



TU 9.1
weight

$$i) \quad w_1 = \begin{pmatrix} 1 \\ -0.5 \\ 0.1 \end{pmatrix} \quad w_2 = \begin{pmatrix} 0 \\ 2 \\ 0.6 \end{pmatrix} \quad w_3 = \begin{pmatrix} -0.5 \\ 0.6 \end{pmatrix}$$

a) biases

$$b_1 = 0 \quad b_2 = 0.5 \quad b_3 = 0.05 \quad (\text{see } \textcircled{+1})$$

$$w = \begin{pmatrix} -0.5 & 0 \\ 0 & 0.6 \end{pmatrix} \quad b = \begin{pmatrix} 0 \\ 0.5 \end{pmatrix}$$

$$\Rightarrow u_1 = 0.25 = x_1^T w_1 + b_1$$

$$g(u_1) y_1 = 0.531$$

$$u_2 = x_1^T w_2 + b_2 = 0.1$$

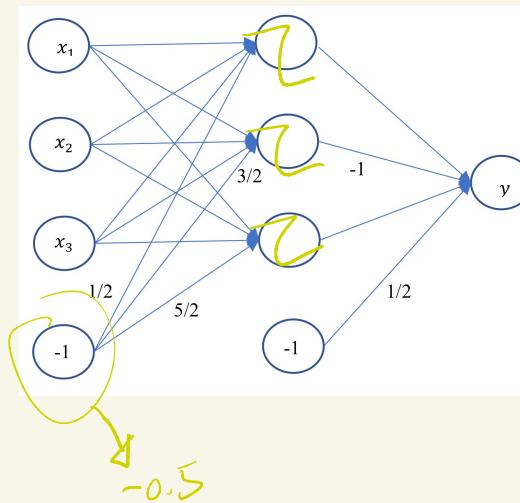
$$g(u_2) y_2 = 0.512$$

$$z = (y_1 \ y_2) \quad u_3 = z^T w_3 + b_3 = 0.092$$

$$y_3 = g(u_3) = \max(0, u_3) = 0.092$$

Q2) (-0.2, 0.2), yes

Q3) logic function



need to list out the truth table!!!

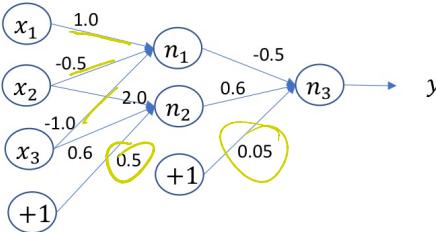


Figure 1

1. Figure 1 shows a two-layer feedforward neural network receiving 3-dimensional inputs $(x_1, x_2, x_3) \in \mathbb{R}^3$. The connection weights and biases of the neurons n_1, n_2 , and n_3 are as indicated in the figure. The hidden-layer neurons have activation functions given by $g(u) = \frac{1.0}{1+e^{-0.5u}}$ where u denotes the synaptic input to the neuron. The activation function $f(u)$ of the output neuron is a ReLU function: $f(u) = \max\{0, u\}$.

- Write weight vectors and biases connected to individual neurons, and the weight matrix and bias vector connected to the hidden layer.
- Find the synaptic inputs and activations of the neurons for the following input signals: $\underline{x} = [x_1 \ x_2 \ x_3]^T$

(i) $(1.0, -0.5, 1.0)$ (ii) $(-1.0, 0.0, -2.0)$ (iii) $(2.0, 0.5, -1.0)$.

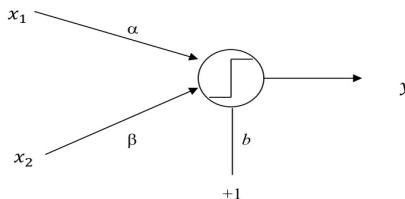


Figure 2

2. Two input binary neuron shown in figure 2 has a unit step activation function with a bias $b = 0.5$ and receives two-dimensional input $(x_1, x_2) \in \mathbb{R}^2$.

- Find the space of possible values of weights (α, β) if the neuron is
 - ON for input $(1.0, 1.0)$
 - ON for input $(0.5, -1.0)$
 - OFF for input $(2.0, -0.5)$.

- Indicate the weight space in 2-D α - β plot and show that $(-0.2, 0.2)$ is in this space.

3. The network shown in figure 3 consists of neurons having threshold activation functions and receives three-bit binary patterns $(x_1, x_2, x_3) \in \{0,1\}^3$. By analyzing the outputs for all possible three-bit input patterns, determine the logic function that the network implements. All unlabeled weights shown in figure 3 are of unity weight.

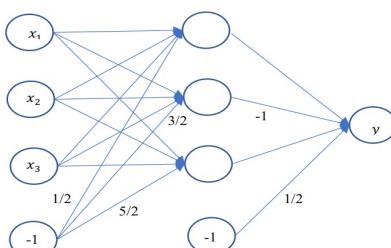


Figure 3

$$(2) g(u) \stackrel{?}{=} \frac{1}{1+e^{-0.5u}}$$

$$\text{a) } n_1 \\ w_1 = \begin{pmatrix} 1 \\ -0.5 \\ -1 \end{pmatrix} \\ b_1 = 0.5 \\ w_2 = \begin{pmatrix} 0 \\ 2 \\ 0.6 \end{pmatrix} \\ b_2 = 0.05 \\ w_3 = \begin{pmatrix} 1 & 0 \\ -1 & 2 \\ -1 & 0.6 \end{pmatrix} \\ b_3 = 0.05$$

$$\text{b) } u_1 = w_1^T x + b_1 \\ = (1 -0.5 -1) \begin{pmatrix} 1 \\ 0.5 \\ 0.6 \end{pmatrix} + 0.5 \\ = 1 + 0.25 - 1 + 0 = 0.25 \\ y_1 = 0.53$$

$$u_2 = w_2^T x + b_2 \\ = (0 \ 2 \ 0.6) \begin{pmatrix} 1 \\ 0.5 \\ 0.6 \end{pmatrix} + 0.05 \\ = 0 - 1 + 0.6 + 0.05 = 0.1 \\ y_2 = \frac{1}{1+e^{-0.5(0.1)}} = 0.512$$

$$u_3 = w_3^T x + b_3 \\ = \begin{pmatrix} 0.5 & 0.6 \\ 0.5 & 0.6 \\ 0.5 & 0.6 \end{pmatrix} \begin{pmatrix} 1 \\ 0.5 \\ 0.6 \end{pmatrix} + 0.05 \\ = 0.0417 + 0.05 = 0.0917$$

$\boxed{\text{DO IT}}$

1) $y = w^T x + b / w^T x = 0$

$g(u) = \frac{1}{1+e^{-0.5u}}$ Relu $f(u) = \max\{0, u\}$

$u > 0 \quad f(u) = u$
 $u < 0 \quad f(u) = 0$

a) weight: $\begin{pmatrix} 1 & -0.5 & -1 \\ 0 & 2 & 0.6 \end{pmatrix} \text{ on, bias } \begin{pmatrix} 0.5 \\ 0.05 \end{pmatrix}$

weight vector = $\begin{pmatrix} 1 & 0 \\ -0.5 & 2 \\ -1 & 0.6 \end{pmatrix}$ weight vector = $\begin{pmatrix} -0.5 \\ 0.6 \end{pmatrix}$

bias vector for $n_1 = 0$
 $n_2 = 0.5$
 $n_3 = 0.05$

for hidden layer !! for outer layer

weight vector $\begin{pmatrix} 1 & 0 \\ -0.5 & 2 \\ -1 & 0.6 \end{pmatrix}$ $\begin{pmatrix} -0.5 \\ 0.6 \end{pmatrix}$

bias vector $\begin{pmatrix} 0 \\ 0.5 \end{pmatrix}$ 0.05

b) Synaptic inputs (u) activation of neurons (y)

for input

i) $x = (-0.5)$ $\begin{pmatrix} u_1 \\ u_2 \end{pmatrix} = w^T x + b$

$\begin{pmatrix} u_1 \\ u_2 \end{pmatrix} = \begin{pmatrix} 1 & -0.5 & -1 \\ 0 & 2 & 0.6 \end{pmatrix} \begin{pmatrix} 1 \\ -0.5 \\ 1 \end{pmatrix} + \begin{pmatrix} 0 \\ 0.5 \end{pmatrix}$

$= \begin{pmatrix} -0.25 - 1 \\ 0 - 1 + 0.6 \end{pmatrix} + \begin{pmatrix} 0 \\ 0.5 \end{pmatrix} \quad g(u) = 0.531$

$\begin{pmatrix} u_1 \\ u_2 \end{pmatrix} = \begin{pmatrix} 0.25 \\ -0.4 \end{pmatrix} + \begin{pmatrix} 0 \\ 0.5 \end{pmatrix} = \begin{pmatrix} 0.25 \\ 0.1 \end{pmatrix} \quad g_1(u) = 0.512$

$u_3 = w_3^T x_3 + b$
 $= (-0.5 \ 0.6) \begin{pmatrix} 0.531 \\ 0.512 \end{pmatrix} + (0.05) < 0.0417 + 0.05 = 0.0917$

$u_3 > 0 \quad f(u_3) = u_3$

ii) $\begin{pmatrix} u_1 \\ u_2 \end{pmatrix} = \begin{pmatrix} 1 & -0.5 & -1 \\ 0 & 2 & 0.6 \end{pmatrix} \begin{pmatrix} -1 \\ 0 \\ -2 \end{pmatrix} + \begin{pmatrix} 0 \\ 0.5 \end{pmatrix}$

$= \begin{pmatrix} -1 + 0 + 2 \\ 0 + 0 - 1.2 \end{pmatrix} + \begin{pmatrix} 0 \\ 0.5 \end{pmatrix} = \begin{pmatrix} 1 \\ -0.7 \end{pmatrix}$

$\begin{pmatrix} g(u) \\ g(u_2) \end{pmatrix} = \begin{pmatrix} \frac{1}{1+e^{-0.5}} \\ \frac{1}{1+e^{-0.5-0.7}} \end{pmatrix} = \begin{pmatrix} 0.622 \\ 0.413 \end{pmatrix}$

$u_3 = (-0.5 \ 0.6) \begin{pmatrix} 0.622 \\ 0.413 \end{pmatrix} + 0.05 = -0.013$

$u_3 < 0 \quad g(u_3) = 0$

iii) $\begin{pmatrix} u_1 \\ u_2 \end{pmatrix} = \begin{pmatrix} 1 & -0.5 & -1 \\ 0 & 2 & 0.6 \end{pmatrix} \begin{pmatrix} 2 \\ 0.5 \\ -1 \end{pmatrix} + \begin{pmatrix} 0 \\ 0.5 \end{pmatrix}$

$= \begin{pmatrix} 2 - 0.25 + 1 \\ 0 + 1 - 0.6 \end{pmatrix} + \begin{pmatrix} 0 \\ 0.5 \end{pmatrix} = \begin{pmatrix} 2.75 \\ 0.9 \end{pmatrix} \quad -0.399 + 0.366 = 0.05$

$\begin{pmatrix} g(u) \\ g(u_2) \end{pmatrix} = \begin{pmatrix} 0.798 \\ 0.611 \end{pmatrix} \quad u_3 = (0.5 \ 0.6) \begin{pmatrix} 0.798 \\ 0.611 \end{pmatrix} + 0.05 = 0.0176$

$u_3 > 0 \quad f(u_3) = u_3 = 0.0176$

2) unit step activation

$$l(x) = \begin{cases} 1 & x \geq 0 \\ 0 & x < 0 \end{cases} \quad w = \begin{pmatrix} \alpha \\ \beta \end{pmatrix}$$

$$b = 0.5$$

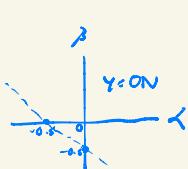
a) $\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$ neuron $y=1$ (assume $y=1 \equiv \text{ON}$) if $y=1$

$$u_1 = w^T x + b$$

$$= (\alpha \beta) \begin{pmatrix} 1 \\ 1 \end{pmatrix} + 0.5$$

$$y = \alpha + \beta + 0.5$$

$$\text{if } u_1 > 0 \Rightarrow \alpha + \beta + 0.5 > 0$$



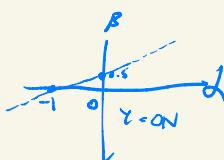
$$u_2 = (\alpha \beta) \begin{pmatrix} 0.5 \\ -1 \end{pmatrix} + 0.5$$

$$u_2 = 0.5\alpha - \beta + 0.5$$

$$\text{if } u_2 > 0 \Rightarrow 0.5\alpha - \beta + 0.5 > 0$$

$$\text{if } \alpha > 0 \Rightarrow \beta < 0.5$$

$$\text{if } \beta > 0 \Rightarrow \alpha < -1$$



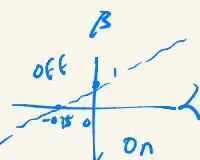
$$u_3 = (\alpha \beta) \begin{pmatrix} 2 \\ -0.5 \end{pmatrix} + 0.5$$

$$u_3 = 2\alpha - 0.5\beta + 0.5$$

$$\text{if } u_3 > 0 \Rightarrow 2\alpha - 0.5\beta + 0.5 > 0$$

$$\text{if } \alpha > 0 \Rightarrow \beta < 1$$

$$\text{if } \beta > 0 \Rightarrow \alpha < -0.25$$



b) $(-0.2, 0.2)$

from picture Yes

OR

use 2 and 3 line

$$u_3 = 2(-0.2) - 0.5(0.2) + 0.5 = -0.4 - 0.1 + 0.5 = 0.8 \rightarrow \text{on} \quad \checkmark$$

$$u_2 = 0.5(-0.2) - 0.2 + 0.5 = -0.1 - 0.2 + 0.5 = 0.2 \rightarrow \text{on}$$

3) weight vector = $\begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{pmatrix}$ hidden layer

output layer $\begin{pmatrix} 1 \\ -1 \end{pmatrix}$

$$\theta = \frac{1}{2}x_1 + \frac{1}{2}x_2 + \frac{1}{2}x_3$$

$$u = w^T x - \theta$$

$$\begin{pmatrix} u_1 \\ u_2 \\ u_3 \end{pmatrix} = \begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} - \begin{pmatrix} \frac{1}{2} \\ \frac{1}{2} \\ \frac{1}{2} \end{pmatrix} = \begin{array}{l} x_1 + x_2 + x_3 - \frac{1}{2} \\ x_1 + x_2 + x_3 - \frac{3}{2} \\ x_1 + x_2 + x_3 - \frac{5}{2} \end{array}$$

$$u = w^T x - \theta$$

$$= (1 \ 1 \ 1) \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} - \frac{1}{2} = u_1 - u_2 + u_3 - \frac{1}{2}$$

x_1	x_2	x_3	u_1	u_2	u_3	y_1	y_2	y_3	u_{final}	y_3
0	0	0	- $\frac{1}{2}$	$-\frac{3}{2}$	$-\frac{5}{2}$	0	0	0	-0.5	0
0	0	1	$\frac{1}{2}$	$-\frac{1}{2}$	$-\frac{1}{2}$	1	0	0	0.5	1
0	1	0	$\frac{1}{2}$	$-\frac{1}{2}$	$-\frac{1}{2}$	1	0	0	0.5	1
0	1	1	1.5	$\frac{1}{2}$	$-\frac{1}{2}$	1	1	0	-0.5	0
1	0	0	$\frac{1}{2}$	$-\frac{1}{2}$	$-\frac{3}{2}$	1	0	0	0.5	1
1	0	1	1.5	$\frac{1}{2}$	$-\frac{1}{2}$	1	1	0	-0.5	0
1	1	0	1.5	$\frac{1}{2}$	$-\frac{1}{2}$	1	1	0	-0.5	0
1	1	1	2.5	1.5	0.5	1	1	1	0.5	1

$$y_1 y_2 y_3 + y_1 y_2 y_3$$

$$y_1 \bar{x}_1 \bar{x}_2 x_3 + x_1 \bar{x}_2 \bar{x}_3 + x_1 \bar{x}_2 x_3 + x_1 x_2 x_3$$

$$\leq \bar{x}_1 (\bar{x}_2 x_3 + x_2 \bar{x}_3) + x_1 (\bar{x}_2 \bar{x}_3 + x_2 x_3)$$

$$= \bar{x}_1 (x_2 \oplus x_3) + x_1 (\bar{x}_2 \oplus x_3)$$

$$= \bar{x}_1 \oplus (x_2 \oplus x_3)$$

$$= x_1 \oplus x_2 \oplus x_3$$

TUT2

I) Linear neuron

$$P = 8$$

a) SGD

$$w_0 = \begin{pmatrix} 0.77 \\ 0.02 \\ 0.63 \end{pmatrix}, b = 0, \alpha = 0.01$$

epoch 1

shuffle the patterns

$$y_p = w^T x_p + b$$

$$w \leftarrow w + \alpha (d_p - y_p) x_p$$

$$b \leftarrow b + \alpha (d_p - y_p)$$

$$Se \leftarrow (d_p - y_p)^2$$

p = 1 choose 3

$$x^T = (0.34 \ 0.65 \ -0.73) \quad d = 0.96$$

$$y_p = (0.77 \ 0.02 \ 0.63) \begin{pmatrix} 0.34 \\ 0.65 \\ -0.73 \end{pmatrix} + 0 = -0.1851$$

$$w \leftarrow \begin{pmatrix} 0.77 \\ 0.02 \\ 0.63 \end{pmatrix} + 0.01 \begin{pmatrix} 0.96 - (-0.1851) \end{pmatrix} \begin{pmatrix} 0.34 \\ 0.65 \\ -0.73 \end{pmatrix} = \begin{pmatrix} 0.774 \\ 0.028 \\ 0.622 \end{pmatrix}$$

$$b \leftarrow 0 + 0.01(-0.1851) = 0.01(-0.1851)$$

p = 2 choose 4

$$x^T = (0.15 \ 0.78 \ -0.58) \quad d = 1.04$$

$$y_p = (0.774 \ 0.028 \ 0.622) \begin{pmatrix} 0.15 \\ 0.78 \\ -0.58 \end{pmatrix} + 0.01(-0.1851) = -0.21161$$

$$w \leftarrow \begin{pmatrix} 0.774 \\ 0.028 \\ 0.622 \end{pmatrix} + 0.01 \begin{pmatrix} 1.04 - (-0.21161) \end{pmatrix} \begin{pmatrix} 0.15 \\ 0.78 \\ -0.58 \end{pmatrix} = \begin{pmatrix} 0.778 \\ 0.037 \\ 0.614 \end{pmatrix}$$

$$b \leftarrow 0.01(-0.1851) + 0.01(-0.21161) = 0.024$$

p = 3 choose 7

$$x^T = (0.63 \ -0.65 \ -0.14) \quad d = -2.39$$

$$y_p = (0.776 \ 0.037 \ 0.614) \begin{pmatrix} 0.63 \\ -0.65 \\ -0.14 \end{pmatrix} + 0.024 = 0.40994$$

$$w \leftarrow \begin{pmatrix} 0.776 \\ 0.037 \\ 0.614 \end{pmatrix} + 0.01 \begin{pmatrix} -2.39 - (0.40994) \end{pmatrix} \begin{pmatrix} 0.63 \\ -0.65 \\ -0.14 \end{pmatrix} = \begin{pmatrix} 0.758 \\ 0.050 \\ 0.618 \end{pmatrix}$$

$$b \leftarrow 0.024 + (-0.027994671) = -0.004$$

p = 4 choose 8

$$x^T = (0.88 \ 0.64 \ -0.83) \quad d = 0.66$$

$$y_p = (0.758 \ 0.05 \ 0.618) \begin{pmatrix} 0.88 \\ 0.64 \\ -0.83 \end{pmatrix} + -0.004 = 0.49095$$

$$w \leftarrow \begin{pmatrix} 0.758 \\ 0.05 \\ 0.618 \end{pmatrix} + 0.01 \begin{pmatrix} 0.66 - (0.49095) \end{pmatrix} \begin{pmatrix} 0.88 \\ 0.64 \\ -0.83 \end{pmatrix} = \begin{pmatrix} 0.75 \\ 0.051 \\ 0.618 \end{pmatrix}$$

$$b \leftarrow -0.004 + = -0.00234$$

$x = (x_1, x_2, x_3)$	y
(0.09 -0.44 -0.15)	-2.57
(0.69 -0.99 -0.76)	-2.97
(0.34 0.65 -0.73)	0.96
(0.15 0.78 -0.34)	1.49
(-0.01 -0.78 -0.56)	0.21
(0.96 0.62 -0.66)	1.05
(0.63 -0.45 -0.14)	-2.39
(0.88 0.64 -0.33)	0.66

shuffle

 $p = 1, 2, 3, 4, 5, 6, 7, 8$

$$p = 1 \quad x_p^T = (0.09 \ -0.44 \ -0.15)$$

$$y_p = (0.77, 0.02, 0.63) \begin{pmatrix} 0.09 \\ -0.44 \\ -0.15 \end{pmatrix} + 0 = -0.034$$

$$w \leftarrow \begin{pmatrix} 0.77 \\ 0.02 \\ 0.63 \end{pmatrix} + 0.01 \begin{pmatrix} -0.034 \end{pmatrix} \begin{pmatrix} 0.09 \\ -0.44 \\ -0.15 \end{pmatrix} = \begin{pmatrix} 0.768 \\ 0.031 \\ 0.634 \end{pmatrix}$$

$$b \leftarrow 0 + -0.02536 = -0.0254$$

$$Se = 1.316$$

$$Se = 1.564$$

$$Se = 7.84$$

$$Se = 0.0286$$

$$P=5 \text{ choose } 1$$

$$x^T = (0.09 \ -0.44 \ -0.15) \quad d = -2.57$$

$$y_p = \begin{pmatrix} 0.76 \\ 0.051 \\ 0.618 \end{pmatrix} \begin{pmatrix} 0.09 \\ -0.44 \\ -0.15 \end{pmatrix} - 0.00234 = -0.099 \quad Se = 6.358$$

$$w \leftarrow \begin{pmatrix} 0.76 \\ 0.051 \\ 0.618 \end{pmatrix} + 0.01 (-2.57 - (-0.099)) \begin{pmatrix} 0.09 \\ -0.44 \\ -0.15 \end{pmatrix} = \begin{pmatrix} 0.757 \\ 0.062 \\ 0.622 \end{pmatrix}$$

$$b \leftarrow -0.00234 + -0.0252 = -0.0276$$

$$P=6 \text{ choose } 5$$

$$(-0.63 \ -0.78 \ -0.56) \quad d = -3.21$$

$$y_p = \begin{pmatrix} 0.757 \\ 0.062 \\ 0.622 \end{pmatrix} \begin{pmatrix} -0.63 \\ -0.78 \\ -0.56 \end{pmatrix} + b = -0.9012$$

$$w \leftarrow \begin{pmatrix} 0.757 \\ 0.062 \\ 0.622 \end{pmatrix} + 0.01 (-3.21 - (-0.9012)) \begin{pmatrix} -0.63 \\ -0.78 \\ -0.56 \end{pmatrix} = \begin{pmatrix} 0.772 \\ 0.08 \\ 0.6345 \end{pmatrix}$$

$$b \leftarrow -0.0276 + -0.02309 = -0.0506$$

$$Se \leftarrow 5.3311$$

$$P=7 \text{ choose } 6$$

$$x^T = 0.96 \ 0.62 \ -0.66 \quad d = 1.05$$

$$y_p = w^T x + b = \begin{pmatrix} 0.772 \\ 0.08 \\ 0.6345 \end{pmatrix} \begin{pmatrix} 0.96 \\ 0.62 \\ -0.66 \end{pmatrix} + (-0.0506) = 0.321$$

$$w \leftarrow \begin{pmatrix} 0.772 \\ 0.08 \\ 0.6345 \end{pmatrix} + 0.01 (1.05 - 0.321) \begin{pmatrix} 0.96 \\ 0.62 \\ -0.66 \end{pmatrix} = \begin{pmatrix} 0.779 \\ 0.085 \\ 0.63 \end{pmatrix}$$

$$b \leftarrow -0.0506 + 0.00729 = -0.0434$$

$$Se \leftarrow 0.5311$$

$$P=8 \quad \text{choose } 2$$

$$X = \begin{pmatrix} 0.69 \\ -0.19 \\ -0.26 \end{pmatrix} \quad d = -2.97 \quad \begin{pmatrix} 0.779 \\ 0.085 \\ 0.63 \\ -0.0434 \end{pmatrix}$$

$$y_p < w^T x + b = (0.779 \ 0.085 \ 0.63) \begin{pmatrix} 0.69 \\ -0.99 \\ -0.76 \end{pmatrix} + b = -0.068$$

$$w \leftarrow w + 0.01 (-2 \cdot q_t - y) x = \begin{pmatrix} 0.759 \\ 0.113 \\ 0.151 \end{pmatrix}$$

$$b \leftarrow b + -0.029$$

$$\text{Set } 0.421 = (d - y_p)^2$$

b) Gd

$$\begin{aligned} y &= u \leftarrow x^T w + b \quad \text{Eq} \\ w &\leftarrow w + \alpha x^T (d - y) \\ b &\leftarrow b + \alpha \quad \text{Eq} \quad (d - y) \end{aligned}$$

$$w = \begin{pmatrix} 0.77 \\ 0.92 \\ 0.83 \end{pmatrix} \quad b < 0 \quad \alpha < 0.0$$

$$X = \begin{pmatrix} 0.09 & -0.44 & -0.15 \\ 0.69 & -0.99 & -0.76 \\ 0.34 & 0.65 & -0.73 \\ 0.15 & 0.78 & -0.58 \\ -0.63 & -0.78 & -0.56 \\ 0.96 & 0.62 & -0.66 \\ 0.63 & -0.45 & -0.14 \\ 0.88 & 0.64 & -0.73 \end{pmatrix} \quad dC = \begin{pmatrix} -2.57 \\ -2.47 \\ 0.92 \\ 1.04 \\ -3.21 \\ 1.05 \\ -2.79 \\ 0.66 \end{pmatrix}$$

$$Y = U \cdot S \cdot V^T = \begin{pmatrix} 0.09 & -0.44 & -0.15 \\ 0.69 & -0.99 & -0.76 \\ 0.34 & 0.65 & -0.73 \\ 0.15 & 0.78 & -0.58 \\ -0.63 & -0.78 & -0.56 \\ 0.46 & 0.61 & -0.66 \\ 0.67 & -0.45 & -0.14 \\ 0.98 & 0.64 & -0.33 \end{pmatrix} \begin{pmatrix} 0.77 \\ 0.02 \\ 0.63 \end{pmatrix} + \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix} = \begin{pmatrix} -0.034 \\ 0.0327 \\ -0.185 \\ -0.234 \\ -0.884 \\ 0.336 \\ 0.388 \\ 0.483 \end{pmatrix}$$

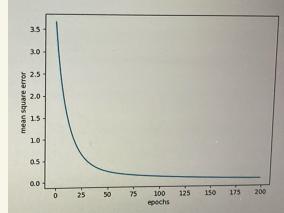
$$W = \begin{pmatrix} 0.77 \\ 0.02 \\ 0.63 \end{pmatrix} + 0.01 \begin{pmatrix} 0.09 & -0.44 & -0.15 \\ 0.69 & -0.91 & -0.76 \\ 0.34 & 0.65 & -0.73 \\ 0.15 & 0.78 & -0.58 \\ -0.63 & -0.78 & -0.56 \\ 0.96 & 0.61 & -0.66 \\ 0.63 & -0.45 & -0.14 \\ 0.98 & 0.11 & -0.27 \end{pmatrix}^T (d - y) = \begin{pmatrix} 0.759 \\ 0.115 \\ 0.653 \end{pmatrix}$$

$$b = 0 + 0.01 \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}^T \begin{pmatrix} d - y \end{pmatrix} = -0.073627$$

$$\text{Output } \mathbf{y} = \mathbf{Xw} + b\mathbf{1}_p = \begin{pmatrix} 0.09 & -0.44 & -0.15 \\ 0.69 & -0.99 & -0.76 \\ 0.34 & 0.65 & -0.73 \\ 0.15 & 0.78 & -0.58 \\ -0.63 & -0.78 & -0.56 \\ 0.96 & 0.62 & -0.66 \\ 0.63 & -0.45 & -0.14 \\ 0.88 & 0.64 & -0.33 \end{pmatrix}_{8 \times 3} \begin{pmatrix} 0.77 \\ 0.02 \\ 0.63 \end{pmatrix}_{3 \times 1} + 0.0 \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{pmatrix}_{8 \times 1} = \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{pmatrix}_{8 \times 1} \begin{pmatrix} -0.03 \\ -0.18 \\ -0.24 \\ -0.85 \\ 0.34 \\ 0.39 \\ 0.48 \\ 0.48 \end{pmatrix}_{8 \times 1} + \mathbf{d} = \begin{pmatrix} -2.57 \\ -2.97 \\ 0.96 \\ 1.04 \\ -3.21 \\ 1.05 \\ -2.39 \\ 0.66 \end{pmatrix}_{8 \times 1}$$

m.s.e. = $\frac{1}{8} \sum_{p=1}^8 (d_p - y_p)^2$
 $= \frac{1}{8} ((-2.57 + 0.03)^2 + (-2.97 - 0.03)^2 + \dots + (0.66 - 0.48)^2)$
 $= 4.02$

epoch	y	mse	w	b
2	$\begin{pmatrix} -0.15 \\ -0.16 \\ -0.22 \\ -0.25 \\ 0.11 \\ 0.29 \\ 0.26 \\ 0.45 \end{pmatrix}$	3.66	$\begin{pmatrix} 0.75 \\ 0.21 \\ 0.67 \end{pmatrix}$	-0.14
200	$\begin{pmatrix} -2.23 \\ -3.29 \\ 1.01 \\ -3.28 \\ 0.92 \\ -2.06 \\ 0.87 \end{pmatrix}$	0.05	$\begin{pmatrix} 0.368 \\ 2.56 \\ -0.21 \end{pmatrix}$	-1.164

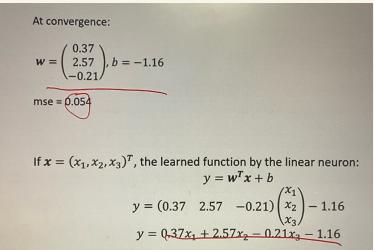


$$\mathbf{w} = \mathbf{w} + \alpha \mathbf{X}^T (\mathbf{d} - \mathbf{y})$$

$$= \begin{pmatrix} 0.77 \\ 0.02 \\ 0.63 \end{pmatrix} + 0.01 \times \begin{pmatrix} 0.09 & 0.69 & 0.34 & 0.15 & -0.63 & 0.96 & 0.63 & 0.88 \\ -0.44 & -0.99 & -0.65 & 0.78 & -0.78 & 0.62 & -0.45 & 0.64 \\ -0.15 & -0.76 & -0.73 & -0.58 & -0.56 & -0.66 & -0.14 & -0.33 \end{pmatrix}_{3 \times 8} \begin{pmatrix} -2.57 \\ -2.97 \\ 0.96 \\ 1.04 \\ -3.21 \\ 1.05 \\ -2.39 \\ 0.66 \end{pmatrix}_{8 \times 1} - \begin{pmatrix} -0.03 \\ -0.18 \\ -0.24 \\ -0.85 \\ 0.34 \\ 0.39 \\ 0.48 \\ 0.48 \end{pmatrix}_{8 \times 1}$$

$$= \begin{pmatrix} 0.76 \\ 0.11 \\ 0.65 \end{pmatrix}$$

$$b = b + \alpha \mathbf{1}_p^T (\mathbf{d} - \mathbf{y}) = 0.0 + 0.01 \times (1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1) \begin{pmatrix} -2.57 \\ -2.97 \\ 0.96 \\ 1.04 \\ -3.21 \\ 1.05 \\ -2.39 \\ 0.66 \end{pmatrix}_{8 \times 1} - \begin{pmatrix} -0.03 \\ -0.18 \\ -0.24 \\ -0.85 \\ 0.34 \\ 0.39 \\ 0.48 \\ 0.48 \end{pmatrix}_{8 \times 1} = -0.07$$



After learning ...				
x	d	y (SGD)	y (GD)	w
(0.09, -0.44, -0.15)	-2.57	-2.23	-2.23	$\begin{pmatrix} 0.369 \\ 2.566 \\ -0.212 \end{pmatrix}$
(0.69, -0.99, -0.76)	-2.97	-3.29	-3.29	
(0.34, 0.65, -0.73)	0.96	0.78	0.78	
(0.15, 0.78, -0.58)	1.04	1.01	1.01	
(-0.63, -0.78, -0.56)	-3.21	-3.28	-3.28	
(0.96, 0.62, -0.66)	1.05	0.92	0.92	
(0.63, -0.45, -0.14)	-2.39	-2.06	-2.06	
(0.88, 0.64, -0.33)	0.66	0.87	0.88	
m.s.e.	0.055	0.054		
w			$\begin{pmatrix} 0.368 \\ 2.567 \\ -0.207 \end{pmatrix}$	
b	-1.165		-1.163	
y	$0.37x_1 + 2.57x_2 - 0.21x_3 - 1.1$		$0.37x_1 + 2.57x_2 - 0.21x_3 - 1.16$	

1b) $\omega = \begin{pmatrix} 0.36871063 \\ 2.56532148 \\ -0.21284826 \end{pmatrix}$

$b = -1.166427293040246^2$

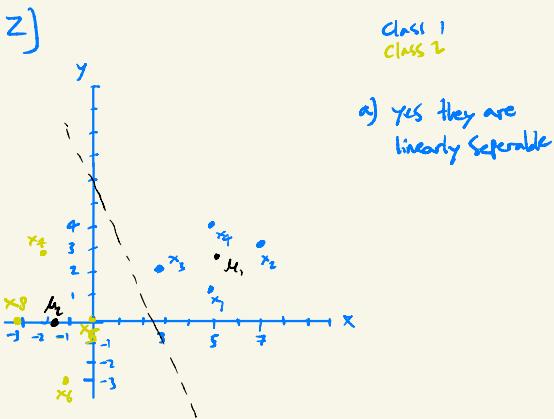
$\text{mse} = 0.0845477$

\downarrow

$y = 0.37x_1 + 2.57x_2 - 0.21x_3 - 1.17$

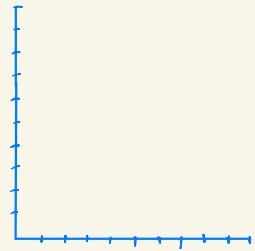
1c) $y_p = \omega^T x + b$

	y
1	-2.23
2	-3.29
3	0.782
4	1.013
5	-3.281
6	0.919
7	-2.059
8	0.87



Class 1
Class 2

a) yes they are
linearly separable



b) $6.5x_1 + 2.5x_2 - 14.5 = 0$ decision boundary?

take Class 1 = 0
Class 2 = 1

converse $w = \begin{pmatrix} 6.5 \\ 2.5 \end{pmatrix}$ $b = -14.5$

$$x = \begin{pmatrix} x_1 & x_2 \end{pmatrix} \quad d = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$$

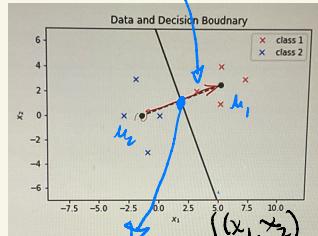
Steps:

① find 2 Centroids

$$\text{Center of Class 1 } (\mu_1) = \frac{1}{4}(x_1 + x_2 + x_3 + x_4) = \frac{1}{4}(20) = (5, 2.5)$$

$$\text{Center of Class 2 } (\mu_2) = \frac{1}{4}(x_5 + x_6 + x_7 + x_8) = \frac{1}{4}(-16) = (-4, -4)$$

② find line joining 2 Centroids



$$\mu_1 - \mu_2 = \begin{pmatrix} 5 \\ 2.5 \end{pmatrix} - \begin{pmatrix} -4 \\ -4 \end{pmatrix} = \begin{pmatrix} 9 \\ 6.5 \end{pmatrix}$$

③ middle point connecting 2 Centroid = $\frac{1}{2}(\mu_1 + \mu_2)$

$$= \frac{1}{2} \left(\begin{pmatrix} 5 \\ 2.5 \end{pmatrix} + \begin{pmatrix} -4 \\ -4 \end{pmatrix} \right) = \begin{pmatrix} 0.75 \\ 1.25 \end{pmatrix}$$

④ line passing perpendicularly through middle point

$$\begin{pmatrix} x_1 - 0.75 \\ x_2 - 1.25 \end{pmatrix}$$

\rightarrow inner product = 0

$$\begin{pmatrix} x_1 - 0.75 \\ x_2 - 1.25 \end{pmatrix} \cdot \begin{pmatrix} 6.5 \\ 2.5 \end{pmatrix} = 0$$

$$\begin{pmatrix} x_1 - 0.75 \\ x_2 - 1.25 \end{pmatrix} \cdot \begin{pmatrix} 6.5 \\ 2.5 \end{pmatrix} = 0$$

$$6.5x_1 - 11.25 + 2.5x_2 - 3.125 = 0$$

$$6.5x_1 + 2.5x_2 - 14.5 = 0$$

$$2c) \text{ use } u = 6.5x_1 + 2.5x_2 - 14.5 \quad \begin{array}{ll} \text{if } u > 0 & \text{class 1} \\ u \leq 0 & \text{class 2} \end{array}$$

2d) \hookrightarrow use $x_1 \dots x_8$ to check (training)

test the 3 input \rightarrow result class 1, 2, 1

2e)

$$\begin{array}{ll} x_1 & 20.5 \\ x_2 & 38.5 \\ x_3 & 10 \\ x_4 & 28 \\ x_5 & -14.5 \\ x_6 & -28.5 \\ x_7 & -20 \\ x_8 & -34 \end{array} \quad \left. \begin{array}{l} \\ \\ \\ \\ \\ \\ \\ \end{array} \right\} \text{Class}_1 \quad \begin{array}{ll} & \text{if } u > 0 \rightarrow \text{class 1} \\ & u \leq 0 \rightarrow \text{class 2} \end{array}$$

$$2d) \begin{pmatrix} 4 \\ 2 \end{pmatrix} \quad u = 16.5 \rightarrow \text{class 1}$$

$$\begin{pmatrix} 0 \\ 5 \end{pmatrix} \quad u = -2 \rightarrow \text{class 2}$$

$$\begin{pmatrix} 3 \\ 1 \\ 3 \\ 0 \end{pmatrix} \quad u = 3.5 \rightarrow \text{class 1}$$

$$\text{Q1) } \text{GD } \alpha = 0.05 \quad w = \begin{pmatrix} 0.88 & 0.08 & -0.34 \\ 0.68 & -0.39 & -0.19 \end{pmatrix} \quad b = 0$$

$$Y = \arg\max_u f(u) \quad f(u) = \frac{e^u}{\sum_{k=1}^2 e^{u_k}}$$

$$w = w - \alpha x^T (-k - f(u))$$

$$b = b - \alpha (f(k) - f(u))^T 1_p$$

1st epoch

$$k = 2$$

Neuron = 9

take class A = 0

class B = 1

class C = 2

$$0 \quad 1 \quad 2$$

$$x = \begin{pmatrix} 0 & 4 \\ -1 & 3 \\ 2 & 3 \\ -2 & 2 \\ 0 & 2 \\ 1 & 2 \\ -1 & 2 \\ -3 & 1 \\ -1 & 1 \end{pmatrix} \quad d = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 1 \\ 0 \\ -1 \\ -1 \\ 1 \end{pmatrix}$$

$$k = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix}$$

$$x = \begin{pmatrix} 2 & 1 \\ 4 & 1 \\ -2 & 0 \\ 1 & 0 \\ 7 & 0 \\ -3 & 1 \\ -2 & -1 \\ 2 & -1 \\ 4 & -1 \end{pmatrix} \quad d = \begin{pmatrix} 2 \\ 2 \\ -1 \\ 2 \\ 2 \\ 1 \\ 1 \\ 2 \\ 2 \end{pmatrix}$$

$$N \times 2 \quad 2 \times 3$$

$$Y = Xw + B = \begin{pmatrix} 0 & 4 \\ 2 & 3 \\ -2 & 2 \\ 0 & 2 \\ 1 & 2 \\ -1 & 2 \\ -3 & 1 \\ -1 & 1 \end{pmatrix} \begin{pmatrix} 0.88 & 0.08 & -0.34 \\ 0.68 & -0.39 & -0.19 \end{pmatrix} + 0$$

$$\text{result} = \begin{pmatrix} 2.72 & -1.56 & -0.76 \\ 1.16 & -1.25 & -0.23 \end{pmatrix}$$

$$\text{diff} = \begin{pmatrix} 2.84 & 0.71 & -1.17 \\ 18.4 & 14 \end{pmatrix}$$

error

$$f(u) = \frac{e^u}{\sum_{k=1}^2 e^{u_k}} = \begin{pmatrix} 0.957 & 0.013 & 0.0215 \end{pmatrix}$$

$$Y = \begin{pmatrix} \bullet \\ \bullet \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 2 \\ 2 \\ 1 \\ 2 \\ 2 \end{pmatrix}$$

cal $\nabla_w J$, w , and b

b , C , d

$$2) \quad x_1 = \begin{pmatrix} 1 \\ 3 \end{pmatrix}, d_1 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}, x_2 = \begin{pmatrix} -2 \\ -2 \end{pmatrix}, d_2 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$

GD of 2 patterns
single feedforward and feedback

a) matrix of
W hidden layer weights
V output layer biases

$$W = \begin{pmatrix} 1 & 2 \\ -2 & 0 \end{pmatrix}, b_h = \begin{pmatrix} 3 & -1 \end{pmatrix}$$

$$V = \begin{pmatrix} 1 & 1 \\ 0 & -1 \end{pmatrix}, b_o = \begin{pmatrix} -2 & 3 \end{pmatrix}$$

$$d = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$$

b) $z = h$ of hidden layer
 $y = v$ of output layer

$$z = xw + b = \begin{pmatrix} 1 & 3 \\ -2 & -2 \end{pmatrix} \begin{pmatrix} 1 & 2 \\ -2 & 0 \end{pmatrix} + \begin{pmatrix} 3 & -1 \end{pmatrix} = \begin{pmatrix} -5 & 2 \\ 2 & -4 \end{pmatrix} + \begin{pmatrix} 3 & -1 \end{pmatrix} = \begin{pmatrix} -2 & 1 \\ 5 & -5 \end{pmatrix}$$

$$h = g(z) = \frac{1}{1+e^{-z}} = \begin{pmatrix} 0.12 & 0.73 \\ 0.99 & 0.01 \end{pmatrix}$$

$$y = Hv + c = \begin{pmatrix} 0.12 & 0.73 \\ 0.99 & 0.01 \end{pmatrix} \begin{pmatrix} 1 & 1 \\ 0 & -2 \end{pmatrix} + \begin{pmatrix} -2 & 3 \end{pmatrix} = \begin{pmatrix} -1.88 & 1.66 \\ -1.01 & 3.97 \end{pmatrix}$$

$$Y = f(v) = \frac{1}{1+e^{-v}} = \begin{pmatrix} 0.13 & 0.84 \\ 0.267 & 0.98 \end{pmatrix}$$

c) MSE = $\frac{1}{2} \sum_{k=1}^K (d_k - Y_k)^2 = \frac{1}{2} \left[(0 - 0.13)^2 + (1 - 0.84)^2 + (1 - 0.267)^2 + (0 - 0.98)^2 \right] = 0.77$

d) $f'(u) = Y(1-Y) = \begin{pmatrix} 0.11 & 0.13 \\ 0.2 & 0.02 \end{pmatrix}$

! $\nabla_u J = -(D-Y)^T = -\left(\begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} - \begin{pmatrix} 0.13 & 0.84 \\ 0.267 & 0.98 \end{pmatrix} \right) = \begin{pmatrix} 0.13 & -0.16 \\ -0.23 & 0.99 \end{pmatrix}, f'(w) = \begin{pmatrix} 0.02 & -0.02 \\ -0.14 & 0.02 \end{pmatrix}$

$$g'(z) = h(1-h) = \begin{pmatrix} 0.1 & 0.2 \\ 0.01 & 0.01 \end{pmatrix}$$

$$\nabla_z J = (\nabla_u J) V^T \cdot g'(z) = \nabla_u J \begin{pmatrix} 1 & 0 \\ 1 & -2 \end{pmatrix} \cdot \begin{pmatrix} 0.1 & 0.2 \\ 0.01 & 0.01 \end{pmatrix} = \begin{pmatrix} 0 & 0.04 \\ 0.02 & -0.01 \end{pmatrix} \begin{pmatrix} 0.1 & 0.2 \\ 0.01 & 0.01 \end{pmatrix} = \begin{pmatrix} -0.001 & 0.01 \\ -0.001 & 0 \end{pmatrix}$$

e) $V \leftarrow V - \alpha H^T \nabla_u J$

$$C \leftarrow C - \alpha (\nabla_u J)^T I_p$$

$$W \leftarrow W - \alpha X^T \nabla_z J$$

$$b \leftarrow b - \alpha (\nabla_z J)^T I_p$$

2) $k \leq 1$

$$Z = \begin{pmatrix} 1 & 1 \\ 0 & 1 \\ 1 & 4 \\ 2 & 2 \\ 2 & -2 \\ -2 & -3 \end{pmatrix} \begin{pmatrix} -0.1 & 0.97 & 0.18 \\ -0.7 & 0.38 & 0.93 \end{pmatrix} + B$$

$$K = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \end{pmatrix}$$

$$H = g(Z)$$

$$U = HV + C$$

$$f(u) = \frac{e^u}{\sum e^u}$$

$$y = \arg \max_k f(u)$$

Perceptron Softmax

do the rest! ↴

$$\text{relu} = \max\{0, u\}$$

3)

0	0	0
0	0	linear neuron = u
0	0	output
hidden		

TUT 6

D) a)

$$I = \begin{pmatrix} 0.7 & 0.1 & 0.2 & 0.3 & 0.3 & 0.5 \\ 0.8 & 0.1 & 0.3 & 0.5 & 0.1 & 0.0 \\ 1.0 & 0.2 & 0.0 & 0.3 & 0.2 & 0.7 \\ 0.8 & 0.1 & 0.5 & 0.6 & 0.3 & 0.4 \\ 0.1 & 0.0 & 0.9 & 0.3 & 0.3 & 0.2 \\ 1.0 & 0.1 & 0.4 & 0.5 & 0.2 & 0.8 \end{pmatrix}$$

padding = 0 strides = (1,1)

$$u_1 = \text{Conv}(I, w_1) + b_1$$

$$u_1(1,1) = I_{1,1} \cdot 0.7 + 0.1 \cdot 0.2 + 0.3 \cdot 0.3 + 0.3 \cdot 0.5 + b_1 = 9$$

$$I_{1,1} = \begin{pmatrix} 0.7 & 0.1 & 0.2 \\ 0.8 & 0.1 & 0.3 \\ 1.0 & 0.2 & 0.0 \end{pmatrix}$$

$$u_2(1,1) = I_{1,1} \cdot 0.1 + 0.2 \cdot 0.3 + 0.3 \cdot 0.5 + 0.3 \cdot 0.0 + b_1 = 6$$

$$I_{1,1} = \begin{pmatrix} 0.1 & 0.2 & 0.3 \\ 0.1 & 0.3 & 0.5 \\ 0.2 & 0.0 & 0.3 \end{pmatrix}$$

$$u_2(1,2) = I_{1,2} \cdot 0.8 + 0.1 \cdot 0.3 + 0.0 \cdot 0.5 + 0.0 \cdot 0.0 + b_1 = 4$$

$$I_{1,2} = \begin{pmatrix} 0.8 & 0.1 & 0.3 \\ 1.0 & 0.2 & 0.0 \\ 0.8 & 0.1 & 0.5 \end{pmatrix}$$

$$u_2(2,1) = I_{2,1} \cdot 0.8 + 0.1 \cdot 0.3 + 0.0 \cdot 0.5 + 0.0 \cdot 0.0 + b_1 = 2$$

$$u_2(2,2) = I_{2,2} \cdot 0.8 + 0.1 \cdot 0.3 + 0.0 \cdot 0.5 + 0.0 \cdot 0.0 + b_1 = 0$$

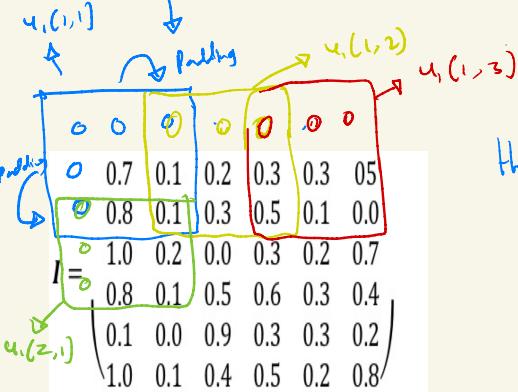
do for u_2 too $u_2 = \text{Conv}(I, w_2) + b_2$

$$Y_1 = f(u_1) = \frac{1}{1+e^{-u_1}}$$

$$Y_2 = f(u_2) = \frac{1}{1+e^{-u_2}}$$

Output spatial size
 $= (\text{Input size} - \text{Filter size} + 2(\text{Padding})) / \text{Stride} + 1$
 $= (6 - 3) / 1 + 1 = 4$

for padding = 1 strides = (2,2)



