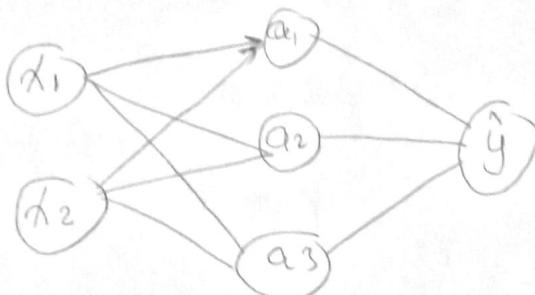


Assignment 1

1



1a) first layer

$$V^{(1)} = \begin{pmatrix} W^{(1)T} X + B^{(1)} \\ (3 \times 2) (2 \times 1) \end{pmatrix} = \begin{bmatrix} 2 & 2 \\ 1 & -1 \\ 3 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 3 \end{bmatrix} + \begin{bmatrix} 1 \\ 0 \\ -1 \end{bmatrix} = \begin{bmatrix} 8 & +1 \\ -2 \\ 6 -1 \end{bmatrix}$$

$$V^{(1)} = \begin{bmatrix} 9 \\ -2 \\ 5 \end{bmatrix}$$

$$\hat{y} = W^{(2)} V^{(1)} + b^{(2)} = [3 \quad 1 \quad 2] \begin{bmatrix} 9 \\ -2 \\ 5 \end{bmatrix} + 1 = 27 - 2 + 10 + 1$$

$$\hat{y} = 36$$

1b)

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

$$\text{ReLU}(V^{(1)}) = \begin{bmatrix} 9 \\ 0 \\ 5 \end{bmatrix}$$

$$\hat{y} = [3 \quad 1 \quad 2] \begin{bmatrix} 9 \\ 0 \\ 5 \end{bmatrix} + 1 = 27 + 10 + 1 = 38$$

$$\hat{y} = 38$$

1c)

$$J = (\hat{y} - y)^2$$

$$\frac{\partial J}{\partial b_1^{(2)}} = \frac{\partial J}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial b_1^{(2)}} = 2(\hat{y} - y)(2)$$

$$\hat{y} = v_1^{(1)} w_1^{(2)} + v_2^{(1)} w_2^{(2)} + v_3^{(1)} w_3^{(2)} + b^{(2)}$$

$$\boxed{\frac{\partial J}{\partial b_1^{(2)}} = 2(36 - 32) = 8}$$

$$\frac{\partial J}{\partial w_{21}^{(2)}} = \frac{\partial J}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial w_{21}^{(2)}} = 2(\hat{y} - y) v_2^{(1)}$$

$$\boxed{\frac{\partial J}{\partial w_{21}^{(2)}} = 2(36 - 32)(-2) = -16}$$

$$\frac{\partial J}{\partial b_2^{(1)}} = \frac{\partial J}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial v_2^{(1)}} \frac{\partial v_2^{(1)}}{\partial b_2^{(1)}}$$

$$\frac{\partial J}{\partial b_2^{(1)}} = 2(\hat{y} - y)(w_{21}^{(2)}) (1) = 2(36 - 32)(1)(1)$$

$$\boxed{\frac{\partial J}{\partial b_2^{(1)}} = 8}$$

$$\begin{aligned} \frac{\partial J}{\partial w_{31}^{(1)}} &= \frac{\partial J}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial v_3^{(1)}} \frac{\partial v_3^{(1)}}{\partial w_{31}^{(1)}} \\ &= 2(\hat{y} - y)(w_{31}^{(2)})(x_1) \\ &= 2 * 4 * 2 * 3 \end{aligned}$$

$$\boxed{\frac{\partial J}{\partial w_{31}^{(1)}} = 48}$$

(2)

(1d)

$$b_2^{(1)}_{n+1} = b_2^{(1)}_n - \gamma \frac{2J}{\alpha b_2^{(1)}_n}$$

$$b_2^{(1)}_{n+1} = 0 - 2(8)$$

$$\boxed{b_2^{(1)}_{n+1} = -16}$$

$$w_{13}^{(1)}_{n+1} = w_{13}^{(1)}_n - \gamma \frac{2J}{\alpha w_{13}^{(1)}_n}$$

$$w_{13}^{(1)}|_{n+1} = 3 - 2(48)$$

$$\boxed{w_{13}^{(1)}|_{n+1} = -93}$$

yes as the test set is uncorrect to our model

Question 1 (e)

