Adadelta and RMSprop, Adam also keeps an exponentially decaying average of past gradients m_t , similar to momentum. The decaying averages of past and past squared gradients m_t and v_t respectively follows.

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t$$

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$$

Problem 2: Learning to implement Neural Network [40 points]

2.2. Please go through the following blog to learn how to recognize handwritten digits using

Neural Network. Here Neural Network is coded using PyTorch Library in Python.

<https://towardsdatascience.com/handwritten-digit-mnist-pytorch-977b5338e627>

Use above code and report your observation based on the following:(15 points)

- (i) Change loss function,
- (ii) Change in learning rate, and
- (iii) Change in Number of hidden layers

Solution:

By updating given code, following observations have generated through python code.

(i) Change loss function:

Training with Loss Function = **NLLLoss()**, Learning Rate = 0.003, Hidden Layer = [128, 64]

Sequential(

(0): Linear(in_features=784, out_features=128, bias=True)

(1): ReLU()

(2): Linear(in_features=128, out_features=64, bias=True)

(3): ReLU()

(4): Linear(in_features=64, out_features=10, bias=True)

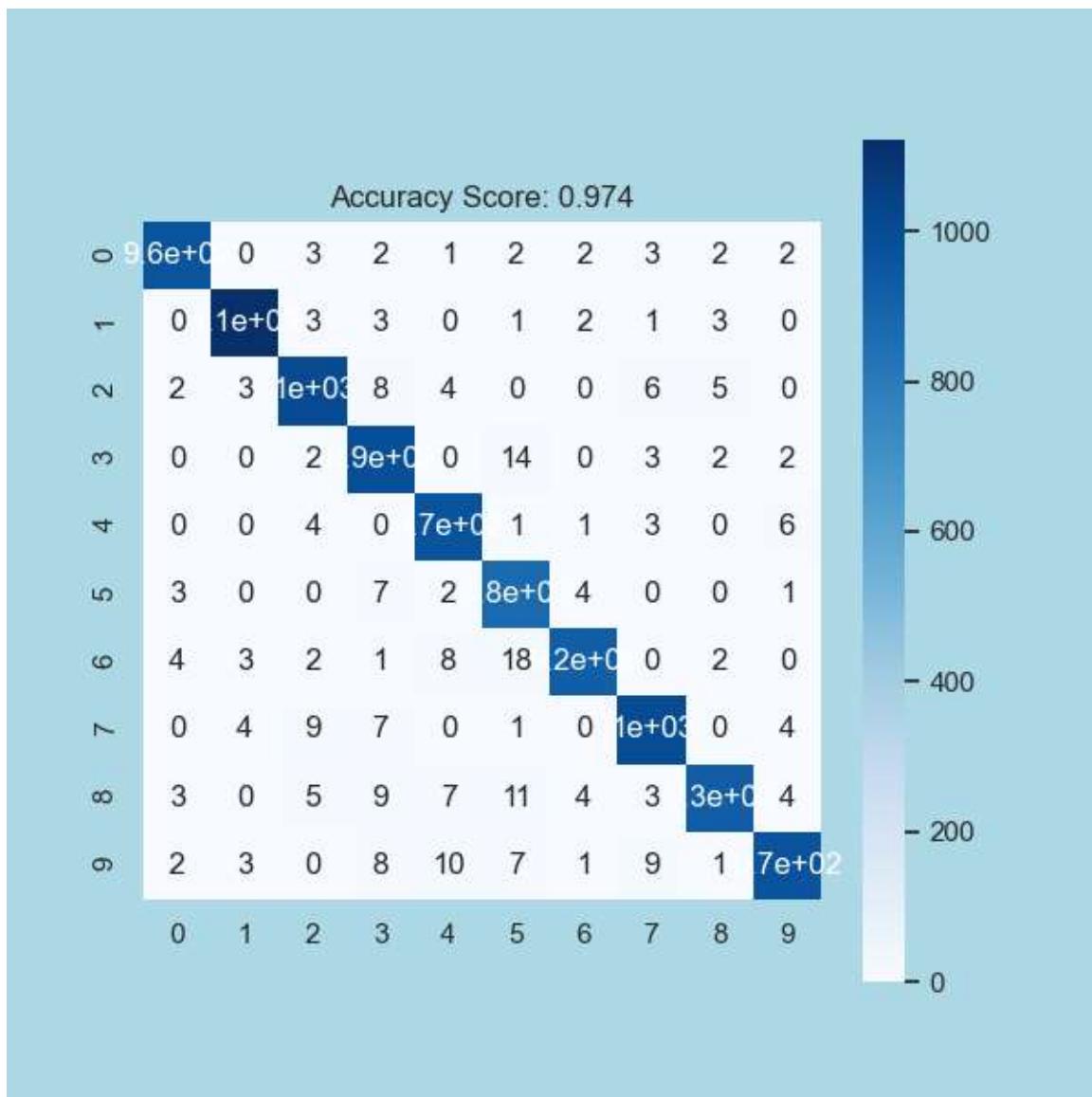
(5): LogSoftmax(dim=1)

)

Training Time (in minutes) = 2.5658482789993284

Overall Accuracy Score: 0.975

Classwise Accuracy Score: [0.9964 0.9973 0.9952 0.9943 0.9948 0.9947 0.9954 0.9945 0.9939 0.9929]



Confusion Matrix for NLLLoss

Training with Loss Function = **MSELoss()**,
Layer = [128, 64]

Learning Rate = 0.003,

Hidden

Sequential(

(0): Linear(in_features=784, out_features=128, bias=True)

(1): ReLU()

(2): Linear(in_features=128, out_features=64, bias=True)

(3): ReLU()

(4): Linear(in_features=64, out_features=10, bias=True)

(5): LogSoftmax(dim=1)

)

Training Time (in minutes) = 7.808586394786834

Overall Accuracy Score: **0.938**

Classwise Accuracy Score: [0.9916 0.9944 0.9846 0.9873 0.9871 0.9863 0.9899 0.9874
0.9839 0.9833]



Training with Loss Function = **KLDivLoss()**, Learning Rate = 0.003, Hidden Layer = [128, 64]

Sequential(

(0): Linear(in_features=784, out_features=128, bias=True)

(1): ReLU()

(2): Linear(in_features=128, out_features=64, bias=True)

(3): ReLU()

(4): Linear(in_features=64, out_features=10, bias=True)

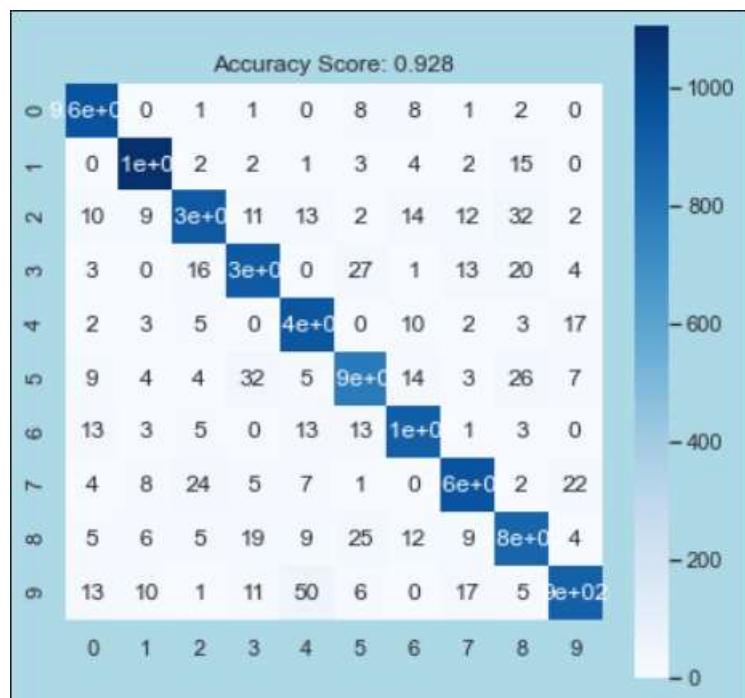
(5): LogSoftmax(dim=1)

)

Training Time (in minutes) = 3.6359249035517376

Overall Accuracy Score: **0.928**

Classwise Accuracy Score: [0.992 0.9928 0.9832 0.9835 0.986 0.9811 0.9886 0.9867 0.9798 0.9831]



Training with Loss Function = **CrossEntropyLoss()**, Learning Rate = 0.003,
Hidden Layer = [128, 64]

Sequential(

(0): Linear(in_features=784, out_features=128, bias=True)

(1): ReLU()

(2): Linear(in_features=128, out_features=64, bias=True)

(3): ReLU()

(4): Linear(in_features=64, out_features=10, bias=True)

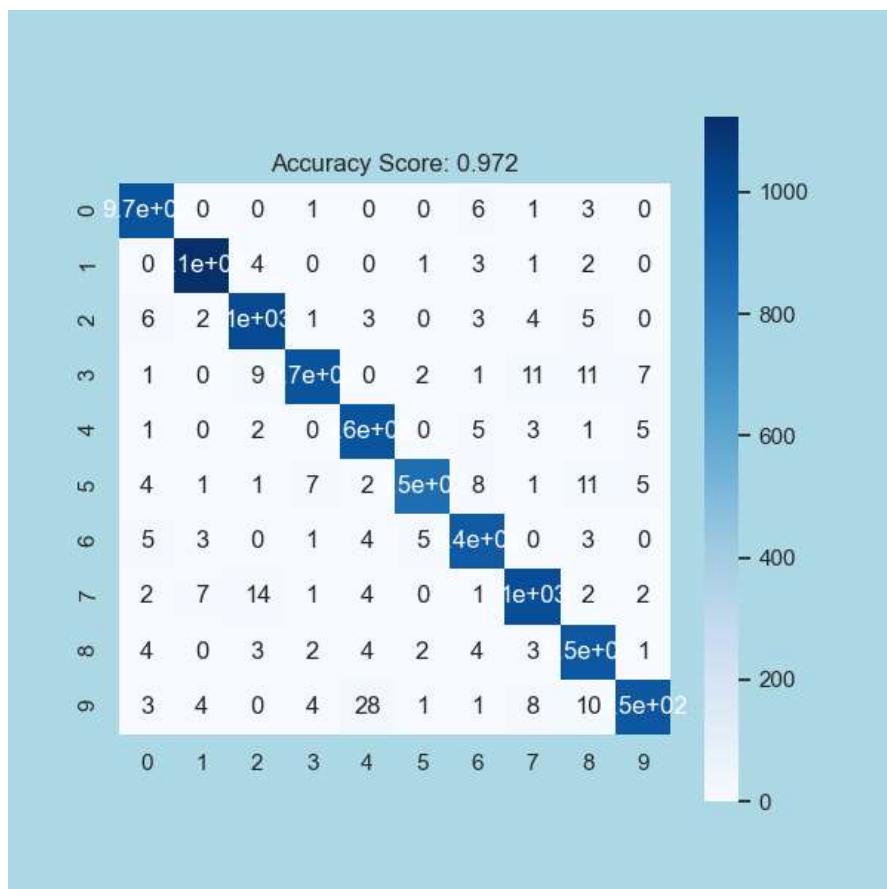
(5): LogSoftmax(dim=1)