The rest same

b) Strides (2,2)

max pooling

Feature maps at the convolutional layer:

$$y = \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} = \begin{pmatrix} 0.94 & 0.8 & 0.8 & 0.87 \\ 0.92 & 0.88 & 0.88 & 0.89 \\ 0.85 & 0.88 & 0.93 & 0.93 \\ 0.94 & 0.88 & 0.96 & 0.88 \\ 0.57 & 0.5 & 0.45 & 0.55 \\ 0.57 & 0.60 & 0.65 & 0.69 \\ 0.48 & 0.69 & 0.75 & 0.43 \\ 0.55 & 0.48 & 0.45 & 0.57 \end{pmatrix}$$

Pooling 2x2 and strides = 2

Max-pooling:

$$o = \begin{pmatrix} 0.94 & 0.89 \\ 0.94 & 0.96 \\ 0.60 & 0.69 \\ 0.69 & 0.75 \end{pmatrix}$$

avg/max

avg
max

Mean-pooling:

$$p_{ave} = \begin{pmatrix} 0.88 & 0.86 \\ 0.89 & 0.92 \\ 0.56 & 0.58 \\ 0.55 & 0.55 \end{pmatrix}$$

Tut 1

TUT 2

TUT 3

TUT 4

$$\begin{aligned}
 & \text{Left side: } \frac{1}{2} \left(\frac{1}{2} \right)^{\frac{1}{2}} = \frac{1}{2} \cdot \frac{1}{\sqrt{2}} = \frac{1}{2\sqrt{2}} = \frac{1}{2\sqrt{2}} \cdot \frac{\sqrt{2}}{\sqrt{2}} = \frac{\sqrt{2}}{4} \\
 & \text{Right side: } \frac{1}{2} \left(\frac{1}{2} \right)^{\frac{1}{2}} + \frac{1}{2} \left(\frac{1}{2} \right)^{\frac{1}{2}} = \frac{1}{2} \left(\frac{1}{2} \right)^{\frac{1}{2}} + \frac{1}{2} \left(\frac{1}{2} \right)^{\frac{1}{2}} = \frac{\sqrt{2}}{4} + \frac{\sqrt{2}}{4} = \frac{2\sqrt{2}}{4} = \frac{\sqrt{2}}{2}
 \end{aligned}$$

108

13

4018

10

1) 4 hidden neurons $\alpha < 0$
 $w, b, f(u)$

$$h = f(Wx + b) \\ = f$$

--

$$\begin{array}{r} 2 \times 5 \\ + 2 \times 2 \\ \hline 0.0005 \end{array} \quad \begin{array}{r} 4 \times 2 = 8 \\ + 4 \times 3 \\ \hline 3 \times 2 = 6 \\ + 3 \times 4 \\ \hline \end{array}$$

1996-1997
年
度
公
司
經
營
成
績
報
告

$$\begin{aligned} & \text{Simplifying the left side:} \\ & \frac{1}{2}x^2 - 2x + 1 = x^2 - 4x + 4 \\ & \text{Subtracting } x^2 - 4x + 4 \text{ from both sides:} \\ & \frac{1}{2}x^2 - 2x + 1 - (x^2 - 4x + 4) = 0 \\ & \frac{1}{2}x^2 - 2x + 1 - x^2 + 4x - 4 = 0 \\ & \frac{1}{2}x^2 - x^2 + 4x - 2x + 1 - 4 = 0 \\ & -\frac{1}{2}x^2 + 2x - 3 = 0 \\ & \text{Multiplying by } -2:} \\ & x^2 - 4x + 6 = 0 \end{aligned}$$

$$B_4 = \begin{pmatrix} -1.91 & \\ -1.61 & \\ & 0.76 \\ & 1.247 \end{pmatrix}$$

$$\begin{aligned} & \left| \begin{array}{cc} 1 & -1 \\ 1 & 1 \end{array} \right| = 1 + 1 = 2 \\ & \text{adj} A = \frac{1}{2} \begin{pmatrix} 1 & -1 \\ 1 & 1 \end{pmatrix} = \begin{pmatrix} \frac{1}{2} & -\frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} \end{pmatrix} \\ & (A^{-1})^{-1} = A \\ & A = \begin{pmatrix} 1 & -1 \\ 1 & 1 \end{pmatrix} \\ & \text{adj } A^{-1} = \frac{1}{2} \begin{pmatrix} 1 & -1 \\ 1 & 1 \end{pmatrix} \\ & \text{adj } A^{-1} = \begin{pmatrix} \frac{1}{2} & -\frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} \end{pmatrix} \end{aligned}$$

$$\begin{aligned}
 & \text{Locus processus } y = 0 \text{ auf } x_1^2 + x_2^2 = 1 \\
 & H_1: x_1 = 0 \quad y = 0 \quad x_2 = \sqrt{1 - x_1^2} = \sqrt{1 - 0^2} = \sqrt{1} = 1 \\
 & H_2: x_1 = 0 \quad y = 0 \quad x_2 = \sqrt{1 - x_1^2} = \sqrt{1 - 0^2} = \sqrt{1} = 1 \\
 & \rightarrow \text{Locus } H_1 \cap H_2 = \{(x_1, x_2) \mid x_1 = 0, x_2 = 1\} = \{(0, 1)\} \quad \checkmark \\
 & \text{Locus } H_1 \cup H_2 = \{(x_1, x_2) \mid x_1 = 0 \text{ oder } x_2 = 1\} = \{(0, 0), (0, 1), (1, 0), (1, 1)\} \quad \checkmark \\
 & \text{Locus } H_1 \cap H_2 = \{(x_1, x_2) \mid x_1 = 0, x_2 = 1\} = \{(0, 1)\} \quad \checkmark \\
 & \text{Locus } H_1 \cup H_2 = \{(x_1, x_2) \mid x_1 = 0 \text{ oder } x_2 = 1\} = \{(0, 0), (0, 1), (1, 0), (1, 1)\} \quad \checkmark \\
 & H_3: x_1 = 0 \quad y = 0 \quad x_2 = \sqrt{1 - x_1^2} = \sqrt{1 - 0^2} = \sqrt{1} = 1 \\
 & \rightarrow \text{Locus } H_1 \cap H_3 = \{(x_1, x_2) \mid x_1 = 0, x_2 = 1\} = \{(0, 1)\} \quad \checkmark \\
 & \text{Locus } H_1 \cup H_3 = \{(x_1, x_2) \mid x_1 = 0 \text{ oder } x_2 = 1\} = \{(0, 0), (0, 1), (1, 0), (1, 1)\} \quad \checkmark \\
 & H_4: x_1 = 0 \quad y = 0 \quad x_2 = \sqrt{1 - x_1^2} = \sqrt{1 - 0^2} = \sqrt{1} = 1 \\
 & \rightarrow \text{Locus } H_1 \cap H_4 = \{(x_1, x_2) \mid x_1 = 0, x_2 = 1\} = \{(0, 1)\} \quad \checkmark
 \end{aligned}$$

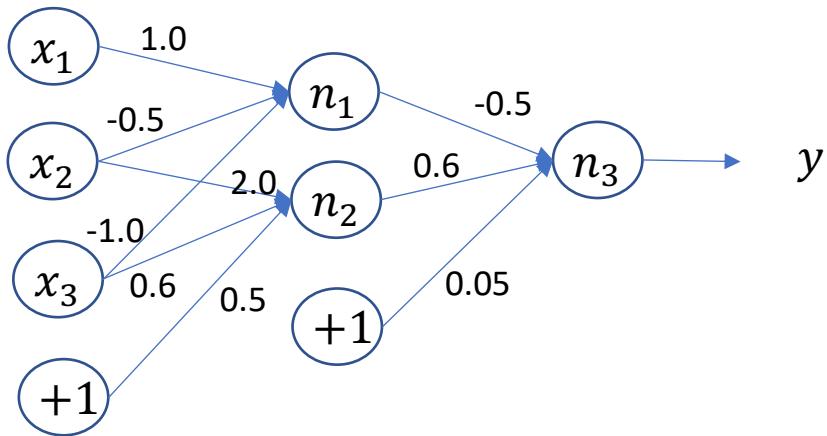


Figure 1

1. Figure 1 shows a two-layer feedforward neural network receiving 3-dimensional inputs $(x_1, x_2, x_3) \in \mathbf{R}^3$. The connection weights and biases of the neurons n_1, n_2 , and n_3 are as indicated in the figure. The hidden-layer neurons have activation functions given by $g(u) = \frac{1.0}{1+e^{-0.5u}}$ where u denotes the synaptic input to the neuron. The activation function $f(u)$ of the output neuron is a ReLU function: $f(u) = \max\{0, u\}$.
 - a. Write weight vectors and biases connected to individual neurons, and the weight matrix and bias vector connected to the hidden layer.
 - b. Find the synaptic inputs and activations of the neurons for the following input signals: u y
 - (i) (1.0, -0.5, 1.0) (ii) (-1.0, 0.0, -2.0) (iii) (2.0, 0.5, -1.0).

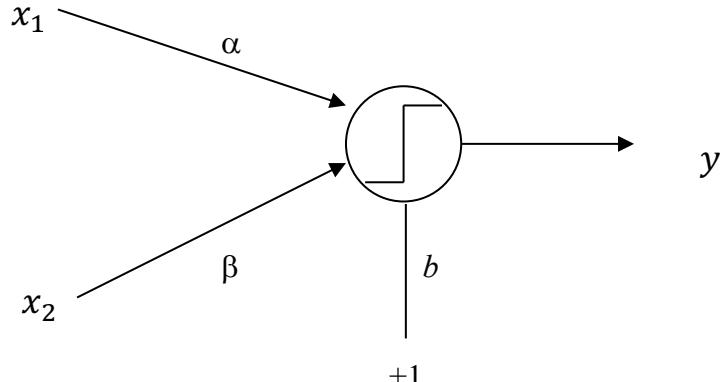


Figure 2

2. Two input binary neuron shown in figure 2 has a unit step activation function with a bias $b = 0.5$ and receives two-dimensional input $(x_1, x_2) \in \mathbf{R}^2$.

(a) Find the space of possible values of weights (α, β) if the neuron is

- (i) ON for input $(1.0, 1.0)$
- (ii) ON for input $(0.5, -1.0)$
- (iii) OFF for input $(2.0, -0.5)$.

(b) Indicate the weight space in 2-D α - β plot and show that $(-0.2, 0.2)$ is in this space.

3. The network shown in figure 3 consists of neurons having threshold activation functions and receives three-bit binary patterns $(x_1, x_2, x_3) \in \{0,1\}^3$. By analyzing the outputs for all possible three-bit input patterns, determine the logic function that the network implements. All unlabeled weights shown in figure 3 are of unity weight.

$Y \text{ or } \bar{Y}$

only take
1/0

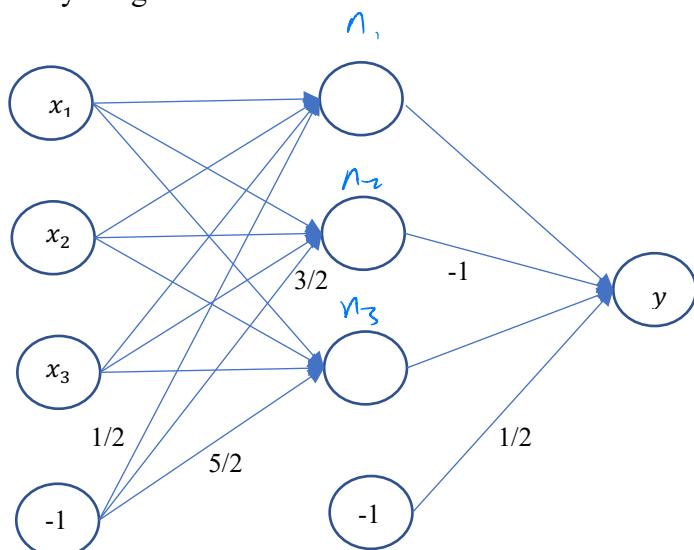


Figure 3

1. Train a linear neuron to learn the mapping from input $x \in \mathbb{R}^3$ to output y from the following examples:

Jpter 34781562

	$x = (x_1, x_2, x_3)$	y
1	(0.09 -0.44 -0.15)	-2.57
2	(0.69 -0.99 -0.76)	-2.97
3	(0.34 0.65 -0.73)	0.96
4	(0.15 0.78 -0.58)	1.04
5	(-0.63 -0.78 -0.56)	-3.21
6	(0.96 0.62 -0.66)	1.05
7	(0.63 -0.45 -0.14)	-2.39
8	(0.88 0.64 -0.33)	0.66

- (a) Show one iteration of learning of the neuron with
- i. Stochastic gradient descent learning
 - ii. Gradient descent learning

Initialize the weights as $\begin{pmatrix} 0.77 \\ 0.02 \\ 0.63 \end{pmatrix}$ and biases to 0.0, and use a leaning factor $\alpha = 0.01$.

- (b) Plot the learning curves (mean square error vs. epochs) until convergence. Determine the learned weights and biases.
- (c) Find the predicted values y of training inputs after the training.

2. Two-dimensional training patterns (inputs) to design a dichotomizer are given as:

$$\mathbf{X}_1 = \begin{bmatrix} 5 \\ 1 \end{bmatrix}; \quad \mathbf{X}_2 = \begin{bmatrix} 7 \\ 3 \end{bmatrix}; \quad \mathbf{X}_3 = \begin{bmatrix} 3 \\ 2 \end{bmatrix}; \quad \mathbf{X}_4 = \begin{bmatrix} 5 \\ 4 \end{bmatrix}; \quad \text{Class 1}$$

$$\mathbf{X}_5 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}; \quad \mathbf{X}_6 = \begin{bmatrix} -1 \\ -3 \end{bmatrix}; \quad \mathbf{X}_7 = \begin{bmatrix} -2 \\ 3 \end{bmatrix}; \quad \mathbf{X}_8 = \begin{bmatrix} -3 \\ 0 \end{bmatrix}; \quad \text{Class 2}$$

- (a) Determine whether the two classes of patterns are linearly separable. (4)
- (b) Find the center of gravity of patterns in each class. Show that a linear decision boundary passing perpendicularly through the middle point of the line joining the two centroids is given by: (3) (2) (1)
- $$6.5x_1 + 2.5x_2 - 14.5 = 0$$

(c) Design a discrete perceptron having the decision boundary as in part (b) for the classification and show that the perceptron separates the points perfectly.

(d) Determine the classes identified by the neuron for following input patterns:

$$\binom{4}{2}, \binom{0}{5}, \binom{36}{13} / \binom{0}{0}$$

3. Use ‘make_blobs’ function from `sklearn.datasets` to create 100 samples of two Gaussian distributed classes for 3-dimensional inputs:

```
from sklearn.datasets import make_blobs
```

Assume each class has a standard deviation = 5.0.

Use `torch.nn.BCELoss()` and `torch.autograd()` functions to train a logistic neuron to separate the two classes.

Find the classification error at convergence and plot the decision boundary.

4. Train a perceptron to learn the following function ϕ :

$$\phi(x, y) = 1.5 + 3.3x - 2.5y + 1.2xy$$

for $0 \leq x, y \leq 1$.

(a) Sample 100 data points randomly from the input space to create a training dataset.

(b) Use the gradient descent algorithm to train the perceptron

(c) Compute the training error and plot the function approximated by the perceptron.

(d) Show how a linear neuron can be used to predict the above function

(e) Compare the results of approximations above by the linear neuron and the perceptron

1. Train a softmax layer of neurons to perform the following classification, given the inputs $x = (x_1, x_2)$ and target class labels d :

(x_1, x_2)	(0 4)	(-1 3)	(2 3)	(-2 2)	(0 2)	(1 2)	(-1 2)	(-3 1)	(-1 1)
d	A	A	A	B	B	A	B	B	B

(x_1, x_2)	(2 1)	(4 1)	(-2 0)	(1 0)	(3 0)	(-3 -1)	(-2 -1)	(2 -1)	(4 -1)
d	C	C	B	C	C	B	B	C	C

- (a) Show one iteration of gradient descent learning at a learning factor 0.05.
 Initialize the weights to $\begin{pmatrix} 0.88 & 0.08 & -0.34 \\ 0.68 & -0.39 & -0.19 \end{pmatrix}$ and biases to zero
- (b) Find the weights and biases at the convergence of learning
 (c) Indicate the probabilities that the network predicts the classes of trained patterns.
 (d) Plot the decision boundaries separating the three classes.
2. Use mini-batch gradient decent learning to train a softmax layer to classify Iris dataset (<https://archive.ics.uci.edu/ml/datasets/Iris>). The dataset contains 150 data points. Use 90 data points for training the classifier and the remaining 60 data points for testing. Plot cross-entropies and classification accuracies against epochs for both train and test data. Set learning rate = 0.1, batch size = 16, and number of epochs = 1000.

You can use the following python commands to load Iris data:

```
from sklearn import datasets
iris = datasets.load_iris()
```

Repeat the classification with batch sizes = 2, 4, 8, 16, 24, 32, and 64, and compare the accuracies and the times taken for an epoch at different batch sizes.

3. Read the Linnerud dataset from `sklearn.datasets` package by using

```
from sklearn.datasets import load_linnerud
```

The Linnerud dataset is a multi-output regression dataset consisting of three exercise (data) and three physiological variables (targets) collected from twenty middle-aged men in a fitness club:

physiological - Weight, Waist and Pulse
exercise - Chins, Situps and Jumps.

Divide the dataset into train and test partitions at 0.75:0.25 ratio and train a perceptron layer to predict physiological data from exercise variables by implementing with

- a) direct gradients
- b) ‘autograd’ functions available in pytorch

Draw learning curves and find mean square error and R^2 values of prediction. Which physiological variable can be better predicted by exercise data?

Compare the time-taken for a weight update and the number of epochs required for convergence for (a) and (b).

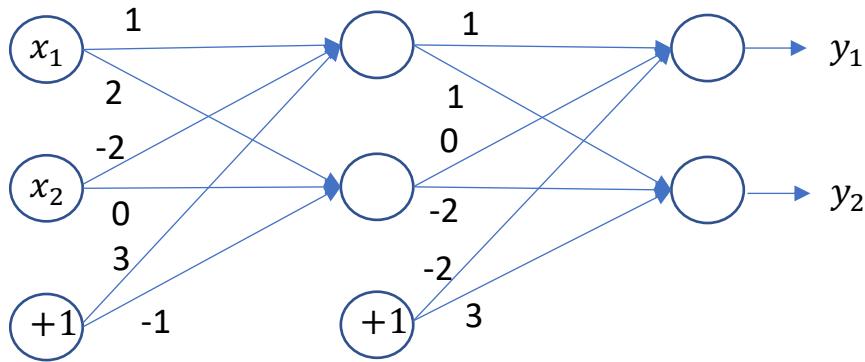


Figure 1

1. The two-layer feedforward perceptron network shown in figure 1 has weights and biases initialized as indicated and receives 2-dimensional inputs (x_1, x_2) . The network is to respond with $\mathbf{d}_1 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$ and $\mathbf{d}_2 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$ for input patterns $x_1 = \begin{pmatrix} 1.0 \\ 3.0 \end{pmatrix}$ and $x_2 = \begin{pmatrix} -2.0 \\ -2.0 \end{pmatrix}$, respectively.

Analyse a single feedforward and feedback step for gradient decent learning of the two patterns by doing the following:

- (a) Find the weight matrix \mathbf{W} to the hidden-layer and weight matrix \mathbf{V} to the output-layer, and the corresponding biases.
- (b) Calculate the synaptic input \mathbf{z} and output \mathbf{h} of the hidden-layer, and the synaptic input \mathbf{u} and output $\mathbf{y} = (y_1, y_2)$ of the output layer.
- (c) Find the mean square error cost J between the outputs and targets.
- (d) Calculate the gradients $\nabla_u J$ and $\nabla_z J$ at the output-layer and the hidden-layer, respectively.
- (e) Compute the new weights and biases.
- (f) Write a program to continue iterations until convergence and find the final weights and biases.

Assume a learning rate of 0.05.

Repeat above (a) – (f) for stochastic gradient decent learning.

2. A feedforward neural network with one hidden layer is to perform the following classification:

class	inputs
A	(1.0, 1.0), (0.0, 1.0)
B	(3.0, 4.0), (2.0, 2.0)
C	(2.0, -2.0), (-2.0, -3.0)

The network has a hidden layer consisting of three perceptrons and a softmax output layer.

Initialize the weights \mathbf{W} and biases \mathbf{b} to the hidden layer, and the weights \mathbf{V} and biases \mathbf{c} to the output layer as follows:

$$\mathbf{W} = \begin{pmatrix} -0.10 & 0.97 & 0.18 \\ -0.70 & 0.38 & 0.93 \end{pmatrix}, \mathbf{b} = \begin{pmatrix} 0.0 \\ 0.0 \\ 0.0 \end{pmatrix}$$

$$\mathbf{V} = \begin{pmatrix} 1.01 & 0.09 & -0.39 \\ 0.79 & -0.45 & -0.22 \\ 0.28 & 0.96 & -0.07 \end{pmatrix}, \mathbf{c} = \begin{pmatrix} 0.0 \\ 0.0 \\ 0.0 \end{pmatrix}$$

Show one iteration of gradient descent learning and plot the learning curves until convergence at a learning rate $\alpha = 0.1$.

Determine the weights and biases at convergence.

Find the class labels predicted by the trained network for patterns:

$$\mathbf{x}_1 = \begin{pmatrix} 2.5 \\ 1.5 \end{pmatrix} \text{ and } \mathbf{x}_2 = \begin{pmatrix} -1.5 \\ 0.5 \end{pmatrix}$$

3. Design a deep neural network consisting of two ReLU hidden layers to approximate the following function:

$$\phi(x, y) = 0.8x^2 - y^3 + 2.5xy$$

for $-1.0 \leq x, y \leq 1.0$.

Use 10 neurons and 5 neurons at first and second hidden layers, respectively, and a linear neuron at the output layer.

- (a) Divide the input space equally into square regions of size 0.25×0.25 and use grid points as data points to learn the function ϕ .
- (b) Train the network using gradient decent learning at learning rate $\alpha = 0.01$ and plot the learning curve (mean square error vs. iterations) and the predicted data points.
- (c) Compare the behaviour of learning by plotting learning curves at rates $\alpha = 0.005, 0.001, 0.01$, and 0.05 .

1. The first hidden layer of a convolution neural network (CNN) has a convolution layer consisting of two feature maps with filters w_1 and w_2 and biases = 0.1, and neurons having sigmoid activation functions, and a pooling layer with a pooling windows of size 2x2:

$$w_1 = \begin{pmatrix} 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{pmatrix} \text{ and } w_2 = \begin{pmatrix} -1 & -1 & -1 \\ 0 & 0 & 0 \\ 1 & 1 & 1 \end{pmatrix}.$$

The input layer is of 6x6 size and receives an input image I :

$$I = \begin{pmatrix} 0.7 & 0.1 & 0.2 & 0.3 & 0.3 & 0.5 \\ 0.8 & 0.1 & 0.3 & 0.5 & 0.1 & 0.0 \\ 1.0 & 0.2 & 0.0 & 0.3 & 0.2 & 0.7 \\ 0.8 & 0.1 & 0.5 & 0.6 & 0.3 & 0.4 \\ 0.1 & 0.0 & 0.9 & 0.3 & 0.3 & 0.2 \\ 1.0 & 0.1 & 0.4 & 0.5 & 0.2 & 0.8 \end{pmatrix}$$

- a. Find the outputs at the first convolution layer if
 - i. padding = 0 and strides = (1, 1).
 - ii. padding = 1 and strides = (2, 2).
- b. Find the outputs at the first pooling layer for Part (a), assuming strides of (2, 2) and pooling is
 - i. max pooling
 - ii. mean pooling

- 2. Given ‘3wolfmoon.jpg’ color image of size 639×516 .
 - a. Initialize weights and biases of a convolutional layer with two kernels of size 9×9 . Note that the input image is in color and has three channels.
 - b. Display the feature maps at the convolution layer, assuming sigmoid activation functions. Use VALID padding and strides = 1.
 - c. Display the outputs of a mean pooling layer with a pooling window size 5×5 and strides = 5.

- 3. Design a CNN with one hidden layer to recognize digit images in MNIST database:
<http://yann.lecun.com/exdb/mnist/>

The convolution layer consists of 25 filters of dimensions 9×9 and the pooling layer has a pooling window size 4×4 . Assume VALID padding and default strides for both convolution

and pooling layer. Train the network with mini batch gradient decent learning with learning factor $\alpha = 10^{-3}$ and batch size = 128.

Plot

- a. The training and test errors against learning epochs.
- b. Final filter weights
- c. Feature maps at the convolution and pooling layers for a representative test pattern
- d. Repeat training by introducing decay parameter $\beta = 10^{-6}$ and momentum term with $\gamma = 0.5$, and compare the learning curves

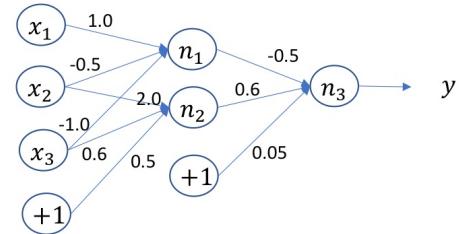


Figure 1

1. Figure 1 shows a two-layer feedforward neural network receiving 3-dimensional inputs $(x_1, x_2, x_3) \in \mathbb{R}^3$. The connection weights and biases of the neurons n_1, n_2 , and n_3 are as indicated in the figure. The hidden-layer neurons have activation functions given by $g(u) = \frac{1.0}{1+e^{-0.5u}}$ where u denotes the synaptic input to the neuron. The activation function $f(u)$ of the output neuron is a ReLU function: $f(u) = \max\{0, u\}$.

- Write weight vectors and biases connected to individual neurons, and the weight matrix and bias vector connected to the hidden layer.
- Find the synaptic inputs and activations of the neurons for the following input signals: $\begin{matrix} u \\ y \end{matrix}$

- (1.0, -0.5, 1.0)
- (-1.0, 0.0, -2.0)
- (2.0, 0.5, -1.0)

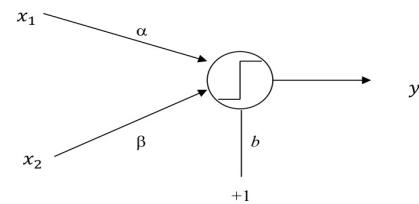


Figure 2

2. Two input binary neuron shown in figure 2 has a unit step activation function with a bias $b = 0.5$ and receives two-dimensional input $(x_1, x_2) \in \mathbb{R}^2$.

- Find the space of possible values of weights (α, β) if the neuron
 - ON for input (1.0, 1.0)
 - ON for input (0.5, -1.0)
 - OFF for input (2.0, -0.5).
- Indicate the weight space in 2-D α - β plot and show that $(-0.2, 0.2)$ is in this space.

3. The network shown in figure 3 consists of neurons having threshold activation functions and receives three-bit binary patterns $(x_1, x_2, x_3) \in \{0,1\}^3$. By analyzing the outputs for all possible three-bit input patterns, determine the logic function that the network implements. All unlabeled weights shown in figure 3 are of unity weight.

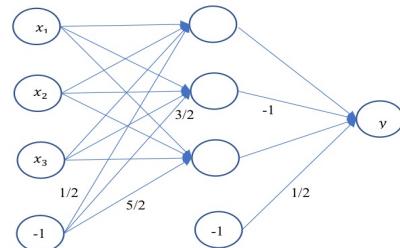


Figure 3

1. Train a linear neuron to learn the mapping from input $x \in R^3$ to output y from the following examples:

	$x = (x_1, x_2, x_3)$	y
1	(0.09 -0.44 -0.15)	-2.57
2	(0.69 -0.99 -0.76)	-2.97
3	(0.34 0.65 -0.73)	0.96
4	(0.15 0.78 -0.58)	1.04
5	(-0.63 -0.78 -0.56)	-3.21
6	(0.96 0.62 -0.66)	1.05
7	(0.63 -0.45 -0.14)	-2.39
8	(0.88 0.64 -0.33)	0.66

- (a) Show one iteration of learning of the neuron with
- i. Stochastic gradient descent learning
 - ii. Gradient descent learning
- Initialize the weights as $\begin{pmatrix} 0.77 \\ 0.02 \\ 0.63 \end{pmatrix}$ and biases to 0.0, and use a leaning factor $\alpha = 0.01$.
- (b) Plot the learning curves (mean square error vs. epochs) until convergence. Determine the learned weights and biases.
- (c) Find the predicted values y of training inputs after the training.

2. Two-dimensional training patterns (inputs) to design a dichotomizer are given as:

$$\mathbf{X}_1 = \begin{bmatrix} 5 \\ 1 \end{bmatrix}; \quad \mathbf{X}_2 = \begin{bmatrix} 7 \\ 3 \end{bmatrix}; \quad \mathbf{X}_3 = \begin{bmatrix} 3 \\ 2 \end{bmatrix}; \quad \mathbf{X}_4 = \begin{bmatrix} 5 \\ 4 \end{bmatrix}; \quad \text{Class 1}$$

$$\mathbf{X}_5 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}; \quad \mathbf{X}_6 = \begin{bmatrix} -1 \\ -3 \end{bmatrix}; \quad \mathbf{X}_7 = \begin{bmatrix} -2 \\ 3 \end{bmatrix}; \quad \mathbf{X}_8 = \begin{bmatrix} -3 \\ 0 \end{bmatrix}; \quad \text{Class 2}$$

- (a) Determine whether the two classes of patterns are linearly separable.
- (b) Find the center of gravity of patterns in each class. Show that a linear decision boundary passing perpendicularly through the middle point of the line joining the two centroids is given by:
- $$6.5x_1 + 2.5x_2 - 14.5 = 0$$
- (c) Design a discrete perceptron having the decision boundary as in part (b) for the classification and show that the perceptron separates the points perfectly.
- (d) Determine the classes identified by the neuron for following input patterns:

$$\binom{4}{2}, \binom{0}{5}, \binom{36/13}{0}$$

3. Use 'make_blobs' function from `sklearn.datasets` to create 100 samples of two Gaussian distributed classes for 3-dimensional inputs:

```
from sklearn.datasets import make_blobs
```

Assume each class has a standard deviation = 5.0.

Use `torch.nn.BCELoss()` and `torch.autograd()` functions to train a logistic neuron to separate the two classes.

Find the classification error at convergence and plot the decision boundary.

4. Train a perceptron to learn the following function ϕ :

$$\phi(x, y) = 1.5 + 3.3x - 2.5y + 1.2xy$$

for $0 \leq x, y \leq 1$.

- (a) Sample 100 data points randomly from the input space to create a training dataset.
- (b) Use the gradient descent algorithm to train the perceptron
- (c) Compute the training error and plot the function approximated by the perceptron.
- (d) Show how a linear neuron can be used to predict the above function
- (e) Compare the results of approximations above by the linear neuron and the perceptron

1. Train a softmax layer of neurons to perform the following classification, given the inputs $x = (x_1, x_2)$ and target class labels d :

(x_1, x_2)	(0 4)	(-1 3)	(2 3)	(-2 2)	(0 2)	(1 2)	(-1 2)	(-3 1)	(-1 1)
d	A	A	A	B	B	A	B	B	B

(x_1, x_2)	(2 1)	(4 1)	(-2 0)	(1 0)	(3 0)	(-3 -1)	(-2 -1)	(2 -1)	(4 -1)
d	C	C	B	C	C	B	B	C	C

- (a) Show one iteration of gradient descent learning at a learning factor 0.05.
 Initialize the weights to $\begin{pmatrix} 0.88 & 0.08 & -0.34 \\ 0.68 & -0.39 & -0.19 \end{pmatrix}$ and biases to zero
 (b) Find the weights and biases at the convergence of learning
 (c) Indicate the probabilities that the network predicts the classes of trained patterns.
 (d) Plot the decision boundaries separating the three classes.
2. Use mini-batch gradient decent learning to train a softmax layer to classify Iris dataset (<https://archive.ics.uci.edu/ml/datasets/Iris>). The dataset contains 150 data points. Use 90 data points for training the classifier and the remaining 60 data points for testing. Plot cross-entropies and classification accuracies against epochs for both train and test data. Set learning rate = 0.1, batch size = 16, and number of epochs = 1000.

You can use the following python commands to load Iris data:

```
from sklearn import datasets
iris = datasets.load_iris()
```

Repeat the classification with batch sizes = 2, 4, 8, 16, 24, 32, and 64, and compare the accuracies and the times taken for an epoch at different batch sizes.

3. Read the Linnerud dataset from sklearn.datasets package by using

```
from sklearn.datasets import load_linnerud
```

The Linnerud dataset is a multi-output regression dataset consisting of three exercise (data) and three physiological variables (targets) collected from twenty middle-aged men in a fitness club:

physiological - Weight, Waist and Pulse
 exercise - Chins, Situps and Jumps.

Divide the dataset into train and test partitions at 0.75:0.25 ratio and train a perceptron layer to predict physiological data from exercise variables by implementing with

- a) direct gradients
- b) ‘autograd’ functions available in pytorch

Draw learning curves and find mean square error and R² values of prediction. Which physiological variable can be better predicted by exercise data?

Compare the time-taken for a weight update and the number of epochs required for convergence for (a) and (b).

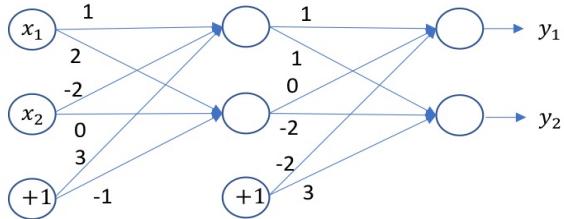


Figure 1

1. The two-layer feedforward perceptron network shown in figure 1 has weights and biases initialized as indicated and receives 2-dimensional inputs (x_1, x_2) . The network is to respond with $d_1 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$ and $d_2 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$ for input patterns $x_1 = \begin{pmatrix} 1.0 \\ 3.0 \end{pmatrix}$ and $x_2 = \begin{pmatrix} -2.0 \\ -2.0 \end{pmatrix}$, respectively.

Analyse a single feedforward and feedback step for gradient decent learning of the two patterns by doing the following:

- Find the weight matrix W to the hidden-layer and weight matrix V to the output-layer, and the corresponding biases.
- Calculate the synaptic input z and output h of the hidden-layer, and the synaptic input u and output $y = (y_1, y_2)$ of the output layer.
- Find the mean square error cost J between the outputs and targets.
- Calculate the gradients $\nabla_u J$ and $\nabla_z J$ at the output-layer and the hidden-layer, respectively.
- Compute the new weights and biases.
- Write a program to continue iterations until convergence and find the final weights and biases.

Assume a learning rate of 0.05.

Repeat above (a) – (f) for stochastic gradient decent learning.

2. A feedforward neural network with one hidden layer is to perform the following classification:

class	inputs
A	(1.0, 1.0), (0.0, 1.0)
B	(3.0, 4.0), (2.0, 2.0)
C	(2.0, -2.0), (-2.0, -3.0)

The network has a hidden layer consisting of three perceptrons and a softmax output layer.

Initialize the weights W and biases b to the hidden layer, and the weights V and biases c to the output layer as follows:

$$W = \begin{pmatrix} -0.10 & 0.97 & 0.18 \\ -0.70 & 0.38 & 0.93 \end{pmatrix}, b = \begin{pmatrix} 0.0 \\ 0.0 \\ 0.0 \end{pmatrix}$$

$$V = \begin{pmatrix} 1.01 & 0.09 & -0.39 \\ 0.79 & -0.45 & -0.22 \\ 0.28 & 0.96 & -0.07 \end{pmatrix}, c = \begin{pmatrix} 0.0 \\ 0.0 \\ 0.0 \end{pmatrix}$$

Show one iteration of gradient descent learning and plot the learning curves until convergence at a learning rate $\alpha = 0.1$.

Determine the weights and biases at convergence.

Find the class labels predicted by the trained network for patterns:

$$x_1 = \begin{pmatrix} 2.5 \\ 1.5 \end{pmatrix} \text{ and } x_2 = \begin{pmatrix} -1.5 \\ 0.5 \end{pmatrix}$$

3. Design a deep neural network consisting of two ReLU hidden layers to approximate the following function:

$$\phi(x, y) = 0.8x^2 - y^3 + 2.5xy$$

for $-1.0 \leq x, y \leq 1.0$.

Use 10 neurons and 5 neurons at first and second hidden layers, respectively, and a linear neuron at the output layer.

- (a) Divide the input space equally into square regions of size 0.25×0.25 and use grid points as data points to learn the function ϕ .
- (b) Train the network using gradient decent learning at learning rate $\alpha = 0.01$ and plot the learning curve (mean square error vs. iterations) and the predicted data points.
- (c) Compare the behaviour of learning by plotting learning curves at rates $\alpha = 0.005, 0.001, 0.01$, and 0.05 .

1. The first hidden layer of a convolution neural network (CNN) has a convolution layer consisting of two feature maps with filters w_1 and w_2 and biases = 0.1, and neurons having sigmoid activation functions, and a pooling layer with a pooling windows of size 2x2:

$$w_1 = \begin{pmatrix} 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{pmatrix} \text{ and } w_2 = \begin{pmatrix} -1 & -1 & -1 \\ 0 & 0 & 0 \\ 1 & 1 & 1 \end{pmatrix}.$$

The input layer is of 6x6 size and receives an input image I :

$$I = \begin{pmatrix} 0.7 & 0.1 & 0.2 & 0.3 & 0.3 & 0.5 \\ 0.8 & 0.1 & 0.3 & 0.5 & 0.1 & 0.0 \\ 1.0 & 0.2 & 0.0 & 0.3 & 0.2 & 0.7 \\ 0.8 & 0.1 & 0.5 & 0.6 & 0.3 & 0.4 \\ 0.1 & 0.0 & 0.9 & 0.3 & 0.3 & 0.2 \\ 1.0 & 0.1 & 0.4 & 0.5 & 0.2 & 0.8 \end{pmatrix}$$

- a. Find the outputs at the first convolution layer if
 - i. padding = 0 and strides = (1, 1).
 - ii. padding = 1 and strides = (2, 2).
- b. Find the outputs at the first pooling layer for Part (a), assuming strides of (2, 2) and pooling is
 - i. max pooling
 - ii. mean pooling

2. Given '3wolfmoon.jpg' color image of size 639×516 .

- a. Initialize weights and biases of a convolutional layer with two kernels of size 9×9 . Note that the input image is in color and has three channels.
- b. Display the feature maps at the convolution layer, assuming sigmoid activation functions. Use VALID padding and strides = 1.
- c. Display the outputs of a mean pooling layer with a pooling window size 5×5 and strides = 5.

3. Design a CNN with one hidden layer to recognize digit images in MNIST database:

<http://yann.lecun.com/exdb/mnist/>

The convolution layer consists of 25 filters of dimensions 9×9 and the pooling layer has a pooling window size 4×4 . Assume VALID padding and default strides for both convolution

and pooling layer. Train the network with mini batch gradient decent learning with learning factor $\alpha = 10^{-3}$ and batch size = 128.

Plot

- a. The training and test errors against learning epochs.
- b. Final filter weights
- c. Feature maps at the convolution and pooling layers for a representative test pattern
- d. Repeat training by introducing decay parameter $\beta = 10^{-6}$ and momentum term with $\gamma = 0.5$, and compare the learning curves

Neuron Layers

SC4001 – Tutorial 3

1. Design a softmax layer of neurons to perform the following classification, given the inputs $x = (x_1, x_2)$ and target class labels d :

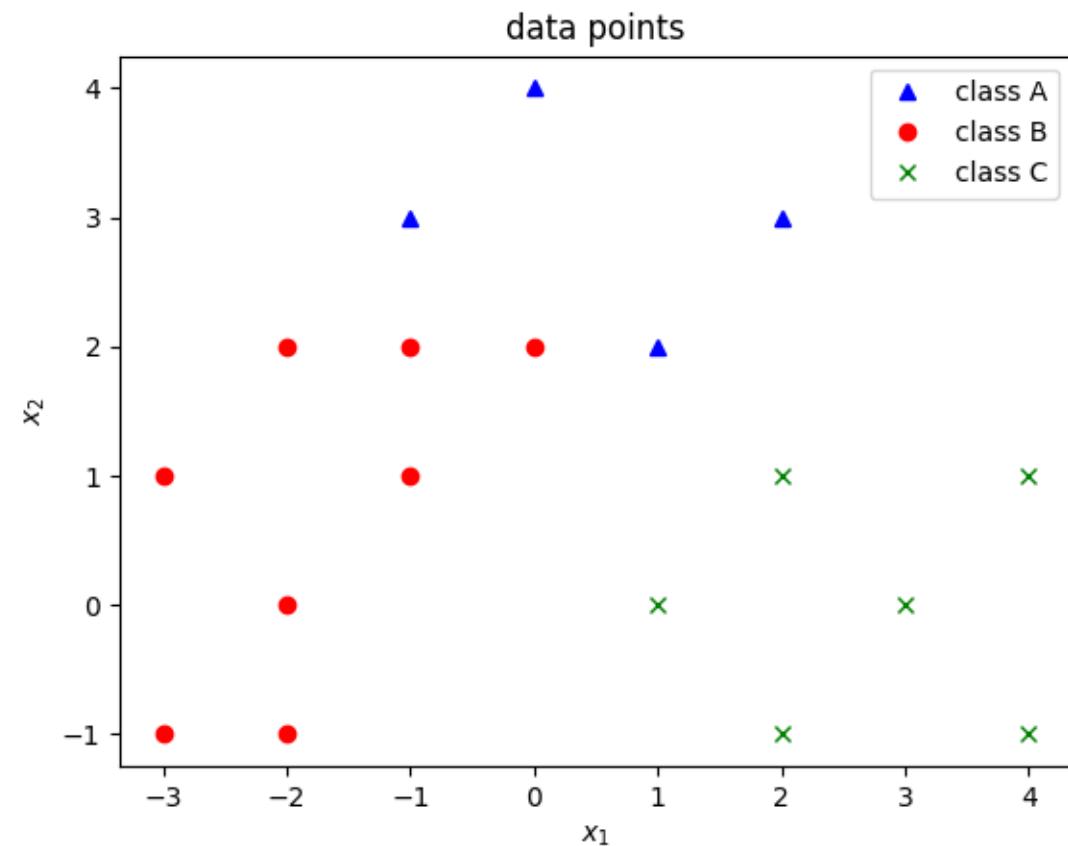
(x_1, x_2)	(0 4)	(-1 3)	(2 3)	(-2 2)	(0 2)	(1 2)	(-1 2)	(-3 1)	(-1 1)
d	A	A	A	B	B	A	B	B	B

(x_1, x_2)	(2 1)	(4 1)	(-2 0)	(1 0)	(3 0)	(-3 -1)	(-2 -1)	(2 -1)	(4 -1)
d	C	C	B	C	C	B	B	C	C

- (a) Show one iteration of gradient descent learning at a learning factor 0.05.

Initialize the weights to $\begin{pmatrix} 0.88 & 0.08 & -0.34 \\ 0.68 & -0.39 & -0.19 \end{pmatrix}$ and biases to zero

- (b) Find the weights and biases at convergence of learning
(c) Indicate the probabilities that the network predicts the classes of trained patterns.
(d) Plot the decision boundaries separating the three classes.



GD for Softmax layer

Given training set (X, d)

Set learning rate α

Initialize W and b

Iterate until convergence:

$$U = XW + B$$

$$f(U) = \frac{e^U}{\sum_{k=1}^K e^{U_k}}$$

$$\nabla_U J = -(\mathbf{K} - f(U))$$

$$W \leftarrow W - \alpha X^T \nabla_U J$$

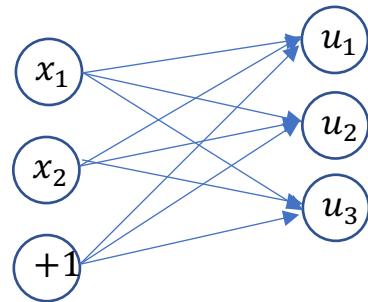
$$b \leftarrow b - \alpha (\nabla_U J)^T \mathbf{1}_P$$

Labels for classes:

Class A → 1, Class B → 2, Class C → 3

The data matrix and target vector:

$$\mathbf{X} = \begin{pmatrix} 0 & 4 \\ -1 & 3 \\ 2 & 3 \\ -2 & 2 \\ 0 & 2 \\ 1 & 2 \\ -1 & 2 \\ \vdots & \vdots \\ 4 & -1 \end{pmatrix}, \quad \mathbf{d} = \begin{pmatrix} 1 \\ 1 \\ 1 \\ 2 \\ 2 \\ 1 \\ 2 \\ \vdots \\ 3 \end{pmatrix}, \quad \mathbf{K} = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ \vdots & \vdots & \vdots \\ 0 & 0 & 1 \end{pmatrix}$$



$$f(\mathbf{U}) = \frac{e^{\mathbf{U}}}{\sum_{k=1}^K e^{\mathbf{U}_k}} = P(y = k|x)$$

$$\mathbf{Y} = \underset{k}{\operatorname{argmax}} f(\mathbf{U})$$

Learning rate $\alpha = 0.05$.