Some how it will be good to use test set as an indicator to measure out of sample error, but data set could be correlated some how, sow we need to test on more data sets to validate our model

$$(2) \quad f: \mathbb{R}^2 \rightarrow \mathbb{R}$$

$$f = \sin g_1 + g_2^2$$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial g_1} \frac{\partial g_1}{\partial x}$$

(1×2)

$$g: \mathbb{R}^2 \rightarrow \mathbb{R}$$

$$g(x) = \begin{bmatrix} x_1 e^{x_2} \\ x_1 + x_2^2 \end{bmatrix}$$

$$\begin{bmatrix} \frac{\partial f}{\partial g_1} & \frac{\partial f}{\partial g_2} \end{bmatrix} \begin{bmatrix} \frac{\partial g_1}{\partial x_1} & \frac{\partial g_1}{\partial x_2} \\ \frac{\partial g_2}{\partial x_1} & \frac{\partial g_2}{\partial x_2} \end{bmatrix}$$

$$= \begin{bmatrix} \frac{\partial f}{\partial g_1} & \frac{\partial g_1}{\partial x_1} \end{bmatrix} + \frac{\partial f}{\partial g_2} \begin{bmatrix} \frac{\partial g_2}{\partial x_1} & \frac{\partial g_2}{\partial x_2} \end{bmatrix}$$

$$= \begin{bmatrix} \frac{\partial f}{\partial x_1} & \frac{\partial f}{\partial x_2} \end{bmatrix}$$

$$\boxed{\frac{\partial f}{\partial x_1} = (\cos g_1)(e^{x_2}) + (2g_2)(1)}$$

$$\boxed{\frac{\partial f}{\partial x_2} = (\cos g_1)(x_1 e^{x_2}) + (2g_2)(2x_2)}$$

$$\frac{\partial f}{\partial x_i} = \sum_{j=1}^m \frac{\partial f}{\partial g_j} \frac{\partial g_j}{\partial x_i}$$

④

(3)

(3.1)

$$f = \frac{1}{1+e^{-z}}$$

$$\frac{df}{dz} = \frac{-1}{(1+e^{-z})^2} * (-e^{-z})$$

$$= \frac{e^{-z}}{(1+e^{-z})^2} = \left(\frac{e^{-z}}{1+e^{-z}}\right) \left(\frac{1}{1+e^{-z}}\right) = \left(\frac{1+e^{-z}-1}{1+e^{-z}}\right) f(z)$$

$$\boxed{\frac{df}{dz} = (1-f(z))(f(z))}$$

(3.2)

$$f(w) = \frac{1}{1-e^{\overset{T}{w^T}x}}$$

(1×1) $\overset{T}{w^T} \in (1 \times D)$ $x \in \mathbb{R}^D$

$$w^T x = w_1 x_1 + w_2 x_2 + \dots + w_D x_D$$

$$\boxed{\frac{\partial f}{\partial w_i} = \frac{\partial f}{\partial z} \frac{\partial z}{\partial w_i} = [(f(w))(1-f(w))] [x_i]}$$

$$\frac{df}{dw} = \left[\frac{\partial f}{\partial w_1} \quad \frac{\partial f}{\partial w_2} \quad \dots \quad \frac{\partial f}{\partial w_D} \right]$$

$$\boxed{\frac{df}{dw} = \overset{T}{x^T} [(f(w))(1-f(w))]} \quad \begin{matrix} \nearrow \\ (1 \times D) \end{matrix}$$

(5)

$$\textcircled{3.3} \quad J(w) = \frac{1}{2} \sum_{i=1}^m \left| w^T x^{(i)} - y^{(i)} \right|$$

(1×1) $\begin{matrix} w^T \\ x^{(i)} \end{matrix}$ $(1 \times D)$ $(D \times 1)$

$$w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_D \end{bmatrix} \quad x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix}$$

$$w^T x = w_1 x_1 + w_2 x_2 + \dots + w_D x_D$$

$$\frac{\partial J}{\partial w} = \begin{bmatrix} \frac{\partial J}{\partial w_1} & \frac{\partial J}{\partial w_2} & \dots & \frac{\partial J}{\partial w_D} \end{bmatrix}$$

$(1 \times D)$

$$\frac{\partial J}{\partial w_1} = \frac{1}{2} \sum_{i=1}^m \text{Sign}(w^T x^{(i)} - y^{(i)}) (x^{(i)}_1)$$

$$\frac{\partial J}{\partial w} = \frac{1}{2} \sum_{i=1}^m \text{Sign}(w^T x^{(i)} - y^{(i)}) (x^{(i)T})$$

