

$$x_0 = 3$$

$$\begin{aligned}x_1 &= x_0 - 0.2 \times f'(x_0) \\&= 3 - 0.2 \times 2.4 \\&= 3 - 0.8 \\&= 2.2\end{aligned}$$

$$\begin{aligned}x_2 &= x_1 - 0.2 \times f'(x_1) \\&= 2.2 - 0.2 \times f'(2.2) \\&\quad \rightarrow 2 \times 2 - 2 \\&= 2.2 - 0.2(2.4) \quad \begin{matrix}4-4-2 \\= 2.4\end{matrix} \\&= 2.2 - 0.48 \\&= 1.72\end{aligned}$$

$$\begin{aligned}x_3 &= x_2 - 0.2 \times f'(x_2) \\&= 1.72 - 0.2 \times (2 \times 1.72 - 2) \\&= 1.72 - 0.2(1.44) \\&= 1.72 - 0.288 \\&= 1.432\end{aligned}$$

$$\begin{aligned}x_4 &= x_3 - 0.2 \times f'(x_3) \\&= 1.432 - 0.2 \times f'(1.432) \\&= 1.432 - 0.2(2 \times 1.432 - 2) \\&= 1.432 - 0.2(0.864) \\&= 1.2592\end{aligned}$$

$$\begin{aligned}
 x_5 &= x_4 - 0.2 \times f'(x_4) \\
 &= 1.2592 - 0.2 \times f'(1.2592) \\
 &= 1.2592 - 0.2(2 \times 1.2592 - 2) \\
 &= 1.2592 - 0.2 \times 0.864 = 0.5184 \\
 &= 1.1556
 \end{aligned}$$

$$\begin{aligned}
 x_6 &= x_5 - 0.2 \times f'(x_5) \\
 &= 1.1556 - 0.2 \times f'(1.1556) \\
 &= 1.1556 - 0.2(2 \times 1.1556 - 2) \\
 &= 1.1556 - 0.2(0.3112) \\
 &= 1.09336
 \end{aligned}$$

$$\begin{aligned}
 x_7 &= x_6 - 0.2 \times f'(x_6) \\
 &= 1.09336 - 0.2(2 \times 1.09336 - 2) \\
 &= 1.09336 - 0.2 \times 0.18 \\
 &= \underline{\text{0.73}} \quad 1.05
 \end{aligned}$$

$$\begin{aligned}
 x_8 &= x_7 - 0.2 \times f'(x_7) \\
 &= \underline{\text{0.73}} \quad 1.05 - 0.2(2 \times \underline{\text{0.73}} \quad 1.05 - 2) \\
 &\Rightarrow \underline{\text{0.73}} \quad 1.05 - 0.2 \\
 &= 1.03
 \end{aligned}$$

$$\begin{aligned}
 x_9 &= x_8 - 0.2 \times f'(x_8) \\
 &= 1.03 - 0.2 \times (2 \times 1.03 - 2) \\
 &= 1.03 - 0.012 \\
 &= 1.01
 \end{aligned}$$

$$x_{10} = x_9 - 0.2 \times f'(x_9)$$

$$= 1.01 - 0.2 \times (2 \times 1.01 - 2)$$

$$= 1.01 - 0.004$$

$$= 1.006$$

$$x_{11} = x_{10} - 0.2 \times f'(x_{10})$$

$$= 1.006 - 0.2 \times (2 \times 1.006 - 2)$$

$$= 1.006 - 0.0024$$

$$= 1.003$$

$$x_{12} = x_{11} - 0.2 \times f'(x_{11})$$

$$= 1.003 - 0.2 \times (2 \times 1.003 - 2)$$

$$= 1.003 - 0.002$$

$$= 1.001$$

$$x_{13} = x_{12} - 0.2 \times f'(x_{12})$$

$$= 1.001 - 0.2 \times (2 \times 1.001 - 2)$$

$$= 1.000$$

- Cost function provides with a quantification of how incorrect the model is in estimating the ideal \hat{y} .
- The primary approach for minimizing the cost in deep learning paradigms is to apply an approach called gradient descent with another approach called back propagation
 - ↳ coming back adjusting the weights

Batch gradient descent:-

Batch gradient descent is also known as \rightarrow Vanilla gradient descent.

- In this it divides the large data into some batches.