)

Training Time (in minutes) = 2.653095289071401

Overall Accuracy Score: **0.975**

Classwise Accuracy Score: [0.9966 0.9973 0.9949 0.9948 0.9943 0.9955 0.9957 0.9936
0.994 0.9933]



(ii) Change in learning rate:

Training with Loss Function = NLLLoss(), Learning Rate = 0.5, Hidden Layer = [128, 64]

Sequential(

(0): Linear(in_features=784, out_features=128, bias=True)

(1): ReLU()

(2): Linear(in_features=128, out_features=64, bias=True)

(3): ReLU()

(4): Linear(in_features=64, out_features=10, bias=True)

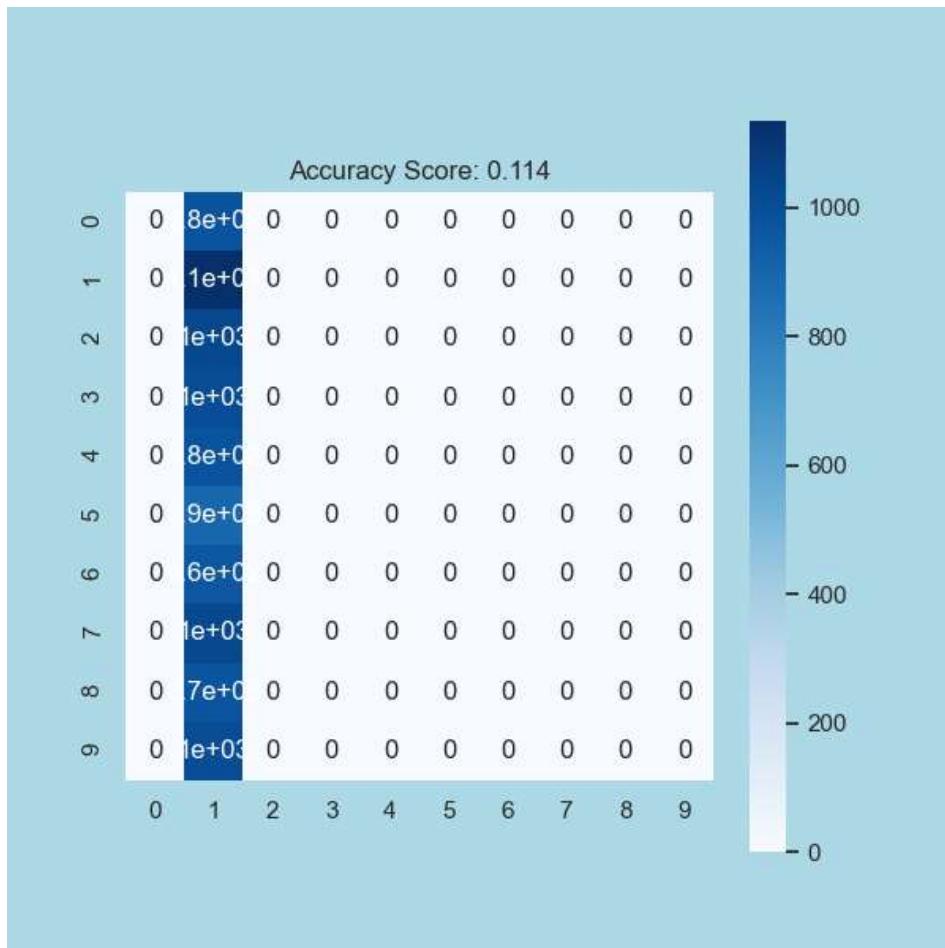
(5): LogSoftmax(dim=1)

)

Training Time (in minutes) = 2.7291388670603434

Overall Accuracy Score: **0.114**

Classwise Accuracy Score: [0.902 0.1135 0.8968 0.899 0.9018 0.9108 0.9042 0.8972 0.9026 0.8991]



Training with Loss Function = NLLLoss(), Learning Rate = 0.01, Hidden Layer = [128, 64]

Sequential(

(0): Linear(in_features=784, out_features=128, bias=True)

(1): ReLU()

(2): Linear(in_features=128, out_features=64, bias=True)

(3): ReLU()

(4): Linear(in_features=64, out_features=10, bias=True)

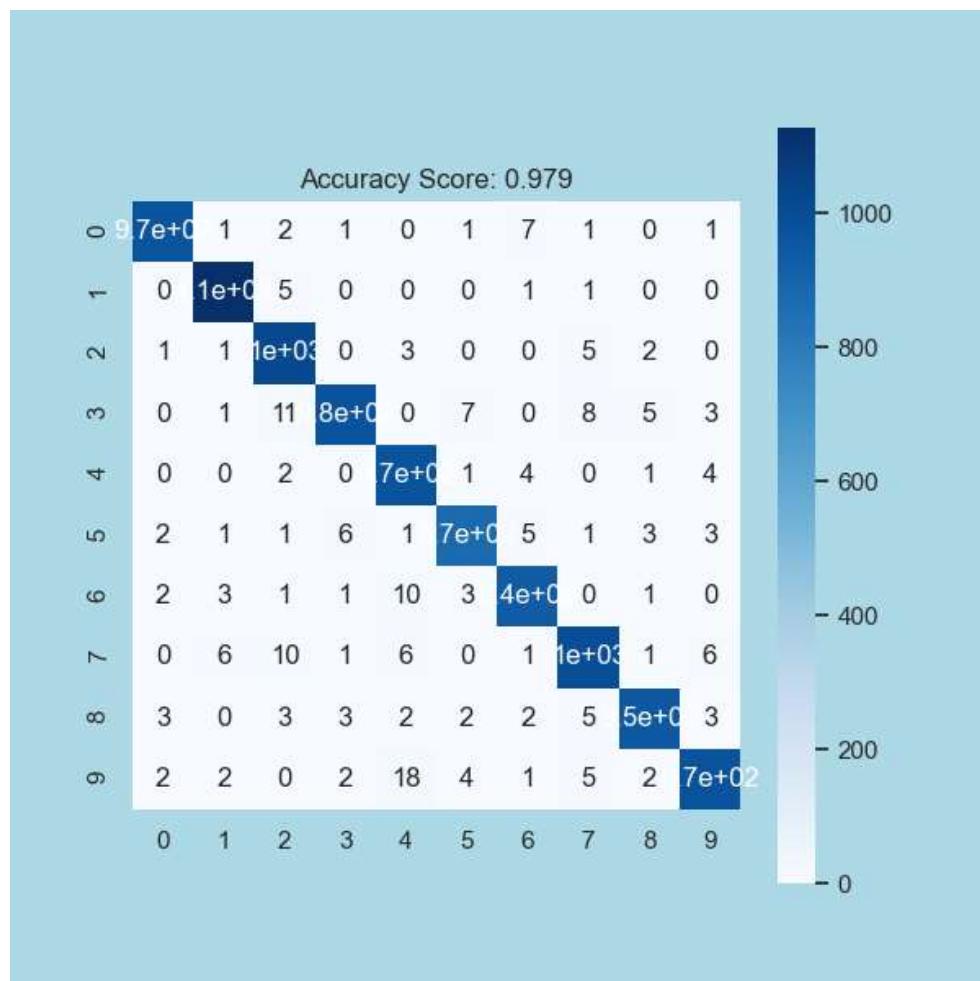
(5): LogSoftmax(dim=1)

)

Training Time (in minutes) = 3.002375260988871

Overall Accuracy Score: **0.979**

Classwise Accuracy Score: [0.9976 0.9978 0.9953 0.9951 0.9948 0.9959 0.9958 0.9943 0.9962 0.9944]



(iii) Change in Number of hidden layers:

Training with Loss Function = NLLLoss(), Learning Rate = 0.01, Hidden Layer = [128]

Sequential(

(0): Linear(in_features=784, out_features=128, bias=True)

(1): ReLU()

(2): Linear(in_features=128, out_features=10, bias=True)

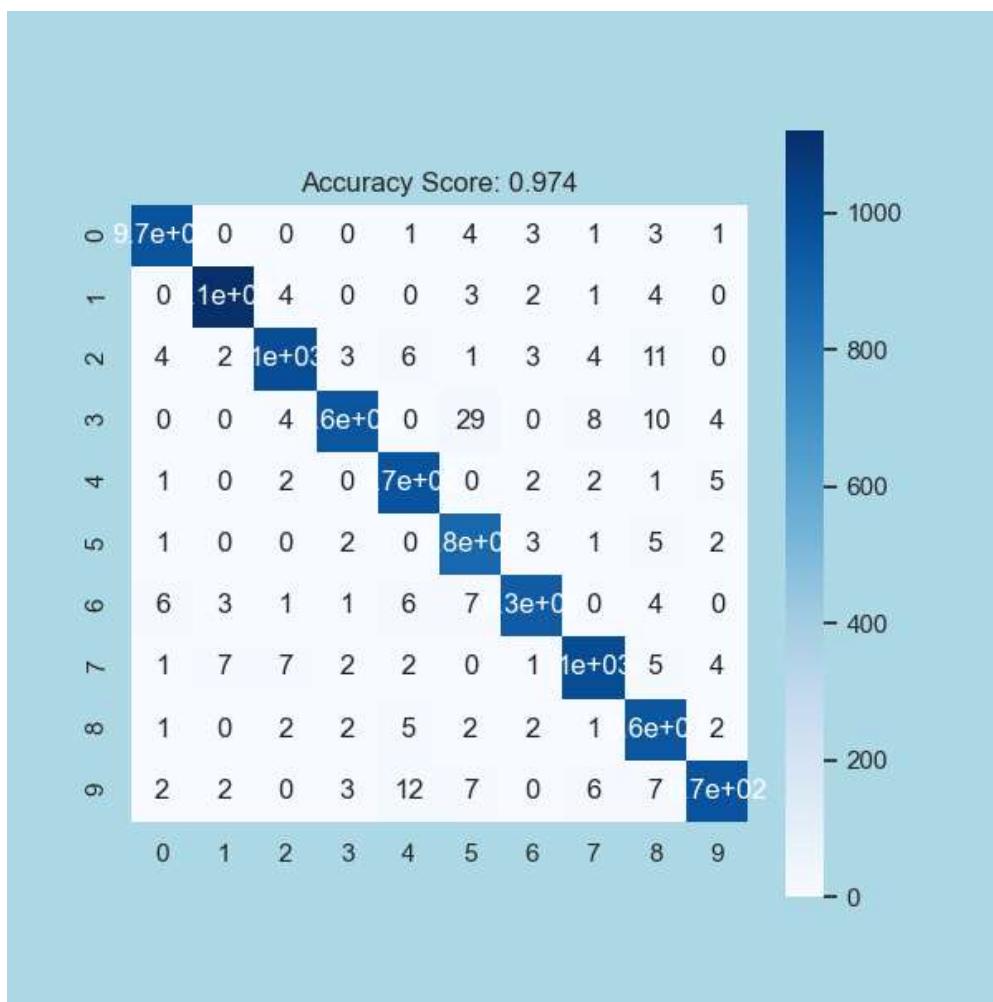
(3): LogSoftmax(dim=1)