(1×1) $(1 \times D)$

side note

$$f(x) = |x| = \sqrt{x^2} \quad f(x)$$

$$\frac{df}{dx} = \begin{cases} 1 & x > 0 \\ -1 & x < 0 \end{cases}$$

$$f(x) = \sqrt{x^2}$$

$$\frac{df}{dx} = \frac{1}{2} (x^2)^{-1/2} \cdot 2x$$

$$\frac{df}{dx} = \frac{x}{\sqrt{x^2}} = \frac{x}{|x|}$$

$$\frac{df}{dx} = \text{Sign}(x)$$

$$x \neq 0$$

$$\text{Sign}(x) = \begin{cases} 1, x > 0 \\ -1, x < 0 \end{cases}$$

⑥

(3.4)

$$J(w) = \frac{1}{2} \left[\sum_{i=1}^m \left(\underbrace{w^T x^{(i)}}_{(1 \times D)} - \underbrace{y^{(i)}}_{(1 \times 1)} \right)^2 \right] + \lambda \|w\|_2^2$$

$$\frac{\partial J}{\partial w} = \left[\frac{\partial J}{\partial w_1} \quad \frac{\partial J}{\partial w_2} \quad + \dots \dots + \frac{\partial J}{\partial w_D} \right]_{(1 \times D)}$$

$$\frac{\partial J}{\partial w_K} = \left[\sum_{i=1}^m \left(w^T x^{(i)} - y^{(i)} \right) x_K^{(i)} \right] + 2\lambda w_K$$

$$\boxed{\frac{\partial J}{\partial w} = \underbrace{\sum_{i=1}^m \left(\underbrace{(w^T x^{(i)} - y^{(i)})}_{(1 \times D)} \underbrace{x^{(i)T}}_{(1 \times D)} \right)}_{(1 \times D)} + 2\lambda w^T}$$

(7)

3.5

$$J(w) = \sum_{i=1}^m y^{(i)} \log \left(\frac{1}{1 + e^{-w^T x^{(i)}}} \right) + (1 - y^{(i)}) \log \left(1 - \frac{1}{1 + e^{-w^T x^{(i)}}} \right)$$

\uparrow
 (1×1)

$$\frac{\partial J}{\partial w} = \begin{bmatrix} \frac{\partial J}{\partial w_1} & \frac{\partial J}{\partial w_2} & \dots & \frac{\partial J}{\partial w_D} \end{bmatrix}$$

$(1 \times D)$

Let $\delta(x) = \frac{1}{1 + e^{-x}}$ $\rightarrow \frac{\partial \delta(x)}{\partial x} = (1 - \delta(x))\delta'(x)$

Side note

$$\log_e \equiv \log$$

$$\frac{d}{dx} (\log x) = \frac{1}{x}$$

$$J(w) = \sum_{i=1}^m y^{(i)} \log(\delta(w^T x^{(i)})) + (1 - y^{(i)}) \log(1 - \delta(w^T x^{(i)}))$$

$$\frac{\partial J}{\partial w_1} = \sum_{i=1}^m \left(y^{(i)} \right) \left(\frac{1}{\delta(w^T x^{(i)})} \right) (1 - \delta(w^T x^{(i)})) \delta'(w^T x^{(i)}) (x_1^{(i)}) \\ + (1 - y^{(i)}) \left(\frac{1}{1 - \delta(w^T x^{(i)})} \right) (-1) (1 - \delta(w^T x^{(i)})) (1 - \delta(w^T x^{(i)})) \delta'(w^T x^{(i)}) (x_1^{(i)})$$

$$\frac{\partial J}{\partial w_1} = \sum_{i=1}^m \left[y^{(i)} (1 - \delta(w^T x^{(i)})) - (1 - y^{(i)}) \delta(w^T x^{(i)}) \right] x_1^{(i)}$$

$$\frac{\partial J}{\partial w} = \sum_{i=1}^m \left[y^{(i)} (1 - \delta(w^T x^{(i)})) - (1 - y^{(i)}) \delta(w^T x^{(i)}) \right] x^{(i)T}$$