Initialize weights and biases:

$$\mathbf{W} = \begin{pmatrix} 0.88 & 0.08 & -0.34 \\ 0.68 & -0.39 & -0.19 \end{pmatrix} \text{ and } \mathbf{b} = \begin{pmatrix} 0.0 \\ 0.0 \\ 0.0 \end{pmatrix}$$

Epoch 1:

$$\mathbf{U} = \mathbf{XW} + \mathbf{B} = \begin{pmatrix} 0 & 4 \\ -1 & 3 \\ 2 & 3 \\ -2 & 2 \\ \vdots & \vdots \\ 4 & -1 \end{pmatrix} \begin{pmatrix} 0.88 & 0.08 & -0.34 \\ 0.68 & -0.39 & -0.19 \end{pmatrix} + \begin{pmatrix} 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \\ \vdots & \vdots & \vdots \\ 0.0 & 0.0 & 0.0 \end{pmatrix} = \begin{pmatrix} 2.72 & -1.54 & -0.75 \\ 1.17 & -1.23 & -0.23 \\ 3.8 & -1.0 & -1.23 \\ -0.39 & -0.93 & 0.30 \\ \vdots & \vdots & \vdots \\ 2.82 & 0.71 & -1.16 \end{pmatrix}$$

$$\mathbf{U} = \begin{pmatrix} 2.72 & -1.54 & -0.75 \\ 1.17 & -1.23 & -0.23 \\ 3.8 & -1.0 & -1.23 \\ -0.39 & -0.93 & 0.30 \\ \vdots & \vdots & \vdots \\ 2.82 & 0.71 & -1.16 \end{pmatrix}$$

$$f(u_i) = \frac{e^{u_i}}{\sum_{k=1}^K e^{u_k}}$$

$$f(\mathbf{U}) == \begin{pmatrix} \frac{e^{2.72}}{e^{2.72} + e^{-1.54} + e^{-0.75}} & \frac{e^{-1.54}}{e^{2.72} + e^{-1.54} + e^{-0.75}} & \frac{e^{-0.75}}{e^{2.72} + e^{-1.54} + e^{-0.75}} \\ \frac{e^{1.17}}{e^{1.17} + e^{-1.23} + e^{-0.23}} & \frac{e^{-1.23}}{e^{1.17} + e^{-1.23} + e^{-0.23}} & \frac{e^{-0.23}}{e^{1.17} + e^{-1.23} + e^{-0.23}} \\ \vdots & \vdots & \vdots \\ \frac{e^{2.82}}{e^{2.82} + e^{0.71} + e^{-1.16}} & \frac{e^{0.71}}{e^{2.82} + e^{0.71} + e^{-1.16}} & \frac{e^{-1.16}}{e^{2.82} + e^{0.71} + e^{-1.16}} \end{pmatrix} = \begin{pmatrix} 0.96 & 0.01 & 0.03 \\ 0.75 & 0.07 & 0.18 \\ \vdots & \vdots & \vdots \\ 0.88 & 0.11 & 0.02 \end{pmatrix}$$

$$\mathbf{y} = \operatorname{argmax}_k f(\mathbf{U}) = \operatorname{argmax}_k \begin{pmatrix} 0.96 & 0.01 & 0.03 \\ 0.75 & 0.07 & 0.18 \\ 0.99 & 0.01 & 0.01 \\ 0.28 & 0.16 & 0.56 \\ \vdots & \vdots & \vdots \\ 0.88 & 0.11 & 0.02 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \\ 1 \\ 3 \\ \vdots \\ 1 \end{pmatrix}$$

$$\mathbf{d} = \begin{pmatrix} 1 \\ 1 \\ 1 \\ 2 \\ \vdots \\ 3 \end{pmatrix}$$

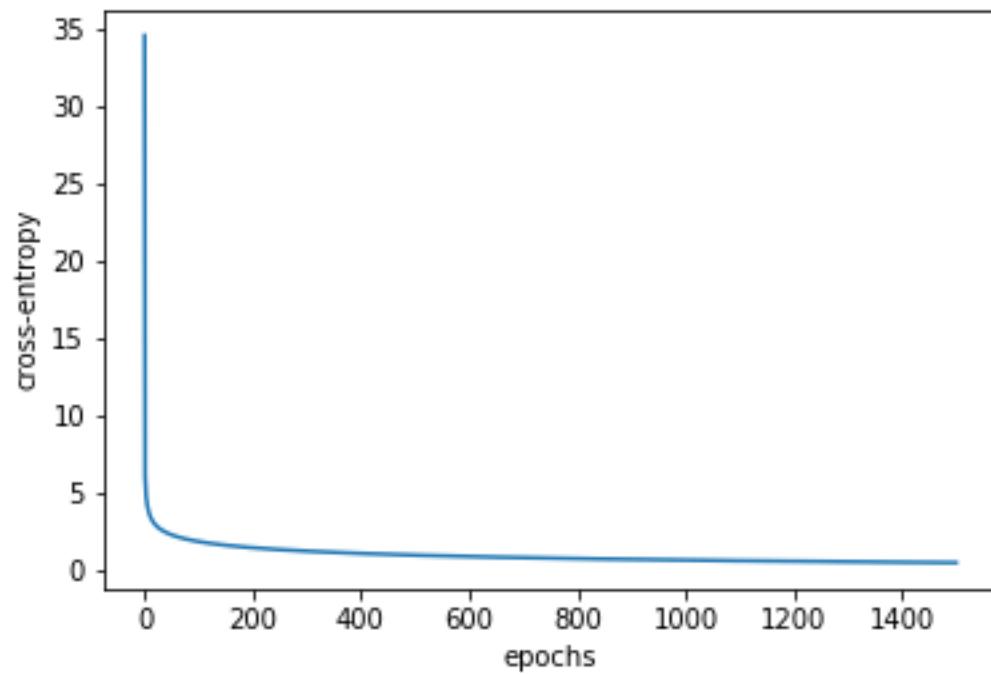
$$\text{Classification error} = \sum_{p=1}^P \mathbb{1}(\mathbf{y} \neq \mathbf{d}) = 14$$

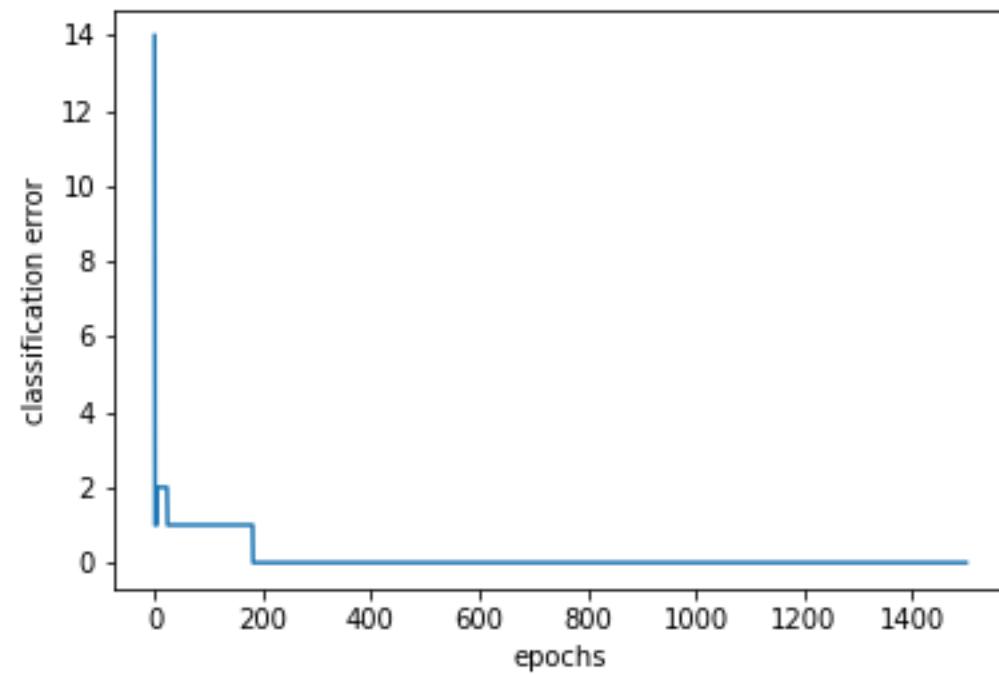
$$\begin{aligned} \text{entropy} &= - \sum_{p=1}^P \log \left(f(u_{pd_p}) \right) \\ &= -\log(0.96) - \log(0.75) - \log(0.99) - \log(0.16) \cdots - \log(0.02) \\ &= 34.36 \end{aligned}$$

$$\nabla_{\mathbf{U}} J = -(\mathbf{K} - f(\mathbf{U})) = - \left(\begin{pmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ \vdots & \vdots & \vdots \\ 0 & 0 & 1 \end{pmatrix} - \begin{pmatrix} 0.96 & 0.01 & 0.03 \\ 0.75 & 0.07 & 0.18 \\ \vdots & \vdots & \vdots \\ 0.88 & 0.11 & 0.02 \end{pmatrix} \right) = \begin{pmatrix} -0.04 & 0.01 & 0.03 \\ -0.25 & 0.07 & 0.18 \\ \vdots & \vdots & \vdots \\ 0.88 & 0.11 & -0.98 \end{pmatrix}$$

$$\begin{aligned} \mathbf{W} &= \mathbf{W} - \alpha \mathbf{X}^T \nabla_{\mathbf{U}} J = \begin{pmatrix} 0.88 & 0.08 & -0.34 \\ 0.68 & -0.39 & -0.19 \end{pmatrix} - 0.05 \begin{pmatrix} 0 & -1 & \cdots & 4 \\ 4 & 3 & \cdots & -1 \end{pmatrix} \begin{pmatrix} -0.04 & 0.01 & 0.03 \\ -0.25 & 0.07 & 0.18 \\ \vdots & \vdots & \vdots \\ 0.88 & 0.11 & -0.98 \end{pmatrix} \\ &= \begin{pmatrix} 0.28 & -0.54 & 0.89 \\ 0.54 & -0.12 & -0.31 \end{pmatrix} \end{aligned}$$

$$\mathbf{b} = \mathbf{b} - \alpha (\nabla_{\mathbf{U}} J)^T \mathbf{1}_P = \begin{pmatrix} 0.0 \\ 0.0 \\ 0.0 \end{pmatrix} + 0.05 \begin{pmatrix} -0.04 & 0.01 & 0.03 \\ -0.25 & 0.07 & 0.18 \\ \vdots & \vdots & \vdots \\ 0.88 & 0.11 & -0.98 \end{pmatrix}^T \begin{pmatrix} 1 \\ 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} = \begin{pmatrix} -0.32 \\ 0.27 \\ 0.06 \end{pmatrix}$$





At convergence:

$$\mathbf{w} = \begin{pmatrix} -0.15 & -3.41 & 4.18 \\ 5.27 & -1.02 & -4.15 \end{pmatrix} \text{ and } \mathbf{b} = \begin{pmatrix} -7.82 \\ 5.81 \\ 2.02 \end{pmatrix}$$

Entropy = 0.58

Classification Error = 0

At convergence:

$$X = \begin{pmatrix} -1 & 2 \\ 0 & 4 \\ -1 & 3 \\ 0 & 2 \\ 3 & 0 \\ -2 & -1 \\ 4 & 1 \\ 1 & 2 \\ 2 & -1 \\ 2 & 3 \\ 2 & 1 \\ -2 & 0 \\ -3 & -1 \\ 1 & 0 \\ -1 & 1 \\ 4 & -1 \\ -3 & 1 \\ -2 & 2 \end{pmatrix}, f(U) = \begin{pmatrix} 1.0 & 0.0 & 0.0 \\ 0.88 & 0.12 & 0.0 \\ 1.0 & 0.0 & 0.0 \\ 0.0 & 1.0 & 0.0 \\ 0.26 & 0.74 & 0.0 \\ 0.89 & 0.1 & 0.0 \\ 0.01 & 0.99 & 0.0 \\ 0.0 & 1.0 & 0.0 \\ 0.0 & 1.0 & 0.0 \\ 0.0 & 0.0 & 1.0 \\ 0.0 & 0.0 & 1.0 \\ 0.0 & 1.0 & 0.0 \\ 0.0 & 0.02 & 0.98 \\ 0.0 & 0.0 & 1.0 \\ 0.0 & 1.0 & 0.0 \\ 0.0 & 1.0 & 0.0 \\ 0.0 & 0.0 & 1.0 \end{pmatrix}, Y = \begin{pmatrix} 1 \\ 1 \\ 1 \\ 2 \\ 2 \\ 1 \\ 2 \\ 2 \\ 2 \\ 3 \\ 3 \\ 3 \\ 3 \\ 3 \\ 2 \\ 2 \\ 3 \\ 3 \end{pmatrix}$$

Probabilities that input patterns belong to target classes are given in RED.

At convergence:

$$\mathbf{w} = \begin{pmatrix} -0.15 & -3.41 & 4.18 \\ 5.27 & -1.02 & -4.15 \end{pmatrix} \text{ and } \mathbf{b} = \begin{pmatrix} -7.82 \\ 5.81 \\ 2.02 \end{pmatrix}$$

Synaptic inputs at the softmax layer for an input $\mathbf{x} = (x_1, x_2)$:

Neuron of class A, $u_1 = \mathbf{w}_1^T \mathbf{x} + b_1 = -0.15x_1 + 5.27x_2 - 7.82$

Neuron of class B, $u_2 = \mathbf{w}_2^T \mathbf{x} + b_2 = -3.41x_1 - 1.02x_2 + 5.81$

Neuron of class C, $u_3 = \mathbf{w}_3^T \mathbf{x} + b_3 = 4.18x_1 - 4.15x_2 + 2.02$

Decision boundaries:

Between class A and class B is given when $u_1 = u_2$

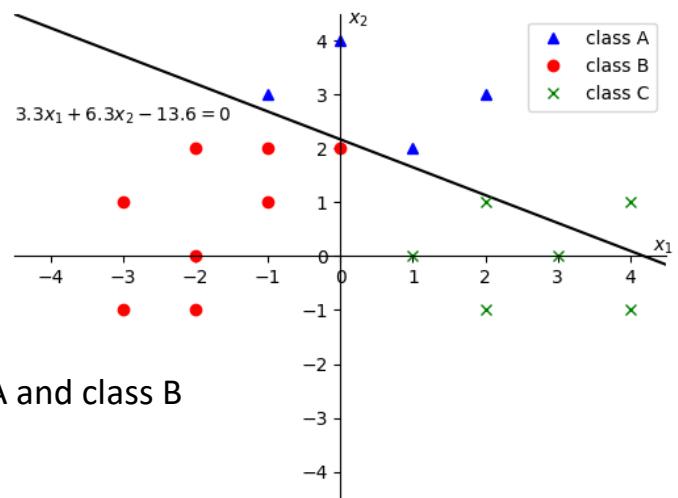
$$\begin{aligned} -0.15x_1 + 5.27x_2 - 7.82 &= -3.41x_1 - 1.02x_2 + 5.81 \\ 3.25x_1 + 6.29x_2 - 13.63 &= 0 \end{aligned}$$

Similarly, between class B and class C $u_2 = u_3$:

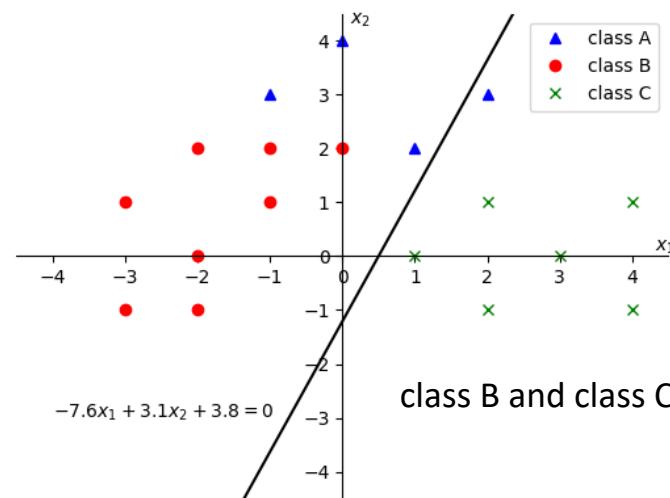
$$-7.59x_1 + 3.13x_2 + 3.79 = 0$$

between class A and class C $u_3 = u_1$:

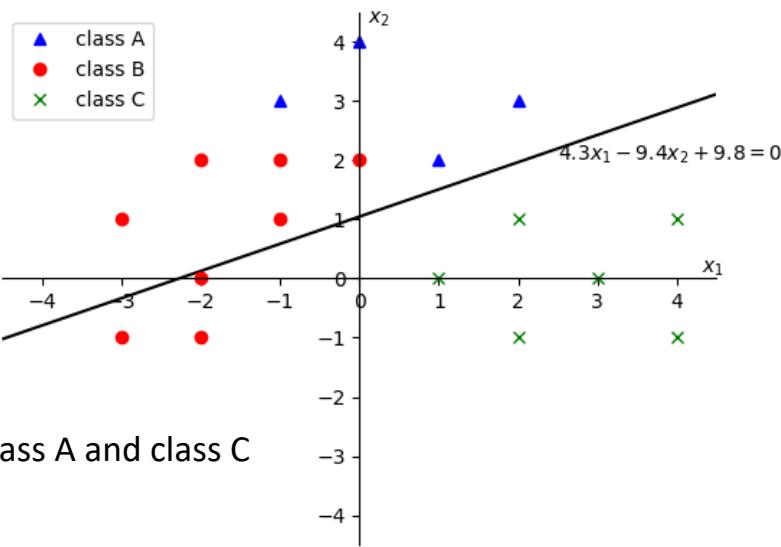
$$4.33x_1 - 9.42x_2 + 9.84 = 0$$



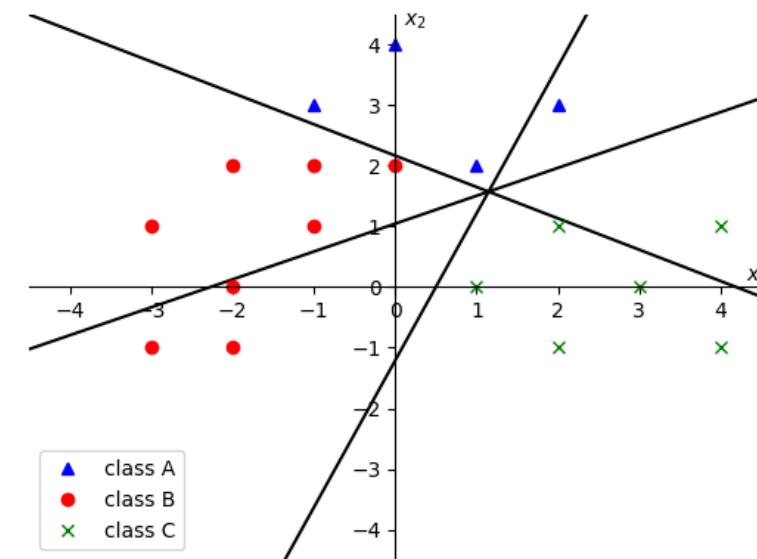
class A and class B



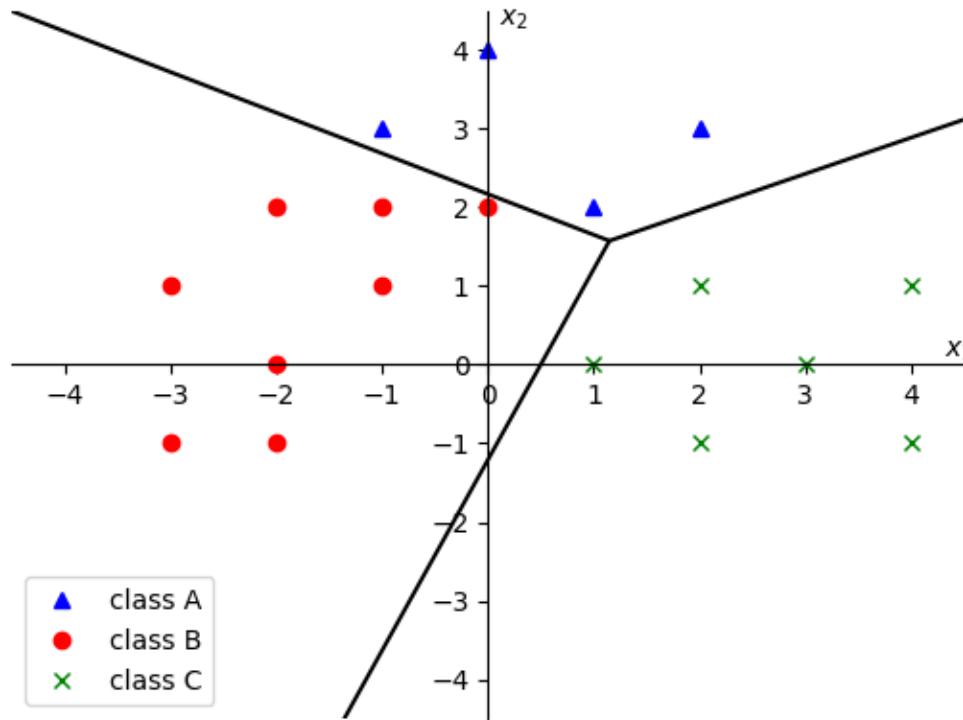
class B and class C



class A and class C



Decision boundaries learnt by the softmax layer



3. Read the Linnerud dataset from `sklearn.datasets` package by using

```
from sklearn.datasets import load_linnerud
```

The Linnerud dataset is a multi-output regression dataset consisting of three exercise (data) and three physiological variables (targets) collected from twenty middle-aged men in a fitness club:

physiological - Weight, Waist and Pulse
exercise - Chins, Situps and Jumps.

Divide the dataset into train and test partitions at 0.75:0.25 ratio and train a perceptron layer to predict physiological data from exercise variables by implementing with

- a) direct gradients
- b) ‘autograd’ functions available in pytorch

Draw learning curves and find mean square error and R^2 values of prediction. Which physiological variable can be better predicted by exercise data?

Compare the time-taken for a weight update and the number of epochs required for convergence for (a) and (b).

```
from sklearn.datasets import load_linnerud
```

```
X, y = load_linnerud(return_X_y=True)
```

```
X  
chins  sit_ups  jumps  
0      5.0       162.0    60.0  
1      2.0       110.0    60.0  
2     12.0       101.0   101.0  
3     12.0       105.0    37.0  
4     13.0       155.0    58.0
```

```
y  
weight  waist  pulse  
0      191.0    36.0     50.0  
1      189.0    37.0     52.0  
2      193.0    38.0     58.0  
3      162.0    35.0     62.0  
4      189.0    35.0     46.0
```

Parameters of the data

	targets			inputs		
--	---------	--	--	--------	--	--

	weight	waist	pulse	chins	sit_ups	jumps
count	20.000000	20.000000	20.000000	20.000000	20.000000	20.000000
mean	9.450000	145.550000	70.30000	178.600000	35.400000	56.100000
std	5.286278	62.566575	51.27747	24.690505	3.201973	7.210373
min	1.000000	50.000000	25.00000	138.000000	31.000000	46.000000
25%	4.750000	101.000000	39.50000	160.750000	33.000000	51.500000
50%	11.500000	122.500000	54.00000	176.000000	35.000000	55.000000
75%	13.250000	210.000000	85.25000	191.500000	37.000000	60.500000
max	17.000000	251.000000	250.00000	247.000000	46.000000	74.000000

Normalizing input variables such that $x' \sim N(0, 1)$.

```
# preprocess input and output data
from sklearn import preprocessing

# normal Gaussian scaling for inputs
X_scaler = preprocessing.StandardScaler().fit(X_train)
X_scaled = X_scaler.transform(X_train)

print(X_scaler.mean_)
print(X_scaler.var_)

[ 9.0666667 161.6 77.8 ]
[ 28.46222222 3795.17333333 2861.62666667]
```

$$x' = \frac{x - \mu}{\sigma}$$

Sclaing output variables such that $y' \sim [0, 1]$.

```
# linear scaling up to [0,1] for outputs
y_scaler = preprocessing.MinMaxScaler().fit(y_train)
y_scaled = y_scaler.transform(y_train)

print(y_scaler.scale_) #
print(y_scaler.min_)
```

```
[0.00917431 0.06666667 0.03571429]
[-1.26605505 -2.06666667 -1.64285714]
```

$$y' = \left(\frac{1}{y_{max} - y_{min}} \right) y + \left(\frac{-y_{min}}{y_{max} - y_{min}} \right)$$

scale *min*

GD for a perceptron layer with direct gradients

Given a training dataset (\mathbf{X}, \mathbf{D})

Set learning parameter α

Initialize \mathbf{W} and \mathbf{b}

Repeat until convergence:

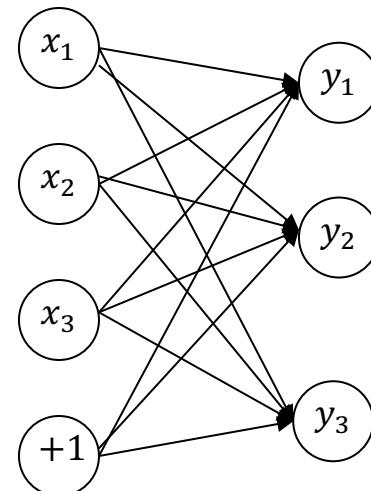
$$\mathbf{U} = \mathbf{X}\mathbf{W} + \mathbf{B}$$

$$\mathbf{Y} = f(\mathbf{U})$$

$$\nabla_{\mathbf{U}} J = -(\mathbf{D} - \mathbf{Y}) \cdot f' (\mathbf{U})$$

$$\mathbf{W} \leftarrow \mathbf{W} - \alpha \mathbf{X}^T \nabla_{\mathbf{U}} J$$

$$\mathbf{b} \leftarrow \mathbf{b} - \alpha (\nabla_{\mathbf{U}} J)^T \mathbf{1}_P$$



Implementing Perceptron Layer using pytorch libraries:

```
from torch import nn # torch nn API provides methods for build custom neural networks
```

nn.Module

- Base class for all neural network modules.
- Your models should also subclass this class.

nn.Linear

- applies a linear transformation $y = w^T x + b$

nn.Sigmoidal

- implements a perceptron layer

nn.Sequential

A sequential container. Modules will be added to it in the order they are passed in the constructor. The value a Sequential provides over manually calling a sequence of modules is that it allows treating the whole container as a single module

Implementing MSE loss and SGD optimizer

```
loss_fn = nn.MSELoss()
```

```
optimizer = torch.optim.SGD(model.parameters(), lr=lr)
```

GD for a perceptron layer with PyTorch nn.Module class

```
# Perceptron layer class
class PerceptronLayer(nn.Module):
    def __init__(self, no_inputs, no_outputs):
        super().__init__()
        self.perceptron_layer = nn.Sequential(
            nn.Linear(no_inputs, no_outputs),
            nn.Sigmoid()
        )

    def forward(self, x):
        logits = self.perceptron_layer(x)
        return logits

# create an instance of the layer
no_inputs, no_outputs = 3, 3
model = PerceptronLayer(no_inputs, no_outputs)
```

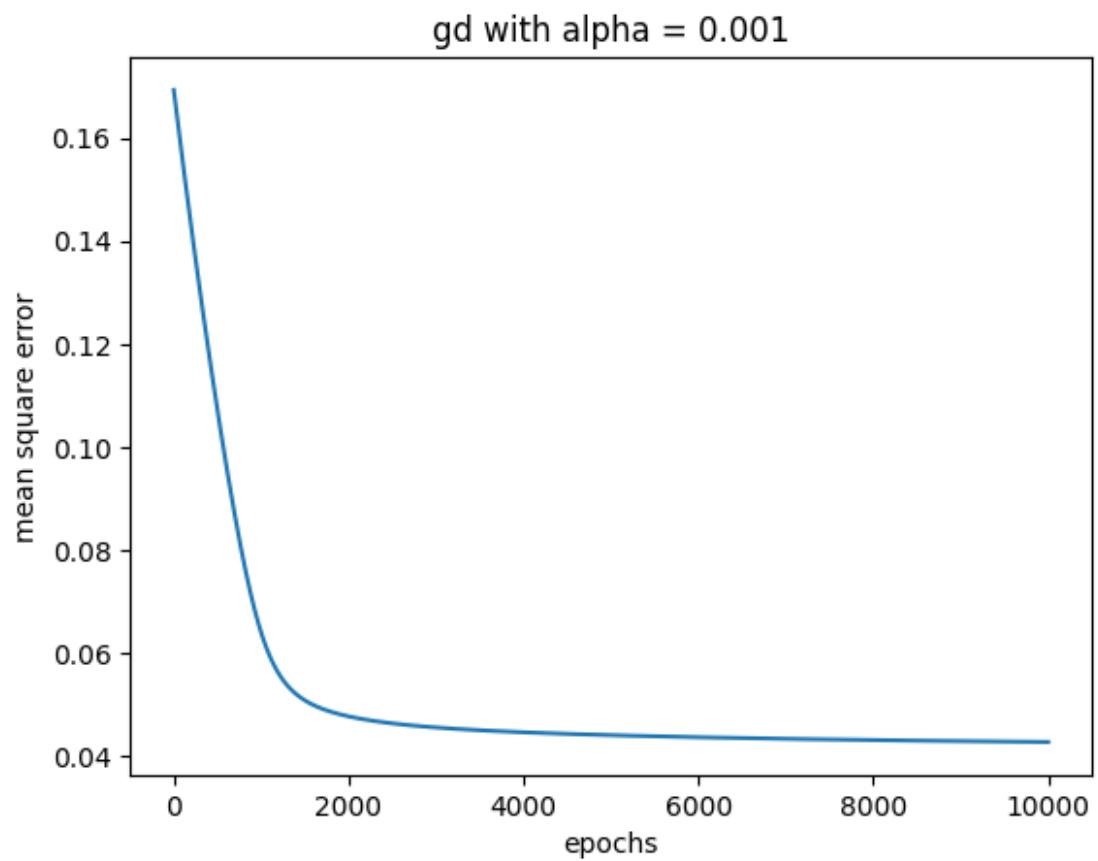
GD for a perceptron layer with PyTorch autograd

```
loss_fn = torch.nn.MSELoss()
optimizer = torch.optim.SGD(model.parameters(), lr=.001)

no_epochs, lr = 50000, 0.001
for epoch in range(no_epochs):
    # Compute prediction and loss
    pred = model(torch.tensor(X_scaled, dtype=torch.float))
    loss = loss_fn(pred, torch.tensor(y_scaled, dtype=torch.float))

    optimizer.zero_grad()
    loss.backward()
    optimizer.step()
```

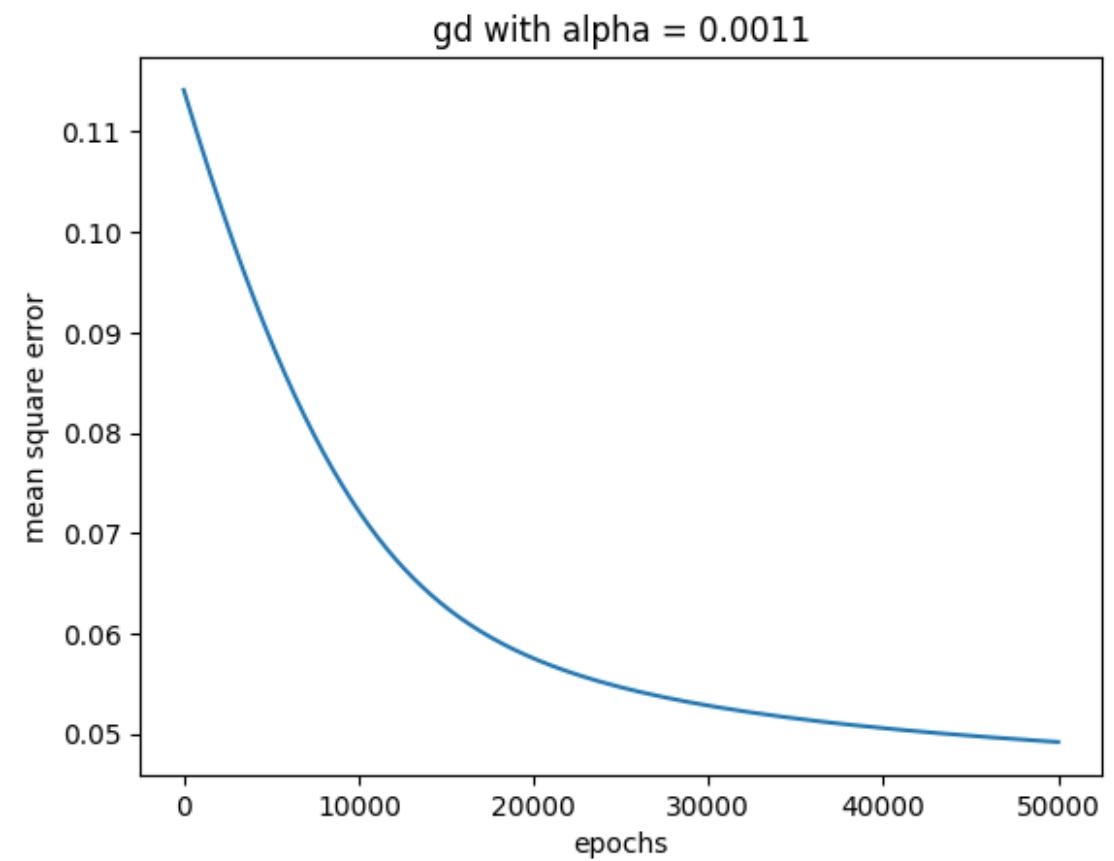
Direct gradients



Time for weight update = 0.074ms

Number of epochs = 10,000

Autograd



Time for weight update = 0.114ms

Number of epochs = 50,000

```
# scaling testing inputs
X_scaled = X_scaler.transform(X_test)

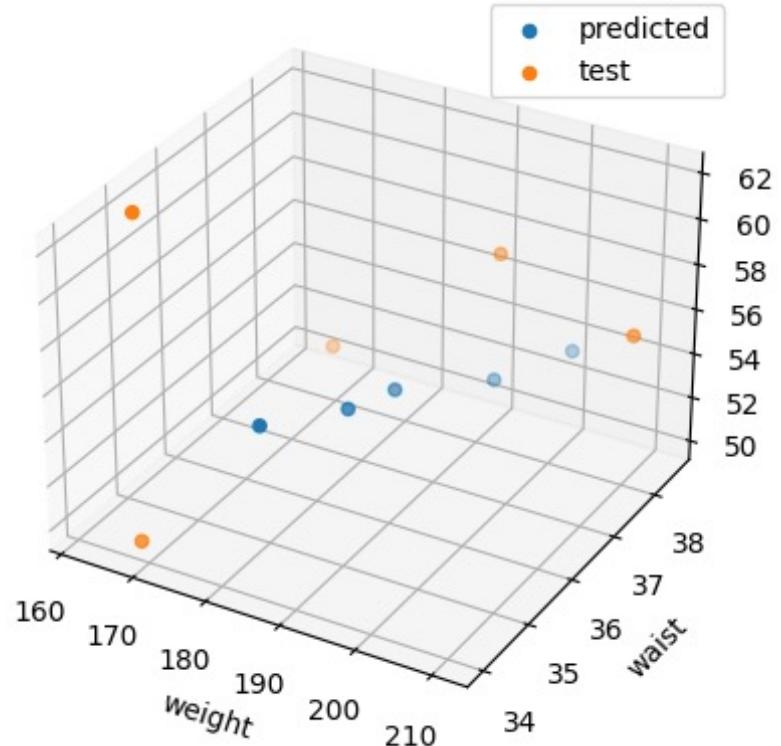
y_pred = model(X_scaled)

# scaling predicted outputs
y_scaled = y_scaler.inverse_transform(y_pred.detach().numpy())

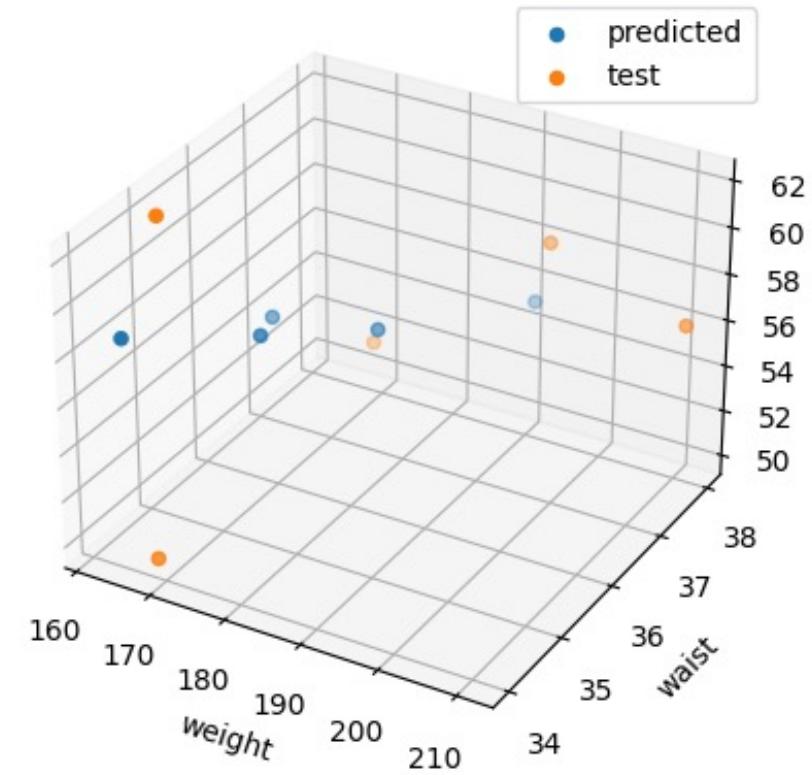
# computing mean square error
from sklearn.metrics import mean_squared_error
rms = mean_squared_error(y_scaled, y_test, squared=True, multioutput='raw_values')

# computing R2
from sklearn.metrics import r2_score
r2 = r2_score(y_scaled, y_test, multioutput='raw_values')
```

Direct gradients



Autograd



MSE = [407.06558087, 1.21179368, 22.40310961]
R² = [-15.16134924, 0.53783585, -11.22642492]

MSE = [356.95227178, 1.61016719, 21.11763412]
R² = [-3.18467906, 0.06741378, -11.90990362]

Exercise can best predict the waist!

R^2 : coefficient of determination

Let \hat{y}_i be the predicted value of data point y_i and the number of data points be n :

The **residual** sum of square error $SS_{res} = \sum_1^n (\hat{y}_i - y_i)^2$

If mean of data, $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$.

The **total** sum of squares (proportional to the variance), $SS_{tot} = \sum_1^n (\bar{y} - y_i)^2$

$$R^2 = 1 - \frac{SS_{res}}{SS_{tot}}$$

In the best case, the modelled value exactly match the observed value $SS_{res}=0$, then $R^2 = 1$.

The baseline model which predicts \bar{y} , then $R^2 = 0$.

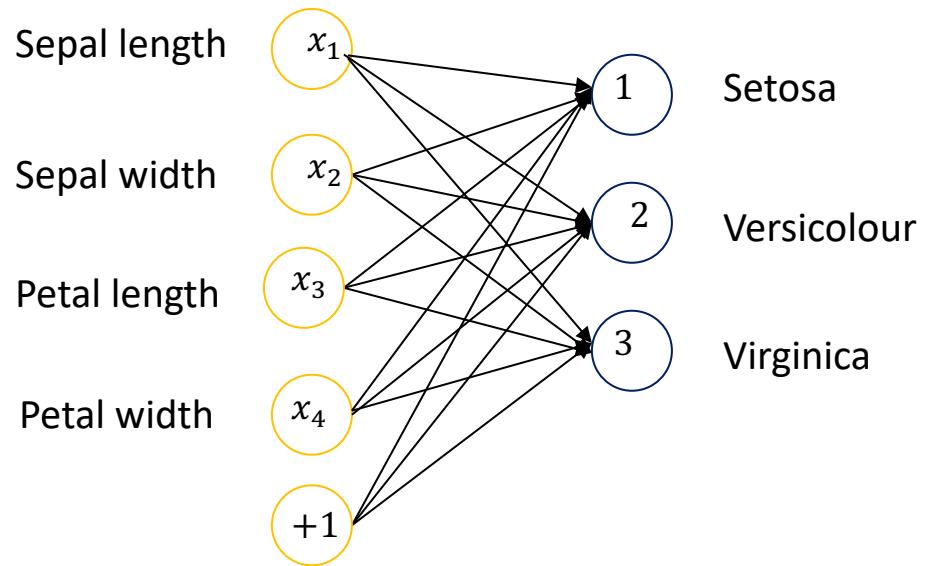
The models that have worse prediction than baseline models have negative R^2 .

2. Use mini-batch gradient decent learning to train a softmax layer to classify Iris dataset (<https://archive.ics.uci.edu/ml/datasets/Iris>). The dataset contains 150 data points. Use 90 data points for training the classifier and the remaining 60 data points for testing. Plot cross-entropies and classification accuracies against epochs for both train and test data. Set learning rate = 0.1, batch size = 16, and number of epochs = 1000.

You can use the following python commands to load Iris data:

```
from sklearn import datasets
iris = datasets.load_iris()
```

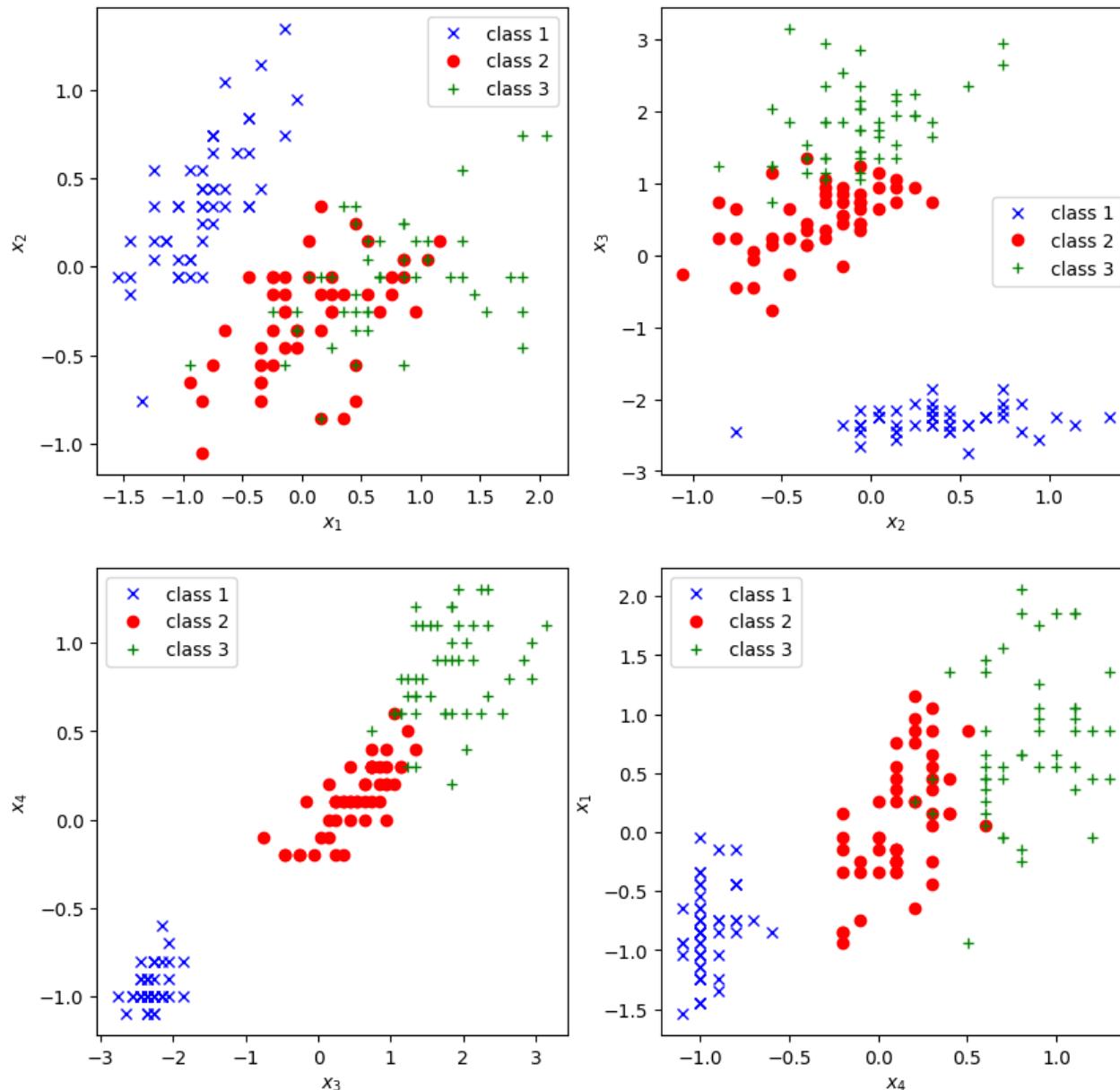
Repeat the classification with batch sizes = 2, 4, 8, 16, 24, 32, and 64, and compare the accuracies and the times taken for an epoch at different batch sizes.