)

Training Time (in minutes) = 5.416209638118744

Overall Accuracy Score: **0.974**

Classwise Accuracy Score: [0.9971 0.9972 0.9946 0.9932 0.9955 0.9933 0.9956 0.9947 0.9933 0.9943]



Training with Loss Function = NLLLoss(), Learning Rate = 0.01, Hidden Layer = [16]

Sequential(

(0): Linear(in_features=784, out_features=16, bias=True)

(1): ReLU()

(2): Linear(in_features=16, out_features=10, bias=True)

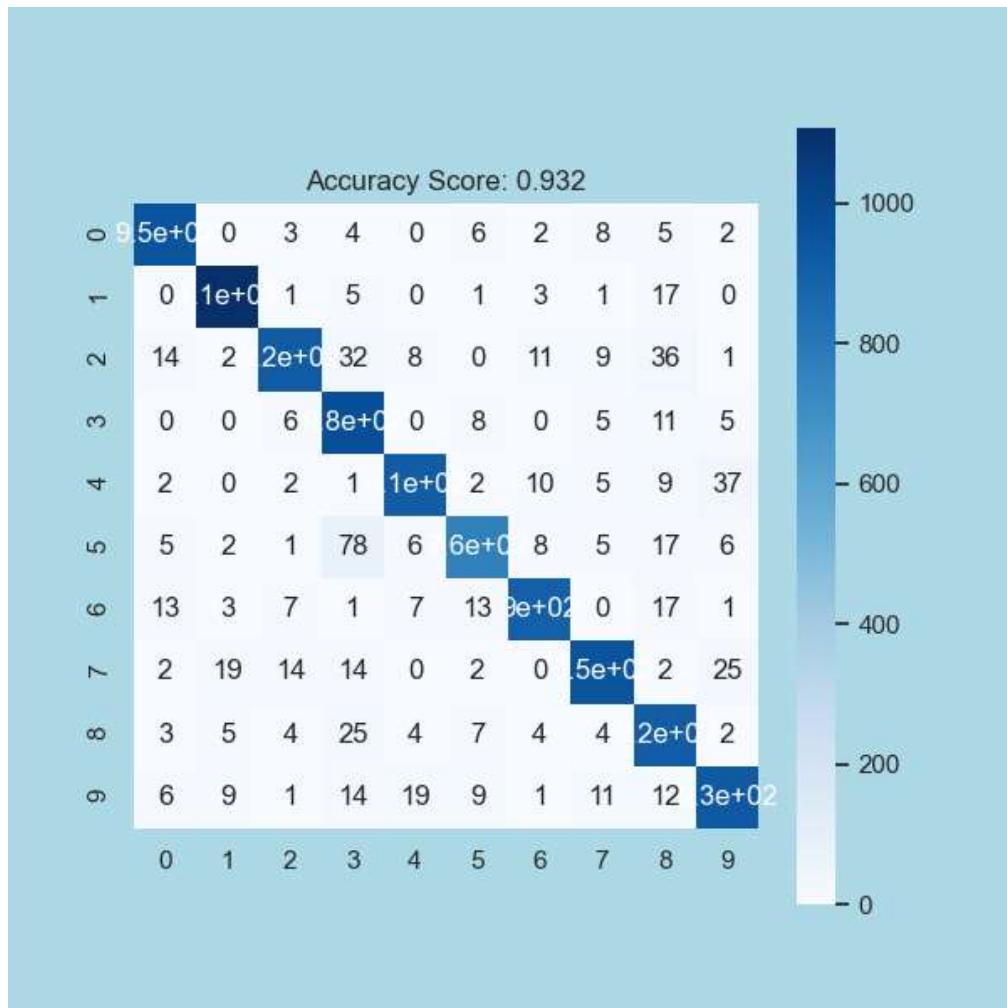
(3): LogSoftmax(dim=1)

)

Training Time (in minutes) = 2.390563261508942

Overall Accuracy Score: **0.932**

Classwise Accuracy Score: [0.9925 0.9932 0.9848 0.9791 0.9888 0.9824 0.9899 0.9874 0.9816 0.9839]



Training with Loss Function = NLLLoss(), Learning Rate = 0.01, Hidden Layer = [256, 128]

Sequential(

(0): Linear(in_features=784, out_features=256, bias=True)

(1): ReLU()

(2): Linear(in_features=256, out_features=128, bias=True)

(3): ReLU()

(4): Linear(in_features=128, out_features=10, bias=True)

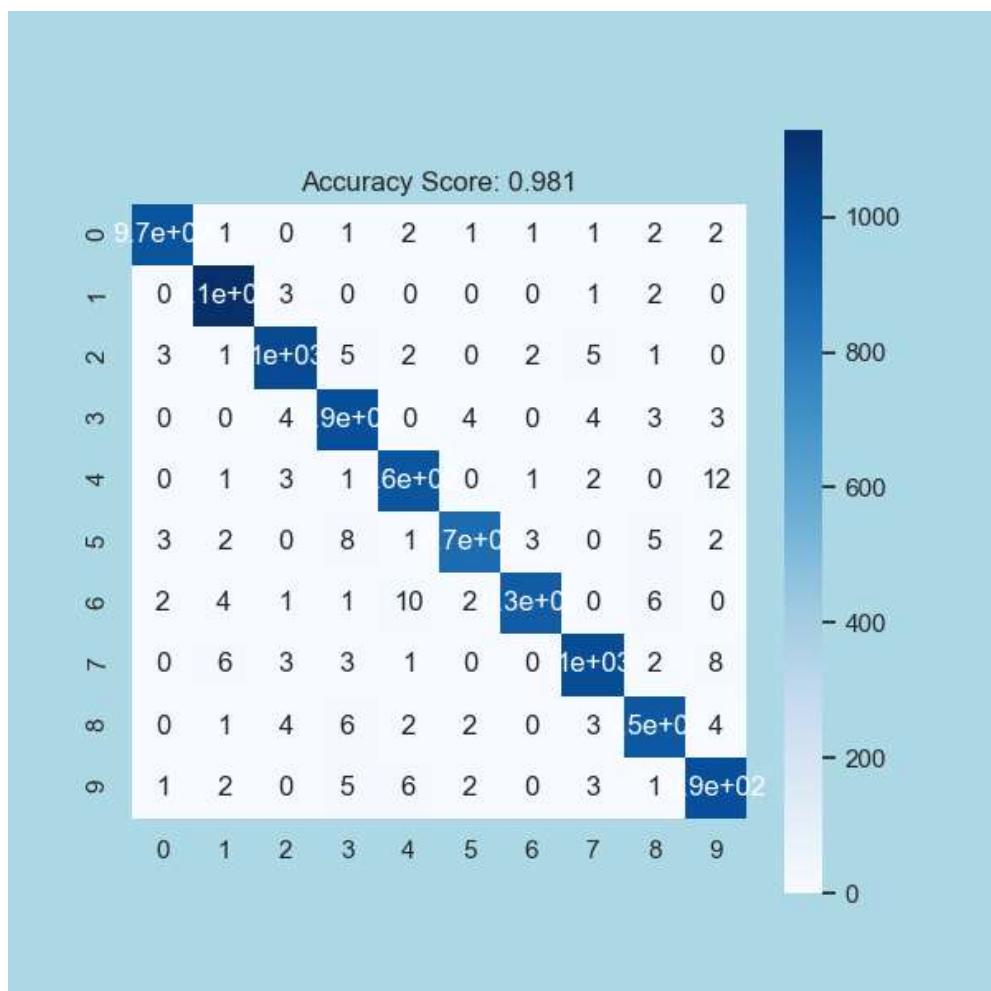
(5): LogSoftmax(dim=1)

)

Training Time (in minutes) = 4.351998011271159

Overall Accuracy Score: **0.981**

Classwise Accuracy Score: [0.998 0.9976 0.9963 0.9952 0.9956 0.9965 0.9967 0.9958 0.9956 0.9949]



Training with Loss Function = NLLLoss(), Learning Rate = 0.01, **Hidden Layer = [64, 16, 16, 16]**

Sequential(

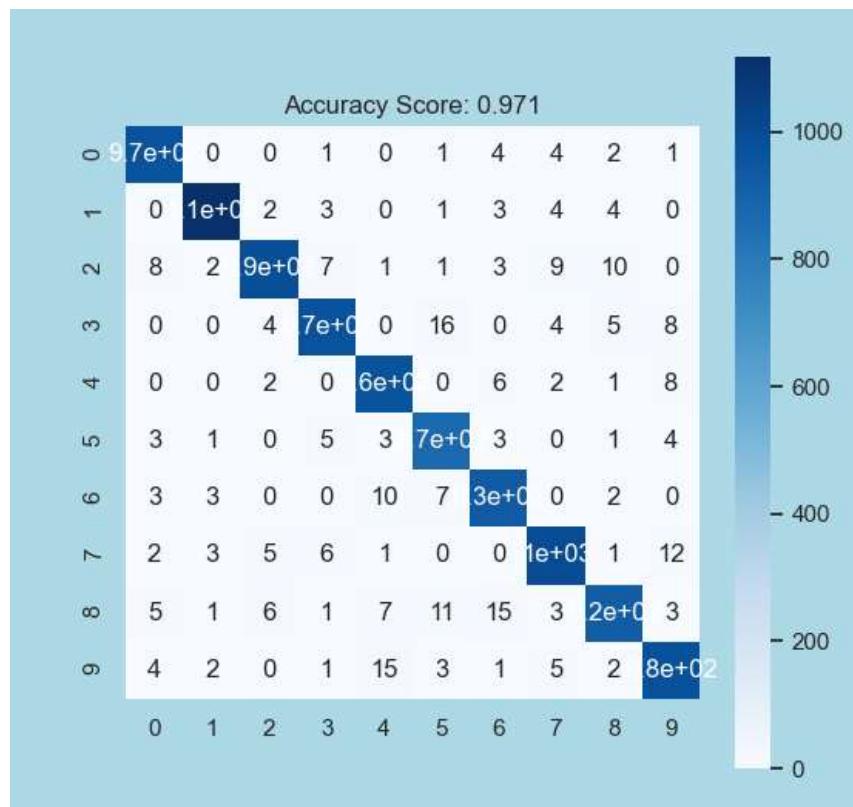
```
(0): Linear(in_features=784, out_features=64, bias=True)
(1): ReLU()
(2): Linear(in_features=64, out_features=16, bias=True)
(3): ReLU()
(4): Linear(in_features=16, out_features=16, bias=True)
(5): ReLU()
(6): Linear(in_features=16, out_features=16, bias=True)
(7): ReLU()
(8): Linear(in_features=16, out_features=10, bias=True)
(9): LogSoftmax(dim=1)
```