$(1 \times D)$

(3.6)

$$f(w) = \tanh [w^T x]$$

(1×1) $(1 \times D)$ $(D \times 1)$

$$\nabla_w f = \begin{bmatrix} \frac{\partial f}{\partial w_1} & \frac{\partial f}{\partial w_2} & \cdots & \frac{\partial f}{\partial w_D} \end{bmatrix}$$

$$\frac{\partial f}{\partial w_1} = 4(1 - g(2w^T x))(g(2w^T x))(x_1)$$

$$\nabla_w f = x^T \left[4(1 - g(2w^T x))(g(2w^T x)) \right]$$

$(1 \times D)$ $(1 \times D)$ (1×1)

Side note

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

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

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

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

$$\tanh(x) = 2g(2x) - 1$$

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

$$= 4(1 - g(2x))(g(2x))$$

(9)

Assignment 1CSE 616

Abdullah Aml, 2101398

[github Ilink for the project](#)

I have trained two data sets: [Iris](#), and [DryBean](#).

The code is:

- customized to run generic number of hidden layers
- creates a new directory to save: plots, scores, lossses, weights in “train-dryBean” directory

Iris

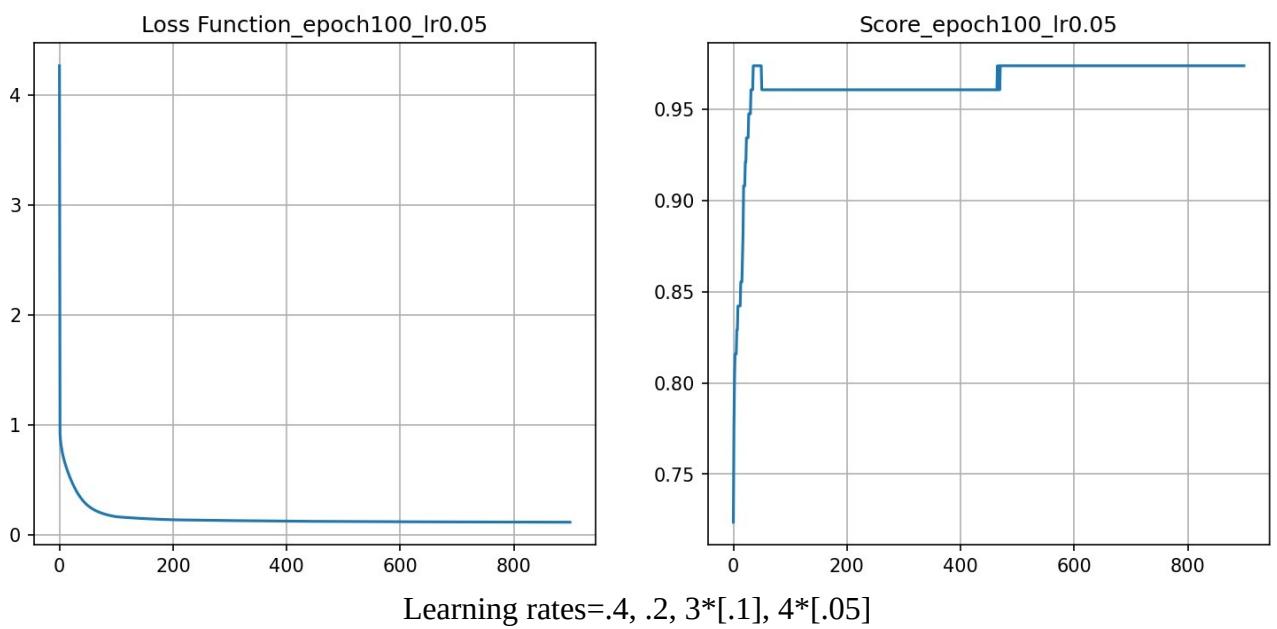
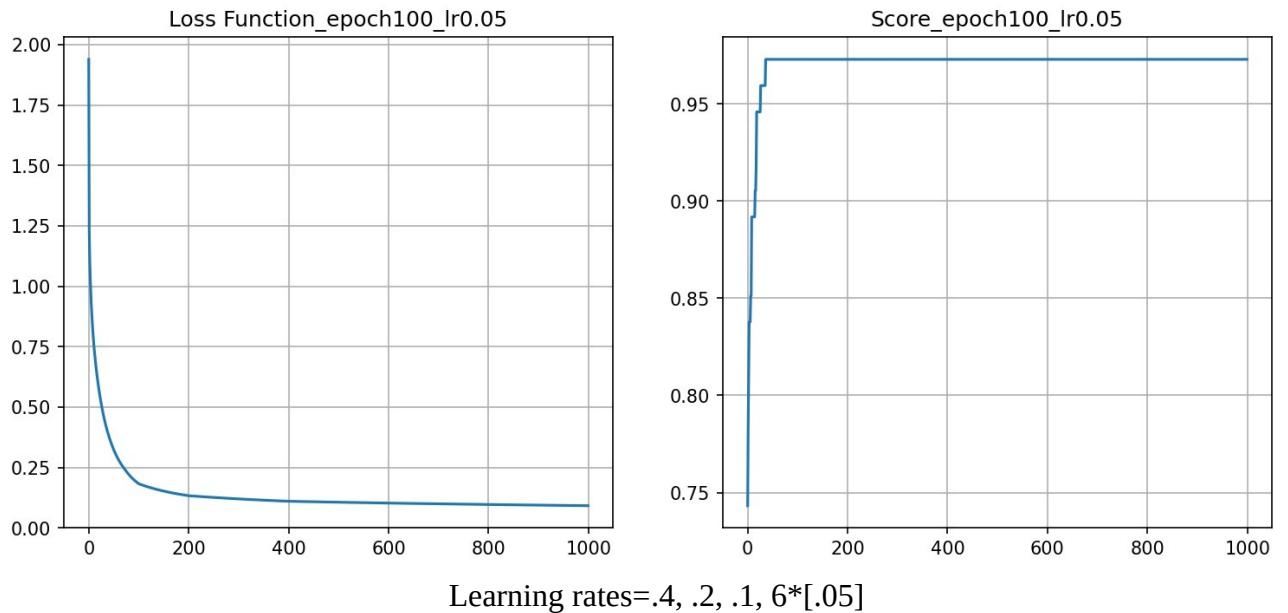
Very simple dataset I used to learn and build working model step by step: it has 4 attributes and 150 sample in total with 3 classes:

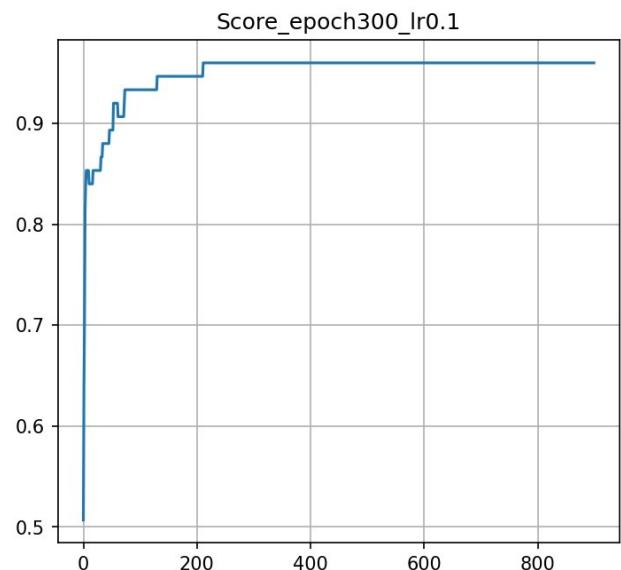
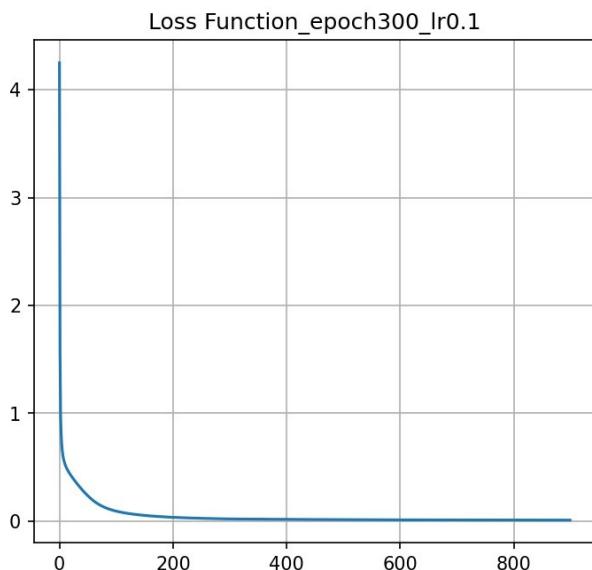
-- Iris Setosa
-- Iris Versicolour
-- Iris Virginica

I used:

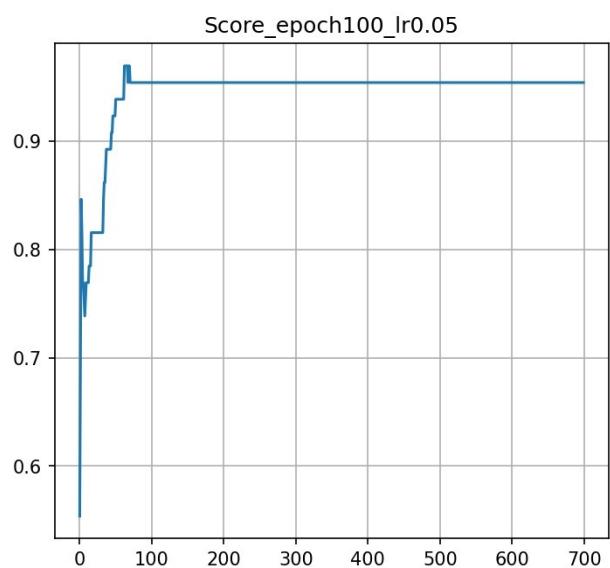
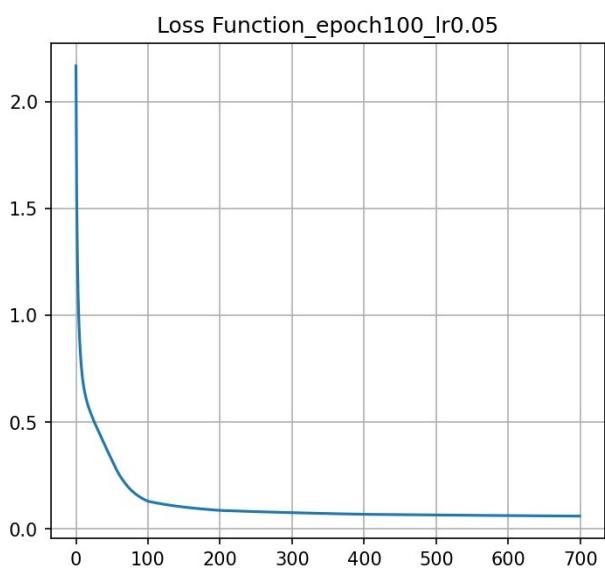
- a single hidden layer with 10 neurons
- LeakyRelu in the middle and sigmoid at the last layer
- cross entropy loss

Learning Rate Sequence	Loss	Epochs	Average Score	Top Score
.4, .2, .1, 6*[.05]	0.091728	1000	0.970635	0.972972
.4, .2, 3*[.1], 4*[.05]	0.112709	900	0.964006	0.973684
.4, .2, .1	.0095086	900	0.949541	0.96
.4, .2, .1, 3*[.05]	0.060299	600	0.94444	0.9692308





Learning rates=.4, .2, .1



Learning rates=.4, .2, .1, 3*[.05]