Four features and three labels

90 data points for training and **60** data points for testing

boundaries in 2-dimensional feature spaces



```
# datasets import
from sklearn import datasets
from sklearn.model_selection import train_test_split

iris = datasets.load_iris()
iris.data -= np.mean(iris.data, axis=0) # mean correct data

# dataset split for train and test
x_train, x_test, y_train, y_test = train_test_split(iris.data, iris.target, test_size=0.2, random_state=2)
```

```
from torch import nn

class SoftmaxLayer(nn.Module):
    def __init__(self, no_inputs, no_outputs):
        super().__init__()
        self.softmax_layer = nn.Sequential(
            nn.Linear(no_inputs, no_outputs),  # applies a linear transformation  $y = w^T x + b$ 
            nn.Softmax(dim=1)      # implements softmax; sum up across rows to 1.0
        )

    def forward(self, x):
        logits = self.softmax_layer(x)
        return logits
```

Mini-batch gradient descent

Batch size = 16, learning factor = 0.1

Torch **Dataset** and **Dataloader** utils provide easy function for mini-batch learning

```
from torch.utils.data import Dataset
```

```
from torch.utils.data import DataLoader
```

We will be creating a subclass of **Dataset** class and using **Dataloaders** to implement
mini-batch gradient descent

```
# create a Dataset class
# A custom Dataset class must implement three functions: __init__, __len__, and __getitem__.
class MyDataset(Dataset):
    def __init__(self, X, y):
        self.X = torch.tensor(X, dtype=torch.float)
        self.y = torch.tensor(y)

    def __len__(self):
        return len(self.y)

    def __getitem__(self, idx):
        return self.X[idx], self.y[idx]

# create Dataset objects for train and test data
train_data = MyDataset(x_train, y_train)
test_data = MyDataset(x_test, y_test)

# create DataLoader objects for train and test data
train_dataloader = DataLoader(train_data, batch_size=batch_size, shuffle=True)
test_dataloader = DataLoader(test_data, batch_size=batch_size, shuffle=True)
```

```
def train_loop(dataloader, model, loss_fn, optimizer):
    size = len(dataloader.dataset)
    num_batches = len(dataloader)

    train_loss, correct = 0, 0
    for batch, (X, y) in enumerate(dataloader):

        # Compute prediction and loss
        pred = model(X)
        loss = loss_fn(pred, y)

        optimizer.zero_grad(). #initialize gradient calculations
        loss.backward()         # compute gradients
        optimizer.step()        # execute one step of SGD

        train_loss += loss.item()
        correct += (pred.argmax(1) == y).type(torch.float).sum().item()

    train_loss /= size
    correct /= size
    return train_loss, correct
```

```
def test_loop(dataloader, model, loss_fn):
    size = len(dataloader.dataset)
    num_batches = len(dataloader)
    test_loss, correct = 0, 0

    with torch.no_grad():
        for X, y in dataloader:
            pred = model(X)
            test_loss += loss_fn(pred, y).item()
            correct += (pred.argmax(dim=1) == y).type(torch.float).sum().item()

        test_loss /= size
        correct /= size

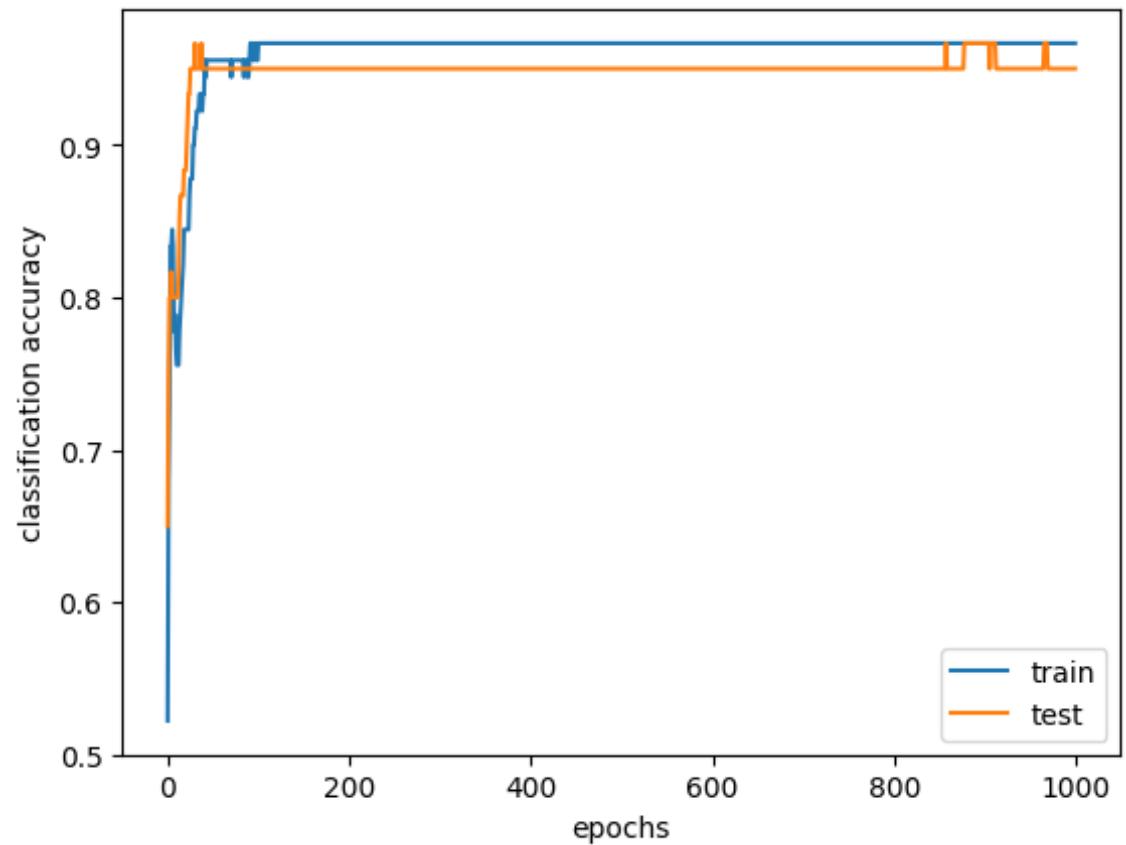
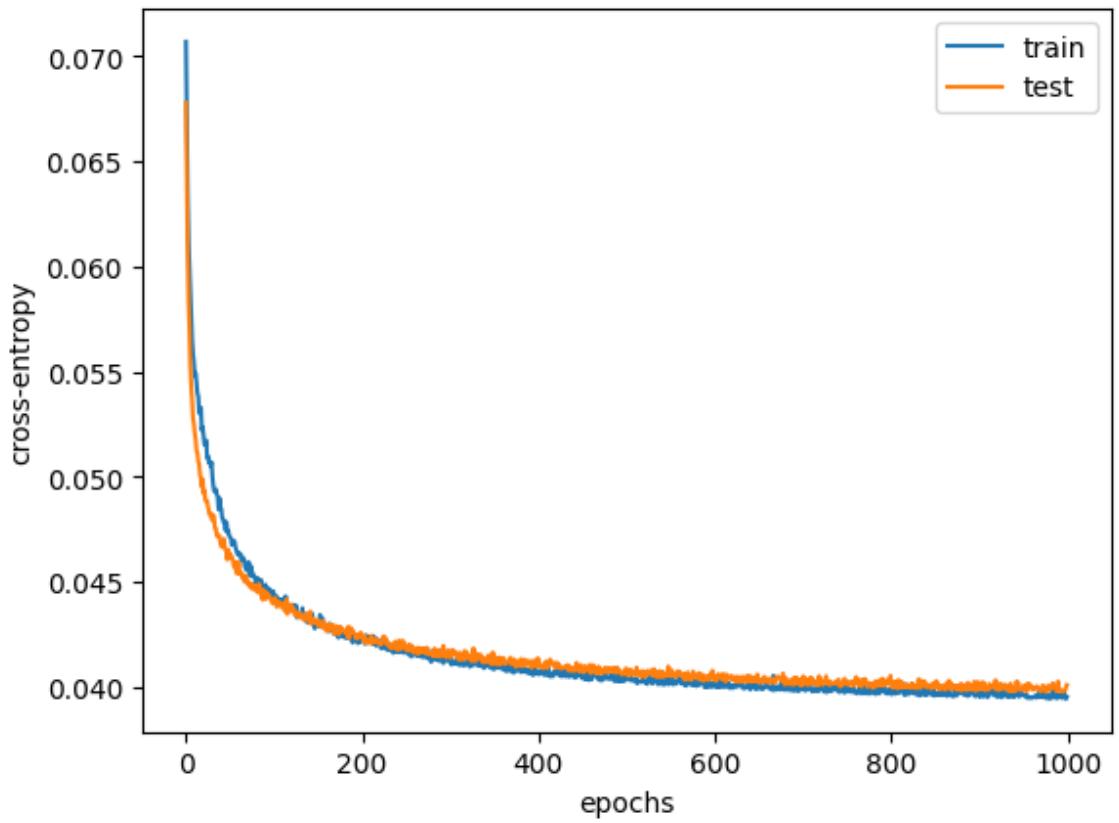
    return test_loss, correct
```

```
train_loss_, train_acc_, test_loss_, test_acc_ = [], [], [], []

for epoch in range(no_epochs):
    train_loss, train_acc = train_loop(train_dataloader, model, loss_fn, optimizer)
    test_loss, test_acc = test_loop(test_dataloader, model, loss_fn)

    train_loss_.append(train_loss), train_acc_.append(train_acc)
    test_loss_.append(test_loss), test_acc_.append(test_acc)
```

Learning curves at batch-size = 16



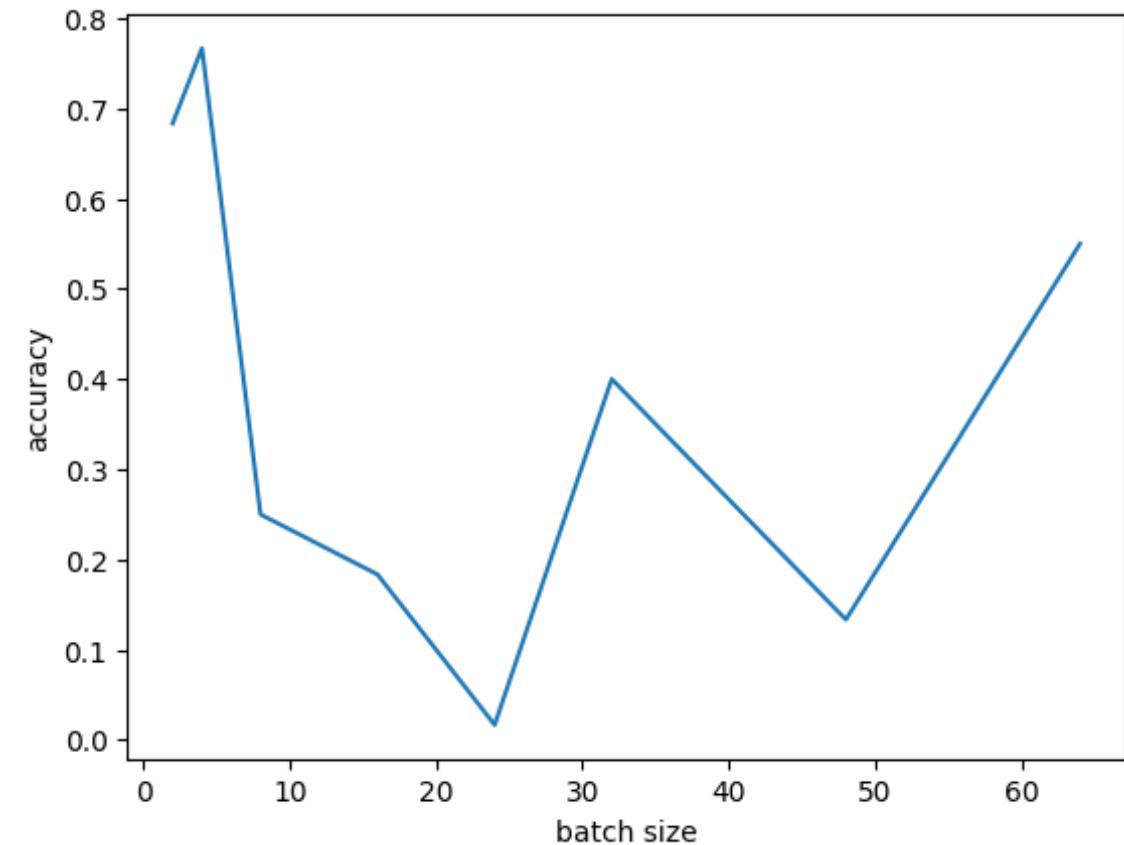
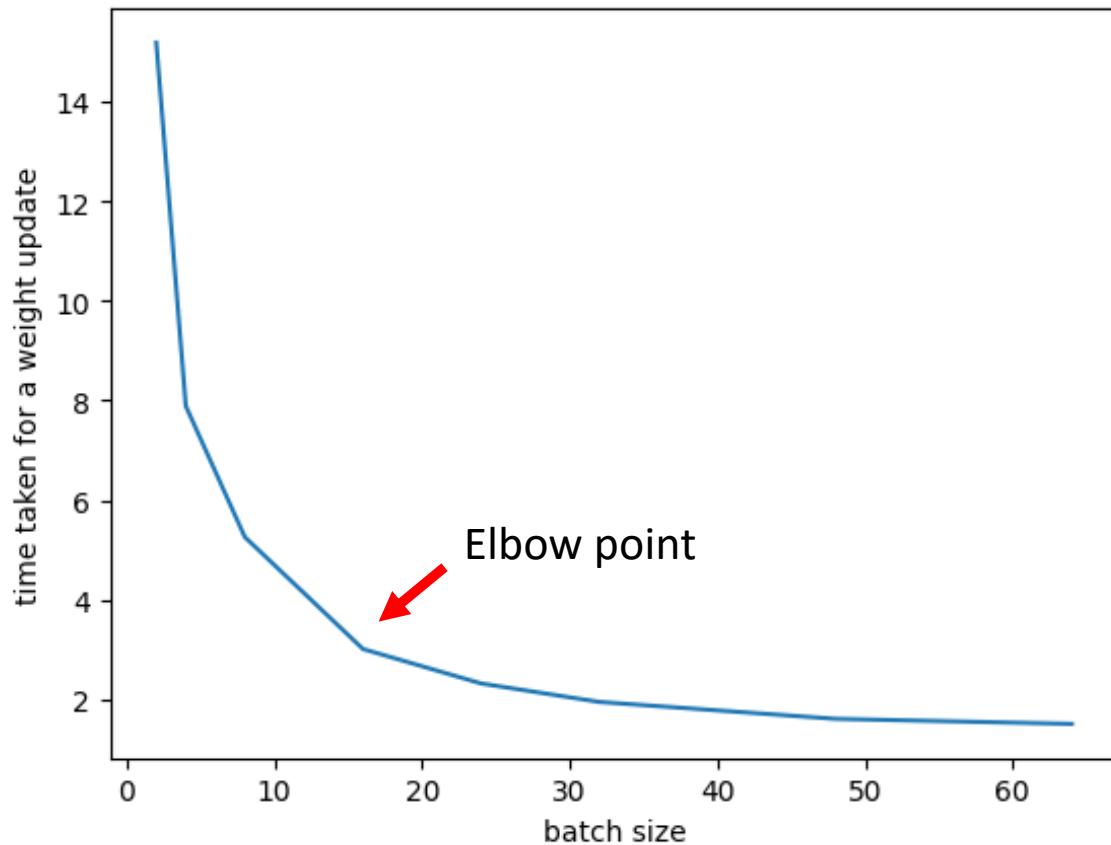
At convergence (1000 epochs):

```
weight = [[-0.875575 1.9825934 -4.016873 -1.6209128 ]  
          [ 0.70730793 -1.0039392 -0.2982792 -2.5409048 ]  
          [-0.3355393 -0.79086286 3.4620962 3.9241226 ]]
```

```
bias = [-0.6561739 3.990793 -2.9169736]
```

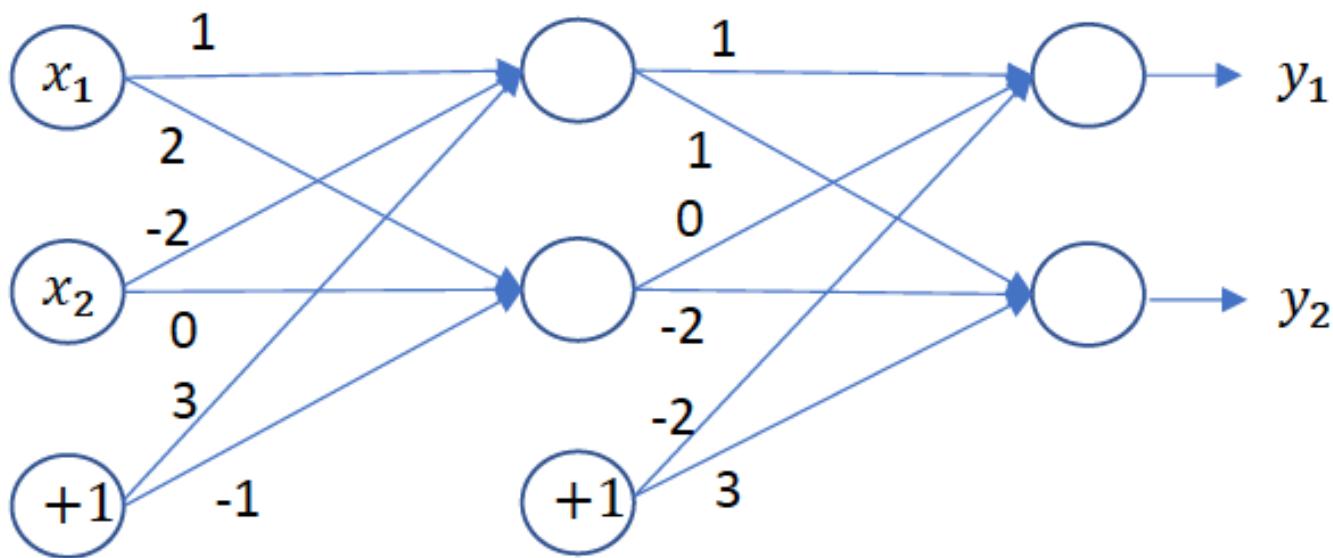
```
train_loss = 0.039542  
train_acc = 0.966667  
test_loss = 0.040103,  
test_acc = 0.950000
```

Accuracies and time-to-update weights against batch-size



Deep neural networks

SC4001 – Tutorial 4



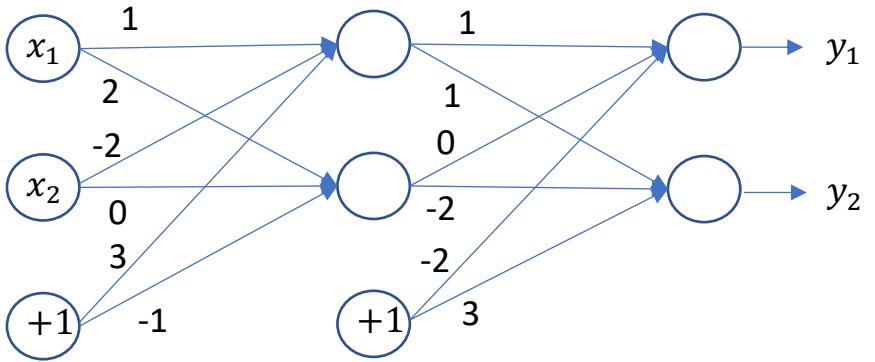
1. The two-layer feedforward perceptron network shown in figure 1 has weights and biases initialized as indicated and receives 2-dimensional inputs (x_1, x_2) . The network is to respond with $d_1 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$ and $d_2 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$ for input patterns $x_1 = \begin{pmatrix} 1.0 \\ 3.0 \end{pmatrix}$ and $x_2 = \begin{pmatrix} -2.0 \\ -2.0 \end{pmatrix}$, respectively.

Analyse a single feedforward and feedback step for gradient decent learning of the two patterns by doing the following:

- (a) Find the weight matrix W to the hidden-layer and weight matrix V to the output-layer, and the corresponding biases.
- (b) Calculate the synaptic input z and output h of the hidden-layer, and the synaptic input u and output $y = (y_1, y_2)$ of the output layer.
- (c) Find the mean square error cost J between the outputs and targets.
- (d) Calculate the gradients $\nabla_u J$ and $\nabla_z J$ at the output-layer and the hidden-layer, respectively.
- (e) Compute the new weights and biases.
- (f) Write a program to continue iterations until convergence and find the final weights and biases.

Assume a learning rate of 0.05.

Repeat above (a) – (f) for stochastic gradient decent learning.



Weight matrix to the hidden layer, $\mathbf{W} = \begin{pmatrix} 1.0 & 2.0 \\ -2.0 & 0.0 \end{pmatrix}$

Bias vector to the hidden-layer $\mathbf{b} = \begin{pmatrix} 3.0 \\ -1.0 \end{pmatrix}$

Weight matrix to the output-layer, $\mathbf{V} = \begin{pmatrix} 1.0 & 1.0 \\ 0.0 & -2.0 \end{pmatrix}$

Bias vector to the output-layer $\mathbf{c} = \begin{pmatrix} -2.0 \\ 3.0 \end{pmatrix}$

GD for 2-layer perceptron network:

Given a training dataset (\mathbf{X}, \mathbf{D})

Set learning parameter α

Initialize $\mathbf{W}, \mathbf{b}, \mathbf{V}, \mathbf{c}$

Repeat until convergence:

$$\mathbf{Z} = \mathbf{X}\mathbf{W} + \mathbf{B}$$

$$\mathbf{H} = g(\mathbf{Z})$$

$$\mathbf{U} = \mathbf{H}\mathbf{V} + \mathbf{C}$$

$$\mathbf{Y} = f(\mathbf{U})$$

Forward propagation
of activation

$$\nabla_{\mathbf{U}} J = -(\mathbf{D} - \mathbf{Y}) \cdot f'(\mathbf{U})$$

$$\nabla_{\mathbf{Z}} J = (\nabla_{\mathbf{U}} J) \mathbf{V}^T \cdot g'(\mathbf{Z})$$

Backward propagation
of gradients

$$\mathbf{V} \leftarrow \mathbf{V} - \alpha \mathbf{H}^T \nabla_{\mathbf{U}} J$$

$$\mathbf{c} \leftarrow \mathbf{c} - \alpha (\nabla_{\mathbf{U}} J)^T \mathbf{1}_P$$

$$\mathbf{W} \leftarrow \mathbf{W} - \alpha \mathbf{X}^T \nabla_{\mathbf{Z}} J$$

$$\mathbf{b} \leftarrow \mathbf{b} - \alpha (\nabla_{\mathbf{Z}} J)^T \mathbf{1}_P$$

Weight updating

$$\mathbf{x}_1 = \begin{pmatrix} 1.0 \\ 3.0 \end{pmatrix} \text{ and } \mathbf{d}_1 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$$

$$\mathbf{x}_2 = \begin{pmatrix} -2.0 \\ -2.0 \end{pmatrix} \text{ and } \mathbf{d}_2 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$

$$\mathbf{X} = \begin{pmatrix} 1.0 & 3.0 \\ -2.0 & -2.0 \end{pmatrix} \text{ and } \mathbf{D} = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$$

Forward propagation:

Synaptic input to hidden-layer, $\mathbf{Z} = \mathbf{XW} + \mathbf{B}$

$$\begin{aligned} &= \begin{pmatrix} 1.0 & 3.0 \\ -2.0 & -2.0 \end{pmatrix} \begin{pmatrix} 1.0 & 2.0 \\ -2.0 & 0.0 \end{pmatrix} + \begin{pmatrix} 3.0 & -1.0 \\ 3.0 & -1.0 \end{pmatrix} \\ &= \begin{pmatrix} -2.0 & 1.0 \\ 5.0 & -5.0 \end{pmatrix} \end{aligned}$$

Output of the hidden layer, $\mathbf{H} = g(\mathbf{Z}) = \frac{\mathbf{1}}{1+e^{-\mathbf{z}}} = \begin{pmatrix} 0.12 & 0.73 \\ 0.99 & 0.01 \end{pmatrix}$

Synaptic input to output-layer, $\mathbf{U} = \mathbf{HV} + \mathbf{C}$

$$\begin{aligned}&= \begin{pmatrix} 0.12 & 0.73 \\ 0.99 & 0.01 \end{pmatrix} \begin{pmatrix} 1.0 & 1.0 \\ 0.0 & -2.0 \end{pmatrix} + \begin{pmatrix} -2.0 & 3.0 \\ -2.0 & 3.0 \end{pmatrix} \\&= \begin{pmatrix} -1.88 & 1.66 \\ -0.99 & 3.98 \end{pmatrix}\end{aligned}$$

Output of the output layer, $\mathbf{Y} = f(\mathbf{U}) = \frac{1}{1+e^{-\mathbf{U}}} = \begin{pmatrix} 0.13 & 0.84 \\ 0.27 & 0.98 \end{pmatrix}$

$$\mathbf{D} = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$$

$$\begin{aligned}m.s.e. &= \frac{1}{2} \sum_{p=1}^2 \sum_{k=1}^2 (d_{pk} - y_{pk})^2 \\&= \frac{1}{2} \left(((0 - 0.13)^2 + (1 - 0.84)^2) + ((1 - 0.27)^2 + (0 - 0.98)^2) \right) \\&= 0.77\end{aligned}$$

Computing gradients (backward propagation):

$$f'(\mathbf{U}) = \mathbf{Y} \cdot (\mathbf{1} - \mathbf{Y}) = \begin{pmatrix} 0.13 & 0.84 \\ 0.27 & 0.98 \end{pmatrix} \cdot \left(\begin{pmatrix} 1.0 & 1.0 \\ 1.0 & 1.0 \end{pmatrix} - \begin{pmatrix} 0.13 & 0.84 \\ 0.27 & 0.98 \end{pmatrix} \right) = \begin{pmatrix} 0.11 & 0.13 \\ 0.20 & 0.02 \end{pmatrix}$$

$$\nabla_{\mathbf{U}} J = -(\mathbf{D} - \mathbf{Y}) \cdot f'(\mathbf{U}) = - \left(\begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} - \begin{pmatrix} 0.13 & 0.84 \\ 0.27 & 0.98 \end{pmatrix} \right) \begin{pmatrix} 0.12 & 0.13 \\ 0.20 & 0.02 \end{pmatrix} = \begin{pmatrix} 0.02 & -0.02 \\ -0.14 & 0.02 \end{pmatrix}$$

$$g'(\mathbf{Z}) = \mathbf{H} \cdot (1 - \mathbf{H}) = \begin{pmatrix} 0.12 & 0.73 \\ 0.99 & 0.01 \end{pmatrix} \cdot \left(\begin{pmatrix} 1.0 & 1.0 \\ 1.0 & 1.0 \end{pmatrix} - \begin{pmatrix} 0.12 & 0.73 \\ 0.99 & 0.01 \end{pmatrix} \right) = \begin{pmatrix} 0.10 & 0.2 \\ 0.01 & 0.01 \end{pmatrix}$$

$$\nabla_{\mathbf{Z}} J = (\nabla_{\mathbf{U}} J) \mathbf{V}^T \cdot g'(\mathbf{Z}) = \begin{pmatrix} 0.02 & -0.02 \\ -0.14 & 0.02 \end{pmatrix} \begin{pmatrix} 1.0 & 0.0 \\ 1.0 & -2.0 \end{pmatrix} \cdot \begin{pmatrix} 0.10 & 0.2 \\ 0.01 & 0.01 \end{pmatrix} = \begin{pmatrix} -0.001 & 0.01 \\ -0.001 & 0.00 \end{pmatrix}$$

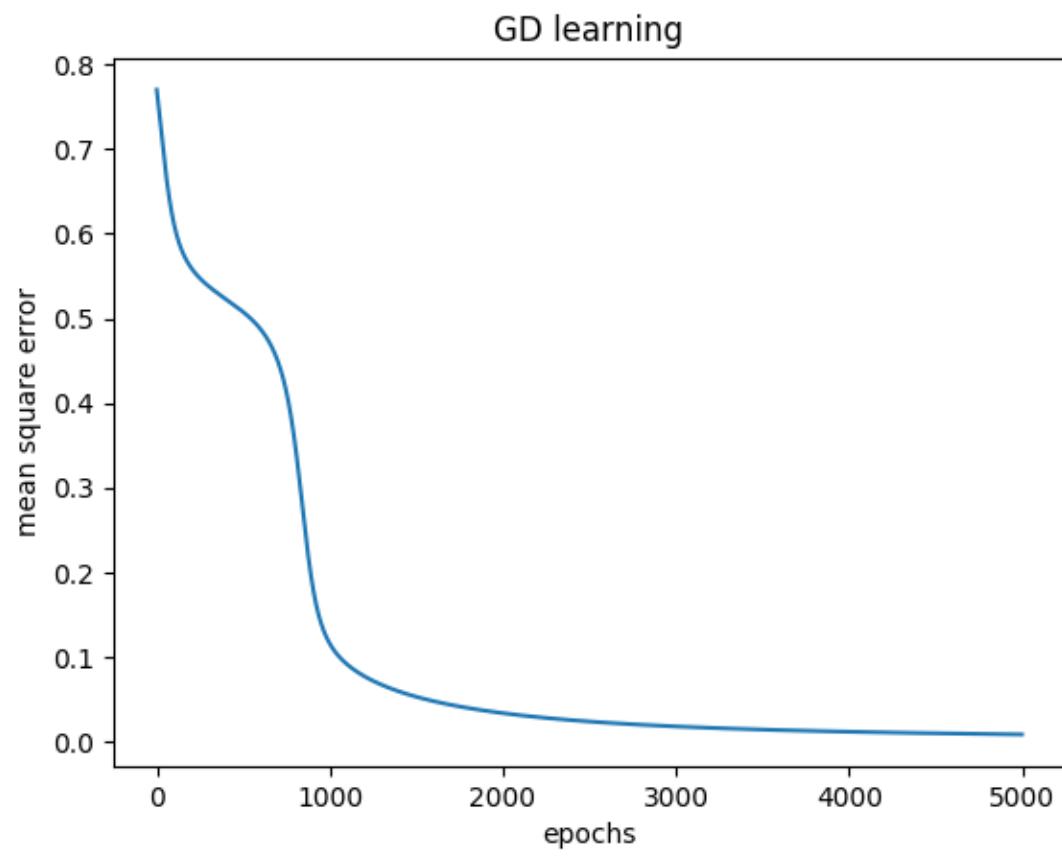
Updating weights:

$$\mathbf{V} \leftarrow \mathbf{V} - \alpha \mathbf{H}^T \nabla_{\mathbf{U}} J = \begin{pmatrix} 1.01 & 1.0 \\ 0.0 & -2.0 \end{pmatrix}$$

$$\mathbf{c} \leftarrow \mathbf{c} - \alpha (\nabla_{\mathbf{U}} J)^T \mathbf{1}_P = \begin{pmatrix} -1.99 \\ 3.00 \end{pmatrix}$$

$$\mathbf{W} \leftarrow \mathbf{W} - \alpha \mathbf{X}^T \nabla_{\mathbf{Z}} J = \begin{pmatrix} 1.0 & 2.0 \\ -2.0 & 0.0 \end{pmatrix}$$

$$\mathbf{b} \leftarrow \mathbf{b} - \alpha (\nabla_{\mathbf{Z}} J)^T \mathbf{1}_P = \begin{pmatrix} 3.0 \\ -1.0 \end{pmatrix}$$



At convergence:

$$\mathbf{W} = \begin{pmatrix} 0.63 & 0.60 \\ -3.0 & -2.0 \end{pmatrix}, \mathbf{b} = \begin{pmatrix} 2.72 \\ -0.74 \end{pmatrix}$$

$$\mathbf{V} = \begin{pmatrix} 4.97 & -3.46 \\ 0.25 & -2.37 \end{pmatrix}, \mathbf{c} = \begin{pmatrix} -2.42 \\ 2.56 \end{pmatrix}$$

Predicted values:

$$\mathbf{y}_1 = \begin{pmatrix} 0.08 \\ 0.93 \end{pmatrix} \text{ and } \mathbf{y}_2 = \begin{pmatrix} 0.94 \\ 0.05 \end{pmatrix}$$

$$\text{m.s.e.} = 0.004$$

SGD learning for 2-layer perceptron network:

Given a training dataset $\{(x, d)\}$

Set learning parameter α

Initialize W, b, V, c

Repeat until convergence:

For every pattern (x, d) :

$$z = W^T x + b$$

$$h = g(z)$$

$$u = V^T h + c$$

$$y = f(u)$$

Forward propagation
of activation

$$\nabla_u J = -(d - y) \cdot f'(z)$$

$$\nabla_z J = V \nabla_u J \cdot g'(z)$$

Backward propagation
of gradients

$$V \leftarrow V - \alpha h (\nabla_u J)^T$$

$$c \leftarrow c - \alpha \nabla_u J$$

$$W \leftarrow W - \alpha x (\nabla_z J)^T$$

$$b \leftarrow b - \alpha \nabla_z J$$

Weight updating

Epoch 1:

Apply **first** pattern $\mathbf{x} = \begin{pmatrix} 1.0 \\ 3.0 \end{pmatrix}$ and $\mathbf{d} = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$:

Synaptic input to the hidden-layer

$$\mathbf{z} = \mathbf{W}^T \mathbf{x} + \mathbf{b} = \begin{pmatrix} 1.0 & -2.0 \\ 2.0 & 0.0 \end{pmatrix} \begin{pmatrix} 1.0 \\ 3.0 \end{pmatrix} + \begin{pmatrix} 3.0 \\ -1.0 \end{pmatrix} = \begin{pmatrix} -2.0 \\ 1.0 \end{pmatrix}$$

Output of the hidden-layer $\mathbf{h} = g(\mathbf{z}) = \frac{1}{1+e^{-\mathbf{z}}} = \begin{pmatrix} 0.12 \\ 0.73 \end{pmatrix}$

Synaptic input to output-layer

$$\mathbf{u} = \mathbf{V}^T \mathbf{h} + \mathbf{c} = \begin{pmatrix} -1.88 \\ 1.66 \end{pmatrix}$$

Output of the output-layer $\mathbf{y} = f(\mathbf{u}) = \frac{1}{1+e^{-\mathbf{u}}} = \begin{pmatrix} 0.13 \\ 0.84 \end{pmatrix}$

$$s.e. = (d_1 - y_1)^2 + (d_2 - y_2)^2 = 0.043$$

Computing gradients:

$$f'(\mathbf{u}) = \mathbf{y} \cdot (1 - \mathbf{y}) = \begin{pmatrix} 0.13 \\ 0.84 \end{pmatrix} \cdot \left(\begin{pmatrix} 1.0 \\ 1.0 \end{pmatrix} - \begin{pmatrix} 0.13 \\ 0.84 \end{pmatrix} \right) = \begin{pmatrix} 0.11 \\ 0.13 \end{pmatrix}$$

$$\nabla_{\mathbf{u}} J = -(\mathbf{d} - \mathbf{y}) \cdot f'(\mathbf{u}) = -\left(\begin{pmatrix} 0 \\ 1 \end{pmatrix} - \begin{pmatrix} 0.13 \\ 0.84 \end{pmatrix} \right) \cdot \begin{pmatrix} 0.12 \\ 0.14 \end{pmatrix} = \begin{pmatrix} 0.02 \\ -0.02 \end{pmatrix}$$

$$g'(\mathbf{z}) = \mathbf{h} \cdot (1 - \mathbf{h}) = \begin{pmatrix} 0.12 \\ 0.73 \end{pmatrix} \cdot \left(\begin{pmatrix} 1.0 \\ 1.0 \end{pmatrix} - \begin{pmatrix} 0.12 \\ 0.73 \end{pmatrix} \right) = \begin{pmatrix} 0.10 \\ 0.20 \end{pmatrix}$$

$$\nabla_{\mathbf{z}} J = \mathbf{V} \nabla_{\mathbf{u}} J \cdot g'(\mathbf{z}) = \begin{pmatrix} 1.0 & 1.0 \\ 0.0 & -2.0 \end{pmatrix} \begin{pmatrix} 0.02 \\ -0.02 \end{pmatrix} \cdot \begin{pmatrix} 0.11 \\ 0.20 \end{pmatrix} = \begin{pmatrix} -0.001 \\ 0.008 \end{pmatrix}$$

Updating weights:

$$\mathbf{V} \leftarrow \mathbf{V} - \alpha \mathbf{h} (\nabla_{\mathbf{u}} J)^T = \begin{pmatrix} 1.0 & 1.0 \\ 0.0 & -2.0 \end{pmatrix} - 0.2 \begin{pmatrix} 0.12 \\ 0.73 \end{pmatrix} (-0.02 \quad 0.022) = \begin{pmatrix} 1.0 & 1.0001 \\ 0.00 & -2.0 \end{pmatrix}$$

$$\mathbf{c} \leftarrow \mathbf{c} - \alpha \nabla_{\mathbf{u}} J = \begin{pmatrix} -2.0 \\ 3.0 \end{pmatrix} + 0.2 \begin{pmatrix} 0.02 \\ -0.02 \end{pmatrix} = \begin{pmatrix} -2.00 \\ 3.001 \end{pmatrix}$$

$$\mathbf{W} \leftarrow \mathbf{W} - \alpha \mathbf{x} (\nabla_{\mathbf{z}} J)^T = \begin{pmatrix} 1.0 & 2.0 \\ -2.00 & -0.001 \end{pmatrix}$$

$$\mathbf{b} \leftarrow \mathbf{b} - \alpha \nabla_{\mathbf{z}} J = \begin{pmatrix} 3.00 \\ -1.00 \end{pmatrix}$$

Apply **second** pattern $x = \begin{pmatrix} -2.0 \\ -2.0 \end{pmatrix}$ and $d = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$:

Synaptic input to the hidden-layer

$$\mathbf{z} = \mathbf{W}^T \mathbf{x} + \mathbf{b} = \begin{pmatrix} 5.0 \\ -5.0 \end{pmatrix}$$

Output of the hidden-layer $\mathbf{h} = g(\mathbf{z}) = \frac{1}{1+e^{-\mathbf{z}}} = \begin{pmatrix} 1.0 \\ 0.007 \end{pmatrix}$

Synaptic input to output-layer

$$\mathbf{u} = \mathbf{V}^T \mathbf{h} + \mathbf{c} = \begin{pmatrix} -0.99 \\ 3.98 \end{pmatrix}$$

Output of the output-layer $\mathbf{y} = f(\mathbf{u}) = \frac{1}{1+e^{-\mathbf{u}}} = \begin{pmatrix} 0.27 \\ 0.98 \end{pmatrix}$

$$s.e. = (d_1 - y_1)^2 + (d_2 - y_2)^2 = 1.5$$

Computing gradients:

$$f'(\mathbf{u}) = \mathbf{y} \cdot (1 - \mathbf{y}) = \begin{pmatrix} 0.195 \\ 0.018 \end{pmatrix}$$

$$\nabla_{\mathbf{u}} J = -(\mathbf{d} - \mathbf{y}) \cdot f'(\mathbf{u}) = \begin{pmatrix} -0.14 \\ 0.018 \end{pmatrix}$$

$$g'(\mathbf{z}) = \mathbf{h} \cdot (1 - \mathbf{h}) = \begin{pmatrix} 0.007 \\ 0.007 \end{pmatrix}$$

$$\nabla_{\mathbf{z}} J = \mathbf{V} \nabla_{\mathbf{u}} J \cdot g'(\mathbf{z}) = \begin{pmatrix} -0.0008 \\ -0.0002 \end{pmatrix}$$

Updating weights:

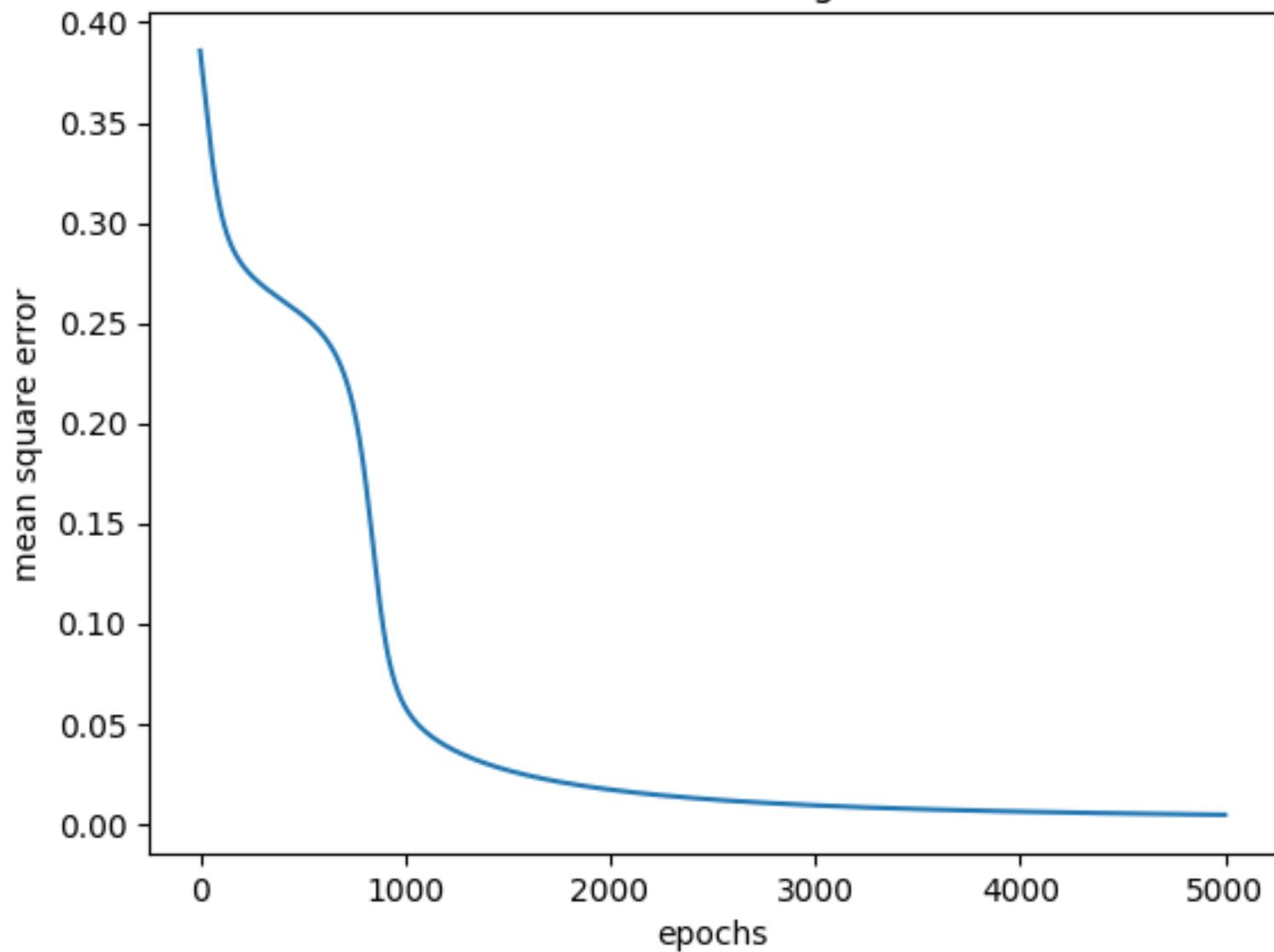
$$\mathbf{V} \leftarrow \mathbf{V} - \alpha \mathbf{h} (\nabla_{\mathbf{u}} J)^T = \begin{pmatrix} 1.007 & 0.99 \\ 0.0 & -2.0 \end{pmatrix}$$

$$\mathbf{c} \leftarrow \mathbf{c} - \alpha \nabla_{\mathbf{u}} J = \begin{pmatrix} -1.99 \\ 3.0 \end{pmatrix}$$

$$\mathbf{W} \leftarrow \mathbf{W} - \alpha \mathbf{x} (\nabla_{\mathbf{z}} J)^T = \begin{pmatrix} 0.999 & 1.99 \\ -1.99 & 0.00 \end{pmatrix}$$

$$\mathbf{b} \leftarrow \mathbf{b} - \alpha \nabla_{\mathbf{z}} J = \begin{pmatrix} 3.00 \\ -1.00 \end{pmatrix}$$

SGD learning



At convergence:

$$\mathbf{W} = \begin{pmatrix} 0.63 & 0.60 \\ -3.0 & -2.0 \end{pmatrix}, \mathbf{b} = \begin{pmatrix} 2.72 \\ -0.74 \end{pmatrix}$$

$$\mathbf{V} = \begin{pmatrix} 4.97 & -3.46 \\ 0.25 & -2.37 \end{pmatrix}, \mathbf{c} = \begin{pmatrix} -2.42 \\ 2.56 \end{pmatrix}$$

Predicted values:

$$\mathbf{y}_1 = \begin{pmatrix} 0.08 \\ 0.93 \end{pmatrix} \text{ and } \mathbf{y}_2 = \begin{pmatrix} 0.94 \\ 0.05 \end{pmatrix}$$

$$\text{m.s.e.} = 0.004$$

2. A feedforward neural network with one hidden layer to perform the following classification:

class	inputs
A	(1.0, 1.0), (0.0, 1.0)
B	(3.0, 4.0), (2.0, 2.0)
C	(2.0, -2.0), (-2.0, -3.0)

The network has a hidden layer consisting of three perceptrons and a softmax output layer.

Show one iteration of gradient descent learning and plot learning curves until convergence at a learning rate $\alpha = 0.1$.

Initialize the weights \mathbf{W} and biases \mathbf{b} to the hidden layer, and the weights \mathbf{V} and biases \mathbf{c} to the output layer as follows:

$$\mathbf{W} = \begin{pmatrix} -0.10 & 0.97 & 0.18 \\ -0.70 & 0.38 & 0.93 \end{pmatrix}, \mathbf{b} = \begin{pmatrix} 0.0 \\ 0.0 \\ 0.0 \end{pmatrix}$$

$$\mathbf{V} = \begin{pmatrix} 1.01 & 0.09 & -0.39 \\ 0.79 & -0.45 & -0.22 \\ 0.28 & 0.96 & -0.07 \end{pmatrix}, \mathbf{c} = \begin{pmatrix} 0.0 \\ 0.0 \\ 0.0 \end{pmatrix}$$

Determine the weights and biases at convergence.

Find the class labels predicted by the trained network for patterns:

$$\mathbf{x}_1 = \begin{pmatrix} 2.5 \\ 1.5 \end{pmatrix} \text{ and } \mathbf{x}_2 = \begin{pmatrix} -1.5 \\ 0.5 \end{pmatrix}$$

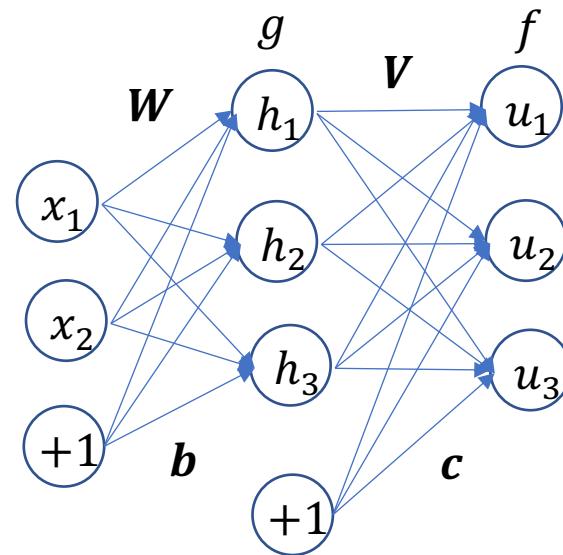
Training examples (patterns):

class	<i>inputs</i>	<i>Target Label</i>
A	(1.0, 1.0), (0.0, 1.0)	1
B	(3.0, 4.0), (2.0, 2.0)	2
C	(2.0, -2.0), (-2.0, -3.0)	3

Feedforward network :

Perceptron hidden layer with 3 neurons

Softmax output layer with 3 neurons



$$g(\mathbf{z}) = \frac{1}{1 + e^{-\mathbf{z}}}$$

$$f(\mathbf{U}) = \frac{e^{\mathbf{U}}}{\sum_{k=1}^K e^{\mathbf{U}_k}}$$

Training inputs and targets

$$X = \begin{pmatrix} 1.0 & 1.0 \\ 0.0 & 1.0 \\ 3.0 & 4.0 \\ 2.0 & 2.0 \\ 2.0 & -2.0 \\ -2.0 & -3.0 \end{pmatrix} \text{ and } D = \begin{pmatrix} 1 \\ 1 \\ 2 \\ 2 \\ 3 \\ 3 \end{pmatrix}$$

Targets as a one hot matrix:

$$K = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \end{pmatrix}$$

Initial weights and biases:

To the hidden layer,

$$\mathbf{W} = \begin{pmatrix} -0.10 & 0.97 & 0.18 \\ -0.70 & 0.38 & 0.93 \end{pmatrix}, \mathbf{b} = \begin{pmatrix} 0.0 \\ 0.0 \\ 0.0 \end{pmatrix}$$

To the output-layer

$$\mathbf{V} = \begin{pmatrix} 1.01 & 0.09 & -0.39 \\ 0.79 & -0.45 & -0.22 \\ 0.28 & 0.96 & -0.07 \end{pmatrix}, \mathbf{c} = \begin{pmatrix} 0.0 \\ 0.0 \\ 0.0 \end{pmatrix}$$

Learning factor $\alpha = 0.1$

Activation functions:

Hidden layer is a continuous perceptron layer: $g(\mathbf{Z}) = \frac{1}{1+e^{-\mathbf{z}}}$

Output layer is a softmax layer: $f(\mathbf{U}) = \frac{e^{\mathbf{U}}}{\sum_{k=1}^K e^{U_k}}$

GD for the feedforward network

Given a training dataset (\mathbf{X}, \mathbf{D})

Set learning parameter α

Initialize $\mathbf{W}, \mathbf{b}, \mathbf{V}, \mathbf{c}$

Repeat until convergence:

$$\mathbf{Z} = \mathbf{X}\mathbf{W} + \mathbf{B}$$

$$\mathbf{H} = g(\mathbf{Z})$$

$$\mathbf{U} = \mathbf{H}\mathbf{V} + \mathbf{C}$$

$$\mathbf{Y} = \arg \max_k f(\mathbf{U})$$

Forward propagation
of activation

$$\nabla_{\mathbf{U}} J = -(\mathbf{K} - f(\mathbf{U}))$$

$$\nabla_{\mathbf{Z}} J = (\nabla_{\mathbf{U}} J) \mathbf{V}^T \cdot g'(\mathbf{Z})$$

Backward propagation
of gradients

$$\mathbf{V} \leftarrow \mathbf{V} - \alpha \mathbf{H}^T \nabla_{\mathbf{U}} J$$

$$\mathbf{c} \leftarrow \mathbf{c} - \alpha (\nabla_{\mathbf{U}} J)^T \mathbf{1}_P$$

$$\mathbf{W} \leftarrow \mathbf{W} - \alpha \mathbf{X}^T \nabla_{\mathbf{Z}} J$$

$$\mathbf{b} \leftarrow \mathbf{b} - \alpha (\nabla_{\mathbf{Z}} J)^T \mathbf{1}_P$$

Weight updating

Synaptic input to hidden-layer,

$$\mathbf{Z} = \mathbf{XW} + \mathbf{B} = \begin{pmatrix} 1.0 & 1.0 \\ 0.0 & 1.0 \\ 3.0 & 4.0 \\ 2.0 & 2.0 \\ 2.0 & -2.0 \\ -2.0 & -3.0 \end{pmatrix} \begin{pmatrix} -0.10 & 0.97 & 0.18 \\ -0.70 & 0.38 & 0.93 \end{pmatrix} + \begin{pmatrix} 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \end{pmatrix} = \begin{pmatrix} -0.80 & 1.35 & 1.10 \\ -0.70 & 0.38 & 0.93 \\ -3.08 & 4.44 & 4.23 \\ -1.59 & 2.70 & 2.21 \\ 1.20 & 1.18 & -1.50 \\ 2.29 & -3.08 & -3.13 \end{pmatrix}$$

Output of the hidden layer, $\mathbf{H} = g(\mathbf{Z}) = \frac{1}{1+e^{-\mathbf{z}}} =$

$$\begin{pmatrix} 0.31 & 0.79 & 0.75 \\ 0.33 & 0.59 & 0.72 \\ 0.04 & 0.99 & 0.99 \\ 0.17 & 0.94 & 0.90 \\ 0.77 & 0.77 & 0.18 \\ 0.91 & 0.04 & 0.04 \end{pmatrix}$$

Synaptic input to output-layer,

$$\mathbf{U} = \mathbf{H}\mathbf{V} + \mathbf{C} = \begin{pmatrix} 0.31 & 0.79 & 0.75 \\ 0.33 & 0.59 & 0.72 \\ 0.04 & 0.99 & 0.99 \\ 0.17 & 0.94 & 0.90 \\ 0.77 & 0.77 & 0.18 \\ 0.91 & 0.04 & 0.04 \end{pmatrix} \begin{pmatrix} 1.01 & 0.09 & -0.39 \\ 0.79 & -0.45 & -0.22 \\ 0.28 & 0.96 & -0.07 \end{pmatrix} + \begin{pmatrix} 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \end{pmatrix} = \begin{pmatrix} 1.15 & 0.40 & -0.34 \\ 1.01 & 0.46 & -0.31 \\ 1.10 & 0.51 & -0.30 \\ 1.16 & 0.47 & -0.33 \\ 1.43 & -0.09 & -0.48 \\ 0.96 & 0.11 & -0.36 \end{pmatrix}$$

Output layer activation $f(\mathbf{U}) = \frac{e^{\mathbf{U}}}{\sum_{k=1}^K e^{U_k}} =$

$$\begin{pmatrix} 0.59 & 0.28 & 0.13 \\ 0.54 & 0.31 & 0.15 \\ 0.56 & 0.31 & 0.14 \\ 0.58 & 0.29 & 0.13 \\ 0.73 & 0.16 & 0.11 \\ 0.59 & 0.25 & 0.16 \end{pmatrix}$$

Output $\mathbf{Y} = \arg \max_k f(\mathbf{U}) =$

$$\begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{pmatrix}$$

$$\mathbf{D} = \begin{pmatrix} 1 \\ 1 \\ 2 \\ 2 \\ 3 \\ 3 \end{pmatrix}, \quad \mathbf{Y} = \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{pmatrix}, \quad f(\mathbf{U}) = \begin{pmatrix} 0.59 & 0.28 & 0.13 \\ 0.54 & 0.31 & 0.15 \\ 0.56 & 0.31 & 0.14 \\ 0.58 & 0.29 & 0.13 \\ 0.73 & 0.16 & 0.11 \\ 0.59 & 0.25 & 0.16 \end{pmatrix}$$

$$\text{Classification error} = \sum 1(\mathbf{D} \neq \mathbf{Y}) = 4$$

$$\begin{aligned} \text{Entropy } J &= -\sum_{p=1}^P \log(f(u_{pd_p})) \\ &= -(log(0.59) + log(0.54) + log(0.31) + log(0.29) + log(0.11) + log(0.16)) \\ &= 7.63 \end{aligned}$$

$$\nabla_{\mathbf{U}} J = -(\mathbf{K} - f(\mathbf{U})) = - \left(\begin{pmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \end{pmatrix} - \begin{pmatrix} 0.59 & 0.28 & 0.13 \\ 0.54 & 0.31 & 0.15 \\ 0.56 & 0.31 & 0.14 \\ 0.58 & 0.29 & 0.13 \\ 0.73 & 0.16 & 0.11 \\ 0.59 & 0.25 & 0.16 \end{pmatrix} \right) = \begin{pmatrix} -0.41 & 0.28 & 0.13 \\ -0.46 & 0.31 & 0.15 \\ 0.56 & -0.69 & 0.14 \\ 0.58 & -0.71 & 0.13 \\ 0.73 & 0.16 & -0.89 \\ 0.59 & 0.25 & -0.84 \end{pmatrix}$$

$$g'(\mathbf{Z}) = \mathbf{H} \cdot (\mathbf{1} - \mathbf{H}) = \begin{pmatrix} 0.21 & 0.16 & 0.19 \\ 0.22 & 0.24 & 0.20 \\ 0.04 & 0.01 & 0.01 \\ 0.14 & 0.06 & 0.09 \\ 0.18 & 0.18 & 0.15 \\ 0.08 & 0.04 & 0.04 \end{pmatrix}$$

$$\nabla_{\mathbf{Z}} J = (\nabla_{\mathbf{U}} J) \mathbf{V}^T \cdot g'(\mathbf{Z}) = \begin{pmatrix} -0.41 & 0.28 & 0.13 \\ -0.46 & 0.31 & 0.15 \\ 0.56 & -0.69 & 0.14 \\ 0.58 & -0.71 & 0.13 \\ 0.73 & 0.16 & -0.89 \\ 0.59 & 0.25 & -0.84 \end{pmatrix} \begin{pmatrix} 1.01 & 0.09 & -0.39 \\ 0.79 & -0.45 & -0.22 \\ 0.28 & 0.96 & -0.07 \end{pmatrix}^T \cdot \begin{pmatrix} 0.21 & 0.16 & 0.19 \\ 0.22 & 0.24 & 0.20 \\ 0.04 & 0.01 & 0.01 \\ 0.14 & 0.06 & 0.09 \\ 0.18 & 0.18 & 0.15 \\ 0.08 & 0.04 & 0.04 \end{pmatrix} = \begin{pmatrix} -0.09 & -0.08 & 0.03 \\ -0.11 & -0.13 & 0.03 \\ 0.02 & 0.01 & -0.01 \\ 0.07 & 0.04 & -0.05 \\ 0.20 & 0.13 & 0.06 \\ 0.08 & 0.02 & 0.02 \end{pmatrix}$$