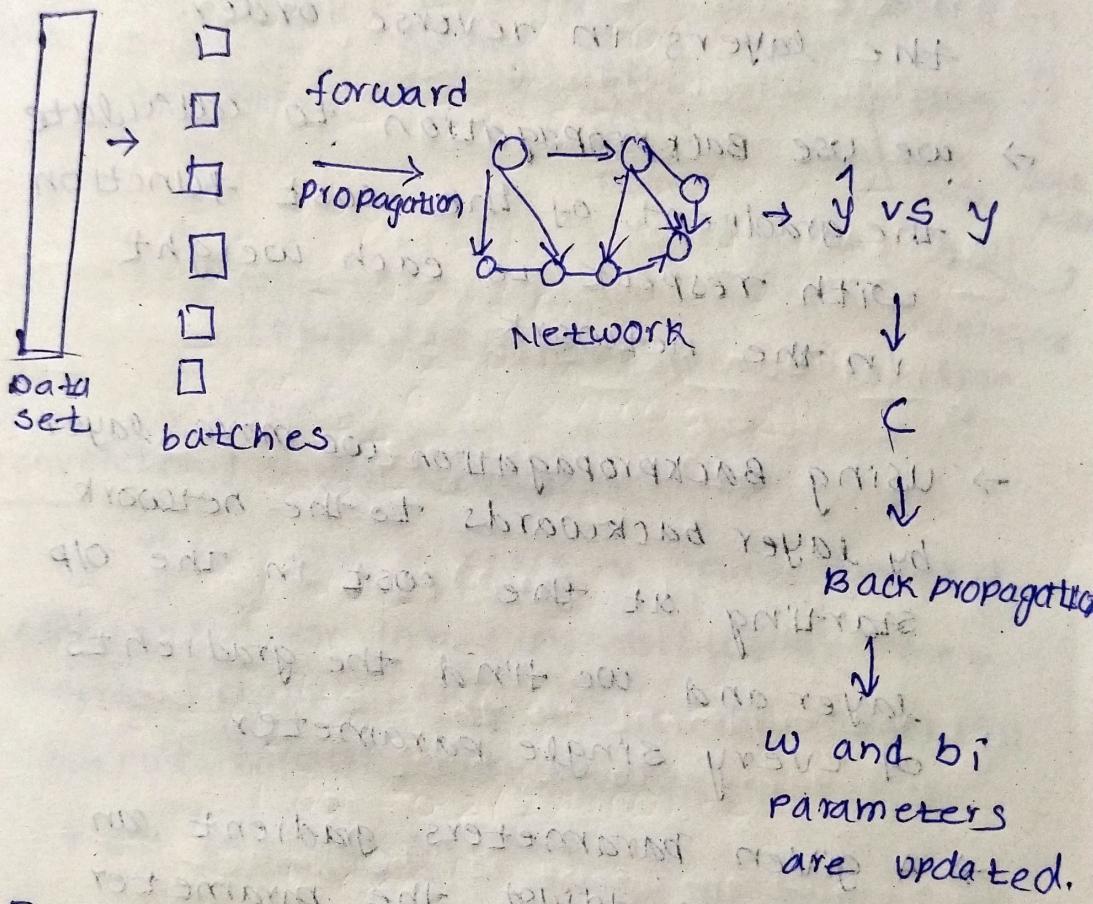
$$\text{No. of batches} = \left\lceil \frac{\text{size of training dataset}}{\text{Batch size}} \right\rceil$$

$$\text{ex} = \left\lceil \frac{60,000}{128} \right\rceil$$

Step-1:- sample a mini-batch of n values.

Step-2:- Forward propagates \vec{x} through network to estimate \hat{y} with \vec{y}

- 3) calculate the cost \hat{C} by comparing y and \hat{y} .
- 4) descent gradient of \hat{C} to adjust w and b , enabling \hat{x} to better predict y .