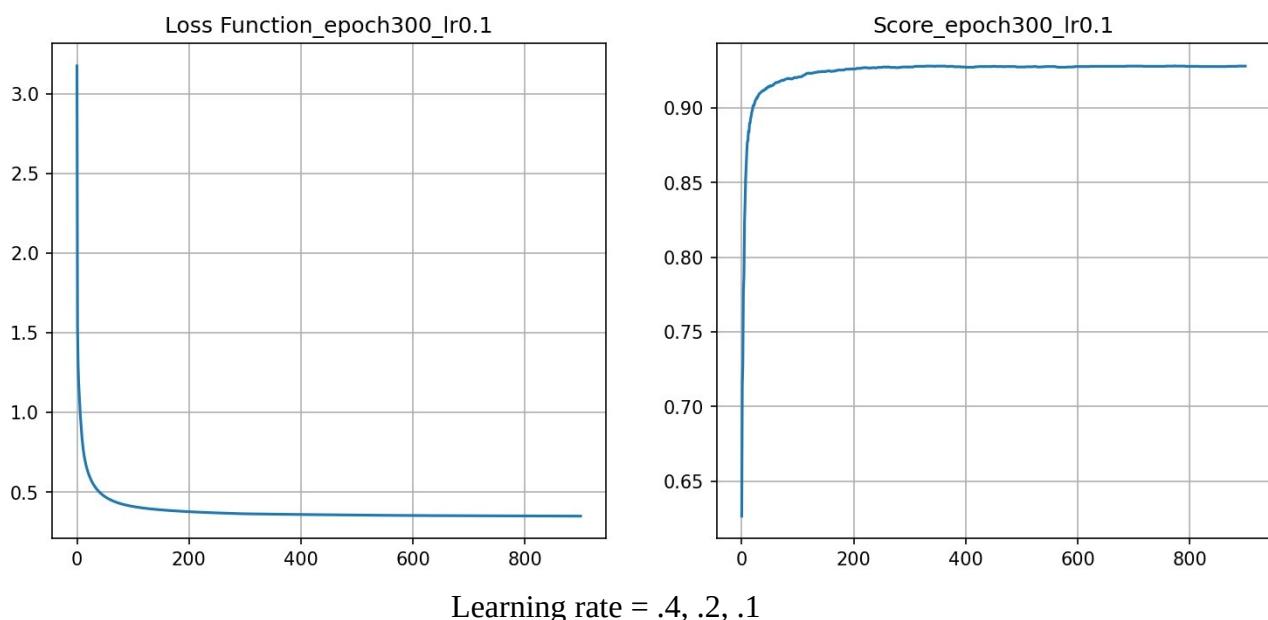
DryBean

The data-set has 16 attributes and 13611 sample in total with 7 classes: (Seker, Barbunya, Bombay, Cali, Dermosan, Horoz and Sira)

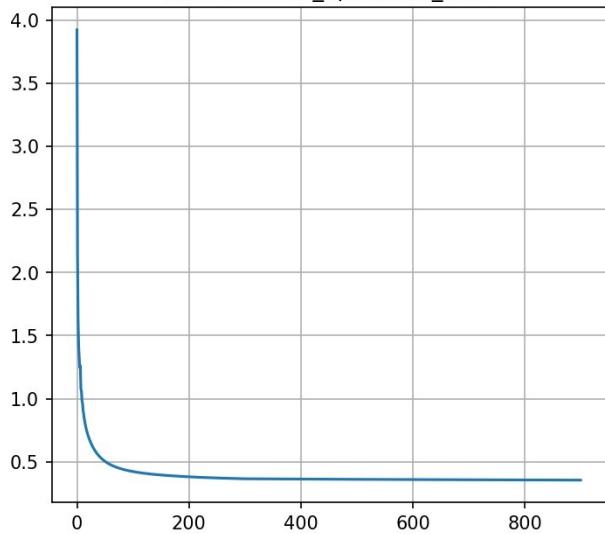
I used:

- a single hidden layer with 50 neurons
- LeakyRelu in the middle and sigmoid at the last layer
- cross entropy loss

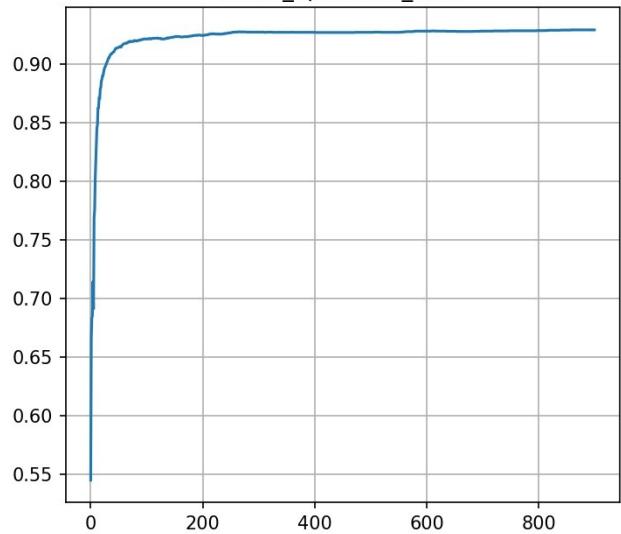
Learning Rate Sequence	Loss	Epochs	Average Score	Top Score
.4, .2, .1	0.348955	900	0.92415	0.92817
.4, .1, .1	0.3558576	900	0.92283	0.92923
.4, .3, .1, .01	0.347872	1200	0.92005	0.927168
.4, .2, .1, .05	0.334702	1200	0.92313	0.926969