)

Training Time (in minutes) = 7.081622072060903

Overall Accuracy Score: **0.971**

Classwise Accuracy Score: [0.9962 0.9971 0.994 0.9939 0.9944 0.994 0.994 0.9939 0.992 0.9931]



Training with Loss Function = NLLLoss(), Learning Rate = 0.01, **Hidden Layer = [256, 128, 64, 32]**

Sequential(

(0): Linear(in_features=784, out_features=256, bias=True)

(1): ReLU()

(2): Linear(in_features=256, out_features=128, bias=True)

(3): ReLU()

(4): Linear(in_features=128, out_features=64, bias=True)

(5): ReLU()

(6): Linear(in_features=64, out_features=32, bias=True)

(7): ReLU()

(8): Linear(in_features=32, out_features=10, bias=True)

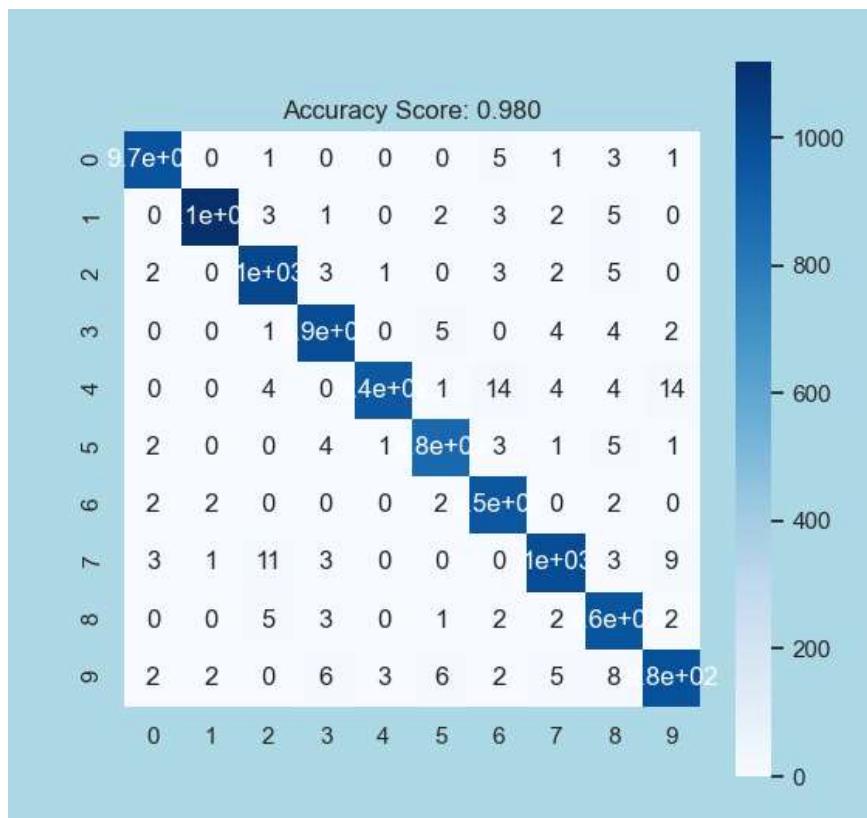
(9): LogSoftmax(dim=1)

)

Training Time (in minutes) = 7.422947756449381

Overall Accuracy Score: **0.980**

Classwise Accuracy Score: [0.9978 0.9979 0.9959 0.9964 0.9954 0.9966 0.996 0.9949 0.9946 0.9937]



2.3. Gurmukhi Handwritten Digit Classification: Gurmukhi is one of the popular Indian scripts widely used in Indian state of Punjab. In this part of the assignment, our goal is to develop a neural network solution for classifying Gurmukhi digits. We provide you Handwritten Gurmukhi digit dataset here:

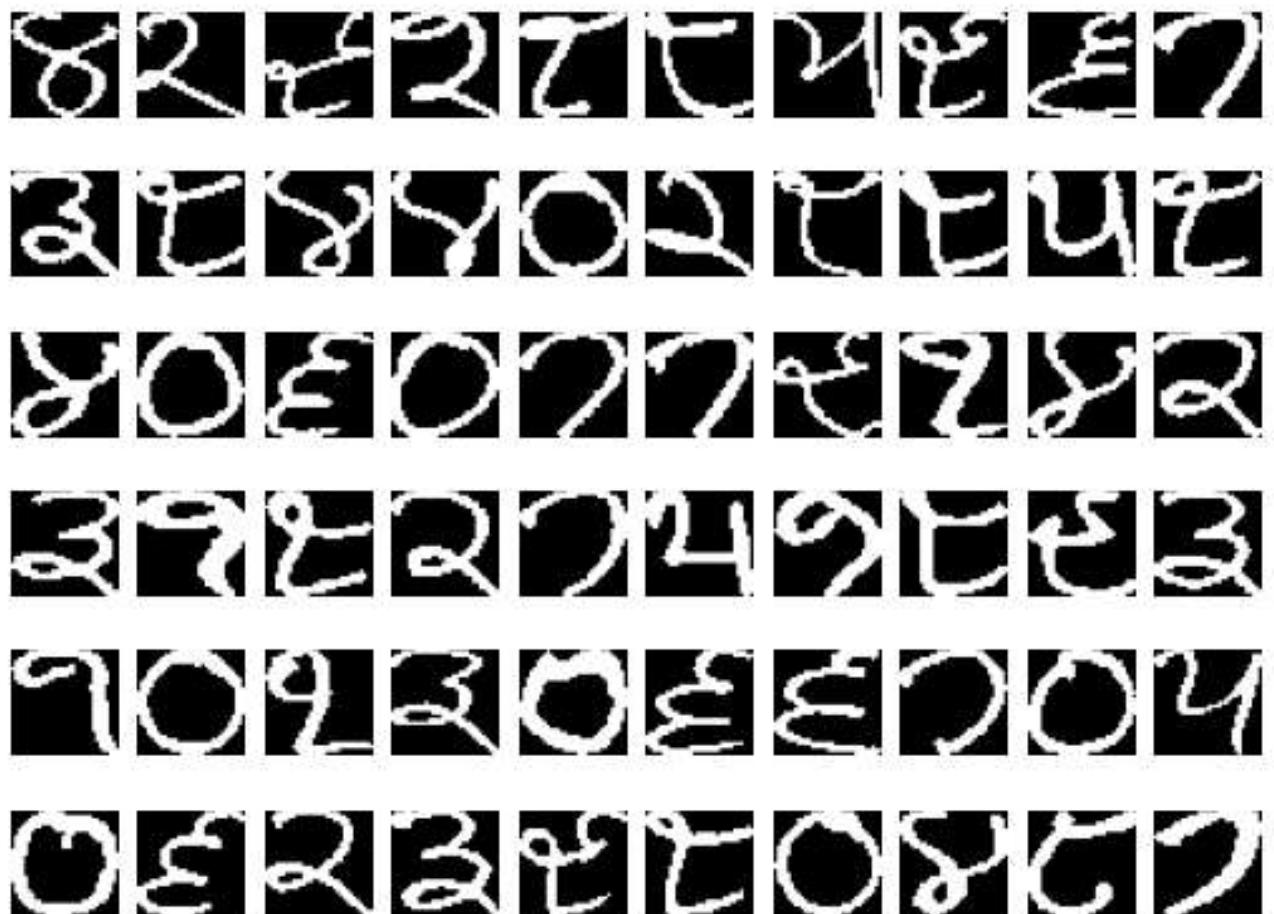
Dataset link

Modifying the code provided in 2, and develop a robust neural network to classify the Gurmukhi digits. Higher performance on test set will have bonus point. Briefly

write your observation and submit your code so that we can evaluate your implementation at our end. (15 points)

Solution:

Sample Data Images:



Model:

```
Sequential(  
    (0): Linear(in_features=1024, out_features=256, bias=True)  
    (1): ReLU()  
    (2): Linear(in_features=256, out_features=64, bias=True)  
    (3): ReLU()  
    (4): Linear(in_features=64, out_features=10, bias=True)  
    (5): LogSoftmax(dim=1)  
)
```

Before backward pass:

None

After backward pass:

```
tensor([[ 3.8962e-04,  3.8962e-04,  3.8962e-04,  ..., -1.3182e-03,  
        -5.4675e-04, -5.9323e-04],  
       [ 6.5492e-04,  6.5492e-04,  8.4178e-04,  ..., -2.3123e-05,  
        2.2393e-04,  4.6235e-05],  
       [ 1.0681e-03,  1.0681e-03,  1.2085e-03,  ..., -9.3416e-04,  
        -5.7025e-04, -3.3054e-04],  
       ...,  
       [-1.1778e-03, -1.1778e-03, -1.1288e-03,  ..., -1.1319e-04,  
        -2.0197e-04,  1.4466e-03],  
       [ 2.4794e-03,  2.3624e-03,  2.0342e-03,  ..., -8.8682e-04,  
        9.1285e-04,  2.0330e-03],  
       [-5.5383e-04, -9.1174e-04, -1.0390e-03,  ..., -2.2249e-03,  
        -1.9248e-03, -6.8608e-04]])
```

