

Stochastic Gradient Descent (SGD)

$$y = w * x + b$$

update rule can be written as :

$$w = w - \text{learning_rate} * \text{gradient}_w$$

$$b = b - \text{learning_rate} * \text{gradient}_b$$

initialize weight parameter w & bias parameter b to small random values $w=0.5, b=0.1$

$$\text{weight gradient} = \frac{\partial L}{\partial w} = -x(y - \hat{y})$$

$$\text{bias gradient} = \frac{\partial L}{\partial b} = -y(y - \hat{y})$$

$$\text{weight update} = w = w - \alpha, \frac{\partial L}{\partial w}$$

$$\text{bias update} = b - \alpha \frac{\partial L}{\partial b}$$

x	y
1	3
2	5
3	7
4	9
5	11

$$\text{MSE} = \frac{1}{2m} \sum_{i=1}^n (y - \hat{y})^2$$

$$= \frac{1}{2(5)} \sum_{i=1}^5 ((3-0.6)^2 + (5-1.6)^2 + (7-1.6)^2 + (9-2.6)^2 + (11-2.6)^2)$$

$$\boxed{\text{MSE} = 16.83}$$

$$y = w * x + b$$

$$\hat{y} = (0.5 * 1) + 0.1$$

$$\hat{y} = (0.5 * 2) + 0.1$$

$$\hat{y} = (0.5 * 3) + 0.1$$

$$\hat{y} = (0.5 * 4) + 0.1$$

$$\hat{y} = (0.5 * 5) + 0.1$$

$$y = w_0 + b$$

for $x=1 \quad w=0.5, b=0.1$

$$\hat{y} = w_0 + b$$

$$\hat{y} = (0.5 \times 1) + 0.1 = 0.6$$

$$\begin{aligned} \text{gradient weight} &= \frac{\partial L}{\partial w} = -(y - \hat{y})x \\ &= -(3 - 0.6)x \\ &= -2.4 \end{aligned}$$

$$\begin{aligned} \text{gradient bias} &= \frac{\partial L}{\partial b} = -\underline{(3 - 0.6)} = -2.4 \end{aligned}$$

$$\begin{aligned} \text{updated weight} &= w - \alpha \frac{\partial L}{\partial w} = 0.5 - (0.01)(-2.4) \\ &= 0.524 \end{aligned}$$

$$\begin{aligned} \text{updated bias} &= b - \alpha \frac{\partial L}{\partial b} = 0.1 - (0.01)(-2.4) \\ &= 0.124 \end{aligned}$$

for $x=2$

$$\hat{y} = w_0 + b$$

$$\hat{y} = 0.524 \times 2 + 0.124 = 1.172$$

$$\begin{aligned} \text{gradient weight} &= \frac{\partial L}{\partial w} = -(y - \hat{y})x \\ &= -(5 - 1.172)x^2 = -7.656 \end{aligned}$$

$$\begin{aligned} \text{gradient bias} &= \frac{\partial L}{\partial b} = -(-5 - 1.172) = -3.828 \end{aligned}$$

$$\begin{aligned} \text{updated weight} &= w - \alpha \frac{\partial L}{\partial w} = 0.524 - 0.01 \\ &= 0.600 \end{aligned}$$

$$\text{updated bias: } b - \alpha \frac{\partial L}{\partial w} = 0.125 - 0.01(3.828) \\ = 0.162$$

for $x_1 = 3$

$$\hat{y} = wx + b$$

$$\hat{y} = (0.6 \times 3) + 0.162 = 1.962$$

$$\text{gradient weight} = \frac{\partial L}{\partial w} = -(y - \hat{y})x \\ = -(7 - 1.962) \times 3 \\ = -15.114$$

$$\text{gradient bias} = -(7 - 1.962) = -5.038$$

$$\text{updated weight} = 0.600 - 0.01(-15.114) \\ = 0.7512$$

$$\text{updated bias} = 0.162 - 0.01(-5.038) = 0.212$$

for $x = 4$

$$\hat{y} = wx + b$$

$$= (0.7512 \times 4) + 0.212 = 3.216$$

$$\text{gradient weight} = \frac{\partial L}{\partial w} = -(y - \hat{y})x \\ = -(9 - 3.216)4 \\ = 3.864$$