Loss Function_epoch300_lr0.1

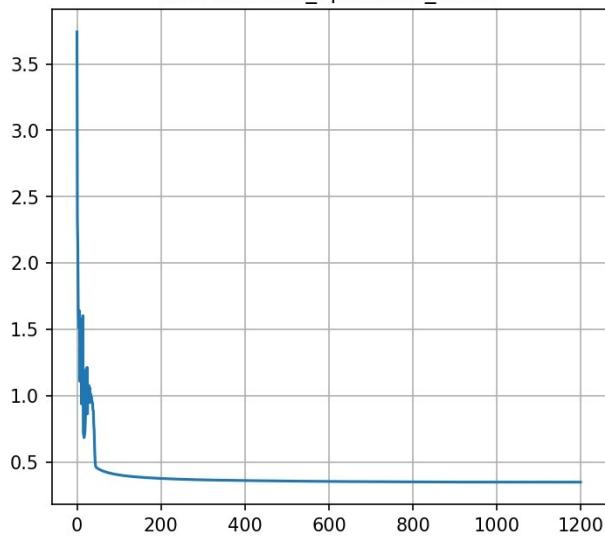


Score_epoch300_lr0.1

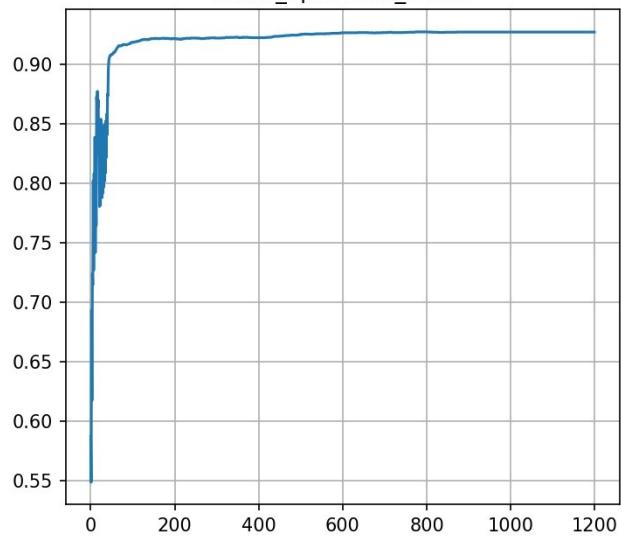


learning rate = .4, .1, .1

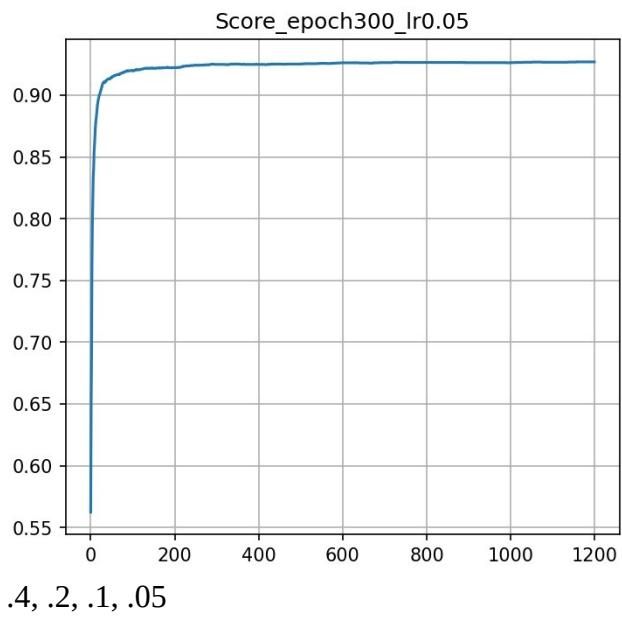
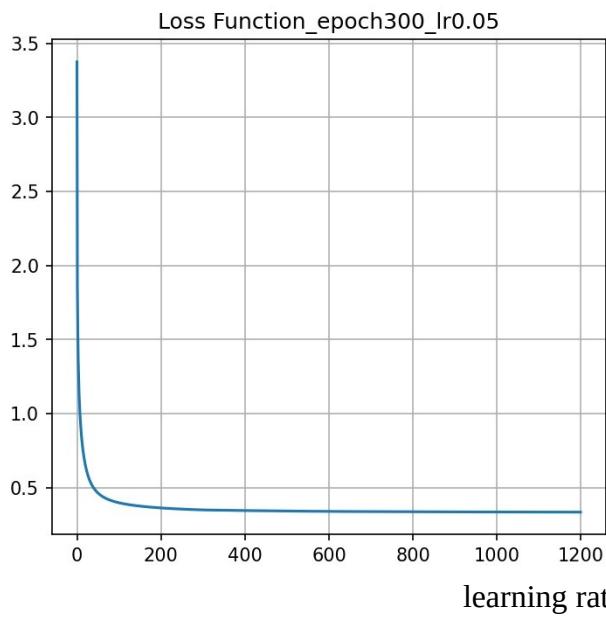
Loss Function_epoch300_lr0.01



Score_epoch300_lr0.01



learning rate = .4, .3, .1, .01



learning rate = .4, .2, .1, .05