Epoch 0 - Training loss: 2.262930393218994 no of images 1000

Epoch 1 - Training loss: 2.079546920955181 no of images 1000

Epoch 2 - Training loss: 1.7204681560397148 no of images 1000

Epoch 3 - Training loss: 1.1766018122434616 no of images 1000

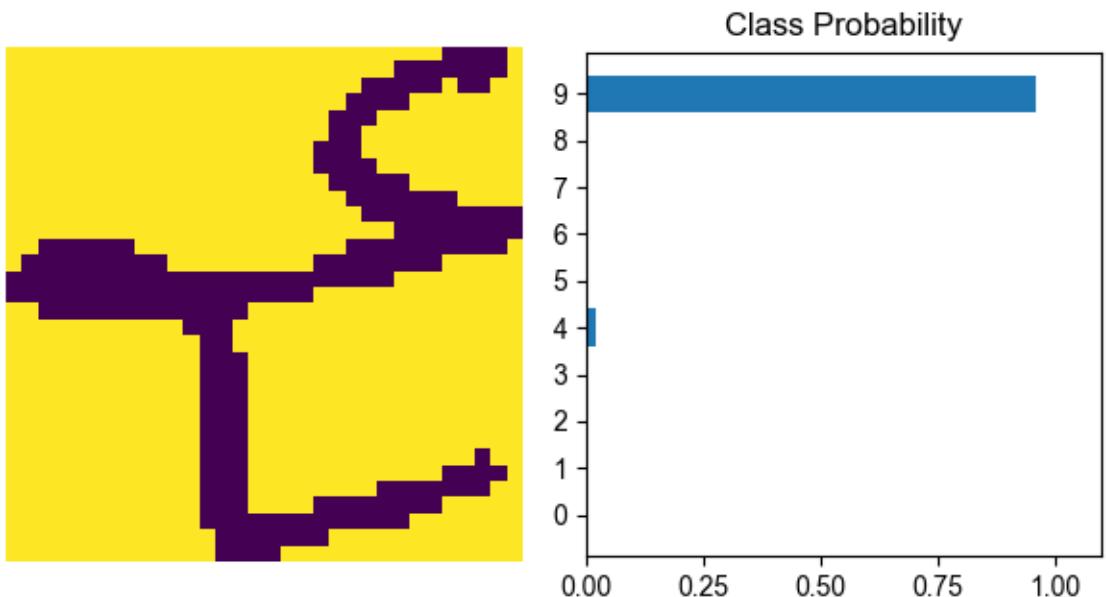
Epoch 4 - Training loss: 0.6834566779434681 no of images 1000

Epoch 5 - Training loss: 0.40679728984832764 no of images 1000

Epoch 6 - Training loss: 0.280624826438725 no of images 1000
Epoch 7 - Training loss: 0.21215492766350508 no of images 1000
Epoch 8 - Training loss: 0.17039840202778578 no of images 1000
Epoch 9 - Training loss: 0.1459644865244627 no of images 1000
Epoch 10 - Training loss: 0.12271467410027981 no of images 1000
Epoch 11 - Training loss: 0.10761574958451092 no of images 1000
Epoch 12 - Training loss: 0.0930116605013609 no of images 1000
Epoch 13 - Training loss: 0.08325718808919191 no of images 1000
Epoch 14 - Training loss: 0.07367376610636711 no of images 1000

Training Time (in minutes) = 0.06675013303756713

Predicted Digit = 9

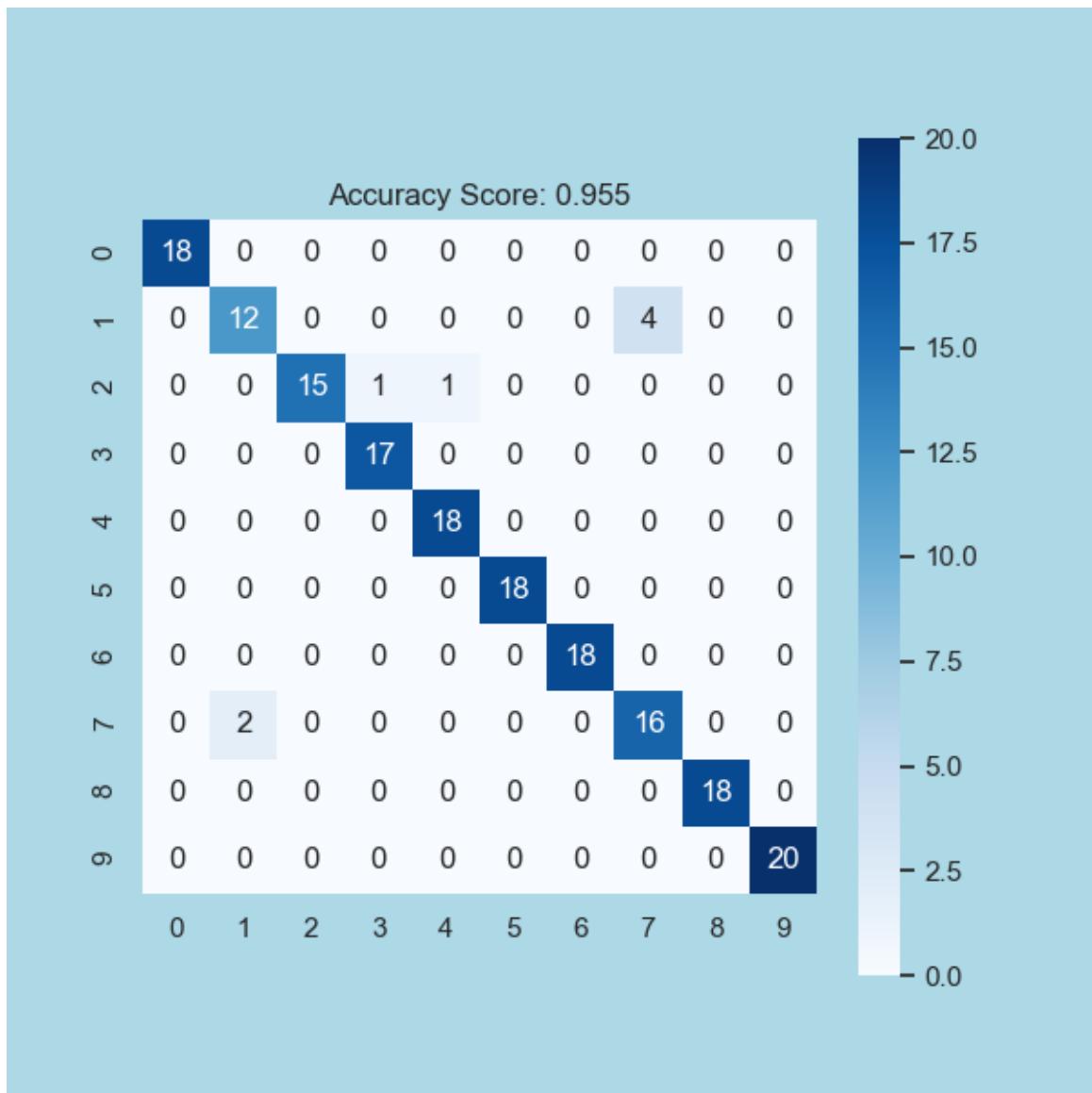


Number Of Images Tested = 178

Model Accuracy = 0.9550561797752809

Overall Accuracy Score: **0.955**

Classwise Accuracy Score: [1. 0.96629213 0.98876404 0.99438202 0.99438202 1.
1. 0.96629213 1. 1.]

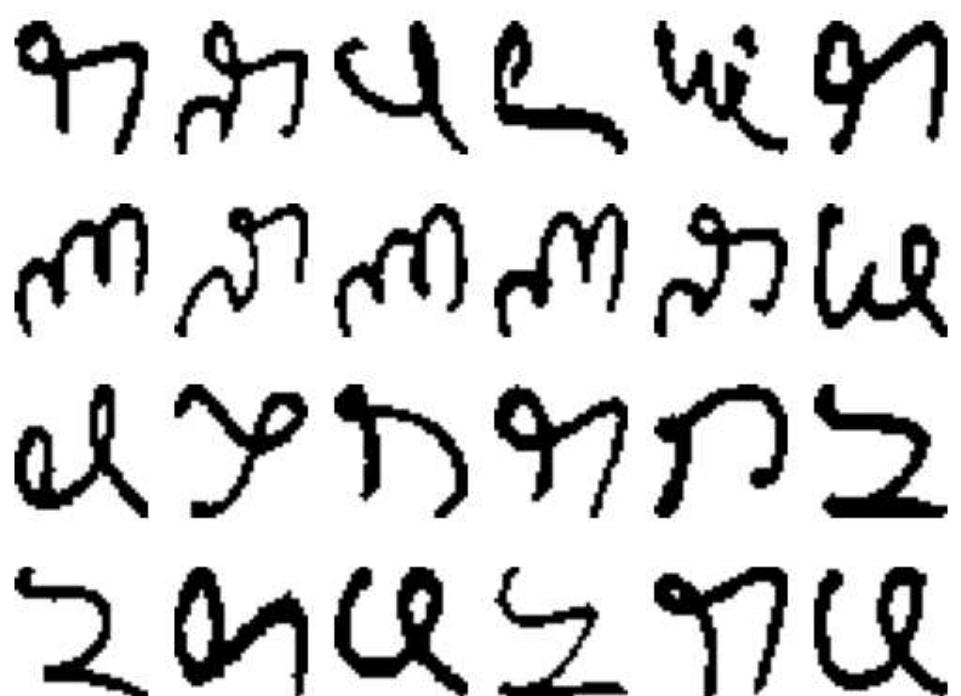


Confusion Matrix for Final prediction of Gurmukhi Image Classification

Problem 4: SVM [20 Points]

1. Solve Subproblem 3 of Problem 2 (i.e. Gurmukhi Handwritten Digit Classification) using SVM. Use raw pixel features and linear, polynomial and RBF Kernel. Resource: <https://scikit-learn.org/stable/modules/svm.html>

Solution:



Accuracy Linear Kernel: 0.9719101123595506

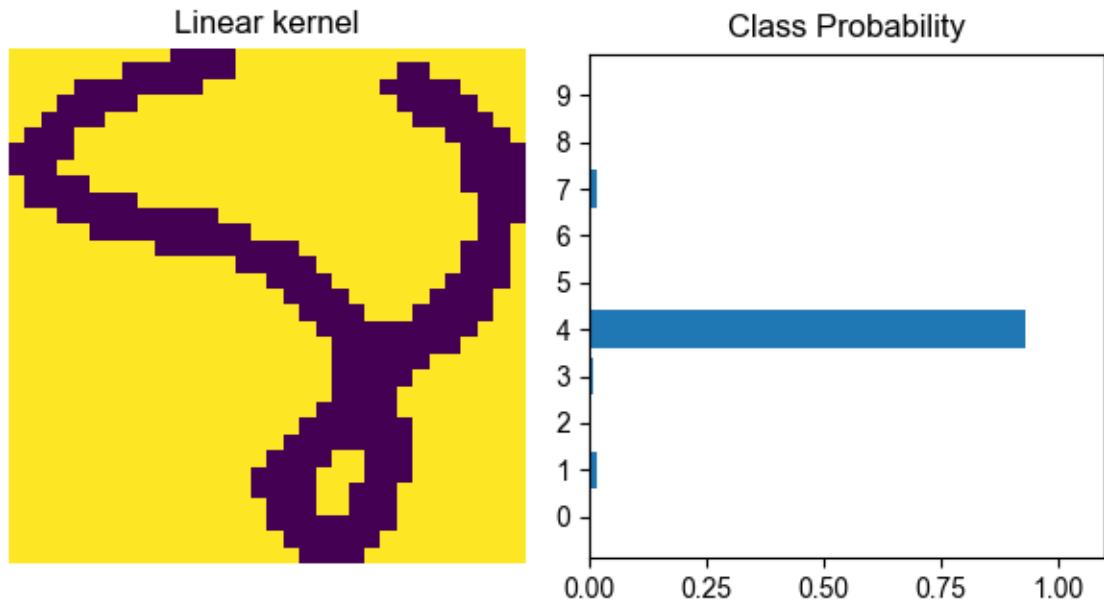
Accuracy Polynomial Kernel: 0.9662921348314607