$$\text{gradient bias} = -(9 - 3.216) = -5.136$$

$$\text{updated weight} = 0.7512 - 0.01(-3.864) \\ =$$

MSE after applying SGD

optimization Algo

$x = 1$	3	$\hat{y} = 1 \times 2.737 + 0.3287$	$(3 - 1.602)^2$
$x = 2$	5	$\hat{y} = 2 \times 2.737 + 0.3287$	$(5 - 2.873)^2$
$x = 3$	7	$\hat{y} = 3 \times 2.737 + 0.3287$	$(7 - 4.146)^2$
$x = 4$	9	$\hat{y} = 4 \times 2.737 + 0.3287$	$(9 - 5.419)^2$
$x = 5$	11	$\hat{y} = 5 \times 2.737 + 0.3287$	$(11 - 6.692)^2$

MSE

$$\text{After SGD} = \frac{1}{2 \times 5} : (1.984 + 4.529 + 8.152 + 12.880 + 18.558)$$

$$= \frac{1}{10} (56.005)$$

$$r_{\text{MSE}} = 5.6005$$

Momentum Gradient Descent
 $\omega = 0.5$, $\alpha = 0.01$, $b = 0.1$, initial velocity $v_w = 0$, $v_b = 0$
momentum coefficient $\beta = 0.9$

Q Given dataset

$$MSE = \frac{1}{2m} \sum_{i=1}^m (y - \hat{y})^2$$

$\hat{y}_1 = (0.5x_1) + 0.1$	$x_1 = 1$	$y_1 = 3$	$= \frac{1}{2 \times 3} ((3 - 0.6)^2 + (5 - 1.1)^2 + (7 - 1.6)^2)$
$\hat{y}_2 = (0.5x_2) + 0.1$	$x_2 = 2$	$y_2 = 5$	$= \frac{1}{6} (5.76 + 15.21 + 29.16)$
$\hat{y}_3 = (0.5x_3) + 0.1$	$x_3 = 3$	$y_3 = 7$	$= 8.35$

velocity updates :

$$v_w = \beta v_w + \alpha \frac{\partial L}{\partial w} \quad v_b = \beta v_b + \alpha \frac{\partial L}{\partial b}$$

new weight update

$$w = w - v_w, b = b - v_b$$

Now for $x=1$

$$\frac{\partial L}{\partial w} = -(y - \hat{y}) \alpha, \frac{\partial L}{\partial b} = -(y - \hat{y}) = -2.4$$

$$= -2.4$$

$$v_w = (0.9 \times 0) + 0.01(-2.4) = -0.024$$

$$v_b = (0.9 \times 0) + 0.01(-2.4) = -0.024$$

$$w = 0.5 + 0.024 = 0.524$$

$$b = 0.1 + 0.024 = 0.124$$

for $x=2$

$$\hat{y} = (0.524 \times 2) + 0.124$$

$$= 1.17$$

$$\frac{\partial L}{\partial w} = -(5 - 1.17) \times 2 = -7.66$$

$$\frac{\partial L}{\partial b} = -(5 - 1.17) = -3.82$$

$$vw = (0.9 \times (-0.025)) + 0.01 (-7.66) = -0.099$$

$$vb = (0.9 \times -0.024) + 0.01 (-3.82) =$$

$$w = 0.525 - (-0.099) = 0.623$$

$$b = 0.12 - (0.000) = 0.184$$

for $x \leq 3$.

$$y = (0.623 \times 3) + 0.184 = 2.053$$

$$\frac{\partial L}{\partial w} = (7 - 2.053) \times 3 = -14.8$$

$$\frac{\partial L}{\partial b} = (7 - 2.053) = -4.9$$

$$vw = (0.9 \times (-0.025)) + 0.01 (-14.9) = -0.237$$

$$vb = (0.9 \times (-0.024)) + 0.01 (-4.9) = -0.103$$

$$\text{final weight } w = 0.623 - (-0.237) = 0.860$$

$$b = 0.184 - (-0.103) = 0.287$$

Calculate MSE from QD

$$\text{MSE} = \frac{1}{2n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad \left| \begin{array}{l} y = (0.860)x_1 + 0.287 = 1.147 \\ \hat{y} = (0.860x_2 + 0.287) = 2.007 \end{array} \right.$$

$$= \frac{1}{2} \times 3 \left((3 - 1.147)^2 + (5 - 2.007)^2 + (7 - 2.867)^2 \right) \quad \left| \begin{array}{l} y = (0.860x_3 + 0.287) = 2.867 \end{array} \right.$$

$$\text{MSE} = 4.912$$

AdaGrad (Adaptive Gradient)