Learning rate $\alpha = 0.1$

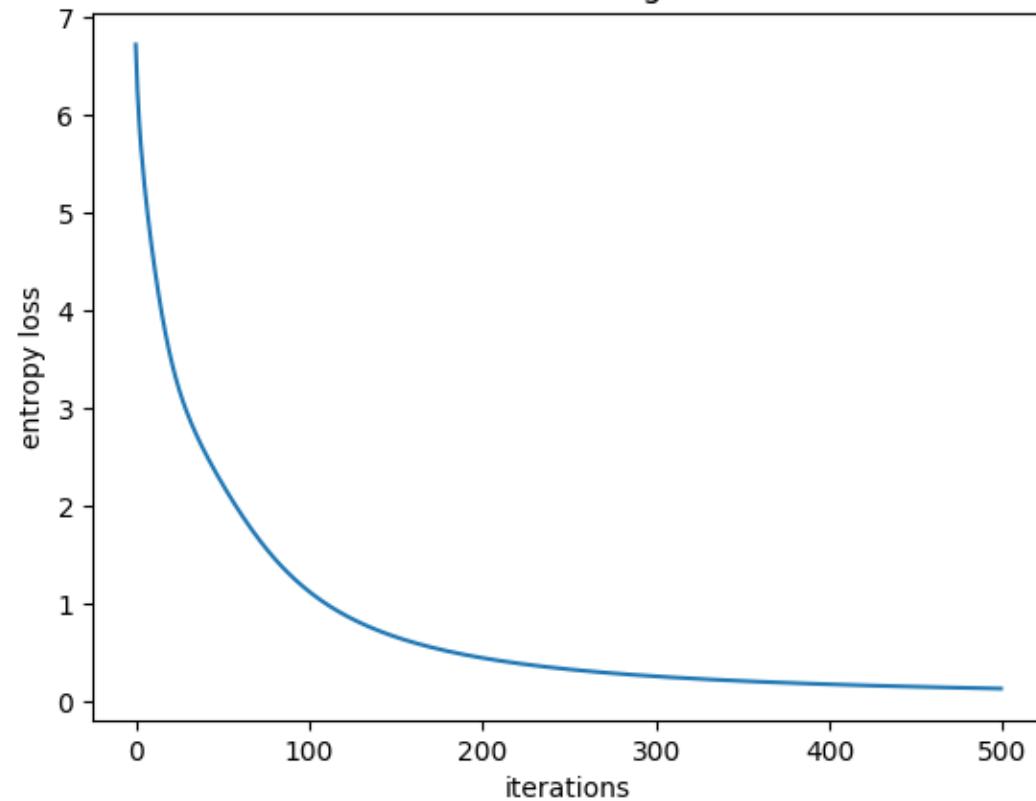
$$\mathbf{V} \leftarrow \mathbf{V} - \alpha \mathbf{H}^T \nabla_{\mathbf{U}} J = \begin{pmatrix} 0.92 & 0.05 & -0.26 \\ 0.68 & -0.36 & -0.19 \\ 0.22 & 1.05 & -0.10 \end{pmatrix}$$

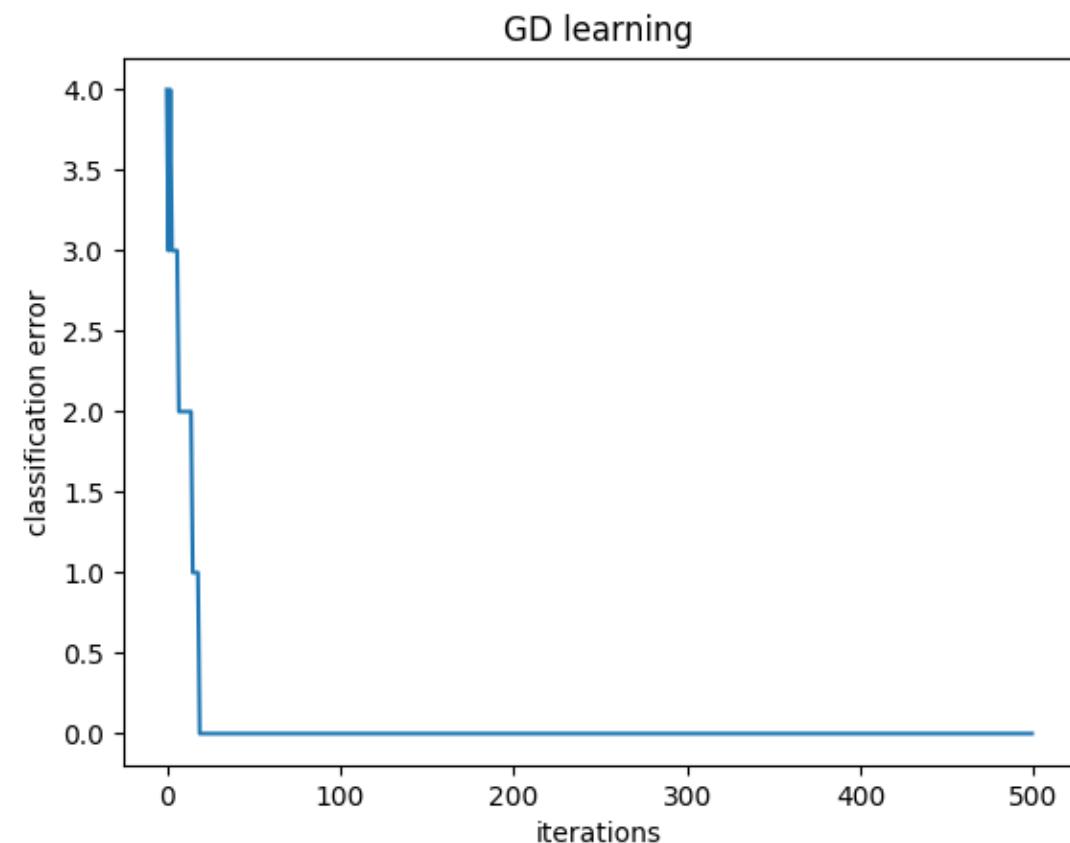
$$\mathbf{c} \leftarrow \mathbf{c} - \alpha (\nabla_{\mathbf{U}} J)^T \mathbf{1}_P = \begin{pmatrix} -0.16 \\ 0.04 \\ 0.12 \end{pmatrix}$$

$$\mathbf{W} \leftarrow \mathbf{W} - \alpha \mathbf{X}^T \nabla_{\mathbf{Z}} J = \begin{pmatrix} -0.13 & 0.95 & 0.18 \\ -0.63 & 0.42 & 0.95 \end{pmatrix}$$

$$\mathbf{b} \leftarrow \mathbf{b} - \alpha (\nabla_{\mathbf{Z}} J)^T \mathbf{1}_P = \begin{pmatrix} -0.02 \\ 0.00 \\ -0.01 \end{pmatrix}$$

GD learning





At convergence:

$$\mathbf{V} = \begin{pmatrix} 2.93 & -5.33 & 3.12 \\ 2.80 & 1.20 & -3.87 \\ 0.09 & 4.55 & -3.47 \end{pmatrix}, \quad \mathbf{c} = \begin{pmatrix} -1.94 \\ -0.06 \\ 2.01 \end{pmatrix}$$

$$\mathbf{W} = \begin{pmatrix} -1.81 & 0.32 & 0.08 \\ -1.40 & 2.92 & 1.91 \end{pmatrix}, \quad \mathbf{b} = \begin{pmatrix} 4.36 \\ 0.73 \\ -1.71 \end{pmatrix}$$

$$\mathbf{Y} = \begin{pmatrix} 1 \\ 1 \\ 2 \\ 2 \\ 3 \\ 3 \end{pmatrix}$$

Entropy = 0.138

Error = 0

Testing patterns:

$$\mathbf{x}_1 = \begin{pmatrix} 2.5 \\ 1.5 \end{pmatrix} \text{ and } \mathbf{x}_2 = \begin{pmatrix} -1.5 \\ 0.5 \end{pmatrix}$$

In batch mode

$$\mathbf{X} = \begin{pmatrix} 2.5 & 1.5 \\ -1.5 & 0.5 \end{pmatrix}$$

Forward propagation of activations

$$\mathbf{Z} = \mathbf{XW} + \mathbf{B} = \begin{pmatrix} -2.26 & 5.91 & 1.36 \\ 6.35 & 1.72 & -0.91 \end{pmatrix}$$

$$\mathbf{H} = f(\mathbf{Z}) = \begin{pmatrix} 0.09 & 1.0 & 0.8 \\ 1.0 & 0.85 & 0.29 \end{pmatrix}$$

$$\mathbf{U} = \mathbf{HV} + \mathbf{C} = \begin{pmatrix} 1.2 & 4.25 & -4.33 \\ 3.38 & -3.006 & 0.82 \end{pmatrix}$$

$$g(\mathbf{U}) = \frac{e^{\mathbf{U}}}{\sum_{k=1}^K e^{\mathbf{U}_k}} = \begin{pmatrix} 0.05 & 0.95 & 0.0 \\ 0.93 & 0.0 & 0.07 \end{pmatrix}$$

$$\mathbf{Y} = \arg \max_k g(\mathbf{U}) = \begin{pmatrix} 2 \\ 1 \end{pmatrix}$$

The class labels: $\mathbf{x}_1 \rightarrow \text{class } B, \mathbf{x}_2 \rightarrow \text{class } A$

3. Design a feedforward neural network consisting of two-hidden layers to approximate the following function:

$$\phi(x, y) = 0.8x^2 - y^3 + 2.5xy$$

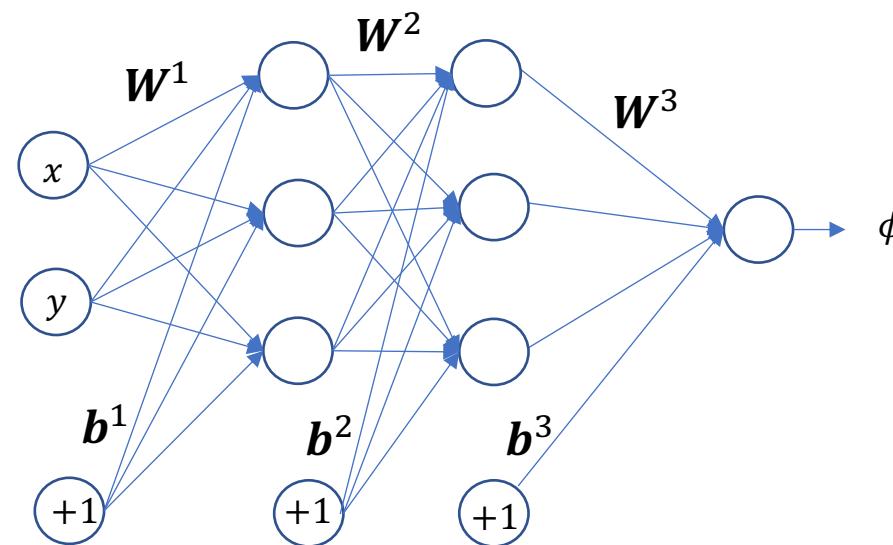
for $-1.0 \leq x, y \leq 1.0$.

Use three ReLU neurons at each hidden layer and a linear neuron at the output layer.

- (a) Divide the input space equally into square regions of size 0.25×0.25 and use grid points as data to learn the function ϕ .
- (b) Train the network using gradient decent learning at learning rate $\alpha = 0.01$ and plot the learning curve (mean square error vs. iterations) and the predicted data points.
- (c) Compare the learning curves when learning the function at learning rates $\alpha = 0.005, 0.01, 0.05$, and 0.1 .

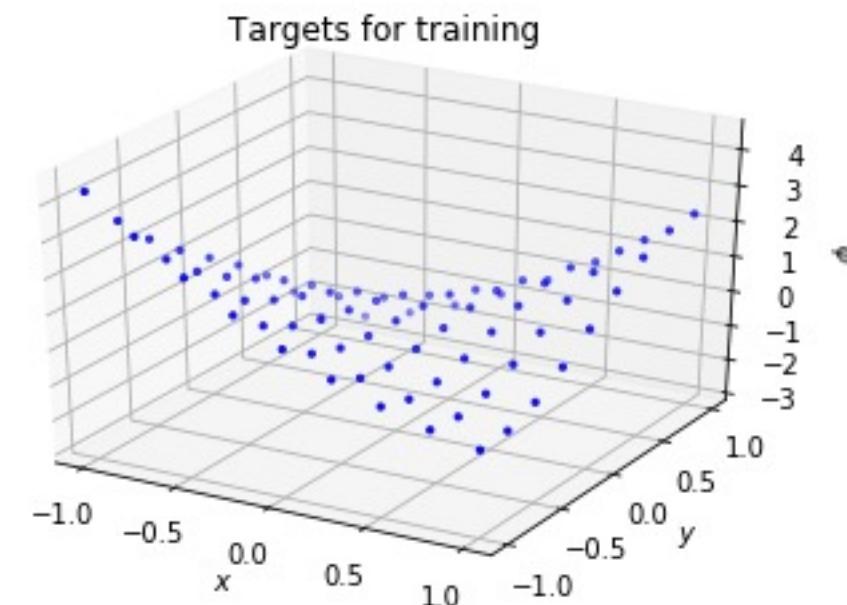
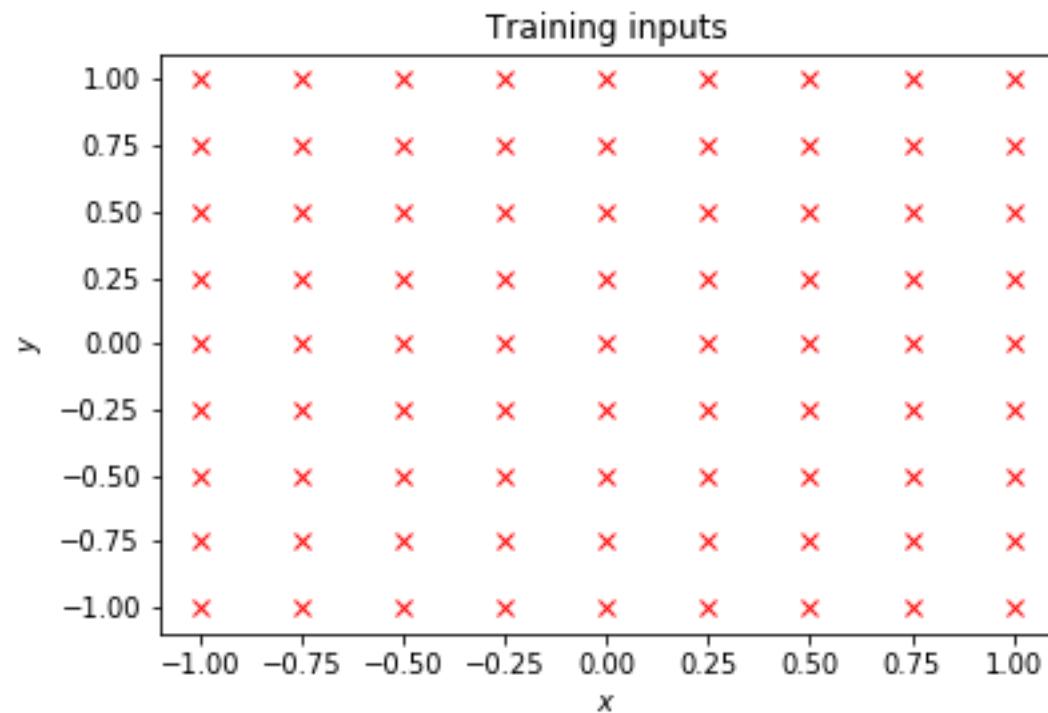
$$\phi(x, y) = 0.8x^2 - x^3 + 2.5xy \quad \text{for } -1.0 \leq x, y \leq 1.0$$

Feedforward neural network with two hidden layers:



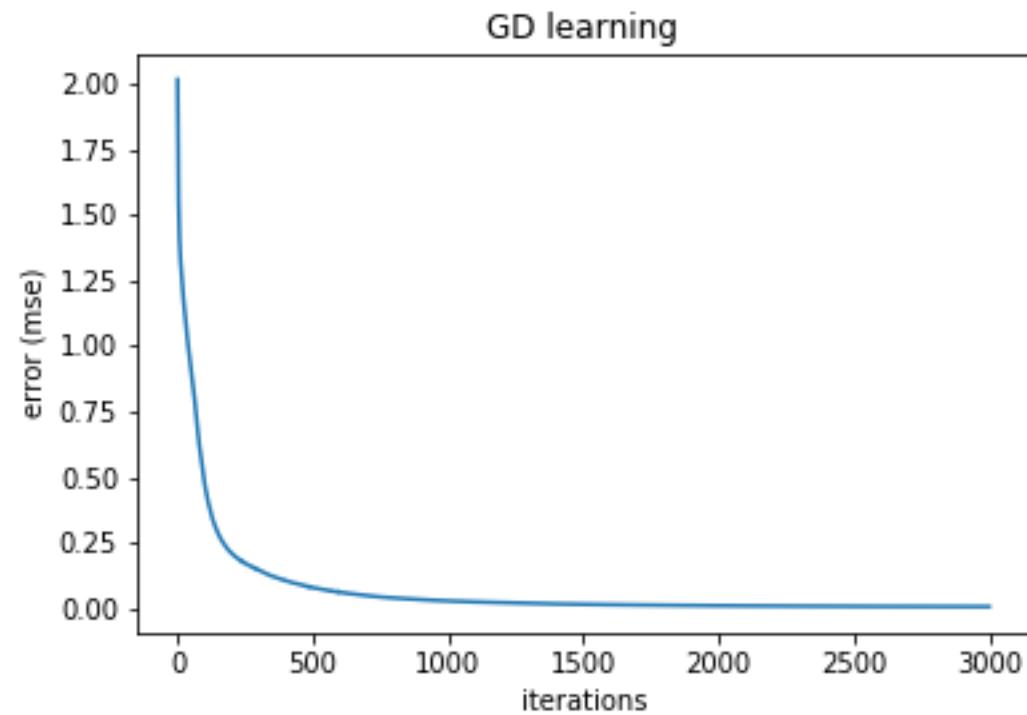
$$\phi(x, y) = 0.8x_1^2 - y^3 + 2.5xy \text{ for } -1.0 \leq x, y \leq 1.0$$

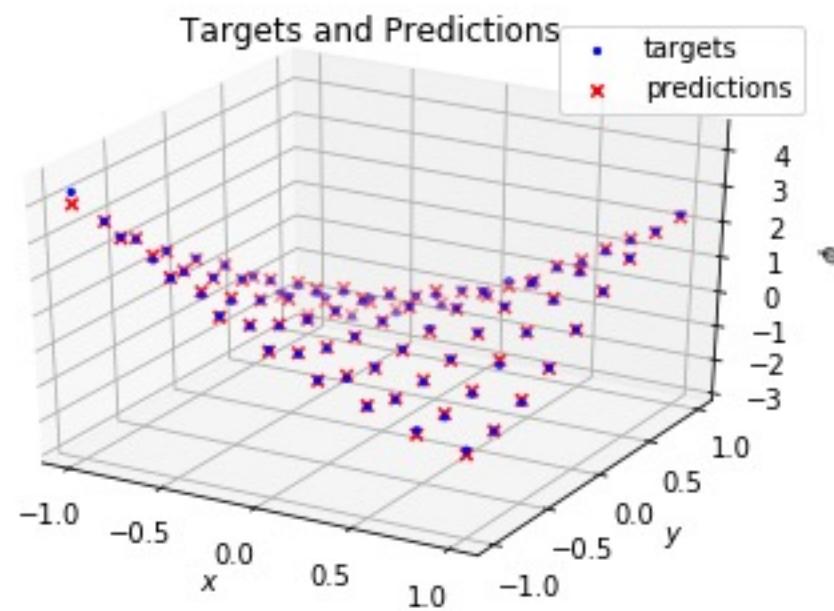
Data points in a grid of size 0.25x0.25:

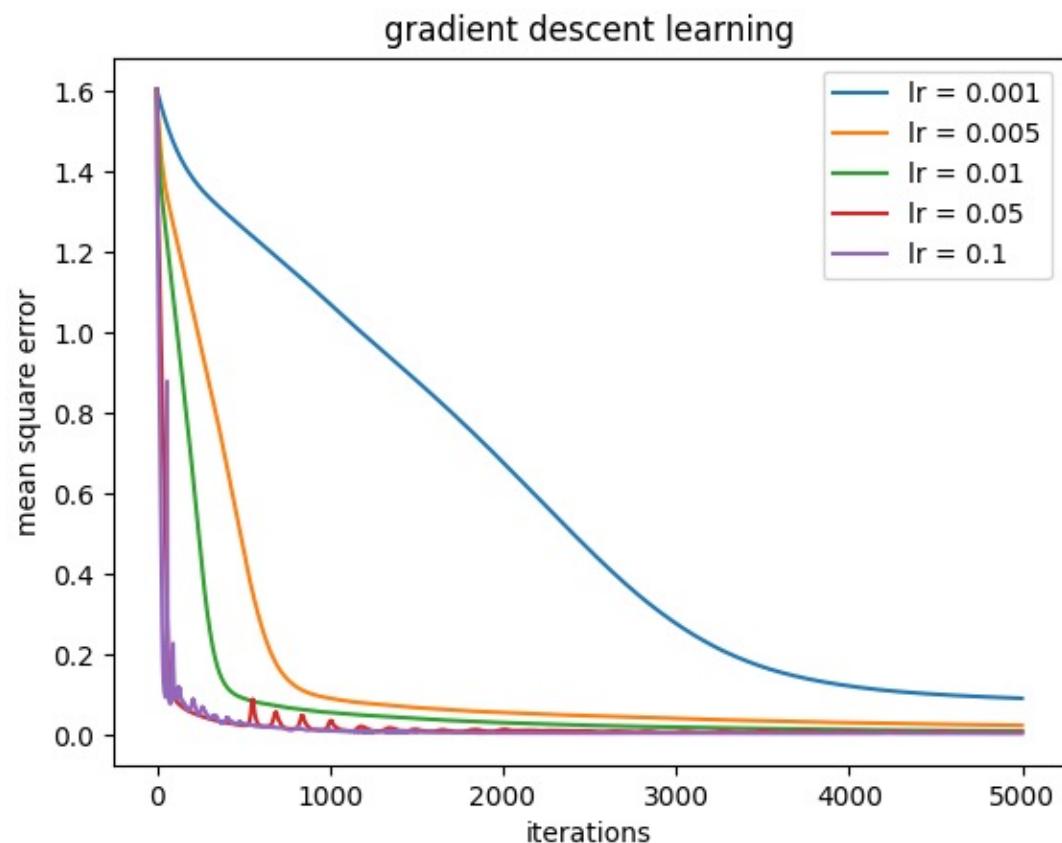


```
class FFN(nn.Module):
    def __init__(self):
        super().__init__()
        self.relu_stack = nn.Sequential(
            nn.Linear(2, 10),
            nn.ReLU(),
            nn.Linear(10, 5),
            nn.ReLU(),
            nn.Linear(5, 1),
        )

    def forward(self, x):
        logits = self.relu_stack(x)
        return logits
```







Tutorial 6

Convolutional Neural

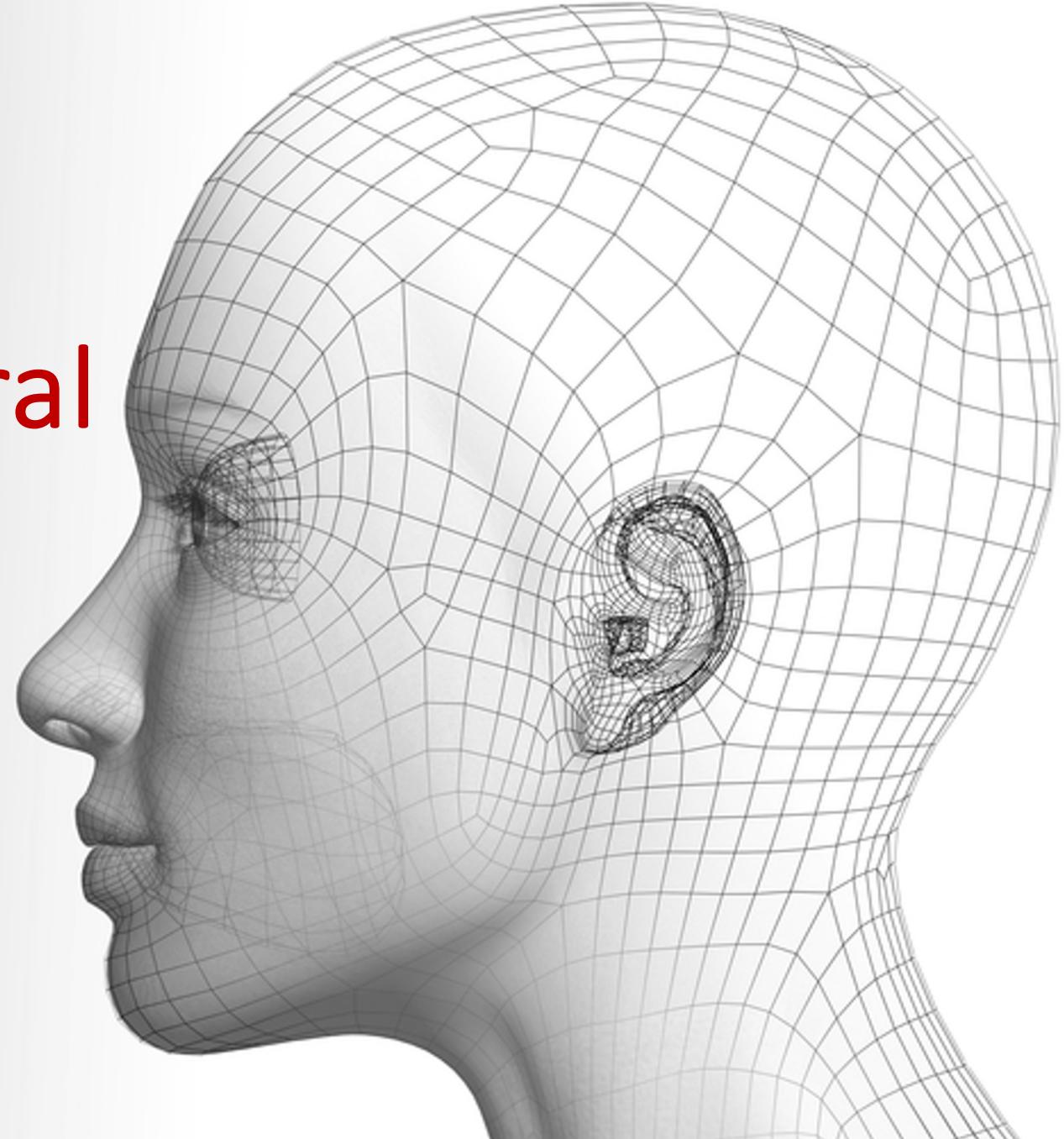
Networks I

Xingang Pan

潘新钢

<https://xingangpan.github.io/>

<https://twitter.com/XingangP>



Question 1

1. The first hidden layer of a convolution neural network (CNN) has a convolution layer consisting of two feature maps with filters w_1 and w_2 and biases = 0.1, and neurons having sigmoid activation functions, and a pooling layer with a pooling windows of size 2x2:

$$w_1 = \begin{pmatrix} 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{pmatrix} \text{ and } w_2 = \begin{pmatrix} -1 & -1 & -1 \\ 0 & 0 & 0 \\ 1 & 1 & 1 \end{pmatrix}.$$

The input layer is of 6x6 size and receives an input image I :

Filters=weights

$$I = \begin{pmatrix} 0.7 & 0.1 & 0.2 & 0.3 & 0.3 & 0.5 \\ 0.8 & 0.1 & 0.3 & 0.5 & 0.1 & 0.0 \\ 1.0 & 0.2 & 0.0 & 0.3 & 0.2 & 0.7 \\ 0.8 & 0.1 & 0.5 & 0.6 & 0.3 & 0.4 \\ 0.1 & 0.0 & 0.9 & 0.3 & 0.3 & 0.2 \\ 1.0 & 0.1 & 0.4 & 0.5 & 0.2 & 0.8 \end{pmatrix}$$

- a. Find the outputs at the first convolution layer if
 - i. padding = 0 and strides = (1,1)
 - ii. padding = 1 and strides = (2,2)
- b. Find the outputs at the first pooling layer for Part (a), assuming strides of (2, 2) and pooling is
 - i. max pooling
 - ii. mean pooling

t6q1.ipynb

Question 1

$$\mathbf{w}_1 = \begin{pmatrix} 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{pmatrix}, \mathbf{w}_2 = \begin{pmatrix} -1 & -1 & -1 \\ 0 & 0 & 0 \\ 1 & 1 & 1 \end{pmatrix} \text{ and } b = 0.1.$$

$$\mathbf{I} = \begin{pmatrix} 0.7 & 0.1 & 0.2 & 0.3 & 0.3 & 0.5 \\ 0.8 & 0.1 & 0.3 & 0.5 & 0.1 & 0.0 \\ 1.0 & 0.2 & 0.0 & 0.3 & 0.2 & 0.7 \\ 0.8 & 0.1 & 0.5 & 0.6 & 0.3 & 0.4 \\ 0.1 & 0.0 & 0.9 & 0.3 & 0.3 & 0.2 \\ 1.0 & 0.1 & 0.4 & 0.5 & 0.2 & 0.8 \end{pmatrix}$$

Synaptic inputs to the feature map with filter \mathbf{w}_1 :

$$\mathbf{u}_1 = \text{Conv}(\mathbf{I}, \mathbf{w}_1) + b_1$$

Padding = 0 and strides = (1,1)

$$\mathbf{u}_1(1,1) = 0.7 \times 0 + 0.1 \times 1 + 0.2 \times 1 + 0.8 \times 1 + 0.1 \times 0 + 0.3 \times 1 + 1.0 \times 1 + 0.2 \times 1 + 0.0 \times 0 + 0.1 = 2.7$$

$$\mathbf{u}_1(1,2) = 0.1 \times 0 + 0.2 \times 1 + 0.3 \times 1 + 0.1 \times 1 + 0.3 \times 0 + 0.5 \times 1 + 0.2 \times 1 + 0.0 \times 1 + 0.3 \times 0 + 0.1 = 1.4$$

$$\mathbf{u}_1(2,1) = 0.8 \times 0 + 0.1 \times 1 + 0.3 \times 1 + 1.0 \times 1 + 0.2 \times 0 + 0.0 \times 1 + 0.8 \times 1 + 0.1 \times 1 + 0.5 \times 0 + 0.1 = 2.4$$

Question 1

$$\mathbf{u}_1 = \begin{pmatrix} 2.7 & 1.4 & 1.4 & 1.9 \\ 2.4 & 2.0 & 2.0 & 2.1 \\ 1.7 & 2.0 & 2.6 & 2.6 \\ 2.8 & 2.0 & 3.1 & 2.0 \end{pmatrix}$$

Similarly, $\mathbf{u}_2 = \text{Conv}(\mathbf{I}, \mathbf{w}_2) + b_2 = \begin{pmatrix} 0.3 & 0.0 & -0.2 & 2.0 \\ 0.3 & 0.4 & 0.6 & 0.8 \\ -0.1 & 0.8 & 1.1 & -0.3 \\ 0.2 & -0.1 & -0.2 & 0.3 \end{pmatrix}$

Feature maps at the convolutional layer

$$\mathbf{y}_1 = f(\mathbf{u}_1) = \frac{1}{1 + e^{-\mathbf{u}_1}} = \begin{pmatrix} 0.94 & 0.8 & 0.8 & 0.87 \\ 0.92 & 0.88 & 0.88 & 0.89 \\ 0.85 & 0.88 & 0.93 & 0.93 \\ 0.94 & 0.88 & 0.96 & 0.88 \end{pmatrix}$$

$$\mathbf{y}_2 = f(\mathbf{u}_2) = \frac{1}{1 + e^{-\mathbf{u}_2}} = \begin{pmatrix} 0.57 & 0.5 & 0.45 & 0.55 \\ 0.57 & 0.60 & 0.65 & 0.69 \\ 0.48 & 0.69 & 0.75 & 0.43 \\ 0.55 & 0.48 & 0.45 & 0.57 \end{pmatrix}$$

Output spatial size

$$= (\text{Input size} - \text{Filter size} + 2(\text{Padding})) / \text{Stride} + 1$$

$$= (6 - 3) / 1 + 1 = 4$$

Question 1

Feature maps at the convolutional layer:

$$\mathbf{y} = \begin{pmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{pmatrix} = \begin{pmatrix} \begin{pmatrix} 0.94 & 0.8 & 0.8 & 0.87 \\ 0.92 & 0.88 & 0.88 & 0.89 \\ 0.85 & 0.88 & 0.93 & 0.93 \\ 0.94 & 0.88 & 0.96 & 0.88 \\ 0.57 & 0.5 & 0.45 & 0.55 \\ 0.57 & 0.60 & 0.65 & 0.69 \\ 0.48 & 0.69 & 0.75 & 0.43 \\ 0.55 & 0.48 & 0.45 & 0.57 \end{pmatrix} \end{pmatrix}$$

Pooling 2×2 and strides = 2

Max-pooling:

$$o = \begin{pmatrix} (0.94 & 0.89) \\ (0.94 & 0.96) \\ (0.60 & 0.69) \\ (0.69 & 0.75) \end{pmatrix}$$

Mean-pooling:

$$p_{ave} = \begin{pmatrix} (0.88 & 0.86) \\ (0.89 & 0.92) \\ (0.56 & 0.58) \\ (0.55 & 0.55) \end{pmatrix}$$

Question 1

$$\mathbf{w}_1 = \begin{pmatrix} 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{pmatrix}, \mathbf{w}_2 = \begin{pmatrix} -1 & -1 & -1 \\ 0 & 0 & 0 \\ 1 & 1 & 1 \end{pmatrix} \text{ and } b = 0.1.$$

$$\mathbf{I} = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0.7 & 0.1 \\ 0 & 0.8 & 0.1 \\ 1.0 & 0.2 & 0.0 \\ 0.8 & 0.1 & 0.5 \\ 0.1 & 0.0 & 0.9 \\ 1.0 & 0.1 & 0.4 \end{pmatrix} \begin{pmatrix} 0.2 & 0.3 & 0.3 & 0.5 \\ 0.3 & 0.5 & 0.1 & 0.0 \\ 0.0 & 0.3 & 0.2 & 0.7 \\ 0.5 & 0.6 & 0.3 & 0.4 \\ 0.3 & 0.3 & 0.3 & 0.2 \\ 0.2 & 0.2 & 0.2 & 0.8 \end{pmatrix}$$

Synaptic inputs to the feature map with filter \mathbf{w}_1 :

$$\mathbf{u}_1 = \text{Conv}(\mathbf{I}, \mathbf{w}_1) + b_1$$

Padding = 1 and strides = (2,2)

$$\mathbf{u}_1(1,1) = 0.0 \times 0 + 0.0 \times 1 + 0.0 \times 1 + 0.0 \times 1 + 0.7 \times 0 + 0.1 \times 1 + 0.0 \times 1 + 0.8 \times 1 + 0.1 \times 0 + 0.1 = 1.0$$

Question 1

$$\mathbf{w}_1 = \begin{pmatrix} 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{pmatrix}, \mathbf{w}_2 = \begin{pmatrix} -1 & -1 & -1 \\ 0 & 0 & 0 \\ 1 & 1 & 1 \end{pmatrix} \text{ and } b = 0.1.$$

$$\mathbf{I} = \begin{pmatrix} 0.7 & 0 & 0 & 0 \\ 0.8 & 0.1 & 0.2 & 0.3 & 0.3 & 0.5 \\ 1.0 & 0.1 & 0.3 & 0.5 & 0.1 & 0.0 \\ 1.0 & 0.2 & 0.0 & 0.3 & 0.2 & 0.7 \\ 0.8 & 0.1 & 0.5 & 0.6 & 0.3 & 0.4 \\ 0.1 & 0.0 & 0.9 & 0.3 & 0.3 & 0.2 \\ 1.0 & 0.1 & 0.4 & 0.5 & 0.2 & 0.8 \end{pmatrix}$$

Synaptic inputs to the feature map with filter \mathbf{w}_1 :

$$\mathbf{u}_1 = \text{Conv}(\mathbf{I}, \mathbf{w}_1) + b_1$$

Padding = 1 and strides = (2,2)

$$\mathbf{u}_1(1,1) = 0.0 \times 0 + 0.0 \times 1 + 0.0 \times 1 + 0.0 \times 1 + 0.7 \times 0 + 0.1 \times 1 + 0.0 \times 1 + 0.8 \times 1 + 0.1 \times 0 + 0.1 = 1.0$$

$$\mathbf{u}_1(1,2) = 0.0 \times 0 + 0.0 \times 1 + 0.0 \times 1 + 0.1 \times 1 + 0.2 \times 0 + 0.3 \times 1 + 0.1 \times 1 + 0.3 \times 1 + 0.5 \times 0 + 0.1 = 0.9$$

Question 1

$$\mathbf{w}_1 = \begin{pmatrix} 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{pmatrix}, \mathbf{w}_2 = \begin{pmatrix} -1 & -1 & -1 \\ 0 & 0 & 0 \\ 1 & 1 & 1 \end{pmatrix} \text{ and } b = 0.1.$$

$$\mathbf{I} = \begin{pmatrix} 0.7 & 0.1 & 0.2 & \boxed{0.3} & 0.3 & 0.5 \\ 0.8 & 0.1 & 0.3 & \boxed{0.5} & 0.1 & 0.0 \\ 1.0 & 0.2 & 0.0 & 0.3 & 0.2 & 0.7 \\ 0.8 & 0.1 & 0.5 & 0.6 & 0.3 & 0.4 \\ 0.1 & 0.0 & 0.9 & 0.3 & 0.3 & 0.2 \\ 1.0 & 0.1 & 0.4 & 0.5 & 0.2 & 0.8 \end{pmatrix}$$

Synaptic inputs to the feature map with filter \mathbf{w}_1 :

$$\mathbf{u}_1 = \text{Conv}(\mathbf{I}, \mathbf{w}_1) + b_1$$

Padding = 1 and strides = (2,2)

$$\mathbf{u}_1(1,1) = 0.0 \times 0 + 0.0 \times 1 + 0.0 \times 1 + 0.0 \times 1 + 0.7 \times 0 + 0.1 \times 1 + 0.0 \times 1 + 0.8 \times 1 + 0.1 \times 0 + 0.1 = 1.0$$

$$\mathbf{u}_1(1,2) = 0.0 \times 0 + 0.0 \times 1 + 0.0 \times 1 + 0.1 \times 1 + 0.2 \times 0 + 0.3 \times 1 + 0.1 \times 1 + 0.3 \times 1 + 0.5 \times 0 + 0.1 = 0.9$$

$$\mathbf{u}_1(1,3) = 0.0 \times 0 + 0.0 \times 1 + 0.0 \times 1 + 0.3 \times 1 + 0.3 \times 0 + 0.5 \times 1 + 0.5 \times 1 + 0.1 \times 1 + 0.0 \times 0 + 0.1 = 1.5$$

Question 1

$$\mathbf{w}_1 = \begin{pmatrix} 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{pmatrix}, \mathbf{w}_2 = \begin{pmatrix} -1 & -1 & -1 \\ 0 & 0 & 0 \\ 1 & 1 & 1 \end{pmatrix} \text{ and } b = 0.1.$$

$$\mathbf{I} = \begin{pmatrix} 0.7 & 0.1 & 0.2 & 0.3 & 0.3 & 0.5 \\ 0.8 & 0.1 & 0.3 & 0.5 & 0.1 & 0.0 \\ 1.0 & 0.2 & 0.0 & 0.3 & 0.2 & 0.7 \\ 0.8 & 0.1 & 0.5 & 0.6 & 0.3 & 0.4 \\ 0.1 & 0.0 & 0.9 & 0.3 & 0.3 & 0.2 \\ 1.0 & 0.1 & 0.4 & 0.5 & 0.2 & 0.8 \end{pmatrix}$$

Synaptic inputs to the feature map with filter \mathbf{w}_1 :

$$\mathbf{u}_1 = \text{Conv}(\mathbf{I}, \mathbf{w}_1) + b_1$$

Padding = 1 and strides = (2,2)

$$\mathbf{u}_1(1,1) = 0.0 \times 0 + 0.0 \times 1 + 0.0 \times 1 + 0.0 \times 1 + 0.7 \times 0 + 0.1 \times 1 + 0.0 \times 1 + 0.8 \times 1 + 0.1 \times 0 + 0.1 = 1.0$$

$$\mathbf{u}_1(1,2) = 0.0 \times 0 + 0.0 \times 1 + 0.0 \times 1 + 0.1 \times 1 + 0.2 \times 0 + 0.3 \times 1 + 0.1 \times 1 + 0.3 \times 1 + 0.5 \times 0 + 0.1 = 0.9$$

$$\mathbf{u}_1(1,3) = 0.0 \times 0 + 0.0 \times 1 + 0.0 \times 1 + 0.3 \times 1 + 0.3 \times 0 + 0.5 \times 1 + 0.5 \times 1 + 0.1 \times 1 + 0.0 \times 0 + 0.1 = 1.5$$

$$\mathbf{u}_1(2,1) = 0.0 \times 0 + 0.8 \times 1 + 0.1 \times 1 + 0.0 \times 1 + 1.0 \times 0 + 0.2 \times 1 + 0.0 \times 1 + 0.8 \times 1 + 0.1 \times 0 + 0.1 = 2.0$$

Question 1

$$\mathbf{u}_1 = \begin{pmatrix} 1.0 & 0.9 & 1.5 \\ 2.0 & 2.0 & 2.1 \\ 2.0 & 2.0 & 2.0 \end{pmatrix}$$

Similarly, $\mathbf{u}_2 = Conv(\mathbf{I}, \mathbf{w}_2) + b_2 = \begin{pmatrix} 1.0 & 1.0 & 0.7 \\ 0.1 & 0.4 & 0.8 \\ 0.3 & -0.1 & 0.3 \end{pmatrix}$

Feature maps at the convolutional layer

$$\mathbf{y}_1 = f(\mathbf{u}_1) = \frac{1}{1 + e^{-\mathbf{u}_1}} = \begin{pmatrix} 0.73 & 0.71 & 0.82 \\ 0.88 & 0.88 & 0.89 \\ 0.88 & 0.88 & 0.88 \end{pmatrix}$$

$$\mathbf{y}_2 = f(\mathbf{u}_2) = \frac{1}{1 + e^{-\mathbf{u}_2}} = \begin{pmatrix} 0.73 & 0.73 & 0.67 \\ 0.52 & 0.60 & 0.69 \\ 0.57 & 0.48 & 0.57 \end{pmatrix}$$

Question 1

Feature maps at the convolutional layer:

$$\mathbf{y} = \begin{pmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{pmatrix} = \begin{pmatrix} \begin{pmatrix} 0.73 & 0.71 \\ 0.88 & 0.88 \\ 0.88 & 0.88 \\ 0.73 & 0.73 \\ 0.52 & 0.60 \end{pmatrix} & 0.82 \\ \begin{pmatrix} 0.88 & 0.89 \\ 0.88 & 0.88 \\ 0.67 & 0.69 \\ 0.57 & 0.57 \end{pmatrix} & 0.89 \end{pmatrix}$$

Pooling VALID 2x2 strides = 2:

Max-pooling:

$$o = \begin{pmatrix} (0.88) \\ (0.73) \end{pmatrix}$$

Mean-pooling:

$$p_{ave} = \begin{pmatrix} (0.80) \\ (0.65) \end{pmatrix}$$

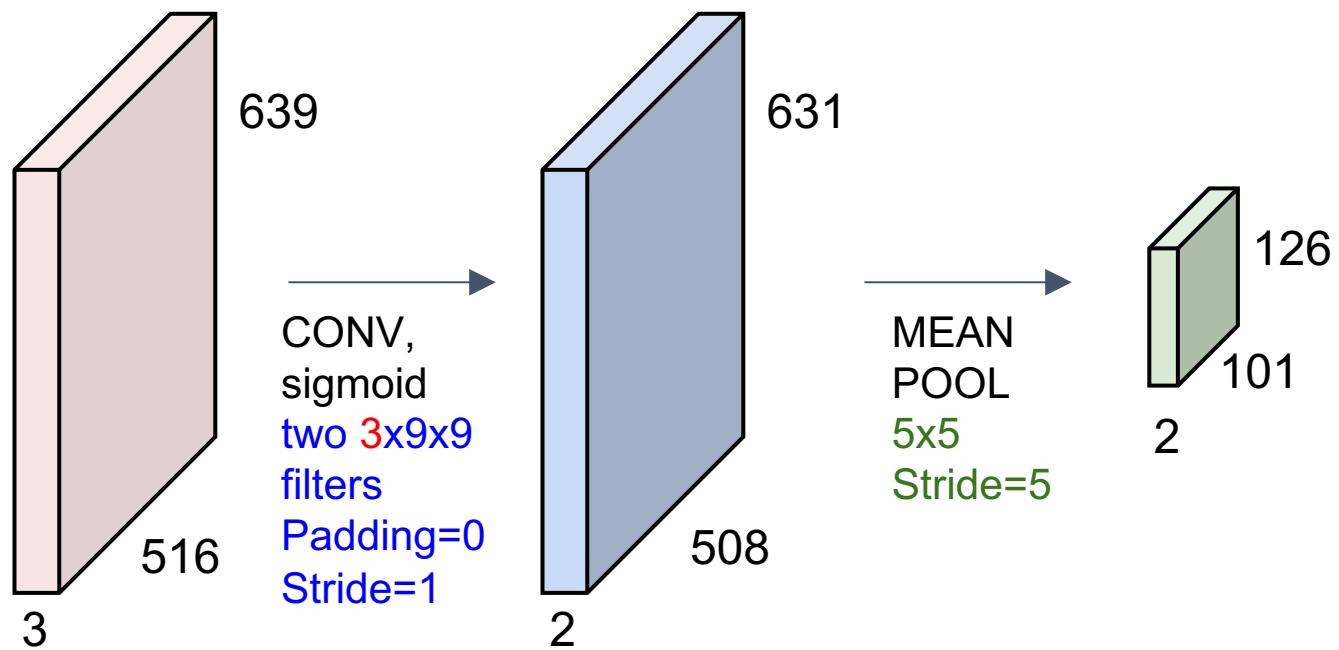
Question 2

2. Given ‘3wolfmoon.jpg’ color image of size 639×516 .

- a. Initialize weights and biases of a convolutional layer with two kernels of size 9×9 .
Note that the input image is in color and has three channels.
- b. Display the feature maps at the convolution layer, assuming sigmoid activation functions. Use VALID padding and strides = 1.
- c. Display the outputs of a mean pooling layer with a pooling window size 5×5 and strides = 5.