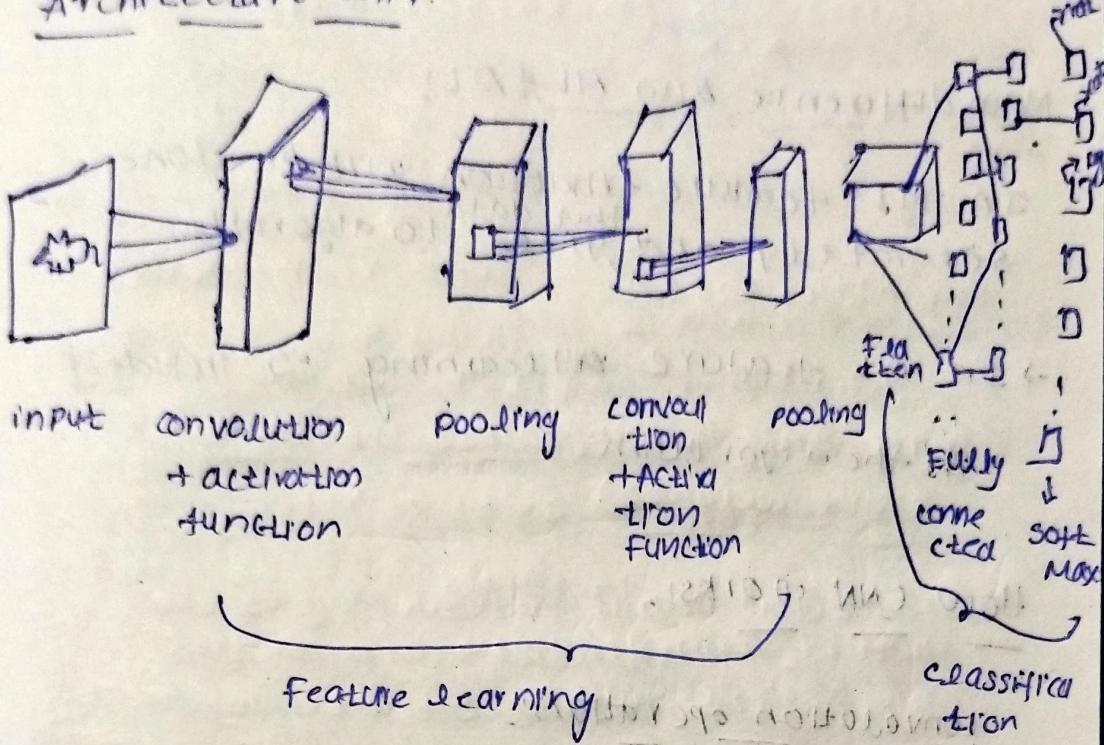
Back propagation!

- to efficiently adjust parameters to multiple layers of artificial neurons stochastic gradient descent is combined with back propagation
- Back propagation is an important application of chain rule in calculus.

- Forward propagation carries information about input x through successive layers of neurons to approximate y with \hat{y} .
- Backpropagation carries information about the cost C backwards through the layers in reverse order.
- we use backpropagation to calculate the gradient of the cost function with respect to each weight in the network.
- Using backpropagation we move layer by layer backwards to the network starting at the cost in the output layer and we find the gradients of every single parameter.
- A given parameters gradient can be used to adjust the parameter up (or) down.

unit - 3

Architecture CNN:-



convolutional Neural Network:-

These are a class of neural network specialized for image processing like other neural networks, they transform input to output through many layers.

- In CNN layers have:
 - 1) convolution step
 - 2) a pooling step (optional)
 - 3) a non-linear activation
- Each layer in a neural network transforms the input tensor into the output tensor through linear and non-linear operations.
- All these intermediate tensors are called activations and they are all

different representation of the input.

Main difference b/w ML & DL:

In ML feature extraction will be done separately, and ^{that data} gives to algorithm.

→ In DL feature learning is included in the algorithm.

How CNN works:

convolution operation:

$$s(t) = \int x(a) w(t-a) \cdot da \rightarrow \text{①}$$

$$s(t) = (x * w)(t) \rightarrow \text{②}$$

→ First argument to the convolutional layer is referred to as input.

→ Second argument is the kernel.

→ The off is referred to the feature map.

→ The data is collected every time ^{is called continuous data}.

→ The data is collected at a specific time with time interval ^{is called discrete data}.

Discrete continuous convolution operation:-

$$s(t) = (x * w)(t) = \sum_{a=-\infty}^{\infty} x(a) \cdot w(t-a)$$

→ (3)

→ the time index t can then take only integer values.

→ the input is usually a multidimensional array of data and the kernel is usually a multidimensional array of parameters that are adapted by the learning algorithm.

- we will refer to these arrays as tensors
- we usually assume that these functions are zero everywhere but the finite set of points for which we store the values.