Accuracy Radial Basis Kernel: 0.20224719101123595

Accuracy Sigmoid Kernel: 0.9550561797752809

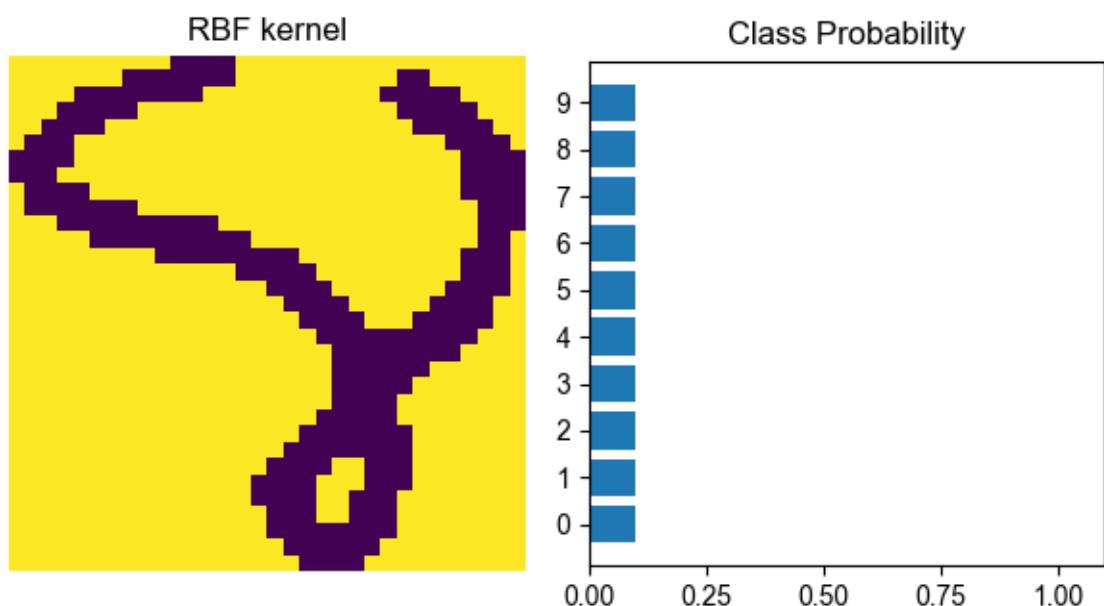
Model = Linear kernel

Predicted Digit = 4



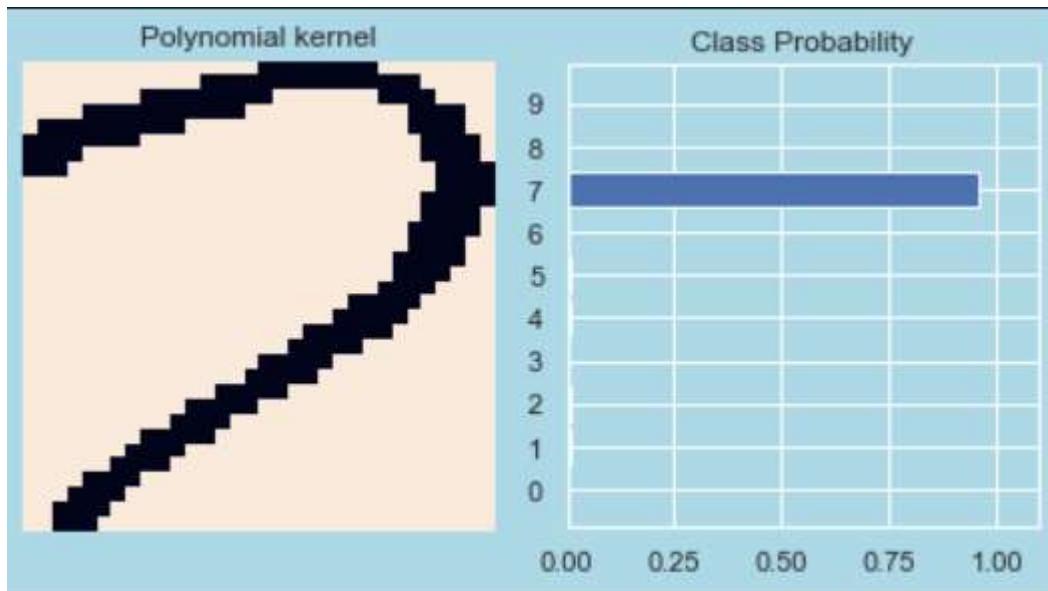
Model = RBF kernel

Predicted Digit = 0



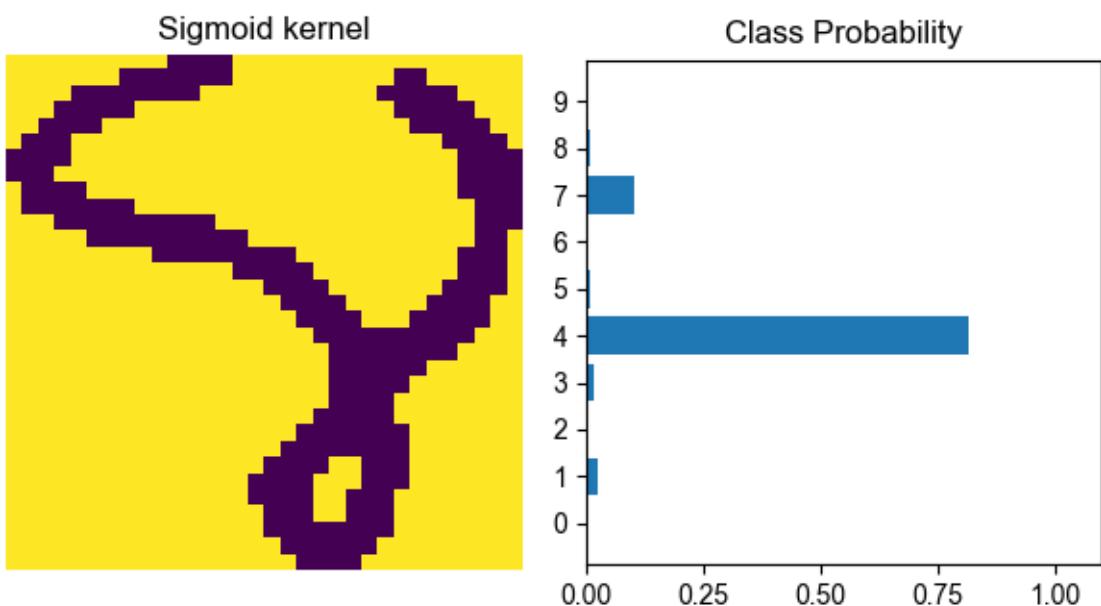
Model = Polynomial kernel

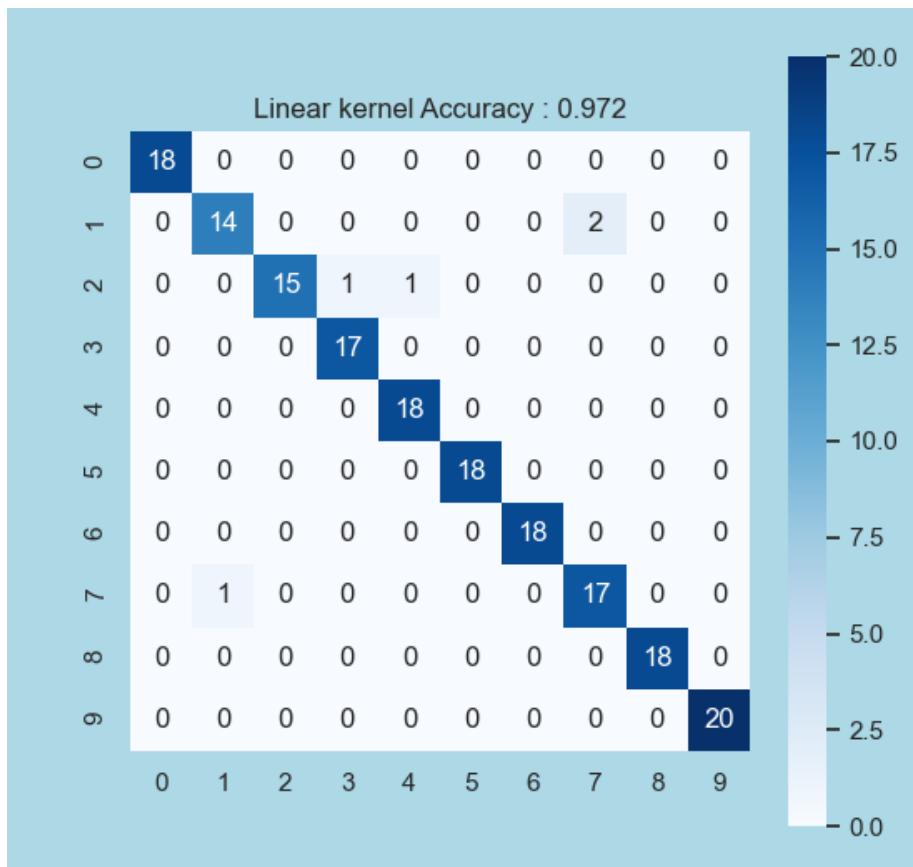
Predicted Digit = 4



Model = Sigmoid kernel

Predicted Digit = 4

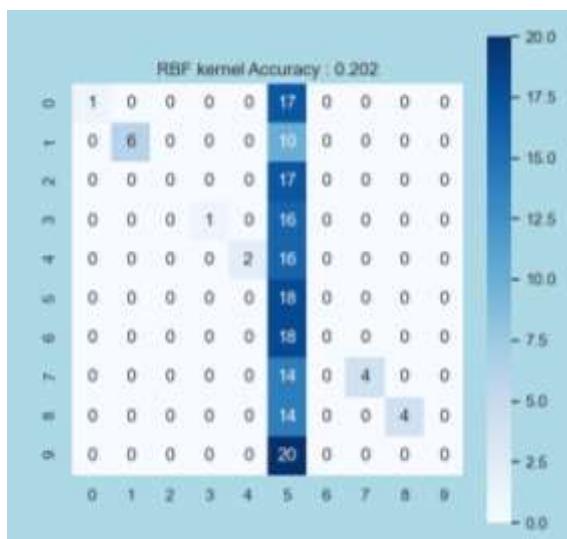




Overall Accuracy Score: 0.972

Classwise Accuracy Score: [1. 0.98314607 0.98876404 0.99438202 0.99438202 1.