Question 2



Question 2



```
# open image and normalize
 img = Image.open('3wolfmoon.jpg')
 img = np.asarray(img, dtype='float32') / 256.
 print(img.shape)
```

```
-----
FileNotFoundError                      Traceback (most recent call last)
<ipython-input-2-becdc6409bf5> in <module>()
      1 # open image and normalize
----> 2 img = Image.open('3wolfmoon.jpg')
      3 img = np.asarray(img, dtype='float32') / 256.
      4 print(img.shape)

/usr/local/lib/python3.7/dist-packages/PIL/Image.py in open(fp, mode)
2841
2842     if filename:
-> 2843         fp = builtins.open(filename, "rb")
2844         exclusive_fp = True
2845

FileNotFoundError: [Errno 2] No such file or directory: '3wolfmoon.jpg'
```

You need to mount Google Drive locally and use the correct path

<https://colab.research.google.com/notebooks/io.ipynb#scrollTo=u22w3BFiOveA>

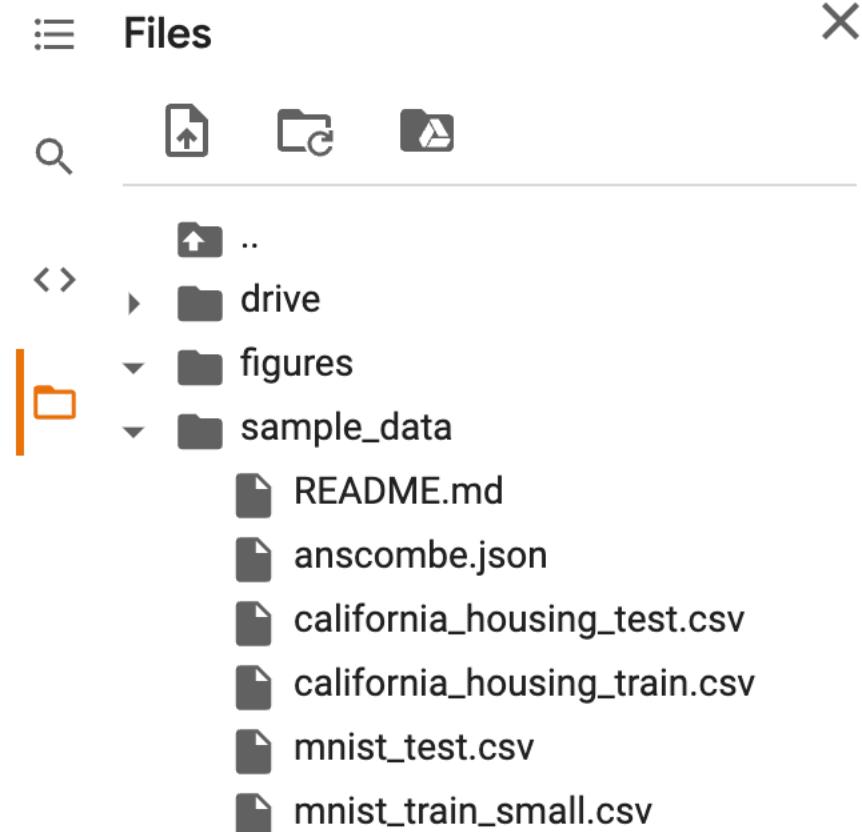
Question 2

```
# mount google drive
from google.colab import drive
drive.mount('/content/drive')
```

Mounted at /content/drive

```
# open image and normalize
img = Image.open('/content/drive/MyDrive/3wolfmoon.jpg')
img = np.asarray(img, dtype='float32') / 256.
print(img.shape)

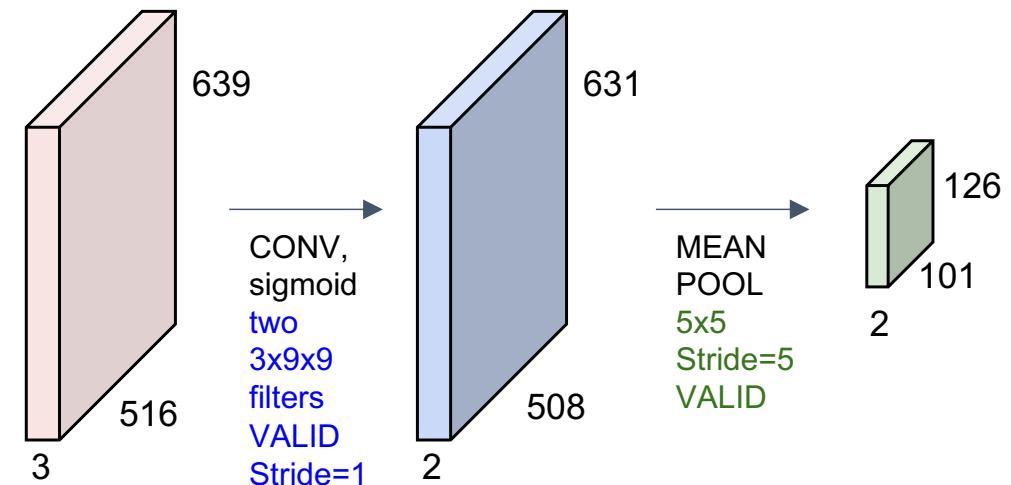
(639, 516, 3)
```



You can copy the path from the browser

Question 2

Plot the three channels of the input image



Original image



Red



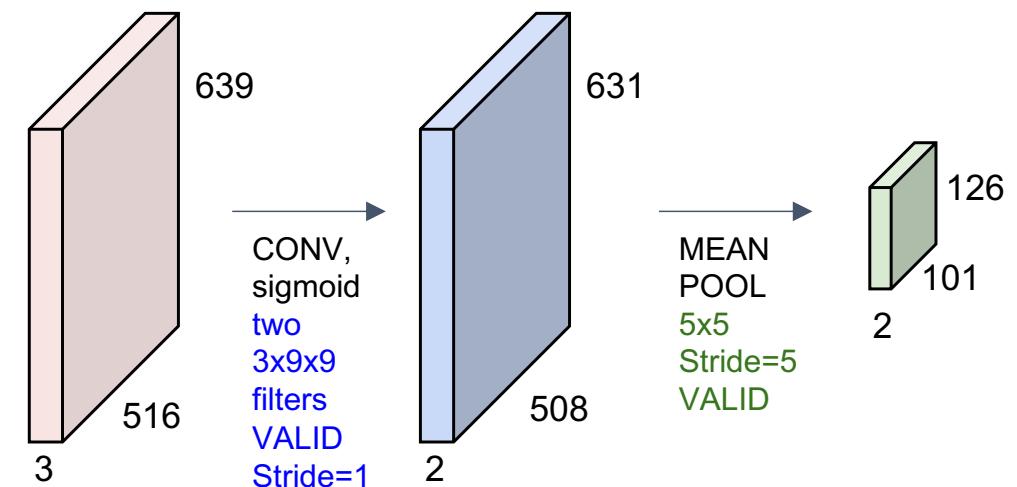
Green



Blue

Question 2

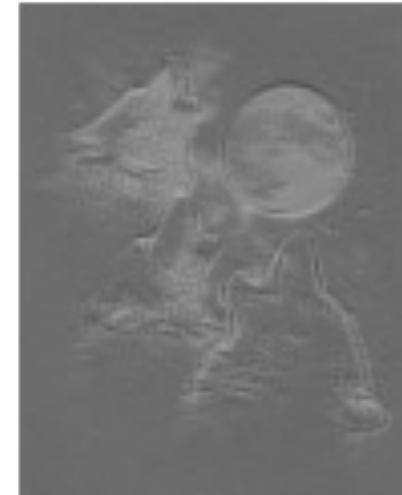
Plot the original image and the first and second components of conv output



Original image



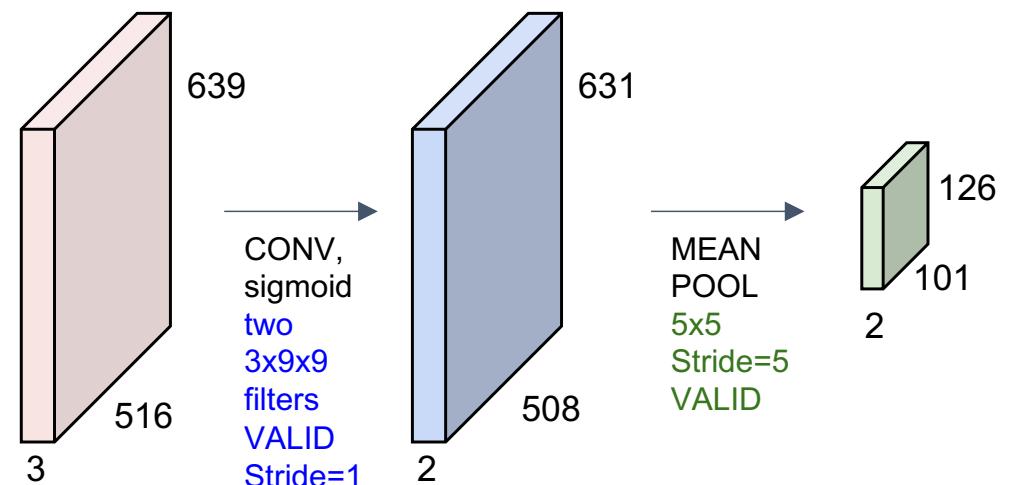
Feature map 1



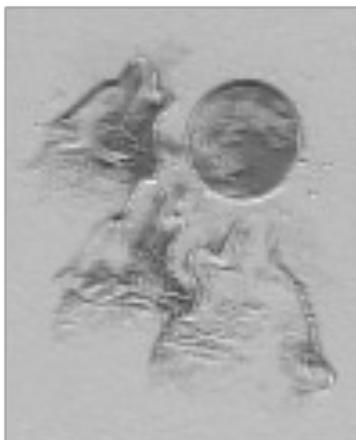
Feature map 2

Question 2

Plot the original image and the first and second components of pooling output



Original image



Pooled feature
map 1



Pooled feature
map 2

Question 3

3. Design a CNN with one hidden layer to recognize digit images in MNIST database:

<http://yann.lecun.com/exdb/mnist/>

The convolution layer consists of 25 filters of dimensions 9x9 and the pooling layer has a pooling window size 4x4. Assume VALID padding and default strides for both convolution and pooling layer. Train the network with mini batch gradient decent learning with learning factor $\alpha = 10^{-3}$ and batch size = 128.

Plot

- a. The training and test errors against learning epochs.
- b. Final filter weights
- c. Feature maps at the convolution and pooling layers for a representative test pattern
- d. Repeat training by introducing decay parameter $\beta = 10^{-6}$ and momentum term with $\gamma = 0.5$, and compare the learning curves

t6q3a.ipynb

t6q3b.ipynb

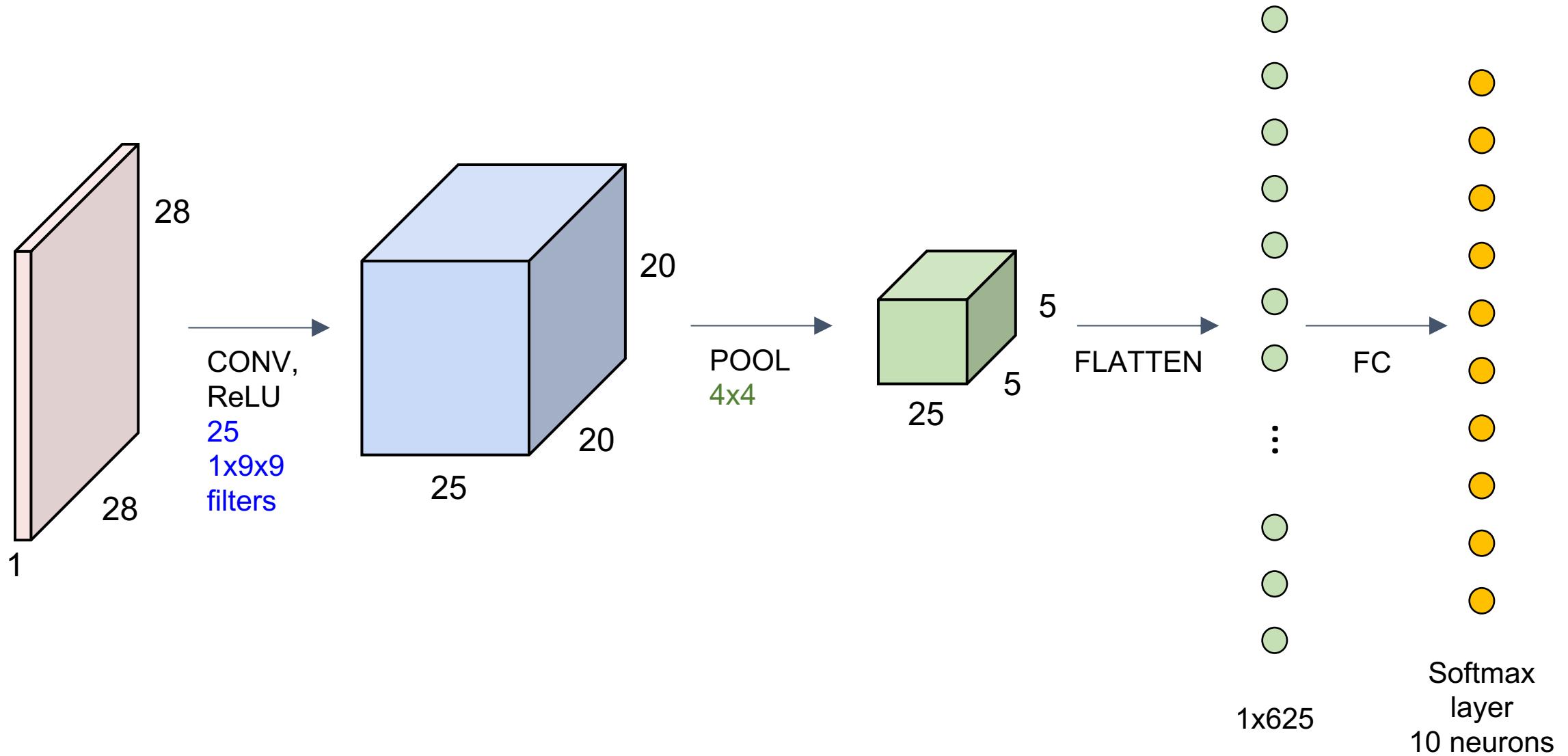
Question 3



MNIST

- Size-normalized and centred 1x28x28 =784 inputs
- Training set = 60,000 images
- Testing set = 10,000 images

Question 3



Question 3

Change runtime type

Runtime type

Python 3

Hardware accelerator

CPU

T4 GPU

A100 GPU

V100 GPU

TPU

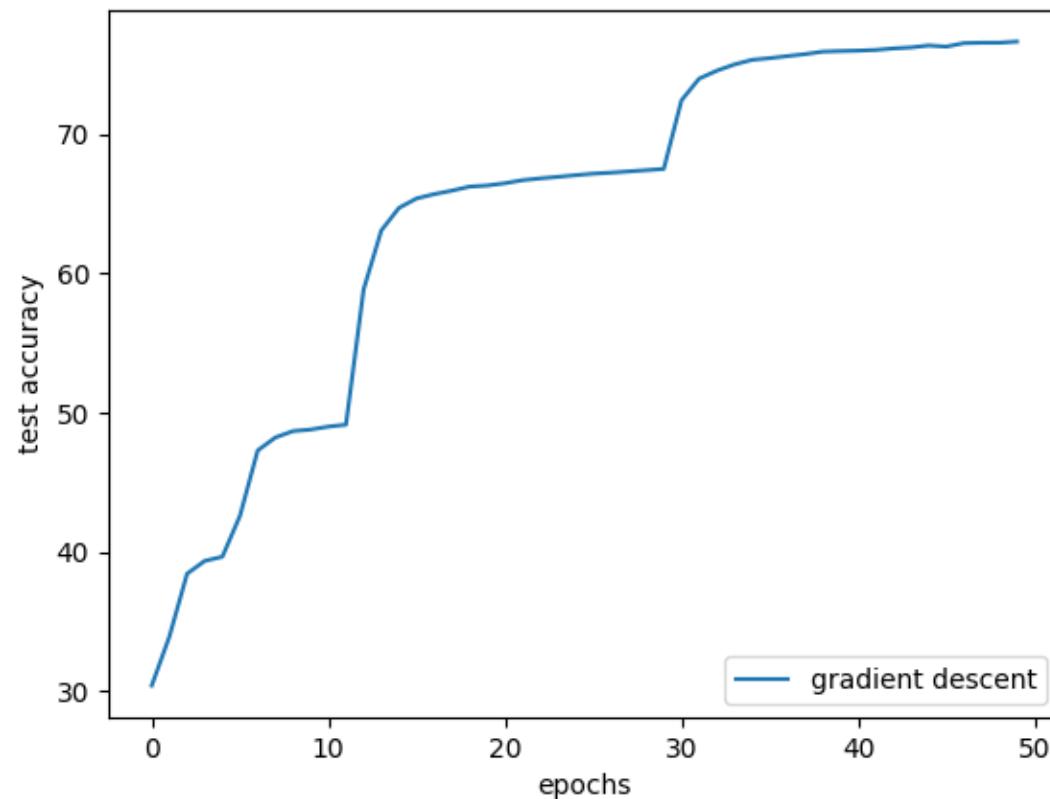
Want access to premium GPUs? [Purchase additional compute units](#)

Cancel

Save

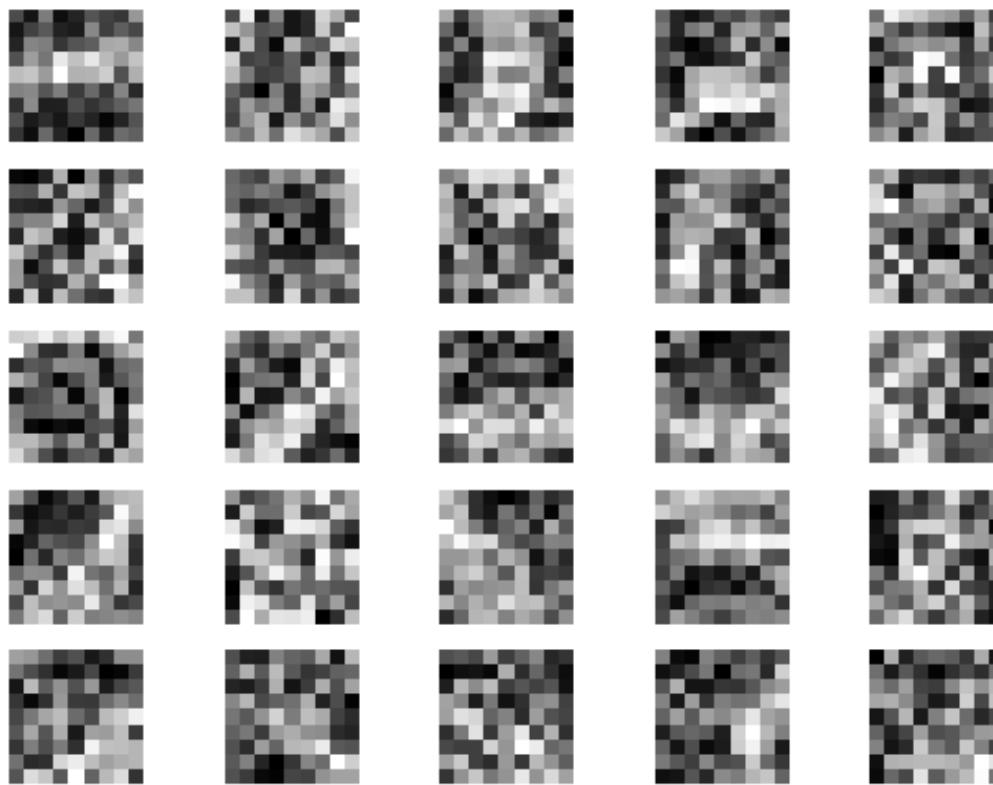
Question 3a

Plot the learning curve (I only show the test curve) - you should be able get better learning curve if you run for more epochs

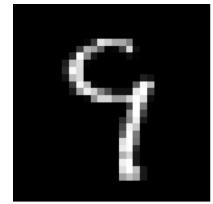


Question 3b

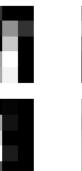
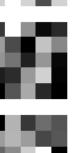
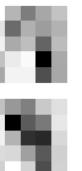
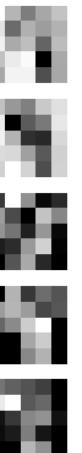
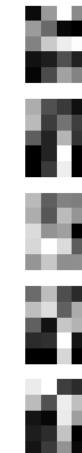
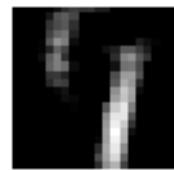
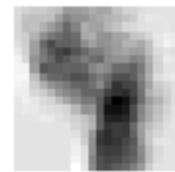
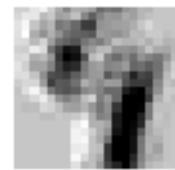
Plot filters learned in the conv layer



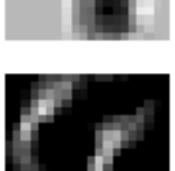
Question 3c



Input 28x28

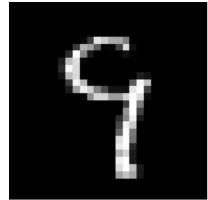


Feature maps of
pooling layer 25x5x5



Feature maps of the
convolution layer 25x20x20

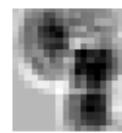
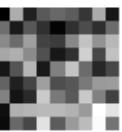
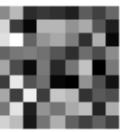
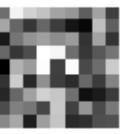
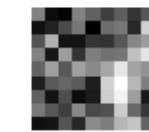
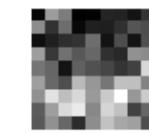
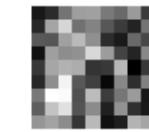
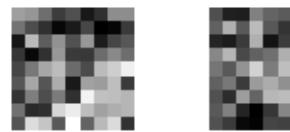
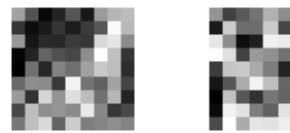
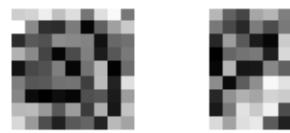
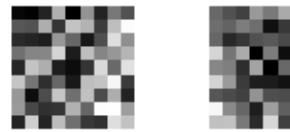
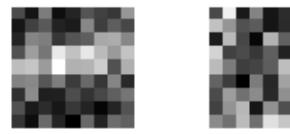
Question 3c



Compare the filters learned and the corresponding feature maps.

What do you see?

Input 28x28



Filters learned in the conv layer

Feature maps of the convolution layer 25x20x20

Question 3d

Gradient descent:

$$\mathbf{W} = \mathbf{W} - \alpha \nabla_{\mathbf{W}} J$$

Weight decay (regularization):

$$J_1(\mathbf{W}, \mathbf{b}) = J(\mathbf{W}, \mathbf{b}) + \beta_2 \sum_{ij} (w_{ij})^2$$

Weight decay:

$$\mathbf{W} = \mathbf{W} - \alpha (\nabla_{\mathbf{W}} J + \beta \mathbf{W})$$

Sgd with momentum:

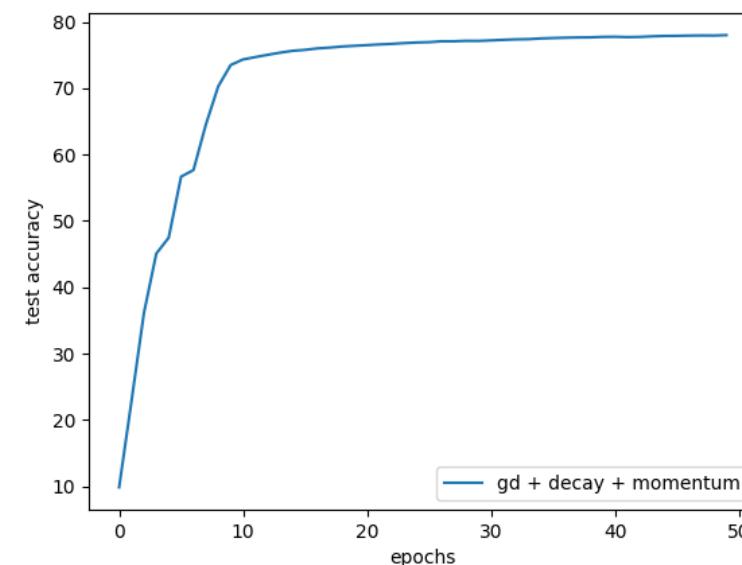
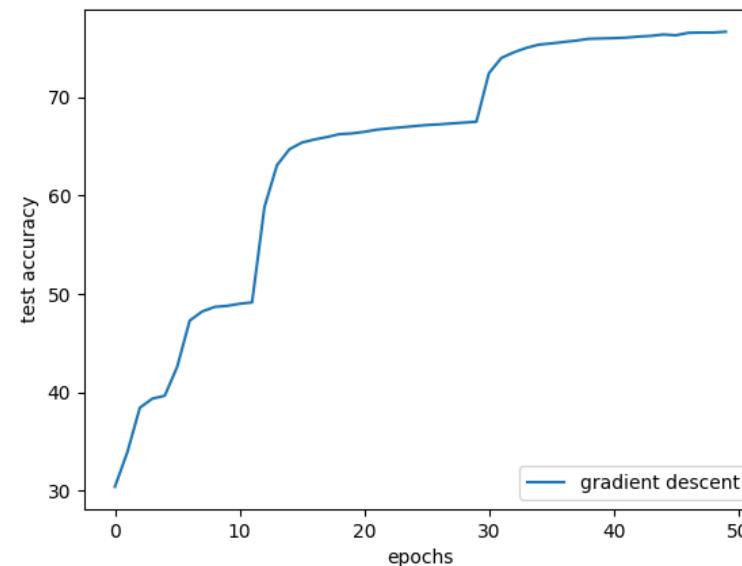
$$\mathbf{V} = \gamma \mathbf{V} - \alpha \nabla_{\mathbf{W}} J$$

$$\mathbf{W} = \mathbf{W} + \mathbf{V}$$

Sgd with decay and momentum:

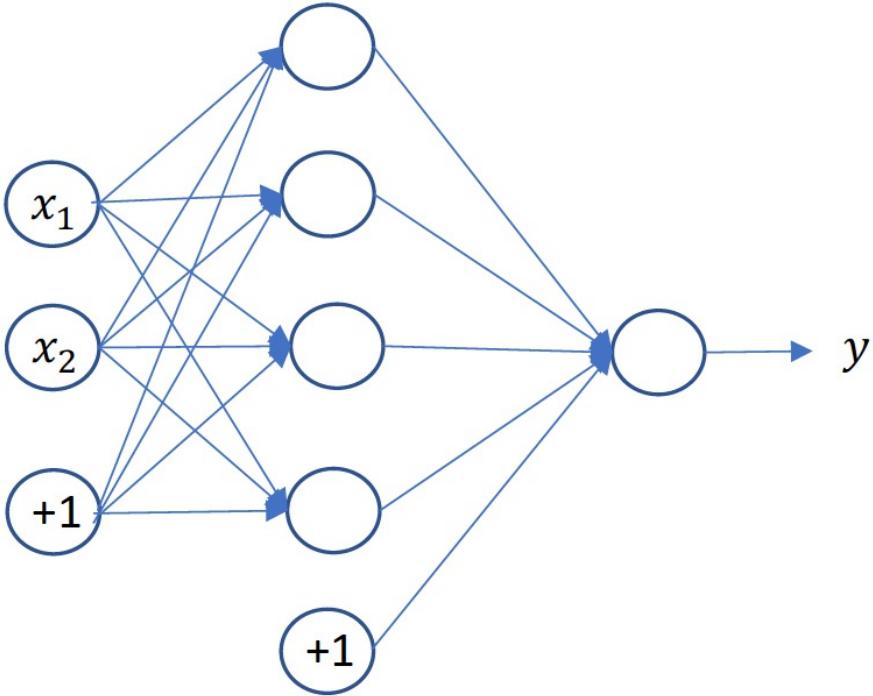
$$\mathbf{V} = \gamma \mathbf{V} - \alpha (\nabla_{\mathbf{W}} J + \beta \mathbf{W})$$

$$\mathbf{W} = \mathbf{W} + \mathbf{V}$$



Model selection and overfitting

CS4001 – Tutorial 5



1. A feedforward network consisting of one hidden layer of perceptrons and a linear output neuron receives inputs two-dimensional inputs (x_1, x_2) . Train the network to predict the following function:

$$y = \sin(\pi x_1) \cos(2\pi x_2)$$

where $-1.0 \leq x_1, x_2 \leq +1.0$.

Use data points distributed in an equally spaced 10×10 grid spanning the input space and the following procedures for training and testing:

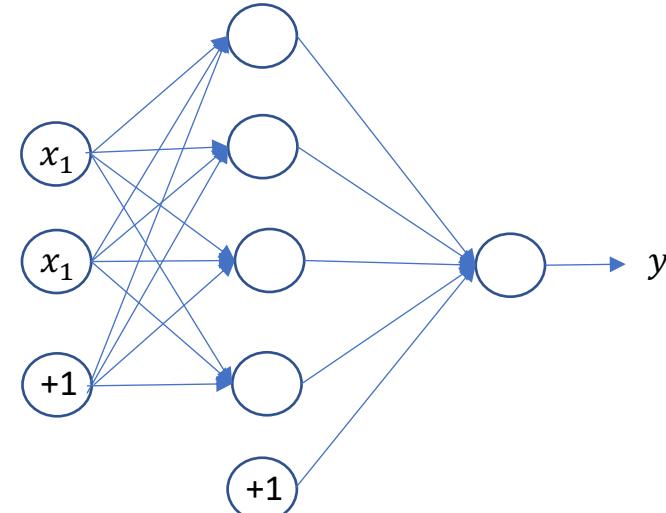
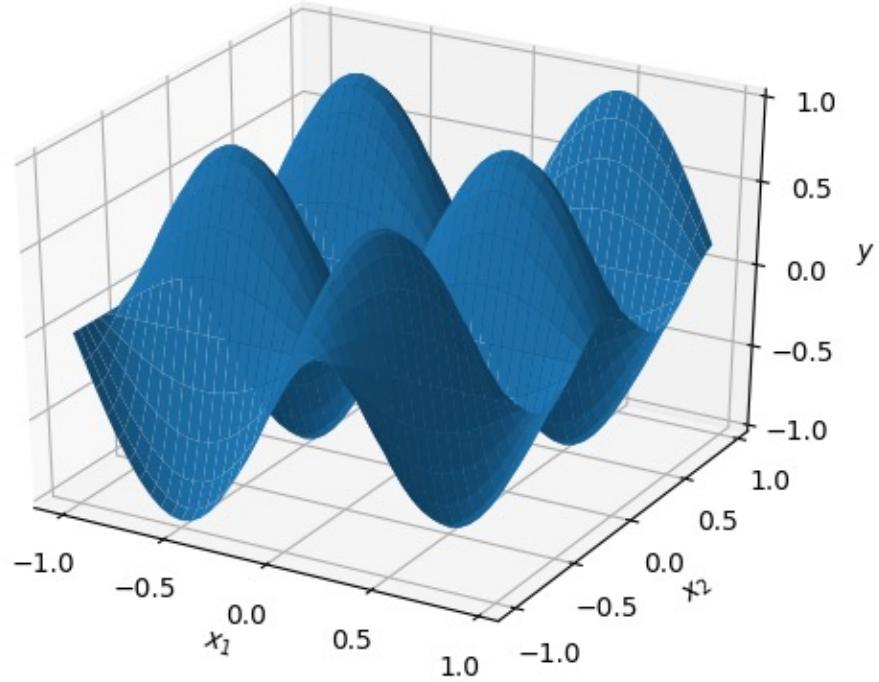
- a) Random subsampling with 70:30 data-split for training and testing.
- b) Five-fold cross validation
- c) Three-way data split

Repeat each procedure in 10 different experiments and determine the optimal number of hidden neurons in the space of $\{2, 4, 6, 8, 10\}$ and the error of the model.

Use a learning factor $\alpha = 0.05$ and early stopping to cease the learning epochs.

$$y = \sin(\pi x_1) \cos(2\pi x_2) \quad \text{where } -1.0 \leq x_1, x_2 \leq +1.0.$$

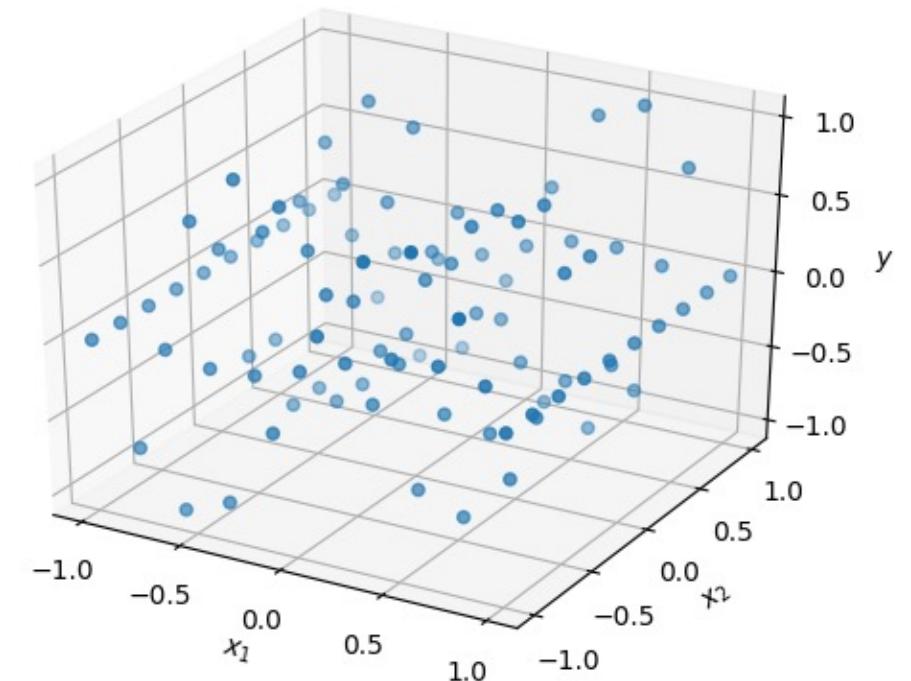
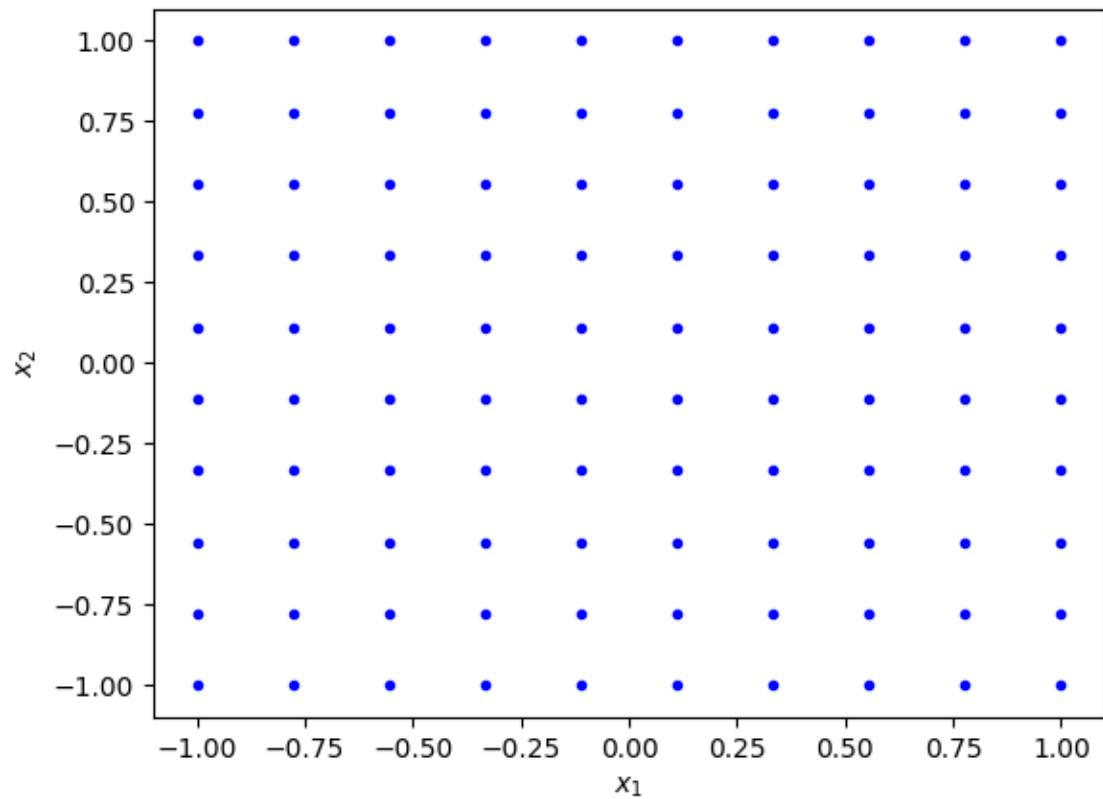
Note that $-1.0 \leq y \leq +1$.



Hidden neurons perceptrons
the output neurons is a linear neuron.

$$y = \sin(\pi x_1) \cos(2\pi x_2) \quad \text{where } -1.0 \leq x_1, x_2 \leq +1.0.$$

Training data is in an input grid of 10x10



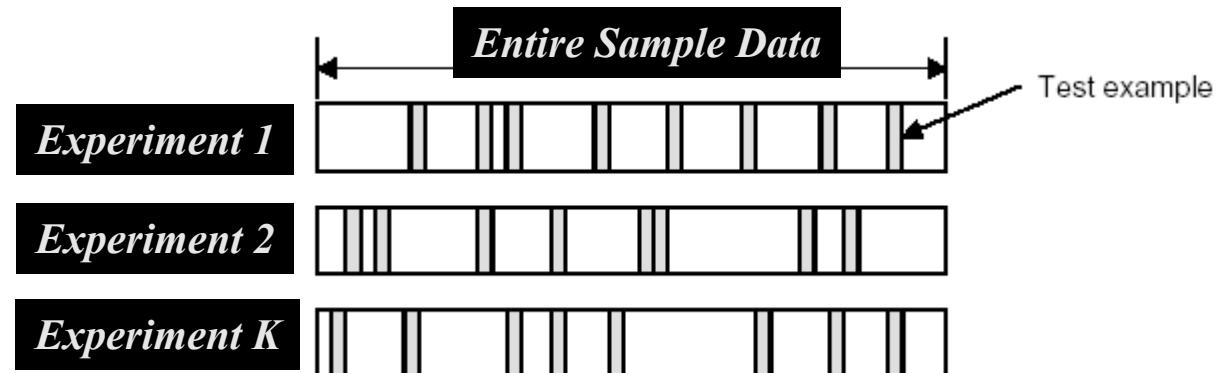
```
no_labels = 1
no_features = 2
no_exps = 10

class MLP(nn.Module):
    def __init__(self, no_features, no_hidden, no_labels):
        super().__init__()
        self.mlp_stack = nn.Sequential(
            nn.Linear(no_features, no_hidden),
            nn.Sigmoid(),
            nn.Linear(no_hidden, no_labels),
        )

    def forward(self, x):
        logits = self.mlp_stack(x)
        return logits
```

```
class EarlyStopper:  
    def __init__(self, patience=10, min_delta=0):  
        self.patience = patience  
        self.min_delta = min_delta  
        self.counter = 0  
        self.min_validation_loss = np.inf  
  
    def early_stop(self, validation_loss):  
        if validation_loss < self.min_validation_loss:  
            self.min_validation_loss = validation_loss  
            self.counter = 0  
        elif validation_loss > (self.min_validation_loss + self.min_delta):  
            self.counter += 1  
            if self.counter >= self.patience:  
                return True  
        return False
```

K Data Splits: Random Subsampling



For each experiment k :

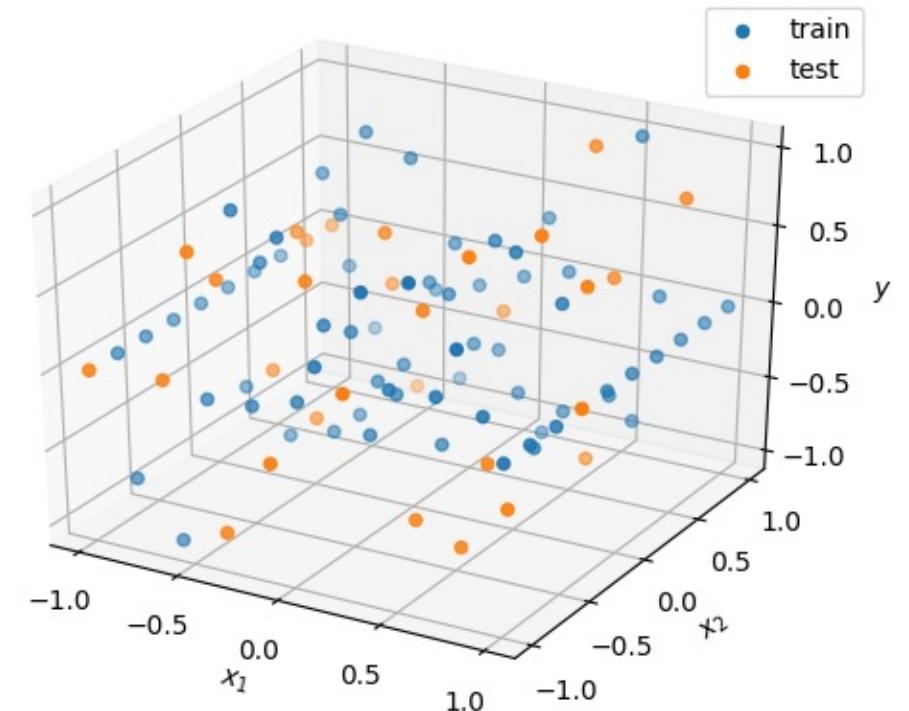
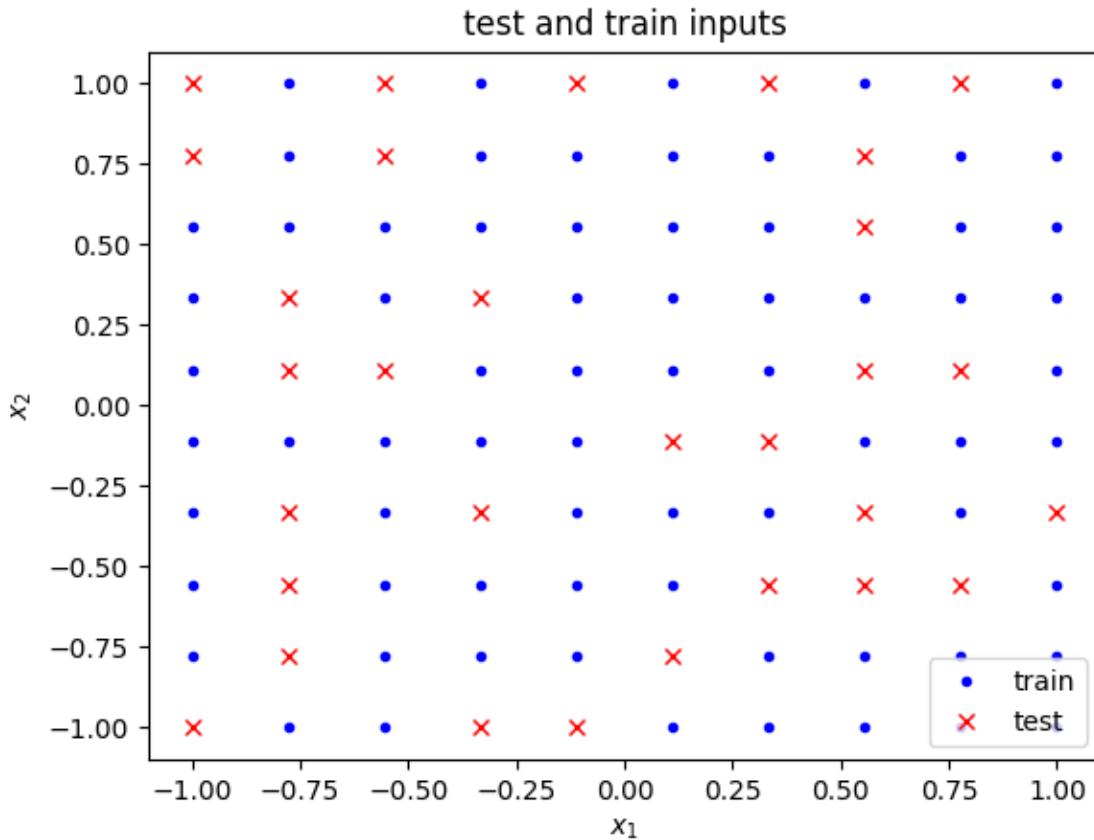
For each model m :

Compute error $e_{k,m}$

Compute mean error $e_m = \frac{1}{K} \sum_{k=1}^K e_{k,m}$

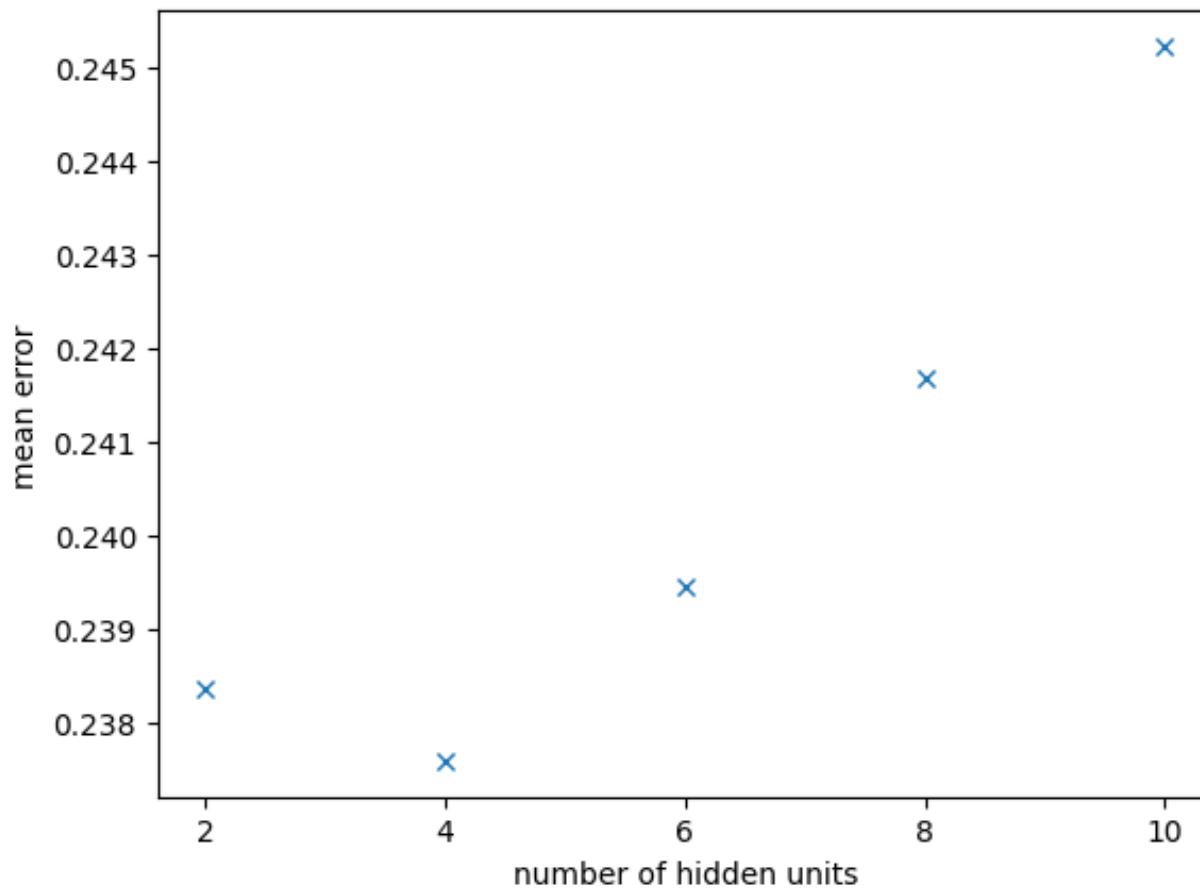
Optimal model $m^* = \operatorname{argmin} e_m$

For each experiment, data is split into train and test data at [30: 70]



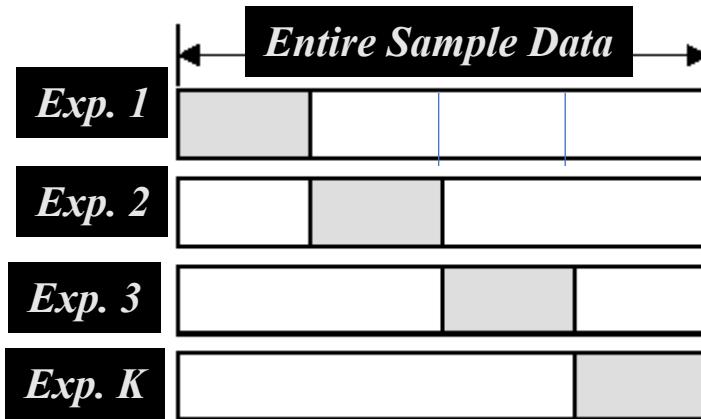
For each experiment, we build models with hidden neurons {2, 4, 6, 8, 10}