Q.1 Given $w = 0.5$ input $(x, y) = (3, 5)$
 \Rightarrow learning rate = 0.1 $\epsilon = 10^{-8}$

$$t = \frac{1}{n} * (y_{\text{pred}} - y)^2 \\ y_{\text{pred}} = 0.5 \times 3 = 1.5$$

$$\frac{\partial L}{\partial w} = (- (y - \hat{y}) x)$$

$$\frac{\partial L}{\partial w} = (- (4 + 10^{-8}) \times 3) = -7.5$$

$$\alpha_t = \sqrt{\sum_{i=1}^t (g_i)^2} = \sqrt{(-7.5)^2} = 7.5$$

$$n't = \frac{n}{\sqrt{\alpha_t + \epsilon}} = \frac{0.1}{\sqrt{7.5 + 10^{-8}}} = 0.013$$

$$w_t = 0.5 - (0.013) (-7.5) = 0.5975$$

$$w_{t+1} = w_t - n't \frac{\partial L}{\partial w(t)}$$

Q2 : $w = 0.597$ input $(4, 5)$

$$\hat{y} = y_{\text{pred}} = 0.597 \times 4 = 2.388$$

$$\frac{\partial L}{\partial w} = -(y - \hat{y}) x = (- (5 - 2.388) \times 4) \\ = -10.448$$

$$\frac{\partial L}{\partial w} = -10.448$$

$$\sigma_t = \sqrt{\sum_{i=1}^t (g_i)^2} = (-10.448)^2 = 109.160$$

$$n't = \frac{n}{\sqrt{\sigma_t + \epsilon}} = \frac{0.1}{\sqrt{109.160 + 10^{-8}}} = 0.0095$$

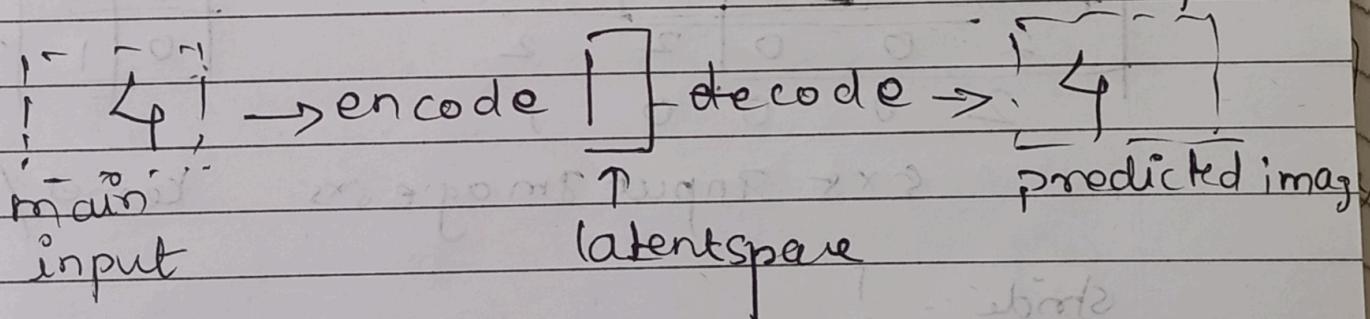
$$wt = 0.597 - (0.0095)(-10.448) \\ = 0.597 + 0.09925$$