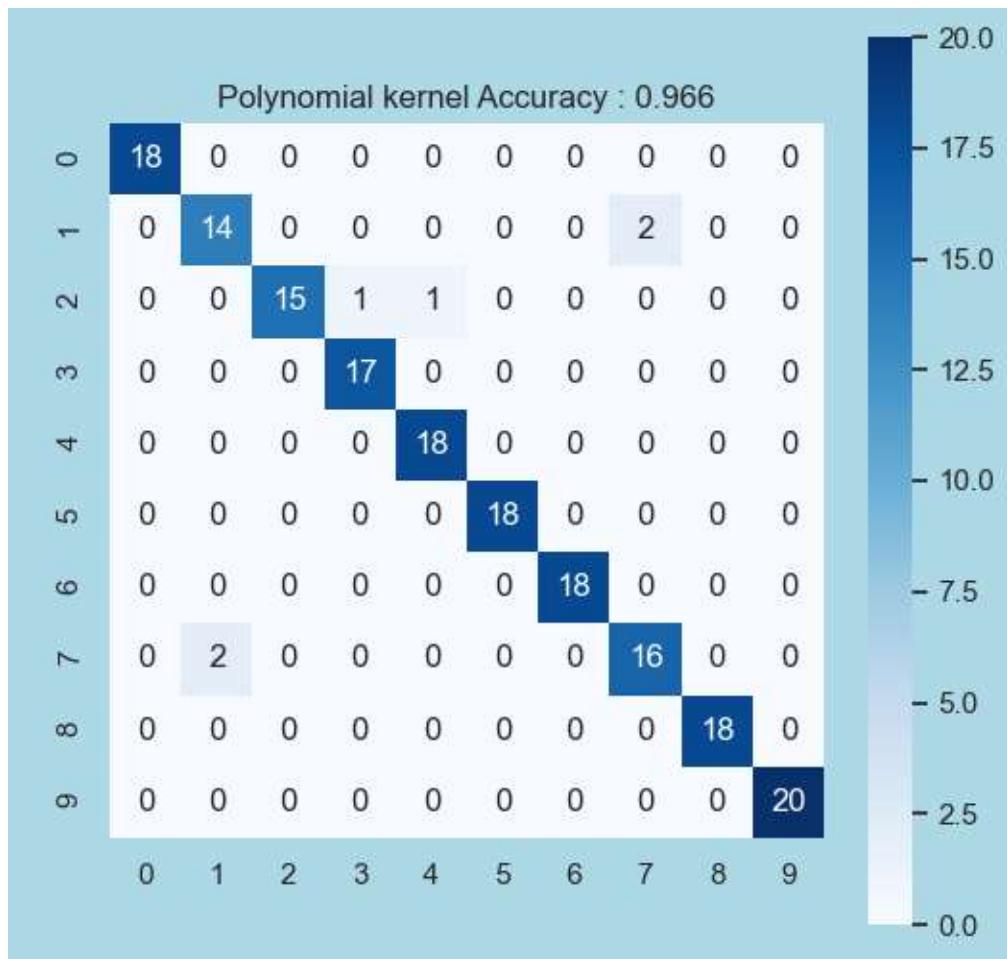
1. 0.98314607 1. 1.]



Overall Accuracy Score: 0.202

Classwise Accuracy Score: [0.90449438 0.94382022 0.90449438 0.91011236 0.91011236
0.20224719

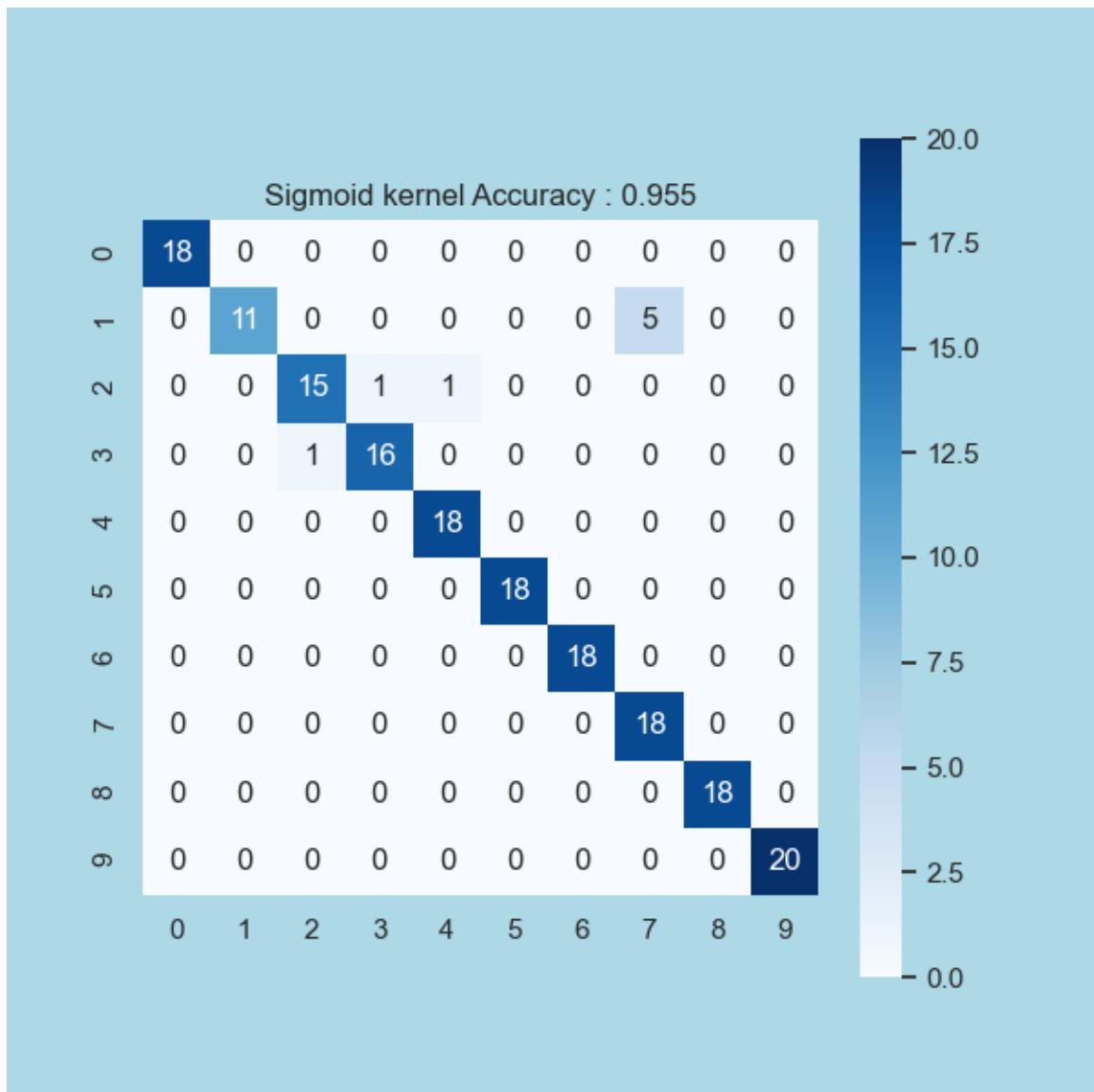
0.8988764 0.92134831 0.92134831 0.88764045]



Overall Accuracy Score: 0.966

Classwise Accuracy Score: [1. 0.97752809 0.98876404 0.99438202 0.99438202 1.

1. 0.97752809 1. 1.]



Overall Accuracy Score: 0.955

Classwise Accuracy Score: [1. 0.97191011 0.98314607 0.98876404 0.99438202 1.

1. 0.97191011 1. 1.]

Problem 2: Learning to implement Neural Network [40 points]

2.1. We have seen in class how backpropagation is used to compute gradients of cost functions with respect to weights. In class example we took Mean Square Error as cost function. Replace this cost function with cross-entropy loss and demonstrate backpropagation. Submit both handwritten explanation as well as modified code. Code link (10 points)

Solution:

Cross-Entropy loss is a popular choice if the problem at hand is a classification problem, and in and of itself it can be classified into either categorical cross-entropy or multi-class cross-entropy (with binary cross-entropy being a special case of the former.)

Let's start with categorical cross-entropy. For this loss function our y 's are one-hot encoded to denote the class our image (or whatever) belongs to. Thus for any x, y is of length equal to the number of classes and the last layer in our model has a neuron for each class. We use Softmax in our last layer to get the probability of x belonging to each of the classes. These probabilities sum to 1.

$$J_{(x,y)} = - \sum_m y_m \cdot \ln(a_m^H)$$

Categorical Cross-Entropy Given One Example. a^H_m is the m th neuron of the last layer (H)

$$\Leftrightarrow z^L = (W^L)^t a^{L-1} + b^L \text{ unless } L = 0 \text{ then } z^L = (W^L)^t x + b^L$$

$$\Leftrightarrow a^L = h(z^L)$$

- $\delta^L = h'(z^L) \odot W^{L+1} \delta^{L+1}$ unless $L = H$ then $\delta^L = (a^L - y) \odot h'(z^L)$
- $\frac{\partial J}{\partial b^L} = \delta^L$
- $\frac{\partial J}{\partial W^L} = a^{L-1} \cdot (\delta^L)^t$ unless $L = 0$ then $\frac{\partial J}{\partial W^L} = x \cdot (\delta^L)^t$

$L=0$ is the first hidden layer, $L=H$ is the last layer. δ is $\partial J / \partial z$

Note that the output (activations vector) for the last layer is a^H and in index notation we would write a^H_n to denote the n th neuron in the last layer. The same applies for the pre-activations vector z^H .

The only equation that we expect to change in the system above is that for δ^H because we used the explicit formula for the loss function to find that. As a result, we'll be only dealing with the **last layer** in the derivation so we might as well drop the superscript for now and keep in mind that we're targeting the last layer whenever we write a , z or δ .