Mean error of 10 experiments



Optimum number of hidden neurons = 4

K-fold Cross Validation



For every fold f :

For every model m

Train the model, using data not in fold f

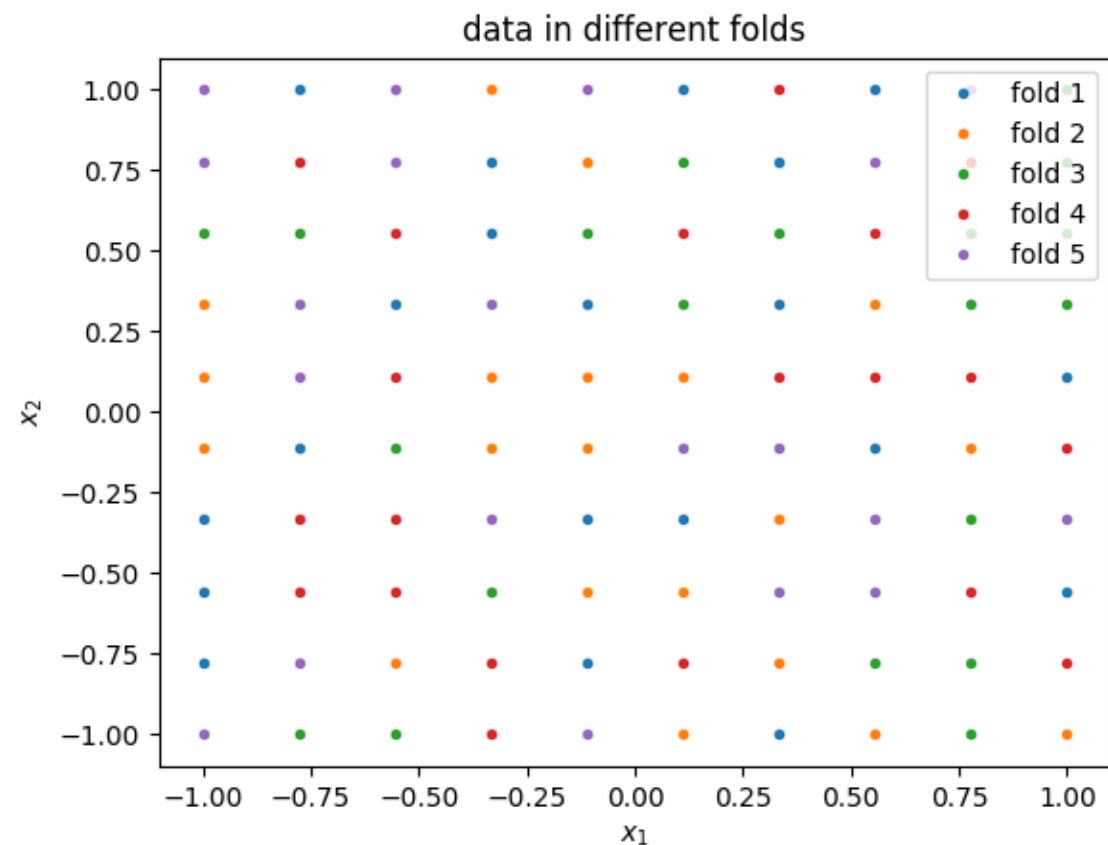
Compute error $e_{m,f}$ on data in fold f

For every model m

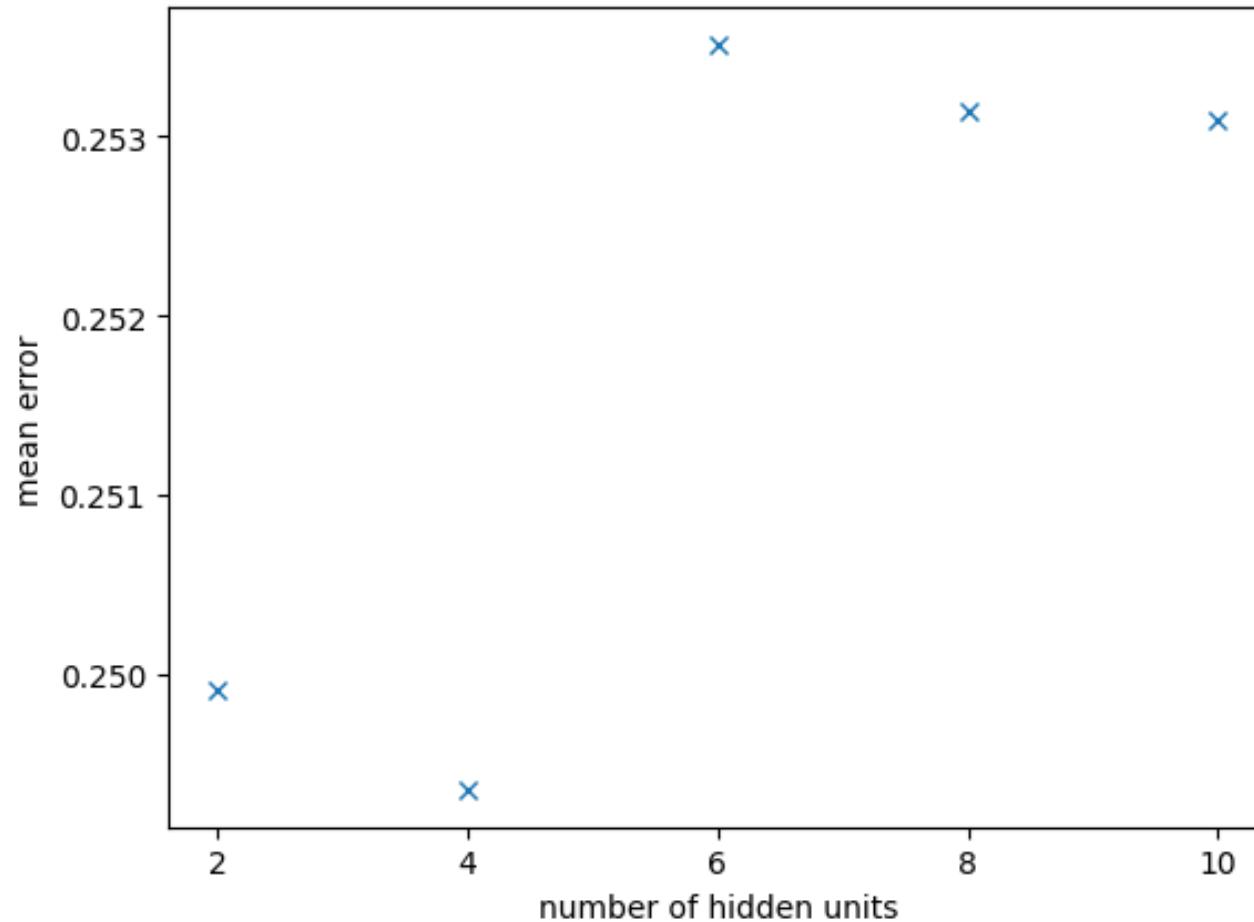
...

$$\text{CV error } e_m = \frac{1}{F} \sum_f e_{m,f}$$

Select the model with minimum CV error, $m^* = \operatorname{argmin} e_m$

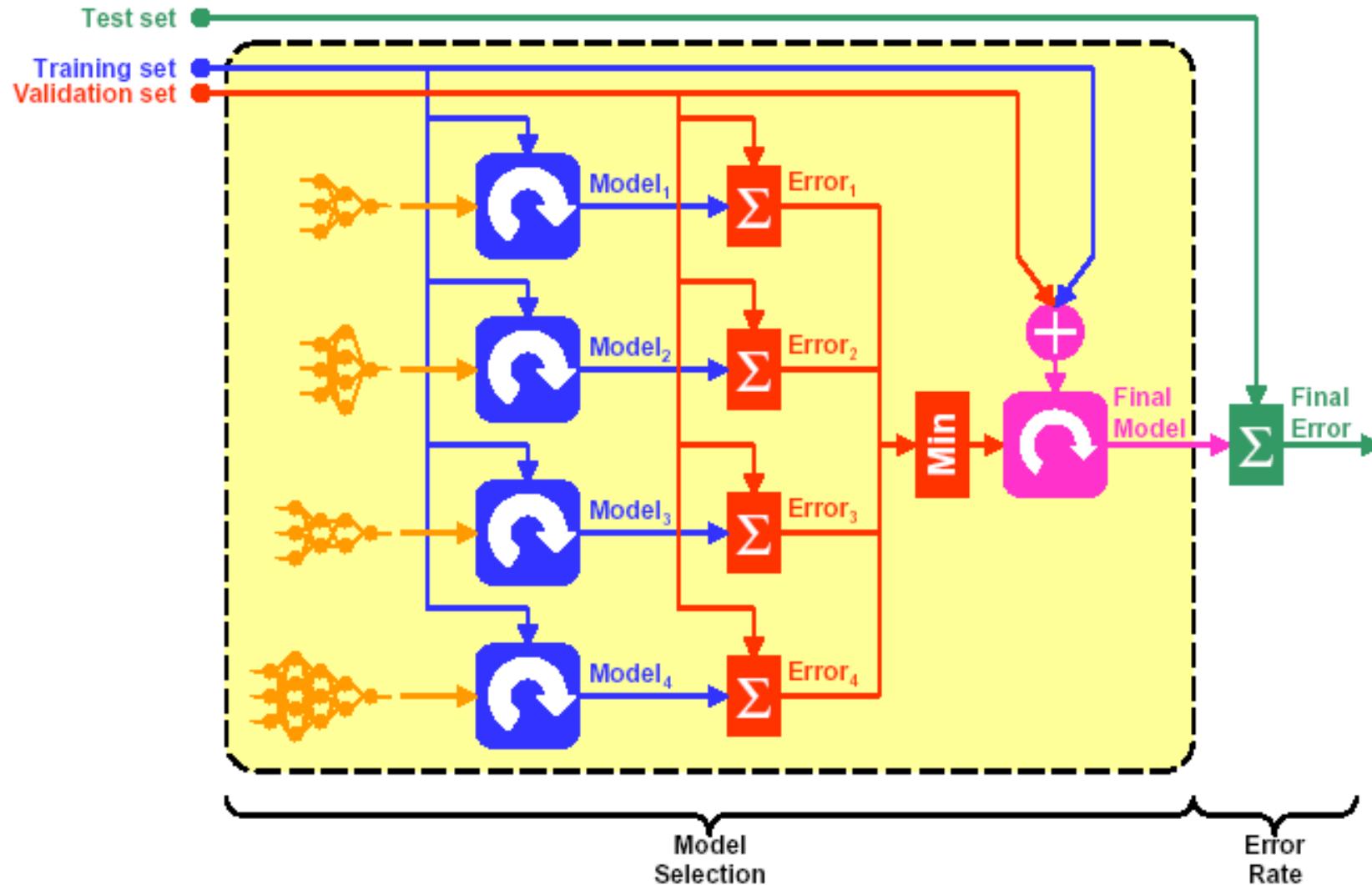


Mean cross-validation error of 10 experiments



Optimum number of hidden neurons = **4**

Three-Way Data Splits Method



Three-Way Data Splits Method

- **Training set:** Each model m , train the network and find the minimum error:

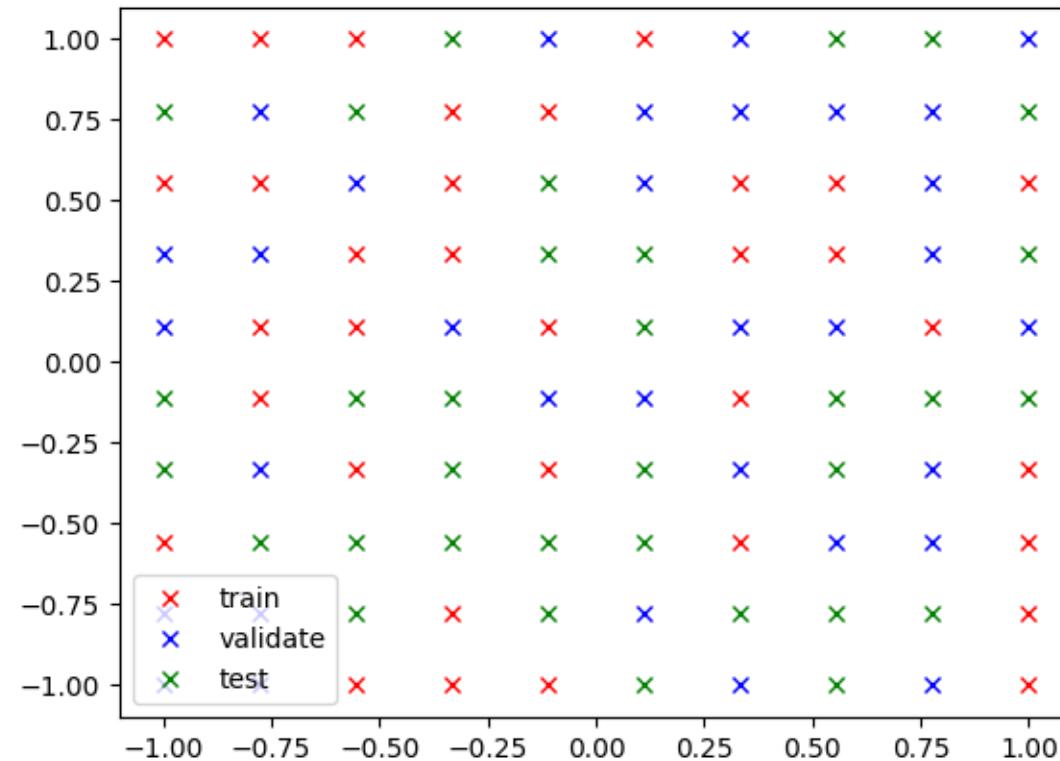
$$\mathbf{W}_m, \mathbf{b}_m = \underset{\mathbf{W}, \mathbf{b}}{\operatorname{argmin}} J$$

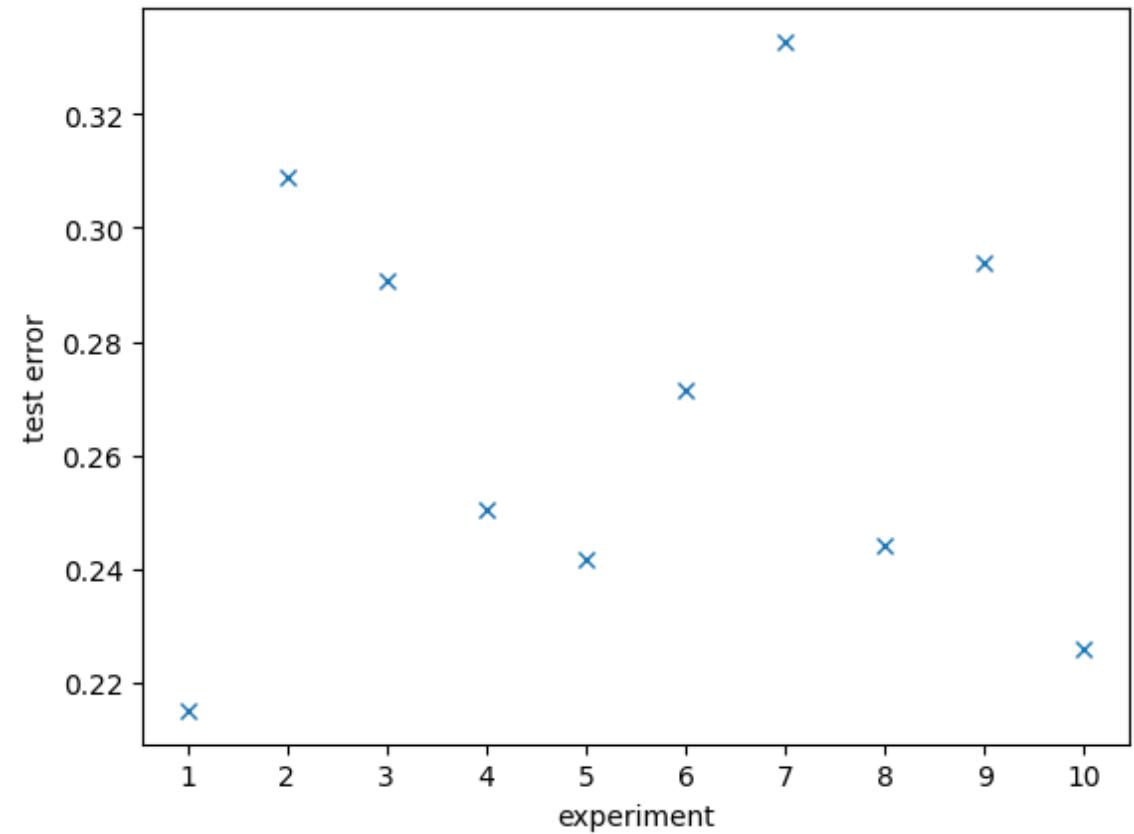
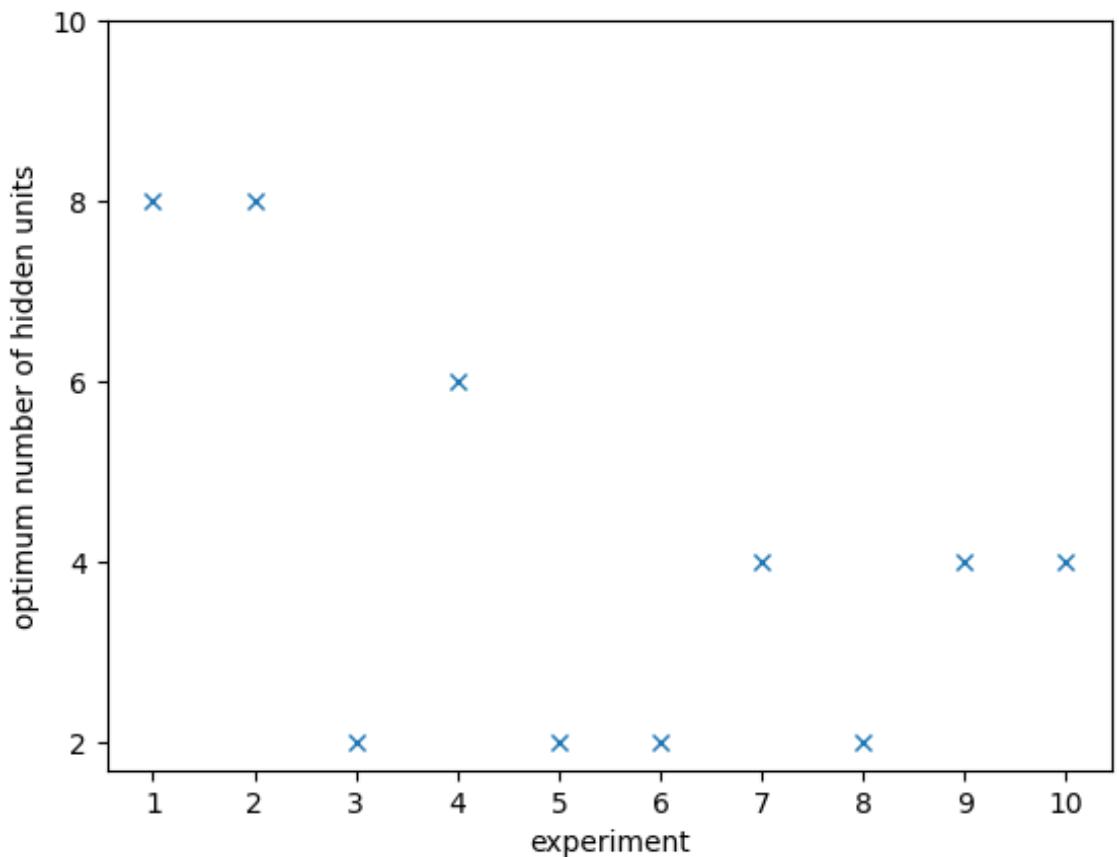
- **Validation set:** The optimal model m^* is given by

$$m^* = \underset{m}{\operatorname{argmin}} J_m$$

- **Training + Validation set:** Train the optimal model m^*
- **Test set:** determine test error on trained m^*

Training, validation and test inputs





Optimal number of hidden neurons is 2 (Max number of times)

Average m.s.e.: 0.215 (by taking average of test error when hidden neurons = 2)

2. The CIFAR-10 dataset consists of 60,000 32x32 colour images from 10 classes, with 6,000 images per class:

<https://www.cs.toronto.edu/~kriz/cifar.html>

There are 50,000 training images and 10,000 test images. Read CIFAR-10 datasets from torchvision.datasets:

```
from torchvision.datasets import CIFAR10
```

- a) Build three DNN of different complexity with three hidden layers to classify the images:
- i. A **small network** with 50 neurons in every hidden layer
 - ii. A **medium network** with 100 neurons in every hidden layer
 - iii. A **large network** with 500 neurons in every layer

Plot cross-entropies against epochs for train and test datasets. Use the Adams optimizer with default parameter for training and early stopping criterion to terminate. Comment on the overfitting and underfitting of the networks if any.

The CIFAR-10 dataset: <https://www.cs.toronto.edu/~kriz/cifar.html>

50000 training and 10000 test images; Images are 32x32 colour images

```
from torchvision.datasets import CIFAR10
```



```
class NeuralNetwork(nn.Module):
    def __init__(self, hidden_size=100):
        super(NeuralNetwork, self).__init__()
        self.flatten = nn.Flatten()
        self.linear_relu_stack = nn.Sequential(
            nn.Linear(32*32*3, hidden_size),
            nn.ReLU(),
            nn.Linear(hidden_size, hidden_size),
            nn.ReLU(),
            nn.Linear(hidden_size, 10),
            nn.Softmax(dim=1)
        )

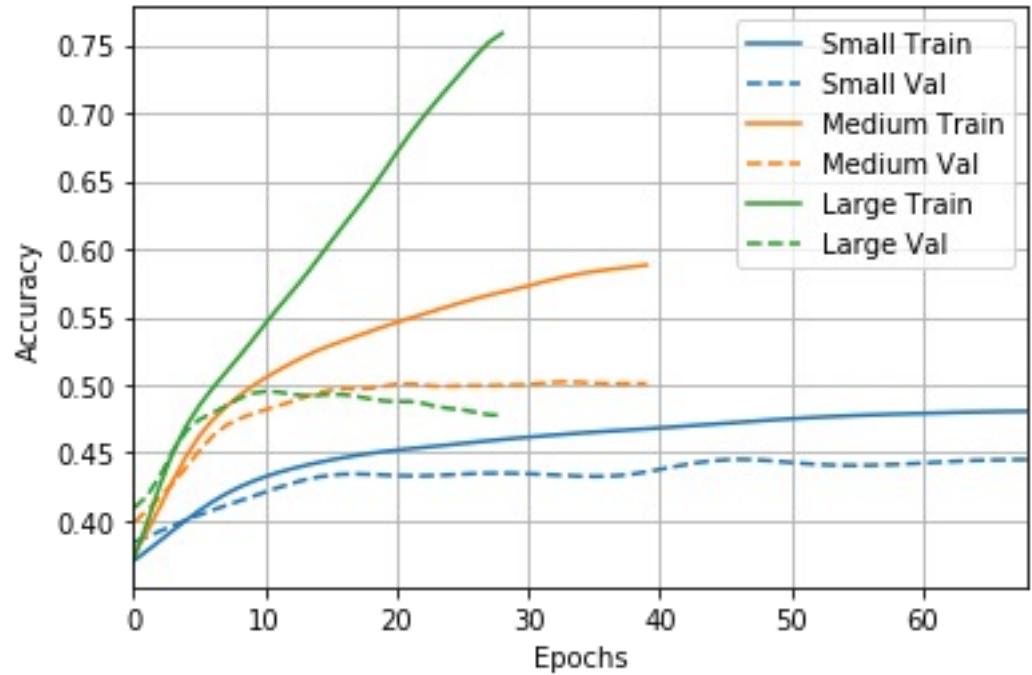
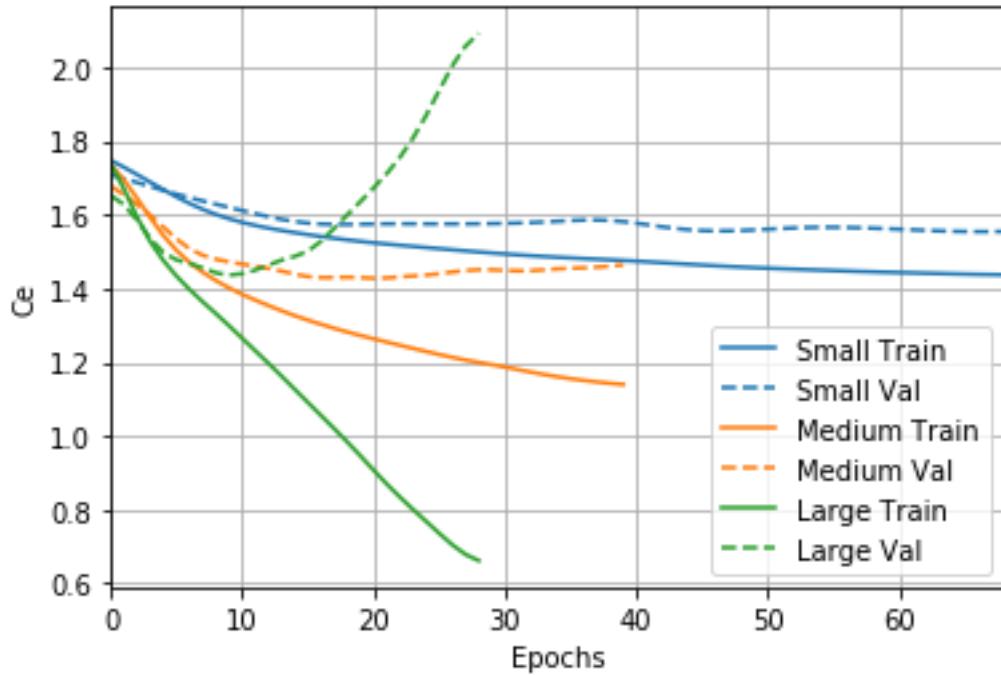
    def forward(self, x):
        x = self.flatten(x)
        logits = self.linear_relu_stack(x)
        return logits
```

```
model = NeuralNetwork(hidden_size=50)

loss_fn = nn.CrossEntropyLoss()
optimizer = torch.optim.Adam(model.parameters())
early_stopper = EarlyStopper(patience=10, min_delta=0)

for t in range(max_epochs):
    train_loss, train_correct = train_loop(train_dataloader, model, loss_fn, optimizer)
    test_loss, test_correct = test_loop(test_dataloader, model, loss_fn)

    if early_stopper.early_stop(test_loss):
        print("Done!")
        break
```

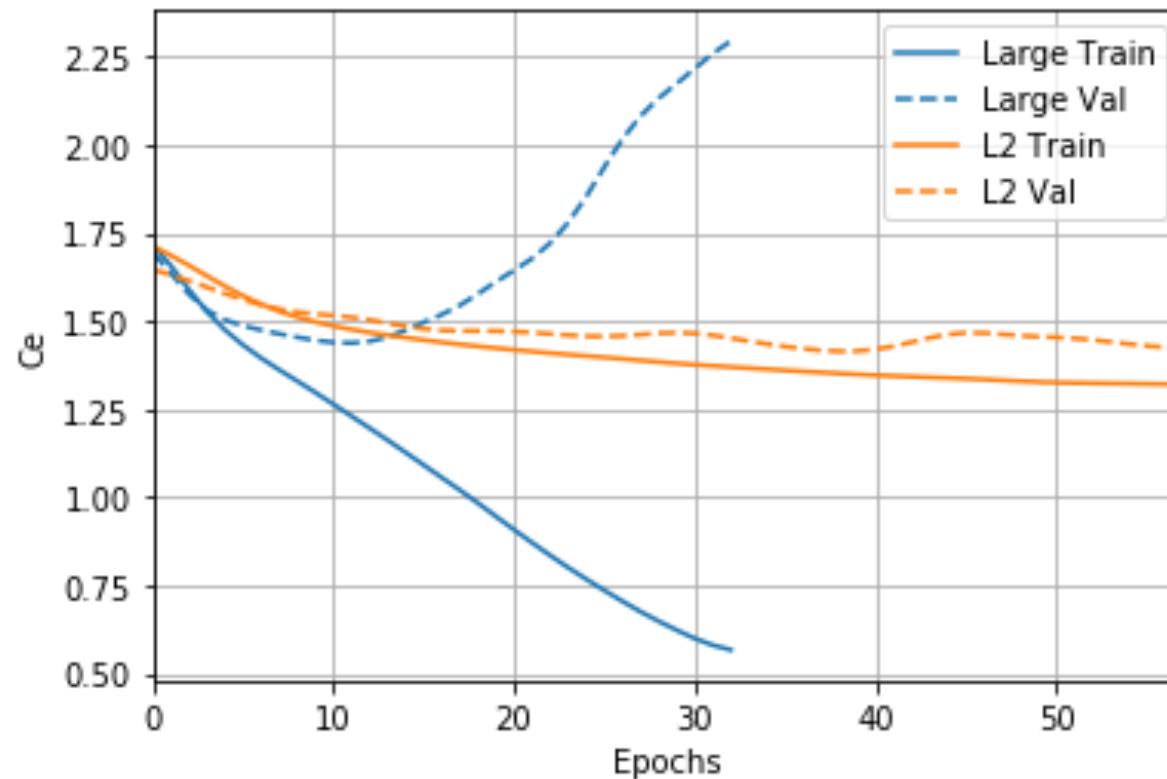


Small model is **underfitting**
Medium model seems just right
Large model is **overfitting**

- b) Using the **large network**, demonstrate how the following methods applied to all the layers could overcome overfitting of the network:
- i. Weight regularization with L2 norm. Use weight decay parameter $\beta = 0.001$.
 - ii. Dropouts at a probability $p = 0.5$.
 - iii. Combining both weight regularization and dropouts from (i) and (ii), respectively.

L2 weight regularization

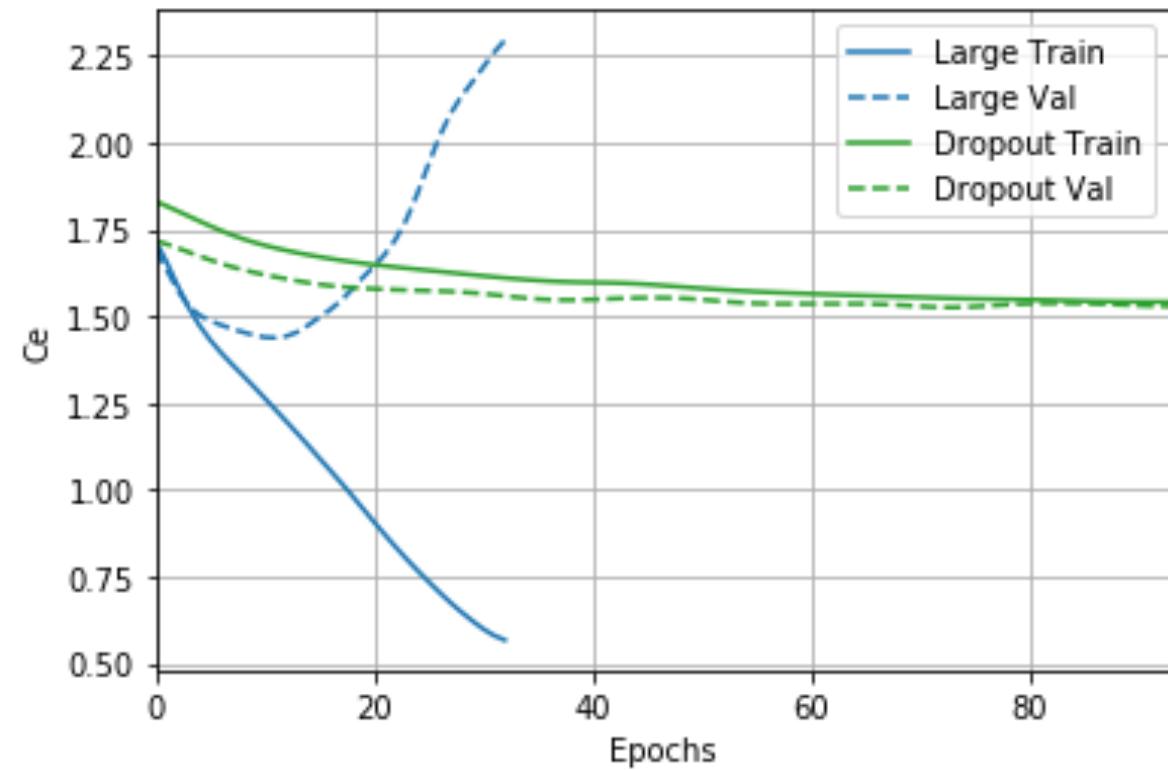
```
optimizer = torch.optim.Adam(model.parameters(), weight_decay=0.001)
```



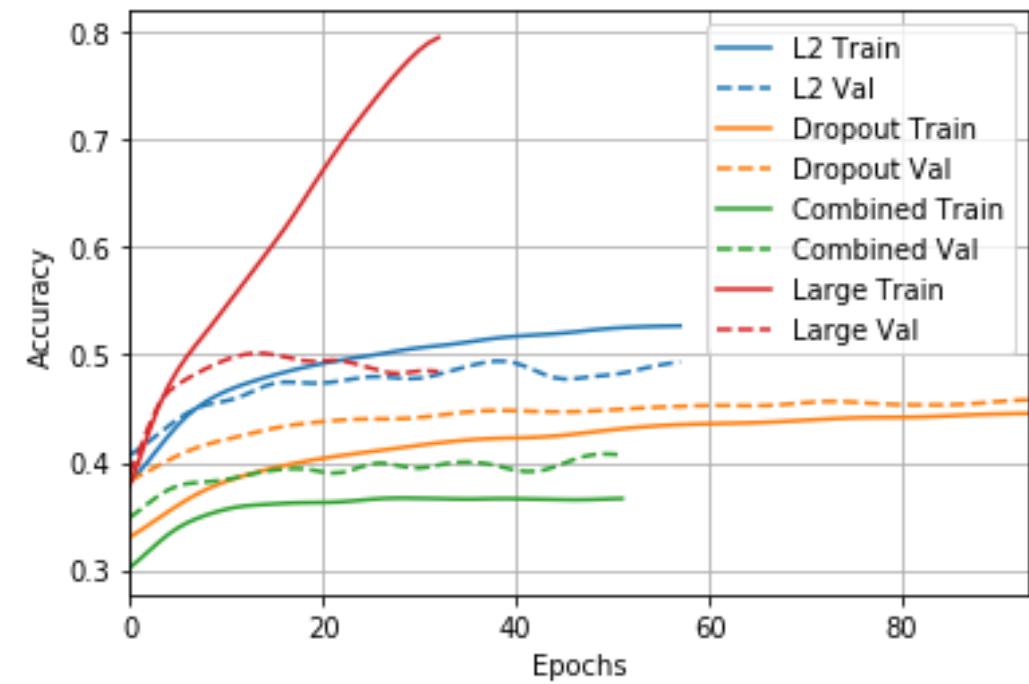
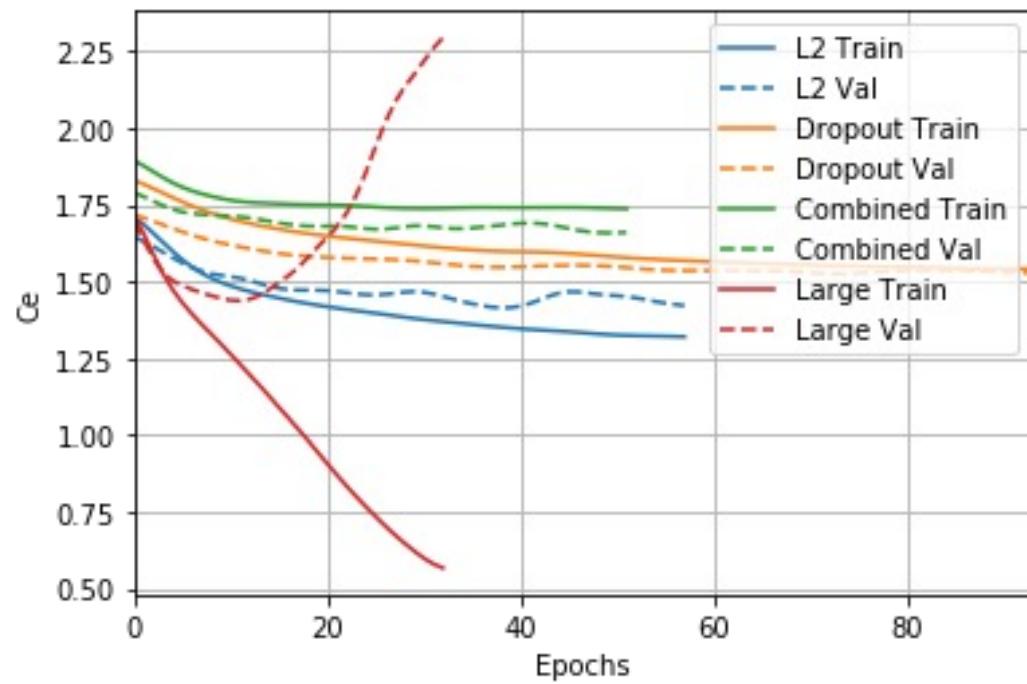
Dropouts

```
class NeuralNetwork_dropout(nn.Module):
    def __init__(self, hidden_size = 500, drop_out=0.5):
        super(NeuralNetwork_dropout, self).__init__()
        self.flatten = nn.Flatten()
        self.linear_relu_stack = nn.Sequential(
            nn.Linear(32*32*3, hidden_size),
            nn.ReLU(),
            nn.Dropout(p=drop_out),
            nn.Linear(hidden_size, hidden_size),
            nn.ReLU(),
            nn.Dropout(p=drop_out),
            nn.Linear(hidden_size, 10),
            nn.Softmax(dim=1)
        )

    def forward(self, x):
        x = self.flatten(x)
        logits = self.linear_relu_stack(x)
        return logits
```



Combined: L2 regularization + dropouts



Tutorial 7

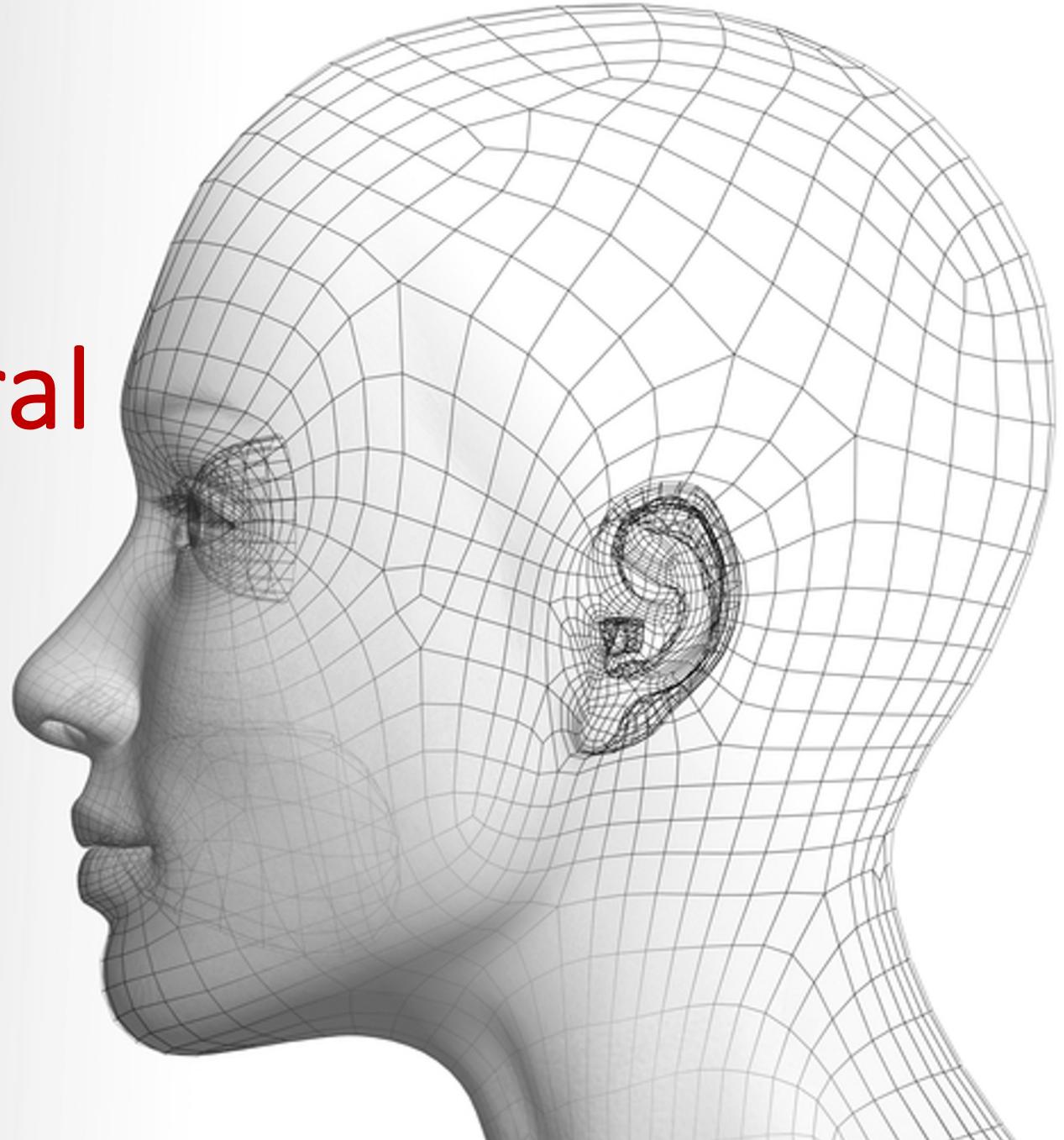
Convolutional Neural Networks II

Xingang Pan

潘新钢

<https://xingangpan.github.io/>

<https://twitter.com/XingangP>



Reminder: Feedback on Teaching

Dear Student,

The Online Faculty Teaching Evaluation for Semester 2 AY23-24 is now open.

Please [click here](#) to provide the feedback. The link is valid from 11-MAR-2024 to 24-MAR-2024.

We would like to strongly encourage you to participate in this exercise. Giving teaching feedback to the School will benefit the student population as we work towards refining and improving our teaching pedagogy and course materials. Together we can develop a more nurturing and effective learning environment for every student.

Thank you and with Best Regards,

SFT Administrator

Question 1

Figure Q1 depicts a block that consists of three convolutional layers. The input volume has a size of $256 \times 32 \times 32$ and the second layer has 32 convolution filters each with a size of $64 \times 3 \times 3$, stride = 1 and padding = 1.

Provide the values of n_1 , d_1 , F_1 , n_2 , d_2 , and F_2 to form a valid block.
Explain your design

Answer is hidden so that you can try the question by yourself first

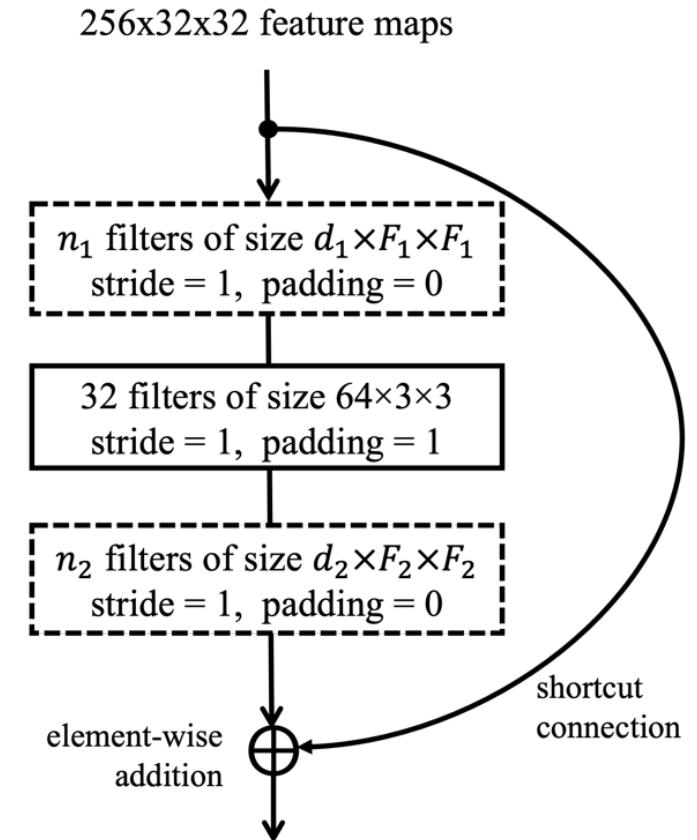


Figure Q1

Question 1

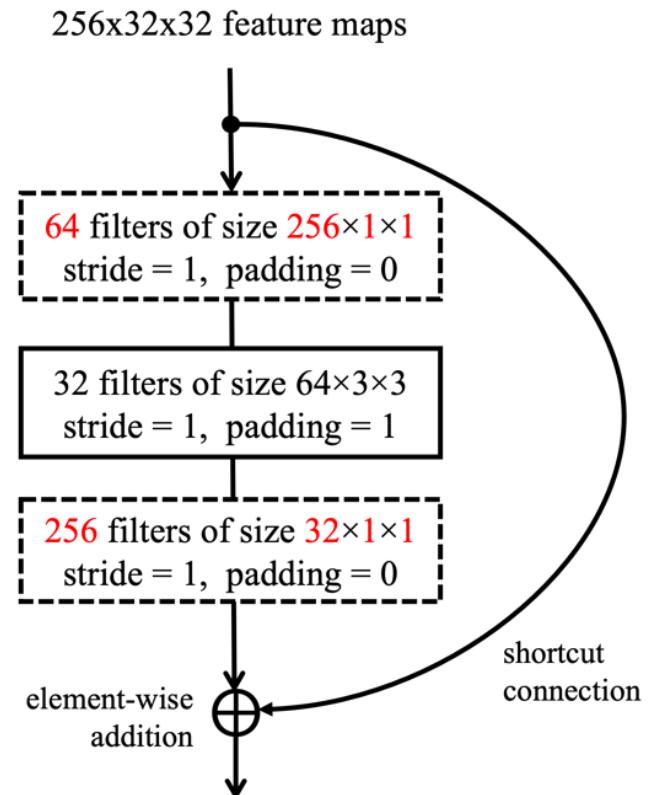
Answer:

$$n_1 = 64, d_1 = 256, F_1 = 1$$

$$n_2 = 256, d_2 = 32, F_2 = 1$$

To form a valid block, the size the of output volume at the residual branch has to be the same size as the input volume, which is $256 \times 32 \times 32$, such that element-wise addition can be performed, thus $n_2 = 256$ and 1×1 is chosen to keep the spatial resolution.

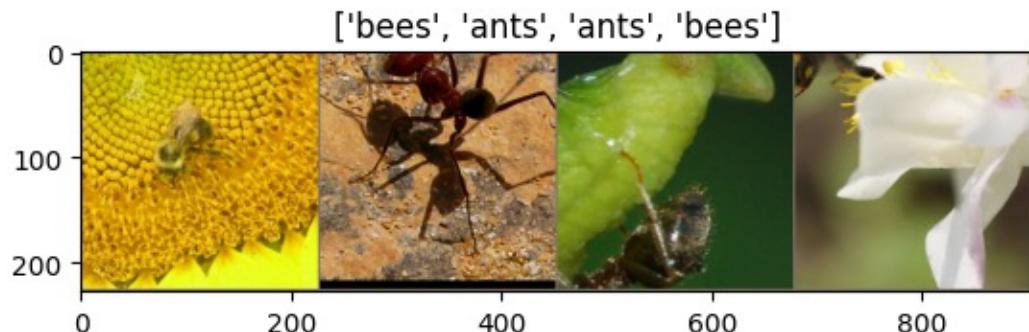
The values of d_1 and d_2 are chosen to match the depth of their corresponding input. n_1 is chosen to match the filter size of the second layer.