$$wt = 0.69625$$

① Auto encoders : $P(A \in \text{dimensional reduction})$

encoder : $h = f(x)$

decoder : $r = g(f(x))$ $\rightarrow r$ as close x as possible



encoding : $z = w_1 * x + b_1$

decoding : $x' = w_2 * z + b_2$

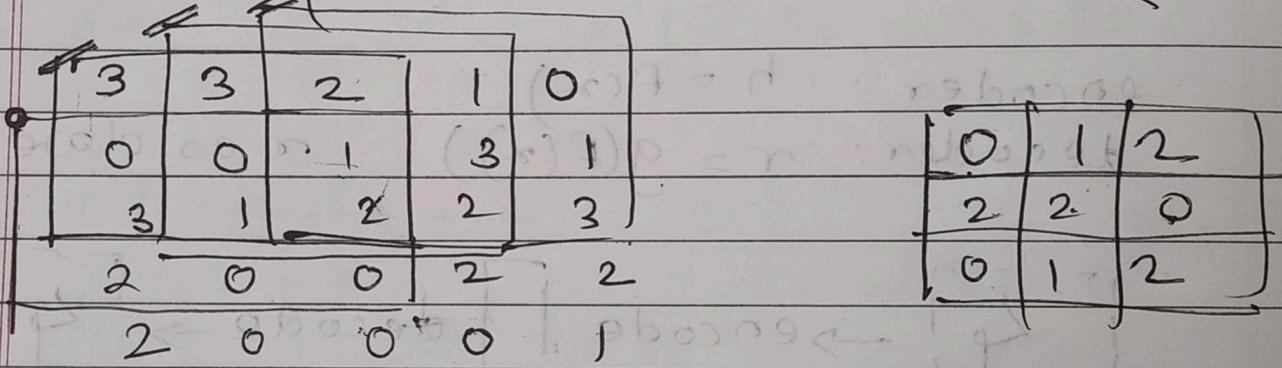
$$\text{loss} = \frac{1}{n} * \sum (x - x')^2$$

- 1) Under complete autoencoder $n_h < d$
- 2) sparse Autoencoder
- 3) overcomplete encoder

Applications

- 1) Dimensionality Reduction
- 2) Anomaly detection
- 3) Denoising Images
- 4) feature learning
- 5) Generative modeling

CNN (Convolutional neural networks)



5x5 input image \times kernel/filter 3x3

$$\text{Stride} = \frac{1}{2}$$

12	12	17
10	17	19
9	6	14

feature map

17	19
17	19

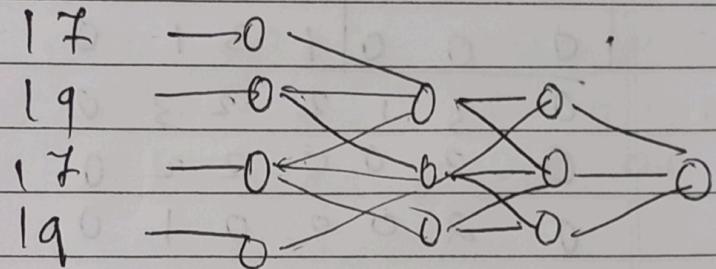
12.5	16.25
10.5	15

Avg max pooling

Avg pooling

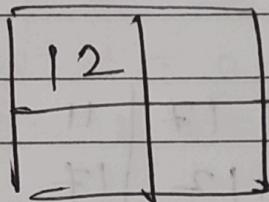
flattening

17 19
17 19

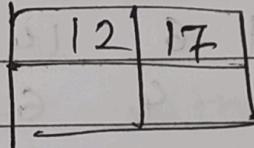


Now stride = 2

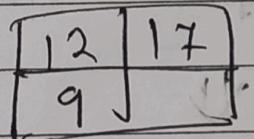
3 3 2
0 0 1
3 1 2



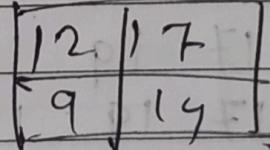
2 1 0
1 3 1
2 2 3



3 1 2
2 0 0
2 0 0

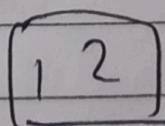


2 2 3
0 2 2
0 0 1



now stride = 3

3 3 2
0 0 1
3 1 2



Solve by max pooling

0	0	0	0	0	0	0
0	3	3	2	1	0	0
0	0	0	1	3	1	0
0	3	1	2	2	3	0
0	2	0	0	2	2	0
0	2	0	0	0	1	0
0	0	0	0	0	0	0