By this, the loss function is

$$J_{(x,y)} = - \sum_m y_m \cdot \ln(a_m)$$

Categorical Cross-Entropy.

with the activation of the n th neuron in the last layer being

$$a_n = h(z_n) = \frac{e^{z_n}}{\sum_m e^{z_m}}$$

Softmax Activation. We'll use this below many times. Keep it in mind.

Notice that the activation of the n th neuron depends on the pre-activations of all other neurons in the layer. This would've not been the case if the last layer involved Sigmoid or ReLU activations. On this account, to find δ for some neuron in the last layer we use the chain-rule by writing

$$\delta_n = \frac{\partial J}{\partial z_n} = \sum_m \frac{\partial J}{\partial a_m} \cdot \frac{\partial a_m}{\partial z_n}$$

Eq 1.0

It will simplify things a bit if we write this as

$$\delta_n = \frac{\partial J}{\partial z_n} = \sum_{m \neq n} \left(\frac{\partial J}{\partial a_m} \cdot \frac{\partial a_m}{\partial z_n} \right) + \frac{\partial J}{\partial a_n} \cdot \frac{\partial a_n}{\partial z_n}$$

considering $m=n$ and $m \neq n$ then adding

Starting with $\partial J / \partial a_n$ we can write

$$\frac{\partial J}{\partial a_n} = \frac{\partial (-\sum_m y_m \cdot \ln(a_m))}{\partial a_n} = \frac{\partial (\textcolor{red}{y_n} \cdot \ln(a_n))}{\partial a_n} = -\frac{y_n}{a_n}$$

The nth neuron is the only survivor in the sum due to ∂a_n

and for $\partial a_n / \partial z_n$ we have

$$\frac{\partial a_n}{\partial z_n} = \frac{\partial (\textcolor{red}{e^{z_n}} / (\sum_m e^{z_m}))}{\partial z_n} = \frac{e^{z_n}}{\sum_m e^{z_m}} \left(\frac{e^{z_n}}{e^{z_n}} - \frac{e^{z_n}}{\sum_m e^{z_m}} \right) = a_n(1 - a_n)$$

You can arrive at the derivative using the quotient rule but there are [other ways](#) that might be worth checking out.

Now multiplying both of our results we and plugging in the original equation we get

$$\delta_n = \frac{\partial J}{\partial z_n} = \sum_{m \neq n} \left(\frac{\partial J}{\partial a_m} \cdot \frac{\partial a_m}{\partial z_n} \right) + -y_n(1 - a_n)$$

By Plugging in Eq. 1

For the remaining sum, let's first compute $\partial J / \partial a_m$ by writing

$$\frac{\partial J}{\partial a_m} = \frac{\partial (-\sum_{m'} y_{m'} \cdot \ln(a_{m'}))}{\partial a_m} = -\frac{y_m}{a_m}$$

If a formula involves a sum over everything you can always change the index to avoid confusion

And then for $\partial a_m / \partial z_n$ ($m \neq n$) we have

$$\begin{aligned}\frac{\partial a_m}{\partial z_n} &= \frac{\partial(\frac{e^{z_m}}{\sum_{m'} e^{z_{m'}}})}{\partial z_n} = \frac{e^{z_m}}{\sum_{m'} e^{z_{m'}}} \left(\frac{\partial e^{z_m}/\partial z_n}{e^{z_m}} - \frac{\partial \sum_{m'} e^{z_{m'}}/\partial z_n}{\sum_{m'} e^{z_{m'}}} \right) = \frac{e^{z_m}}{\sum_{m'} e^{z_{m'}}} (0 - \frac{e^{z_n}}{\sum_{m'} e^{z_{m'}}}) \\ &= -a_m a_n\end{aligned}$$

You can also use the quotient rule to arrive at the same result.

Now multiplying both results, we get

$$\frac{\partial J}{\partial a_m} \cdot \frac{\partial a_m}{\partial z_n} = -\frac{y_m}{a_m} \cdot -a_m a_n = y_m a_n$$

and propagating that back to the original equation we get

$$\delta_n = \frac{\partial J}{\partial z_n} = \sum_{m \neq n} (\frac{\partial J}{\partial a_m} \cdot \frac{\partial a_m}{\partial z_n}) + -y_n(1 - a_n) = \sum_{m \neq n} y_m a_n + -y_n(1 - a_n)$$

which can be simplified to

$$\delta_n = \frac{\partial J}{\partial z_n} = \sum_{m \neq n} y_m a_n - y_n(1 - a_n) = \sum_{m \neq n} y_m a_n + y_n a_n - y_n = \sum_m y_m a_n - y_n = a_n \sum_m y_m - y_n$$

provided that y is a one-hot vector, we know that $\sum_m y_m = 1$; thus, we can write

$$\delta_n = \frac{\partial J}{\partial z_n} = a_n - y_n$$

or in vector form after reimposing the superscript (H to denote the last layer)

$$\delta^H = \frac{\partial J}{\partial z^H} = a^H - y$$

To conclude the proof, let's update the backpropagation equations to

$$\Rightarrow z^L = (W^L)^t a^{L-1} + b^L \text{ unless } L = 0 \text{ then } z^L = (W^L)^t x + b^L$$

$$\Rightarrow a^L = h(z^L)$$

- $\delta^L = h'(z^L) \odot W^{L+1} \delta^{L+1}$ unless $L = H$ then $\delta^L = (a^L - y)$
- $\frac{\partial J}{\partial b^L} = \delta^L$
- $\frac{\partial J}{\partial W^L} = a^{L-1} \cdot (\delta^L)^t$ unless $L = 0$ then $\frac{\partial J}{\partial W^L} = x \cdot (\delta^L)^t$

Only what's in red has changed.

We're not all there yet. Using categorical cross-entropy your model would do good classifying the image below as a dog; however, your dataset might include many images that have both cats and dogs in them. In this case, y is no longer a one-hot vector (and rather looks something like [0 1 0 0 1] if two classes are present within the image.) Thus, we use multi-class cross-entropy and refrain from using Softmax in the last layer; instead, we use Sigmoid. We conclude that a class is present in the image if its Sigmoid activation is greater than some threshold (0.5 for instance.)

The multi-class cross-entropy loss function for one example is given by

$$J_{(x,y)} = - \sum_m y_m \cdot \ln(a_m^H) - \sum_m (1 - y_m) \cdot \ln(1 - a_m^H)$$

a_m^H is the m th neuron in the last layer (H)

If we go back to dropping the superscript we can write

$$J_{(x,y)} = - \sum_m y_m \cdot \ln(a_m) - \sum_m (1 - y_m) \cdot \ln(1 - a_m)$$

Because we're using Sigmoid, we also have

$$a_n = h(z_n) = \frac{1}{1 + e^{-z_n}}$$

Unlike Softmax a_n is only a function in z_n ; thus, to find δ for the last layer, all we need to consider is that

$$\delta_n = \frac{\partial J}{\partial z_n} = \frac{\partial J}{\partial a_n} \cdot \frac{\partial a_n}{\partial z_n}$$

Eq. 2

or more precisely

$$\frac{\partial J}{\partial a_n} = \frac{\partial(-\sum_m y_m \cdot \ln(a_m) - \sum_m (1 - y_m) \cdot \ln(1 - a_m))}{\partial a_n}$$

which is

$$\frac{\partial J}{\partial a_n} = -(\frac{\partial \sum_m y_m \cdot \ln(a_m)}{\partial a_n} + \frac{\partial \sum_m (1 - y_m) \cdot \ln(1 - a_m)}{\partial a_n})$$

by differentiating (like we did earlier) then adding the two fractions we get

$$\frac{\partial J}{\partial a_n} = -\left(\frac{y_n}{a_n} + \frac{1 - y_n}{1 - a_n}\right) = \frac{a_n - y_n}{a_n(1 - a_n)}$$

Now we need to consider $\partial a_n / \partial z_n$ before plugging in equation 2.0

$$\frac{\partial a_n}{\partial z_n} = \frac{\partial h(z_n)}{\partial z_n} = \frac{\partial(1/(1 + e^{-z_n}))}{\partial z_n} = \frac{1}{1 + e^{-z_n}} \left(0 - \frac{\partial(1 + e^{-z_n})/\partial z_n}{1 + e^{-z_n}}\right)$$

this is simply the derivative of the Sigmoid function

$$\frac{\partial a_n}{\partial z_n} = \frac{1}{1 + e^{-z_n}} \left(\frac{e^{-z_n}}{1 + e^{-z_n}} \right) = \frac{1}{1 + e^{-z_n}} \left(\frac{1 + e^{-z_n} - 1}{1 + e^{-z_n}} \right) = \frac{1}{1 + e^{-z_n}} (1 - \frac{1}{1 + e^{-z_n}})$$

by substituting a_n for it's definition we get

$$\frac{\partial a_n}{\partial z_n} = a_n(1 - a_n)$$

Now using both results in the original equation

$$\delta_n = \frac{\partial J}{\partial z_n} = \frac{\partial J}{\partial a_n} \cdot \frac{\partial a_n}{\partial z_n} = \frac{a_n - y_n}{a_n(1 - a_n)} \cdot a_n(1 - a_n)$$

thereby,

$$\delta_n = a_n - y_n$$

and now by reimposing the superscript and writing this in vector form

$$\delta^H = a^H - y$$

which quite marvelously is the same result we got for categorical cross-entropy with Softmax as an activation. So we still use

$$\Rightarrow z^L = (W^L)^t a^{L-1} + b^L \text{ unless } L = 0 \text{ then } z^L = (W^L)^t x + b^L$$

$$\Rightarrow a^L = h(z^L)$$

- $\delta^L = h'(z^L) \odot W^{L+1} \delta^{L+1}$ unless $L = H$ then $\delta^L = (a^L - y)$

- $\frac{\partial J}{\partial b^L} = \delta^L$

- $\frac{\partial J}{\partial W^L} = a^{L-1} \cdot (\delta^L)^t$ unless $L = 0$ then $\frac{\partial J}{\partial W^L} = x \cdot (\delta^L)^t$

as our backpropagation equations.

References:

1. <https://towardsdatascience.com/multiclass-classification-with-support-vector-machines-svm-kernel-trick-kernel-functions-f9d5377d6f02>
2. <https://towardsdatascience.com/deriving-backpropagation-with-cross-entropy-loss-d24811edeaf9>
3. <http://neuralnetworksanddeeplearning.com/chap2.html>
4. <https://towardsdatascience.com/mse-is-cross-entropy-at-heart-maximum-likelihood-estimation-explained-181a29450a0b>
5. Youtube.com
6. NPTEL