Question 2

Study and try the tutorial t7q2.ipynb on transfer learning. In particular,

- Understand how data augmentation is performed
- Review the transfer learning steps
- Try the code to perform transfer learning on the classification of bees vs. ants



Check the code for more details

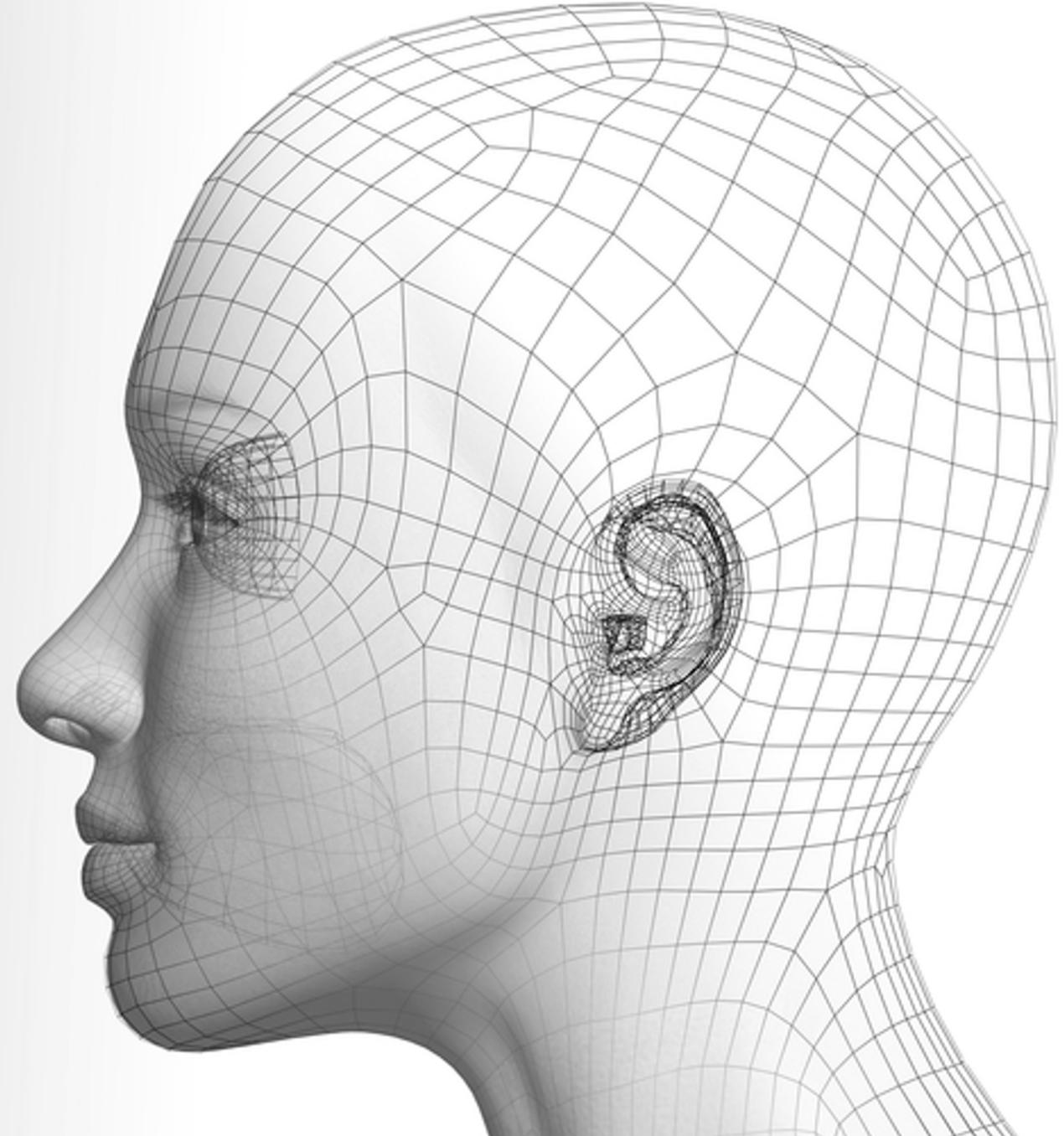
Tutorial 8

Recurrent Neural

Networks

Xingang Pan

潘新钢



Question 1

1. A recurrent neural network (RNN) receives sequences of 3-dimensional inputs and has two hidden neurons and one output neuron. The weight matrix \mathbf{U} connecting the input to the hidden layer, the weight matrix \mathbf{V} connecting the hidden output to the output layer, the hidden layer bias \mathbf{b} and the output layer bias \mathbf{c} are given by

$$\mathbf{U} = \begin{pmatrix} -1.0 & 0.5 \\ 0.5 & 0.1 \\ 0.2 & -2.0 \end{pmatrix}, \mathbf{V} = \begin{pmatrix} 2.0 \\ -1.5 \end{pmatrix}, \mathbf{b} = \begin{pmatrix} 0.2 \\ 0.2 \end{pmatrix}, \text{ and } \mathbf{c} = 0.5.$$

Find the output sequence for an input sequence of $(\mathbf{x}(t))_{t=1}^4$ where

$$\mathbf{x}(1) = \begin{pmatrix} 1 \\ 2 \\ -1 \end{pmatrix}, \mathbf{x}(2) = \begin{pmatrix} -1 \\ 1 \\ -2 \end{pmatrix}, \mathbf{x}(3) = \begin{pmatrix} 0 \\ 3 \\ 1 \end{pmatrix}, \text{ and } \mathbf{x}(4) = \begin{pmatrix} 2 \\ -1 \\ 0 \end{pmatrix} \text{ if}$$

- (a) The RNN is of hidden-recurrence (Elman) type with the recurrence weight matrix \mathbf{W} connecting the previous hidden output to the current hidden layer input is given by

$$\mathbf{W} = \begin{pmatrix} 2.0 & 1.3 \\ 1.5 & 0.0 \end{pmatrix}.$$

Assume that the hidden activations are initialized to zero.

- (b) The RNN is of top-down recurrence (Jordan) type with the weight matrix \mathbf{W} connecting the previous output to the current state of hidden layer is given by

$$\mathbf{W} = (2.0 \quad 1.3).$$

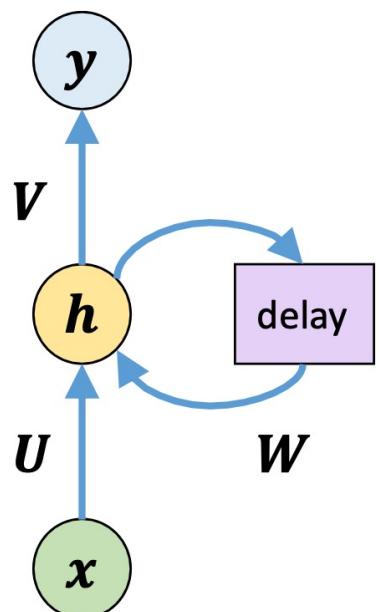
Assume that the output activations are initialized to zero.

Question 1

Hidden Recurrence:

$$\mathbf{U} = \begin{pmatrix} -1.0 & 0.5 \\ 0.5 & 0.1 \\ 0.2 & -2.0 \end{pmatrix}, \mathbf{V} = \begin{pmatrix} 2.0 \\ -1.5 \end{pmatrix} \text{ and } \mathbf{W} = \begin{pmatrix} 2.0 & 1.3 \\ 1.5 & 0.0 \end{pmatrix}.$$

$$\mathbf{b} = \begin{pmatrix} 0.2 \\ 0.2 \end{pmatrix}, \text{ and } \mathbf{c} = 0.5$$



Three input neurons, two hidden neurons, and one output neuron.

$$\mathbf{h}(t) = \phi(\mathbf{U}^T \mathbf{x}(t) + \mathbf{W}^T \mathbf{h}(t-1) + \mathbf{b})$$

$$\mathbf{y}(t) = \sigma(\mathbf{V}^T \mathbf{h}(t) + \mathbf{c})$$

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

$$\sigma(u) = \text{sigmoid}(u) = \frac{1}{1+e^{-u}}.$$

Assume $\mathbf{h}(0) = (0 \quad 0)^T$.

Question 1

At $t=1$, $\mathbf{x}(1) = \begin{pmatrix} 1 \\ 2 \\ -1 \end{pmatrix}$:

$$\begin{aligned}\mathbf{h}(1) &= \phi(\mathbf{U}^T \mathbf{x}(1) + \mathbf{W}^T \mathbf{h}(0) + \mathbf{b}) \\ &= \tanh\left(\begin{pmatrix} -1.0 & 0.5 & 0.2 \\ 0.5 & 0.1 & -2.0 \end{pmatrix} \begin{pmatrix} 1 \\ 2 \\ -1 \end{pmatrix} + \begin{pmatrix} 2.0 & 1.5 \\ 1.3 & 0.0 \end{pmatrix} \begin{pmatrix} 0 \\ 0 \end{pmatrix} + \begin{pmatrix} 0.2 \\ 0.2 \end{pmatrix}\right) = \begin{pmatrix} 0.0 \\ 0.99 \end{pmatrix}\end{aligned}$$

$$\mathbf{y}(1) = \sigma(\mathbf{V}^T \mathbf{h}(1) + \mathbf{c}) = \text{sigmoid}\left((2.0 \quad -1.5) \begin{pmatrix} 0.0 \\ 0.99 \end{pmatrix} + (0.5)\right) = (0.27)$$

Question 1

At $t=2$, $\mathbf{x}(2) = \begin{pmatrix} -1 \\ 1 \\ -2 \end{pmatrix}$:

$$\begin{aligned}\mathbf{h}(2) &= \phi(\mathbf{U}^T \mathbf{x}(2) + \mathbf{W}^T \mathbf{h}(1) + \mathbf{b}) \\ &= \tanh\left(\begin{pmatrix} -1.0 & 0.5 & 0.2 \\ 0.5 & 0.1 & -2.0 \end{pmatrix} \begin{pmatrix} -1 \\ 1 \\ -2 \end{pmatrix} + \begin{pmatrix} 2.0 & 1.5 \\ 1.3 & 0.0 \end{pmatrix} \begin{pmatrix} 0.0 \\ 0.99 \end{pmatrix} + \begin{pmatrix} 0.2 \\ 0.2 \end{pmatrix}\right) = \begin{pmatrix} 0.99 \\ 1.0 \end{pmatrix}\end{aligned}$$

$$\mathbf{y}(2) = \sigma(\mathbf{V}^T \mathbf{h}(2) + \mathbf{c}) = \text{sigmoid}\left((2.0 \quad -1.5) \begin{pmatrix} 0.99 \\ 1.0 \end{pmatrix} + (0.5)\right) = (0.73)$$

Similarly,

$$\text{at } t=3, \mathbf{x}(3) = \begin{pmatrix} 0 \\ 3 \\ 1 \end{pmatrix}; \mathbf{h}(3) = \begin{pmatrix} 1.0 \\ -0.21 \end{pmatrix} \text{ and } \mathbf{y}(3) = (0.94)$$

$$\text{at } t=4, \mathbf{x}(4) = \begin{pmatrix} 2 \\ -1 \\ 0 \end{pmatrix}; \mathbf{h}(4) = \begin{pmatrix} -0.54 \\ 0.98 \end{pmatrix} \text{ and } \mathbf{y}(4) = (0.11)$$

Question 1

The input sequence $(\mathbf{x}(1), \mathbf{x}(2), \mathbf{x}(3), \mathbf{x}(4))$:

$$\left(\begin{pmatrix} 1 \\ 2 \\ -1 \end{pmatrix} \quad \begin{pmatrix} -1 \\ 1 \\ -2 \end{pmatrix} \quad \begin{pmatrix} 0 \\ 3 \\ 1 \end{pmatrix} \quad \begin{pmatrix} 2 \\ -1 \\ 0 \end{pmatrix} \right)$$

The output sequence $(\mathbf{y}(1), \mathbf{y}(2), \mathbf{y}(3), \mathbf{y}(4))$:

$$((0.27) \quad (0.73) \quad (0.94) \quad (0.11))$$

Question 1

1. A recurrent neural network (RNN) receives sequences of 3-dimensional inputs and has two hidden neurons and one output neuron. The weight matrix \mathbf{U} connecting the input to the hidden layer, the weight matrix \mathbf{V} connecting the hidden output to the output layer, the hidden layer bias \mathbf{b} and the output layer bias \mathbf{c} are given by

$$\mathbf{U} = \begin{pmatrix} -1.0 & 0.5 \\ 0.5 & 0.1 \\ 0.2 & -2.0 \end{pmatrix}, \mathbf{V} = \begin{pmatrix} 2.0 \\ -1.5 \end{pmatrix}, \mathbf{b} = \begin{pmatrix} 0.2 \\ 0.2 \end{pmatrix}, \text{ and } \mathbf{c} = 0.5.$$

Find the output sequence for an input sequence of $(\mathbf{x}(t))_{t=1}^4$ where

$$\mathbf{x}(1) = \begin{pmatrix} 1 \\ 2 \\ -1 \end{pmatrix}, \mathbf{x}(2) = \begin{pmatrix} -1 \\ 1 \\ -2 \end{pmatrix}, \mathbf{x}(3) = \begin{pmatrix} 0 \\ 3 \\ 1 \end{pmatrix}, \text{ and } \mathbf{x}(4) = \begin{pmatrix} 2 \\ -1 \\ 0 \end{pmatrix} \text{ if}$$

- (a) The RNN is of hidden-recurrence (Elman) type with the recurrence weight matrix \mathbf{W} connecting the previous hidden output to the current hidden layer input is given by

$$\mathbf{W} = \begin{pmatrix} 2.0 & 1.3 \\ 1.5 & 0.0 \end{pmatrix}.$$

Assume that the hidden activations are initialized to zero.

- (b) The RNN is of top-down recurrence (Jordan) type with the weight matrix \mathbf{W} connecting the previous output to the current state of hidden layer is given by

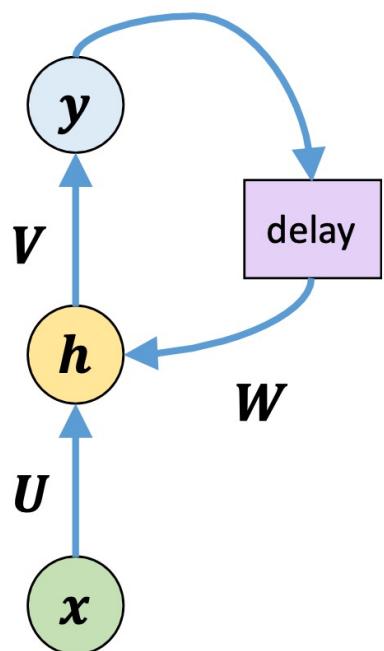
$$\mathbf{W} = (2.0 \quad 1.3).$$

Assume that the output activations are initialized to zero.

Question 1

Top-down Recurrence:

$$\mathbf{U} = \begin{pmatrix} -1.0 & 0.5 \\ 0.5 & 0.1 \\ 0.2 & -2.0 \end{pmatrix}, \mathbf{V} = \begin{pmatrix} 2.0 \\ -1.5 \end{pmatrix}, \mathbf{W} = (2.0 \quad 1.3)$$



$$\mathbf{b} = \begin{pmatrix} 0.2 \\ 0.2 \end{pmatrix}, \text{ and } \mathbf{c} = 0.5$$

Three input neurons, two hidden neurons, and one output neuron.

$$\begin{aligned}\mathbf{h}(t) &= \phi(\mathbf{U}^T \mathbf{x}(t) + \mathbf{W}^T \mathbf{y}(t-1) + \mathbf{b}) \\ \mathbf{y}(t) &= \sigma(\mathbf{V}^T \mathbf{h}(t) + \mathbf{c})\end{aligned}$$

Initially,

$$\mathbf{Y}(0) = (0.0)$$

Question 1

At $t=1$, $\mathbf{x}(1) = \begin{pmatrix} 1 \\ 2 \\ -1 \end{pmatrix}$:

$$\begin{aligned}\mathbf{h}(1) &= \phi(\mathbf{U}^T \mathbf{x}(1) + \mathbf{W}^T \mathbf{y}(0) + \mathbf{b}) \\ &= \tanh\left(\begin{pmatrix} -1.0 & 0.5 & 0.2 \\ 0.5 & 0.1 & -2.0 \end{pmatrix} \begin{pmatrix} 1 \\ 2 \\ -1 \end{pmatrix} + \begin{pmatrix} 2.0 \\ 1.3 \end{pmatrix}(0) + \begin{pmatrix} 0.2 \\ 0.2 \end{pmatrix}\right) = \begin{pmatrix} 0.0 \\ 0.99 \end{pmatrix}\end{aligned}$$

$$\mathbf{y}(1) = \sigma(\mathbf{V}^T \mathbf{h}(1) + \mathbf{c}) = \text{sigmoid}\left((2.0 \quad -1.5 \quad 0.2) \begin{pmatrix} 0.0 \\ 0.99 \end{pmatrix} + (0.5)\right) = (0.27)$$

Question 1

At $t=2$, $\mathbf{x}(2) = \begin{pmatrix} -1 \\ 1 \\ -2 \end{pmatrix}$:

$$\begin{aligned}\mathbf{h}(2) &= \phi(\mathbf{U}^T \mathbf{x}(2) + \mathbf{W}^T y(1) + \mathbf{b}) \\ &= \tanh\left(\begin{pmatrix} -1.0 & 0.5 & 0.2 \\ 0.5 & 0.1 & -2.0 \end{pmatrix} \begin{pmatrix} -1 \\ 1 \\ -2 \end{pmatrix} + \begin{pmatrix} 2.0 \\ 1.3 \end{pmatrix} (0.27) + \begin{pmatrix} 0.2 \\ 0.2 \end{pmatrix}\right) = \begin{pmatrix} 0.95 \\ 1.0 \end{pmatrix}\end{aligned}$$

$$\mathbf{y}(2) = \sigma(\mathbf{V}^T \mathbf{h}(2) + \mathbf{c}) = \text{sigmoid}\left((2.0 \quad -1.5 \quad 0.2) \begin{pmatrix} 0.95 \\ 1.0 \end{pmatrix} + (0.5)\right) = (0.71)$$

Similarly,

$$\text{at } t=3, \mathbf{x}(3) = \begin{pmatrix} 0 \\ 3 \\ 1 \end{pmatrix}; \mathbf{h}(3) = \begin{pmatrix} 1.0 \\ -0.52 \end{pmatrix} \text{ and } \mathbf{y}(3) = (0.96)$$

$$\text{at } t=4, \mathbf{x}(4) = \begin{pmatrix} 2 \\ -1 \\ 0 \end{pmatrix}; \mathbf{h}(4) = \begin{pmatrix} -0.36 \\ 0.98 \end{pmatrix} \text{ and } \mathbf{y}(4) = (0.16)$$

Question 1

The input sequence $(\mathbf{x}(1), \mathbf{x}(2), \mathbf{x}(3), \mathbf{x}(4))$:

$$\left(\begin{pmatrix} 1 \\ 2 \\ -1 \end{pmatrix} \quad \begin{pmatrix} -1 \\ 1 \\ -2 \end{pmatrix} \quad \begin{pmatrix} 0 \\ 3 \\ 1 \end{pmatrix} \quad \begin{pmatrix} 2 \\ -1 \\ 0 \end{pmatrix} \right)$$

The output sequence $(\mathbf{y}(1), \mathbf{y}(2), \mathbf{y}(3), \mathbf{y}(4))$:

$$((0.27) \quad (0.71) \quad (0.96) \quad (0.16))$$

Question 2

2. An RNN receives 3-dimensional input sequence and produces 2-dimensional output sequence. It has 5 hidden neurons and receives sequences of 100 time steps.
- Generate 8 training input sequences by drawing values uniformly between 0 and 1.0.
 - If the sequence $(\mathbf{y}(t))_{t=1}^{100}$ where $\mathbf{y}(t) = (y_k(t))_{k=1}^2 \in \mathbf{R}^2$ denotes the output sequence for an input sequence $(\mathbf{x}(t))_{t=1}^{100}$ where $\mathbf{x}(t) = (x_i(t))_{i=1}^3 \in \mathbf{R}^3$, generate the corresponding output sequences for the input training sequences as follows:

$$\begin{aligned}y_1(t) &= 5x_1(t) - 0.2x_3(t-1) + 0.1\epsilon \\y_2(t) &= 25x_2(t-1)x_3(t-3) + 0.1\epsilon\end{aligned}$$

where ϵ is standard normally distributed random variable.

Train an RNN to learn the above sequences by using gradient descent learning. Use a learning factor $\alpha = 0.01$ and Adam optimizer.

Question 2

▼ Set random seed and initialize hyper-parameters

```
✓ 0s [2] # Set random seed for reproducibility  
    seed = 10  
    np.random.seed(seed)  
    torch.manual_seed(seed)
```

```
<torch._C.Generator at 0x7eeddd747dd0>
```

```
✓ 0s [3] # Initialize hyper-parameters  
    n_in = 3  
    n_hidden = 5  
    n_out = 2  
    n_steps = 100  
    n_seqs = 8  
    n_iters = 10000  
    lr = 0.01
```

▼ Generate training data

```
✓ 0s [4] x_train = np.random.rand(n_seqs, n_steps, n_in)  
    y_train = np.zeros([n_seqs, n_steps, n_out])  
  
    y_train[:, 1:, 0] = 5 * x_train[:, 1:, 0] - 0.2 * x_train[:, :-1, 2]  
    y_train[:, 3:, 1] = 25 * x_train[:, 2:-1, 1] * x_train[:, :-3, 2]  
    y_train += 0.1 * np.random.randn(n_seqs, n_steps, n_out)  
  
    x_train = torch.tensor(x_train, dtype=torch.float32)  
    y_train = torch.tensor(y_train, dtype=torch.float32)
```

$$\left. \begin{array}{l} y_1(t) = 5x_1(t) - 0.2x_3(t-1) + 0.1\epsilon \\ y_2(t) = 25x_2(t-1)x_3(t-3) + 0.1\epsilon \end{array} \right\}$$

Question 2

The basic slice syntax is $i:j:k$ where i is the starting index, j is the stopping index, and k is the step.

`y_train[:,3:,1] = [y2(3), y2(4), y2(5) ..., y2(99)]`

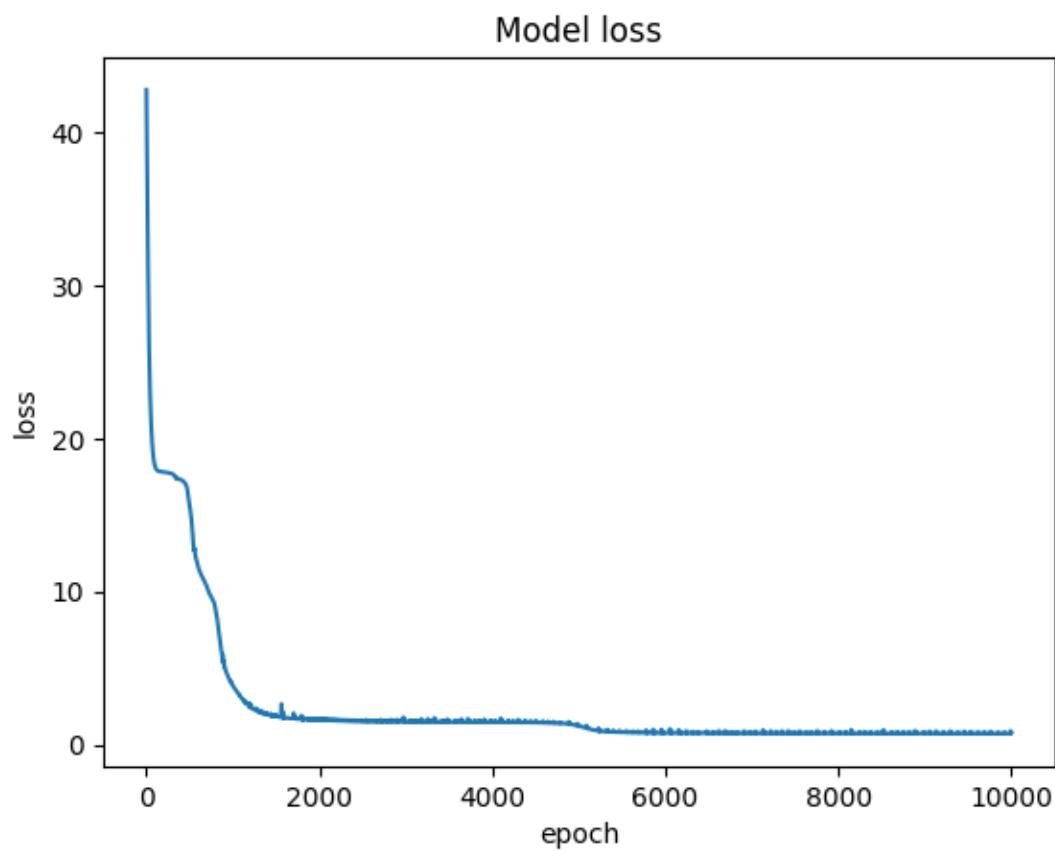
} Assume n is the number of elements in the dimension being sliced
If j is not given it defaults to n
If k is not given it defaults to 1

`x_train[:,2:-1,1] = [x2(2), x2(3), x2(4) ..., x2(98)]`

} Negative i and j are interpreted as $n + i$ and $n + j$ where n is the number of elements in the corresponding dimension

`x_train[:, :-3,2] = [x3(0), x3(1), x3(2) ..., x3(96)]`

Question 2



Question 2

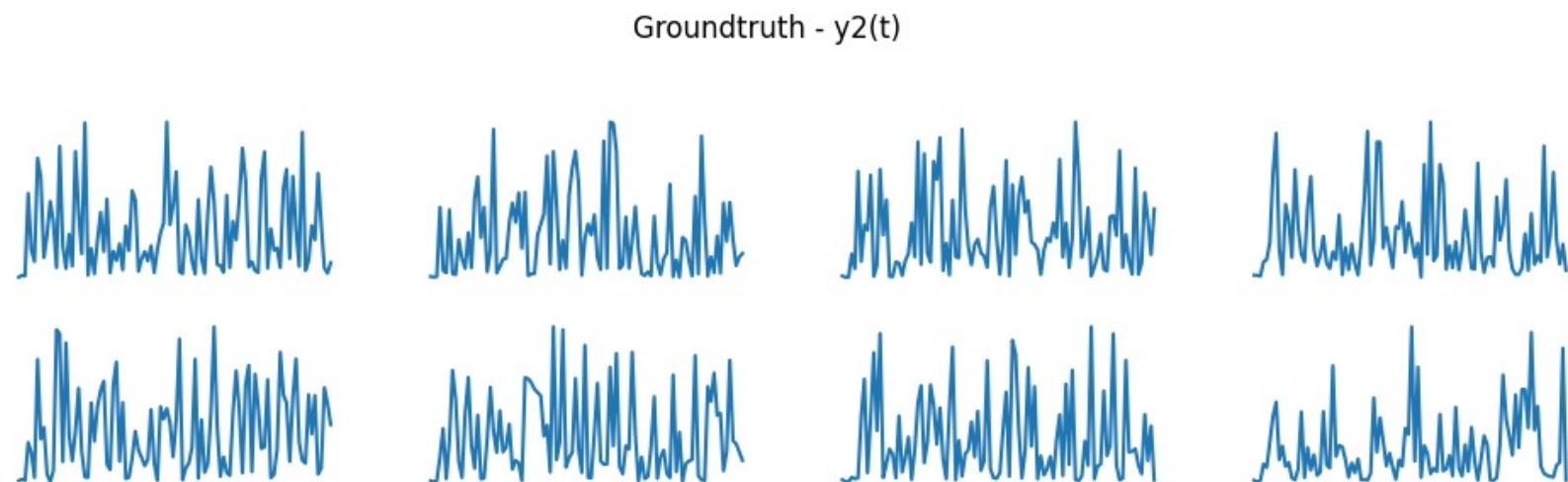
$$y_1(t) = 5x_1(t) - 0.2x_3(t-1) + 0.1\varepsilon$$

Groundtruth - $y_1(t)$



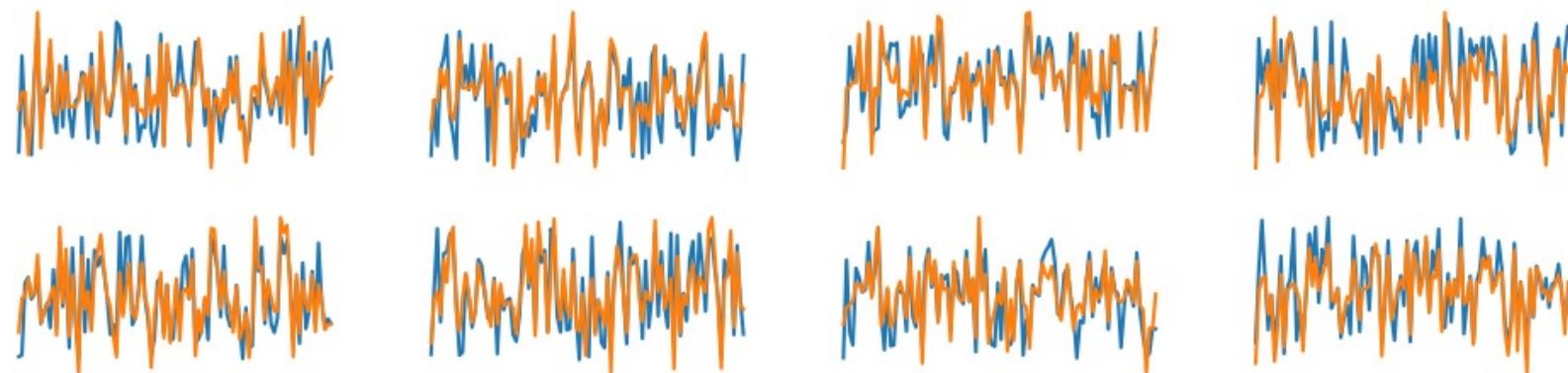
Question 2

$$y_2(t) = 25x_2(t - 1)x_3(t - 3) + 0.1\varepsilon$$



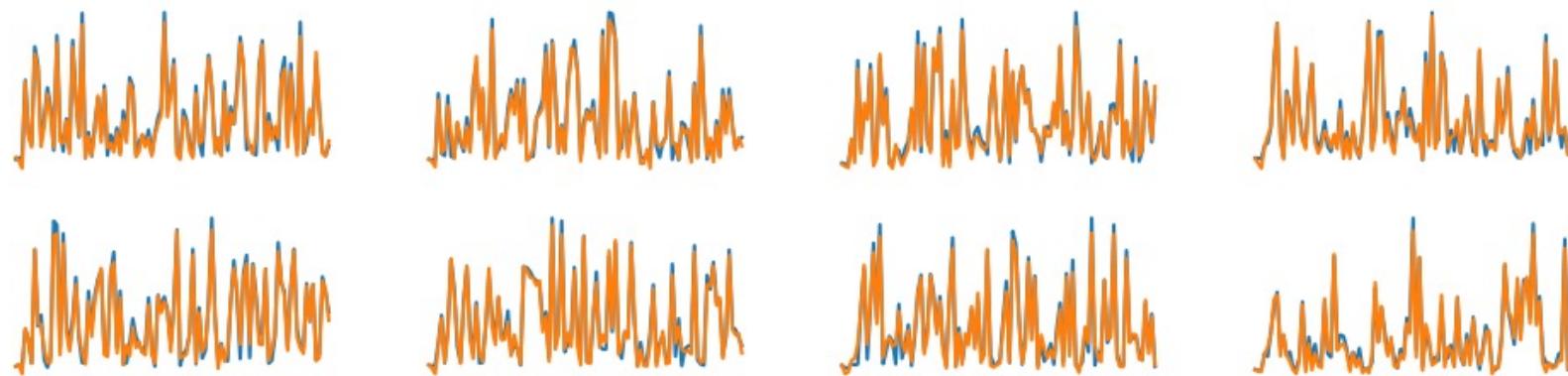
Question 2

Prediction (orange) vs Groundtruth (blue) - $y_1(t)$



Question 2

Prediction (orange) vs Groundtruth (blue) - $y_2(t)$



Question 3

3. Design an LSTM layer with 10 units to map the following input and output sequences:

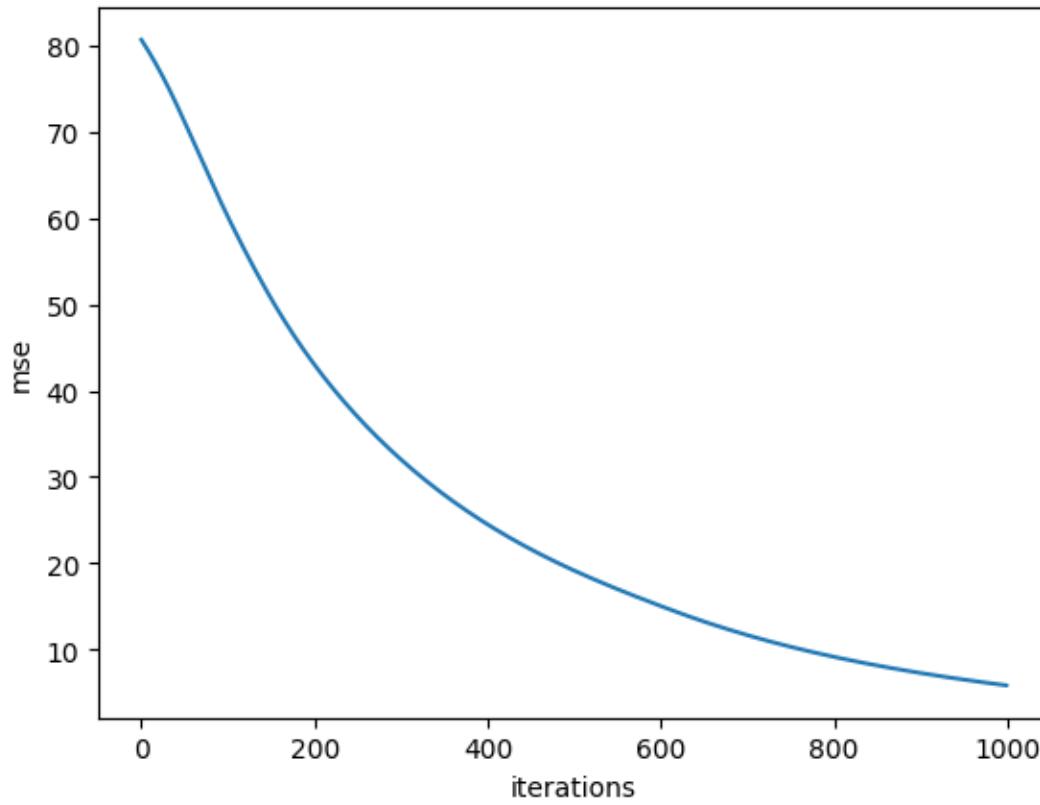
Input x	Output y
(1 2 5 6)	(1 3 7 11)
(5 7 7 8)	(5 12 14 15)
(3 4 5 7)	(3 7 9 12)

Plot the learning curves at a rate $\alpha = 0.001$.

Find the output sequences for the following input sequences:

$$\begin{pmatrix} 1 & 2 & 3 & 4 \\ 4 & 5 & 6 & 7 \end{pmatrix}$$

Question 3



Question 3

Find the output sequences for the following input sequences:

$$\begin{pmatrix} 1 & 2 & 3 & 4 \\ 4 & 5 & 6 & 7 \end{pmatrix}$$

```
test_x = torch.tensor(  
    [[[1], [2], [3], [4]], # 1, 3, 5, 7  
     [[4], [5], [6], [7]]], dtype=torch.float32) # 4, 9, 11, 13  
out = predictor(test_x)  
print(out)
```

```
tensor([[0.9988, 2.7404, 6.3064, 8.8883], [4.6289,  
8.4632, 9.4166, 9.5583]], grad_fn=<SqueezeBackward1>)
```

Question 4

4. Design a character RNN layer with 10 GRU units to learn the sentiments about a movie, given in table 1.

Convert the input text into character ids and set the maximum text length to be 40. Train the network using gradient descent learning with the Adam optimizer and at a learning rate $\alpha = 0.001$. Use dropouts at the hidden layer of GRU at a probability $p = 0.7$.

Plot the cost and the accuracies against epochs during training.

Question 4

Table 1

<i>Input</i>	<i>Sentiment</i>
I did not like the movie	negative
The movie was not good	negative
I watched the movie with great interest	positive
I have seen better movies	negative
Good to see that movie	positive
I am not a fan of movies	negative
I liked the movie great	positive
The movie was of interest to me	positive
I thought they could show interesting scenes	negative
The movie did not have good scenes	negative
Family did not like the movie at all	negative

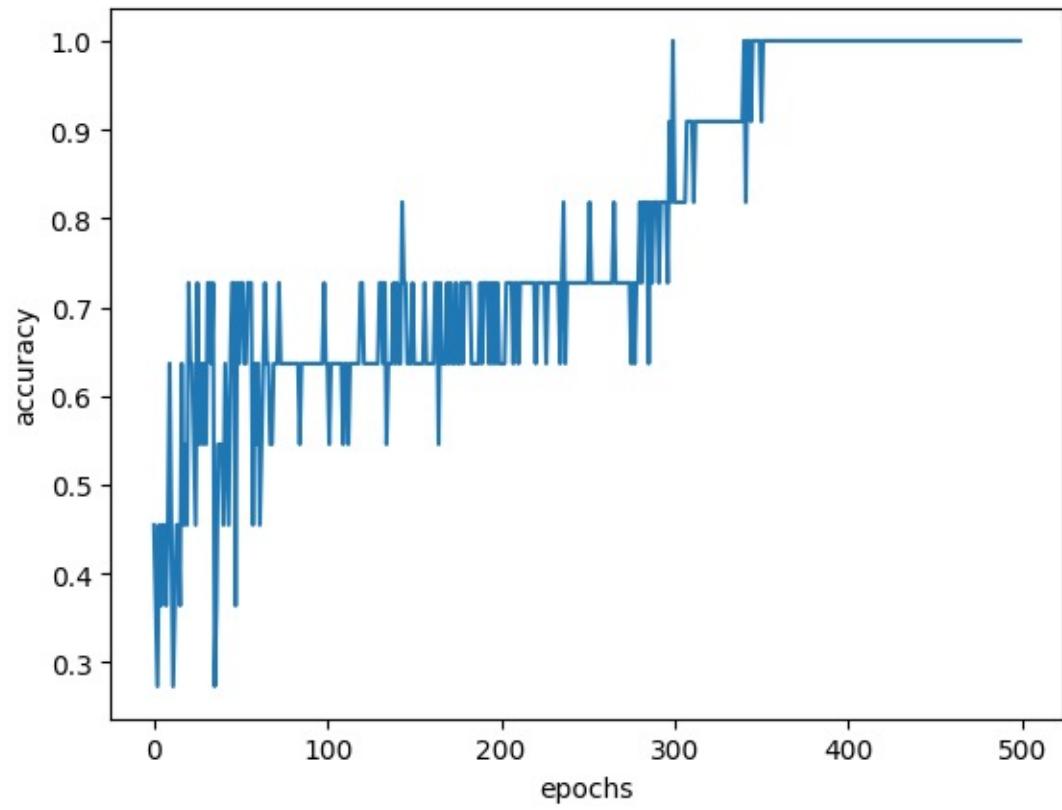
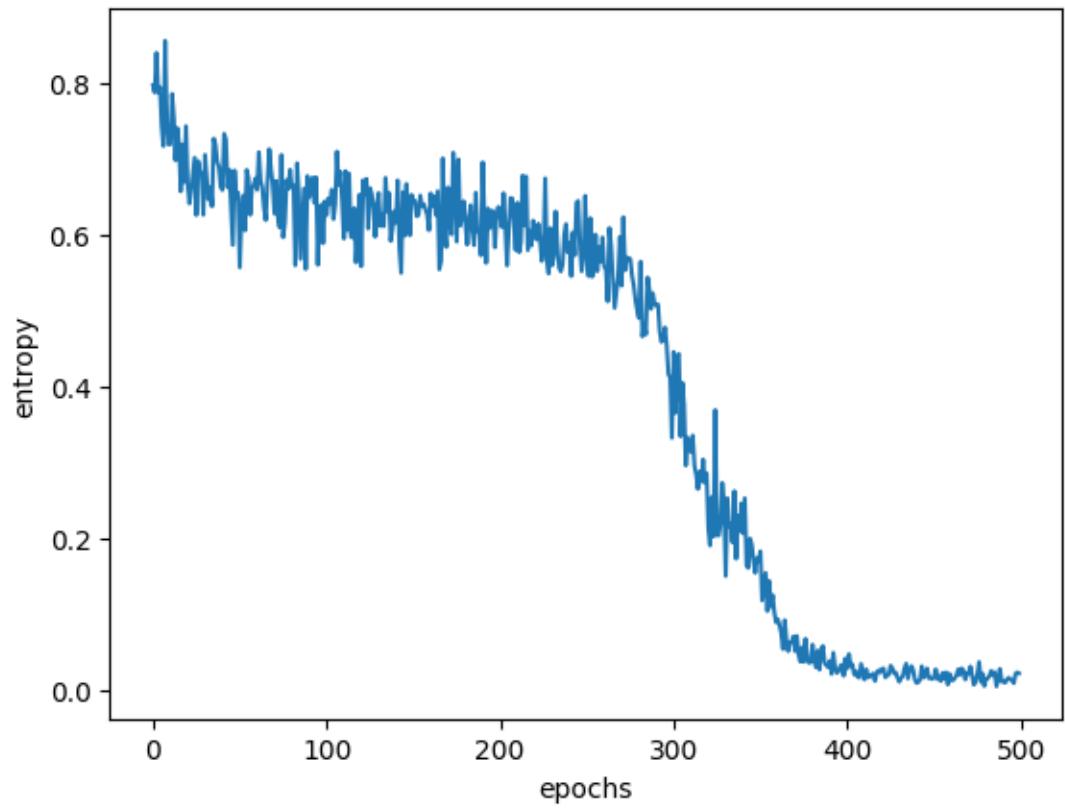
After training, find the likelihoods of the sentiments of the following statements:

The movie was not interesting to me

I liked the movie with great interest

Question 4

Question 4



Question 4

```
test_data = ['The movie was not interesting to me',  
'I liked the movie with great interest']
```

```
test_logits = model(x_test, drop_rate=0)  
probs = nn.functional.softmax(test_logits, dim=1)  
print(probs)
```

```
tensor(  
[[0.0119, 0.9881],  
[0.4359, 0.5641]])
```

'The movie was not interesting to me' is a negative sentiment with a probability of 0.9881

'I liked the movie with great interest' is a positive sentiment with a probability of 0.4359 (hasn't trained well)

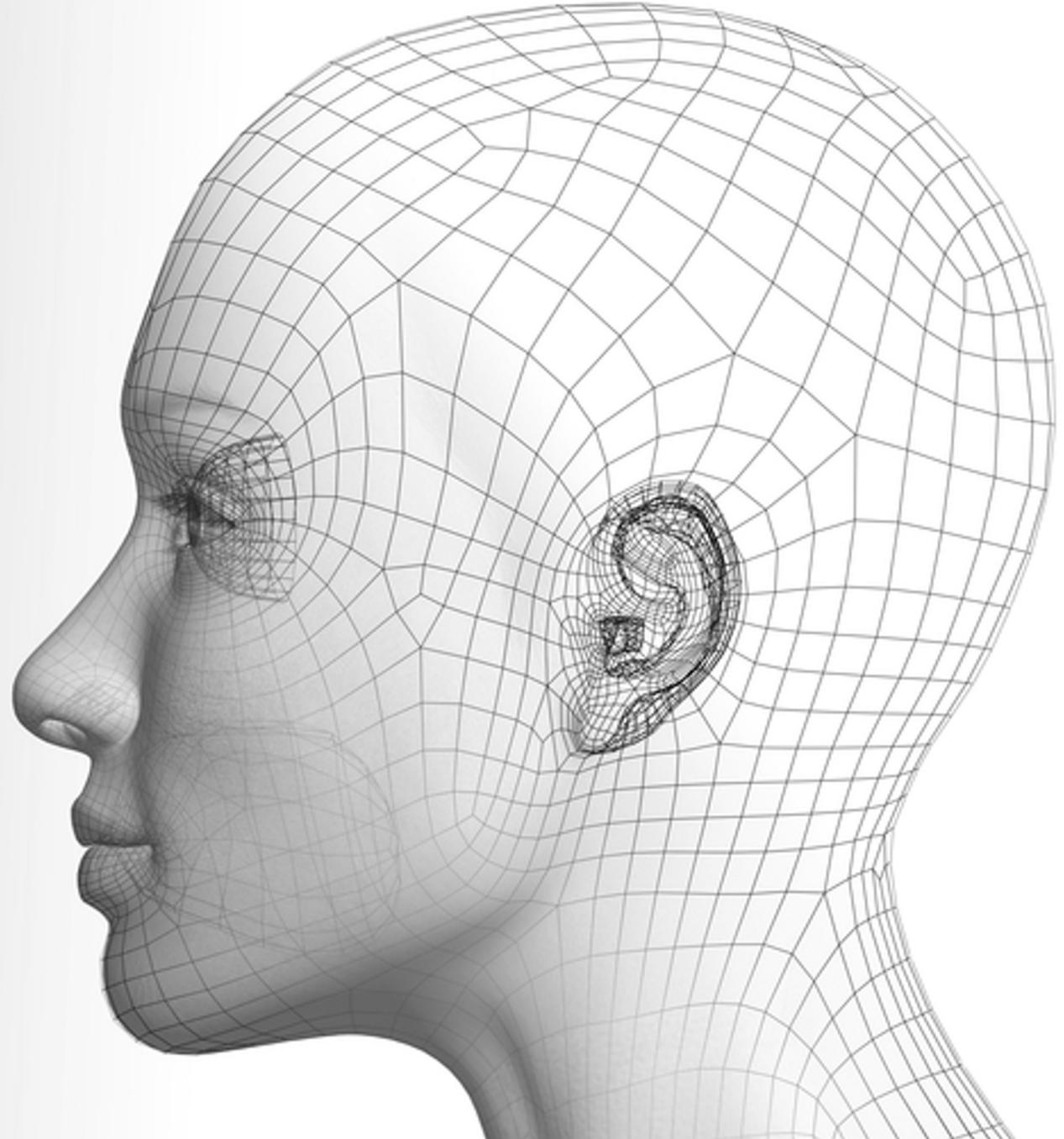
Tutorial 9 Attention

Xingang Pan

潘新钢

<https://xingangpan.github.io/>

<https://twitter.com/XingangP>



Question 1

The tutorial provides a simple walkthrough of the Vision Transformer. We hope you will be able to understand how it works by looking at the actual data flow during inference.

t9q1.ipynb

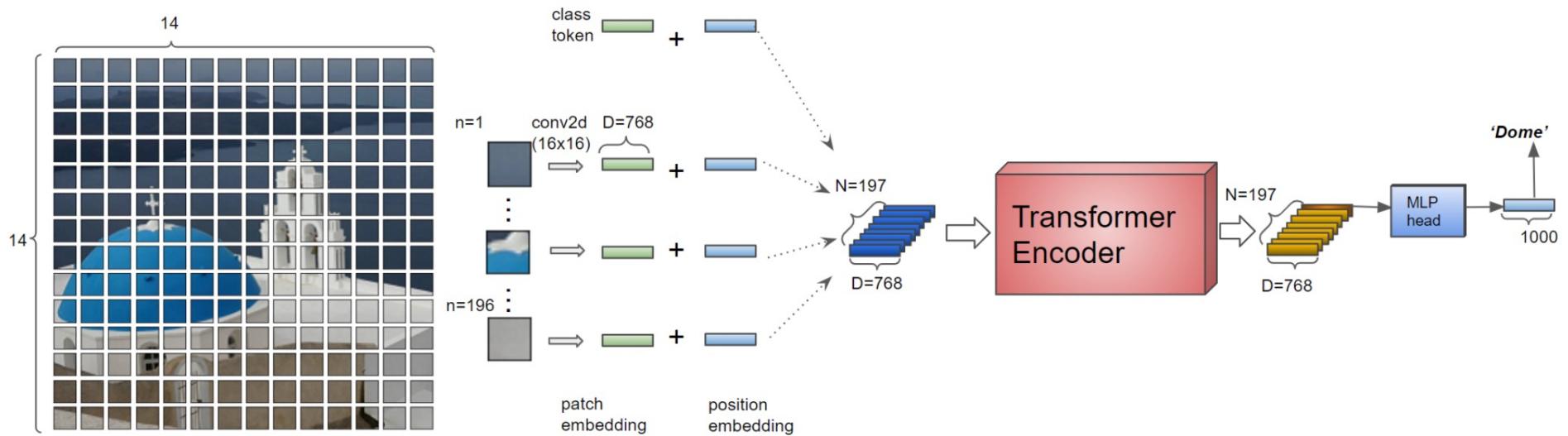


Figure 1. Vision Transformer inference pipeline.

1. Split Image into Patches

The input image is split into 14×14 vectors with dimension of 768 by Conv2d ($k=16 \times 16$) with stride=(16, 16).

2. Add Position Embeddings

Learnable position embedding vectors are added to the patch embedding vectors and fed to the transformer encoder.

3. Transformer Encoder

The embedding vectors are encoded by the transformer encoder. The dimension of input and output vectors are the same. Details of the encoder are depicted in Fig. 2.

4. MLP (Classification) Head

The 0th output from the encoder is fed to the MLP head for classification to output the final classification results.

Question 1