0 1 2

2 2 0

0 1 2

3x3 kernel

0 = stride down

7x7

6	14	17	11	13
14	12	12	17	11
8	10	17	19	13
11	9	6	14	12
6	4	4	6	9

8 8 8

1 0 0

5 1 8

0 1 6

1 8 1

8 5 6

5x5 feature matrix

$\begin{pmatrix} 14 & 17 & 17 \\ 14 & 17 & 19 \\ 11 & 17 & 19 \end{pmatrix}$
max pooling

14	17	17	17
14	17	17	19
11	17	19	19
11	9	14	14

14	0	0
17	0	0
17	0	0
17	0	0
19	0	0

17	0	0
17	0	0
19	0	0
17	0	0
19	0	0

19	0	0
11	0	0
17	0	0
19	0	0
19	0	0

11	0	0
17	0	0
19	0	0
17	0	0
19	0	0

17	0	0
19	0	0
19	0	0
17	0	0
19	0	0

19	0	0
11	0	0
17	0	0
19	0	0
19	0	0

11	0	0
17	0	0
19	0	0
17	0	0
19	0	0

17	0	0
19	0	0
19	0	0
17	0	0
19	0	0

19	0	0
11	0	0
17	0	0
19	0	0
19	0	0

Ch.	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	kernel 1	kernel 2
																										0.1 0.2 0.3	-0.1 -0.1 -0.1
																										0.1 0.2 0.3	0.1 0.1 0.1
																										0.1 0.2 0.3	0.0 0.0 0.0

Input 5×5

1) stride = 1

2) Bias Term for Kernel 1 = 0.1

Bias Term for Kernel 2 = 0.2

3) Activation f^n is ReLU.

4) Max pooling is 2×2

5) Flattening of FC layer with hidden layer with 2 neurons.

\Rightarrow	13.2	15	16.8				
	22.2	24	25.8				

class 1
class 2
class 3

13.2	15	16.8
22.2	24	25.8
31.2	33	34.8

feature map 1

Adding bias for feature

class 1	1.5	1.5	1.5
class 2	1.5	1.5	1.5
class 3	1.5	1.5	1.5

feature map 2

13.3	15.1	16.9
22.3	24.1	25.9
31.3	33.1	34.9

bias

1.7	1.7	1.7
1.7	1.7	1.7
1.7	1.7	1.7

Activation "convert
- to f value

, 1st row , 1st row

PAGE No.	
DATE	5/1/1

13.3 15.1 16.9 1.7 1.7 1.7

22.3 24.1 25.9 1.7 1.7 1.7

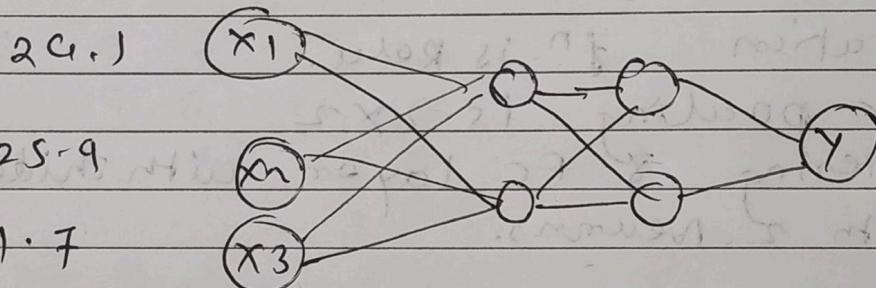
31.3 33.1 34.9 1.7 1.7 1.7

Apply ReLU Activation function

after apply ReLU matrix will be same as above

stride 2 max pooling 2x2

24.1 25.9 1.7



flattening

24.1

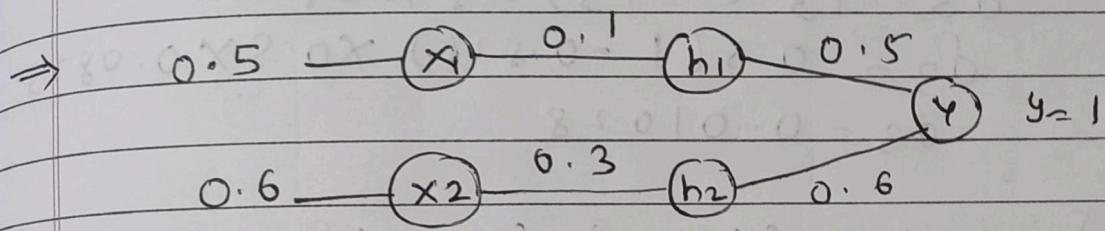
25.9
1.7

LeNet-5 Architecture Autoencoder (theory)
CNN optimization backpropagation (Neural network)