Cross-correlation operations:-

Input

| | | | |
|---|---|---|---|
| a | b | c | d |
| e | f | g | h |
| i | j | k | l |

Kernel

| | |
|---|---|
| w | x |
| y | z |

$$aw + bx + cy + fz$$

$$bw + cx + dy + gz$$

$$ew + dx + gy + hz$$

→ Stride defines the no. of columns & the no. of rows shifted.

$$(n \times n) \xrightarrow{K \times K} (n-k+1)(n-k+1)$$

order of resultant matrix when the data is multiplied

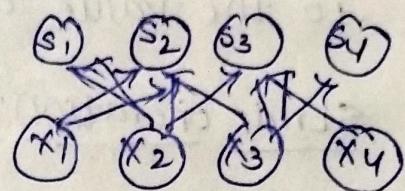
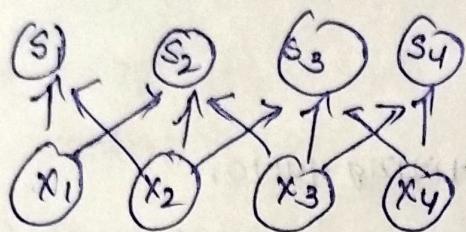
Advantages of CNNs: with kernel ($K \times K$)

slow connection exchange in input doesn't affect the OIP more.

Sparse Interactions (sparse connectivity) or sparse weights

- 1) Traditional neural network layers use matrix multiplications by a matrix of parameters with a separate parameter.
- 2) Every OIP unit interacts with every input unit.
- 3) It is accomplished by making the kernel smaller than the OIP.
- 4) This will reduce the no of parameters that requires less memory and improves statistical efficiency.
- 5) For traditional it needs $m \times n$ output requires $O(mn)$ runtime.
- 6) The sparsely connected approach requires $O(K \times n)$ complexity where $K \ll m$.

conventional NN! - Traditional method!



- n -inputs and s -output then the inputs of 's' are called as receptive fields

The receptive fields of units in the deeper layers of a CNN is larger than the receptive fields of units in shallow layers.

- this effect increases if the network includes architectural features like ~~is strided convolution~~
- pooling
- Even though the direct connections in a convolutional network are very sparse - units in the deeper layers can be indirectly connected to all (or) most of the input units.

Parameter sharing!

- parameter sharing refers to using the same parameter for more than one function in a model.
- In a traditional network, each element of the weight matrix is used exactly once when computing the output of a layer.
- the parameter sharing mechanism uses tied weights which means the value of

the weight applied to one input is tied to the value of weight applied elsewhere.

Stride operation:-

- stride is known as shifting factor.
- stride = 1 means shifting of 1 row & 1 column both horizontally and vertically.
- let us say n : order of input matrix
 k : order of kernel

s :- shifting factor (stride)

$$\left(\frac{n-1}{s} + 1\right) \times \left(\frac{k-1}{s} + 1\right)$$

e.g., $n=6$, $k=3$, $s=2$.

The output matrix would be the size 2×2 .

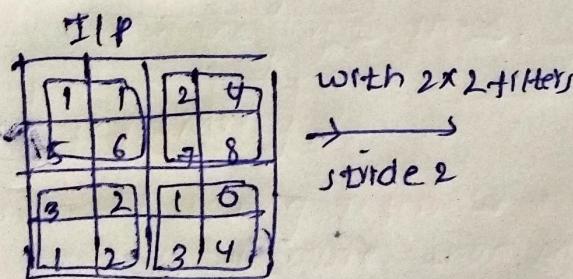
- stride is used to reduce the size of the matrix x and to get a smaller output matrix.

Pooling operation:-

- A pooling function replaces the output of the net at a certain location with a summary statistic of the nearby outputs.
- consider for instance of image size 96×96 , and suppose we have learned 400 features over 8×8 inputs.
- convolution operation results in an output of size $(96 - 8 + 1) \times (96 - 8 + 1)$

- suppose we have 400 features (dimensions)
this results in a vector of size: 792×400
- it is very difficult to process and prone to overfitting.
- formally, after obtaining convolved features as described earlier, we divide the size of region say $m \times n$ to pool our convolved features
we divide our convolved features into disjoint $m \times n$ regions, and take mean of it.
- It is used to reduce the computation complexity.
- for translation and rotational it is invariant.
- Types of Pooling:-
 1) max pooling
 2) mean "
 3) Average "

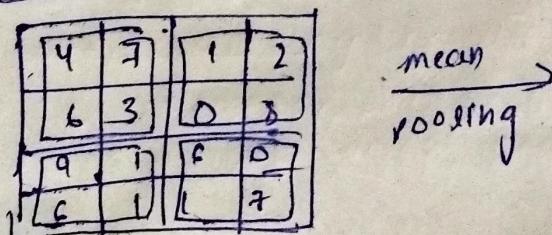
Max Pooling:-



O/P

| | |
|---|---|
| 6 | 8 |
| 3 | 4 |

Mean Pooling:-



| | |
|---|------|
| 5 | 2.75 |
| 5 | 3.5 |