$5 \times 5, 3 \times 3$, max pooling

Assignment Class Test - 8

Q. Perform one step back propagation in simple neuron network with one I/P layer having 2 neurons in hidden layer with $2N = 1 \times 0.1$
 $AF = \text{sigmoid}$ Target O/P $y = 1$
 $x = [0.5 \ 0.6]^T$ $w_{11} = [0.1 \ 0.3]$



$$H_3 = x_1 w_1 = 0.5 \times 0.1 = 0.05$$

$$H_3 = 0.5125$$

$$H_4 = x_2 w_2 = 0.6 \times 0.3 = 0.18$$

$$H_4 = 0.549$$

$$O_S = H_3 w_3 + H_4 w_4$$

$$= 0.5125 \times 0.5 + 0.549 \times 0.6$$

$$= 0.58313$$

$$O_S = 0.64178$$

$$\text{error} = y_1 - 0.5 = 1 - 0.64178 = 0.35822$$

$$S_5 = y(1-y)(y_1 - y)$$

$$S_5 = 0.0823$$

$$\Delta w_{45} = 0.1 \times 0.0823 \times 0.5498$$

$$w_{45} = 0.6 + 0.00448$$

$$= 0.60498$$

$$\Delta w_{35} = 0.1 \times 0.0823 \times 0.5125 \\ = 0.004521$$

$$w_{35} = 0.5 + 0.004521 = 0.504521$$

for hidden layer,

$$\delta_3 = y_3 (1 - y_3) w_5 \times \delta_5$$

$$\delta_3 = 0.5 (1 - 0.518) \times 0.5 \times 0.0823$$

$$\delta_3 = 0.01028$$

$$\delta_4 = y_4 (1 - y_4) w_{45} \times \delta_5$$

$$\delta_4 = 0.545 (1 - 0.545) \times 0.6$$

$$\delta_4 = 0.01222$$

$$\Delta w_{13} = 0.1 \times \delta_3 \times 0.1 = 0.1 \times 0.01028$$

$$w_{13} = 0.1 \times 0.005 \\ = 0.10054$$

$$\Delta w_{24} = 0.1 \times \delta_4 \times 0.2 = 0.1 \times 0.012 \times 0.6 \\ = 0.00073$$

$$w_{24} = 0.0007349$$

Iteration 2

$$H_3 = 0.5 \times 0.1005 \quad H_4 = 0.6 \times 0.300 \\ = 0.05 \quad (x-1) = 0.18$$

$$H_3 = 0.512$$

$$H_4 = 0.5459$$

$$O_5 = H_3 w_{35} + H_4 w_{45}$$

$$= 0.58787$$

$$O_5 = 1.80015$$

a.2 perform CNN output kernel matrix by 5×5
 Input & X Kernel

$$\begin{bmatrix} 1 & 2 & 3 & 4 & 5 \\ 2 & 3 & 4 & 5 & 1 \\ 3 & 4 & 5 & 1 & 2 \\ 4 & 5 & 1 & 2 & 3 \\ 5 & 1 & 2 & 3 & 4 \end{bmatrix} \quad \begin{bmatrix} 0.2 & 0.2 & 0.3 \\ 0.2 & 0.2 & 0.3 \\ 0.2 & 0.2 & 0.3 \end{bmatrix}$$

feature map 1

$$\begin{matrix} 6.6 & 7.2 & 6.8 \\ 7.2 & 6.8 & 5.4 \\ 6.8 & 5.4 & 5.5 \end{matrix}$$

AF = ReLU

As all vals are the same matrix

max pooling

$$\begin{matrix} 7.2 & 7.2 \\ 7.2 & 6.8 \end{matrix}$$

4

flattening

$$[7.2 \ 7.2 \ 7.2 \ 6.8]$$