# Lecture 22 - 19/11/24

## Weight Initialization

- Can we initialize all weights to an identical value?
    - No, because all nodes within the same hidden layer will behave identically
    - this is known as symmetry problem.

## Gradient Vanishing

- In a deep neural network:
    - Vanishing → gradients become very small = W's barely update.
    - Exploding → gradients become very large = W's updating causes instability.
    - Releva elements of the problem : activation function, weights initialization

- Vanishing : mostly related to activation
    - i.e Sigmoid:

        For its activation there is a large difference between variance (range) of input, output.

        This means that output gets quickly close to 0 or 1 where derivatives are near zero

        Many layers - problem

- Exploding : mostly related to large weights.
    - Leads to instability.

- How to initialize weights to reduce vanishing & exploding?
    - Good technique → keep equal variance of activations across layers.
    - There are some heuristics:

        Glorot : used for sigmoid & tanh

        He : used for ReLU

## Pro & Cons of various activation

- Sigmoid

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

    Toward either end of sigmoid, output changes slower. (gradient vanishing)
    Computationally expenses.

- Tanh

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

    Bigger range, derivatives are steeper.
    Suffers from vanishing gradient.
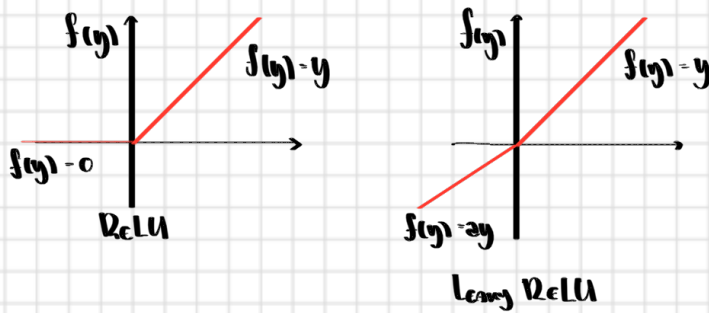    Computationally expensive due to exponential.

- ReLU

    $$f(x) = \max(0, x)$$

    Computationally efficient
    Dying ReLU problem might happen → negative input implies gradient = 0

# Leaky ReLU

- ^ It is an attempt to solve the dying ReLU problem
- ^ Leak increases range output of ReLU $(-\infty, \infty)$
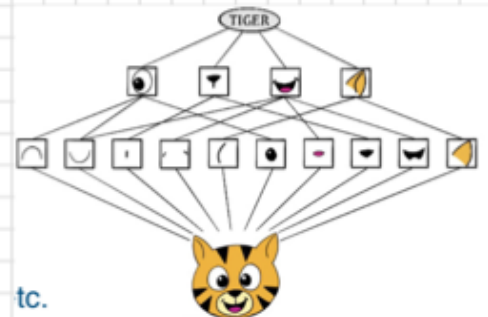- ^ $f$ and its derivatives are monotonic



ReLU

Leaky ReLU

$f(y) = 0$  $f(y) = y$  $f(y) = ay$

## Convolutional Neural Network

- With image flattening we loose **spacial dependencies**.
- When using **fully connected** NN for an image:
  - Lose **spatial dependencies**
  - **Large # of parameters** :
    - ○ Computer needs a lot of memory!
    - ○ too many params lead to **overfitting!**

**IDEA**

- ^ Check for certain **features in different image patches**
- ^ Create a **hierarchy** of features.
- ^ Handle **variations** of the same object.



tc.

## Filters & Convolution Operation

- A **FILTER** is a **matrix that is much smaller than the image.**
- **Filters** are also known as **"KERNEL".**
- **Convolutional sum** is a **linear operation.**
- **Slide the filter over input matrix:**
  - Execute convolution
  - multiply **overlapping values** & **sum them up** → Gives back final pixel



RESULTS IN

$$\begin{bmatrix} 1 & 1 & 1 \\ 0 & 1 & 1 \\ 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} = 1 + 0 + 1 + 0 + 1 + 0 + 0 + 0 + 1 = 4$$

## PADDING

^ **OUTPUT IMAGE IS SMALLER THAN INPUT IMAGE**

? How to apply filters to BORDERS?

? What if we want THE SAME SIZE?

^ Solution:

- ADD EXTRA PIXELS ON RIM OF ORIGINAL, DEPENDING ON KERNEL'S SIZE



| Input | Kernel | Output |
|---|---|---|

## STRIDE
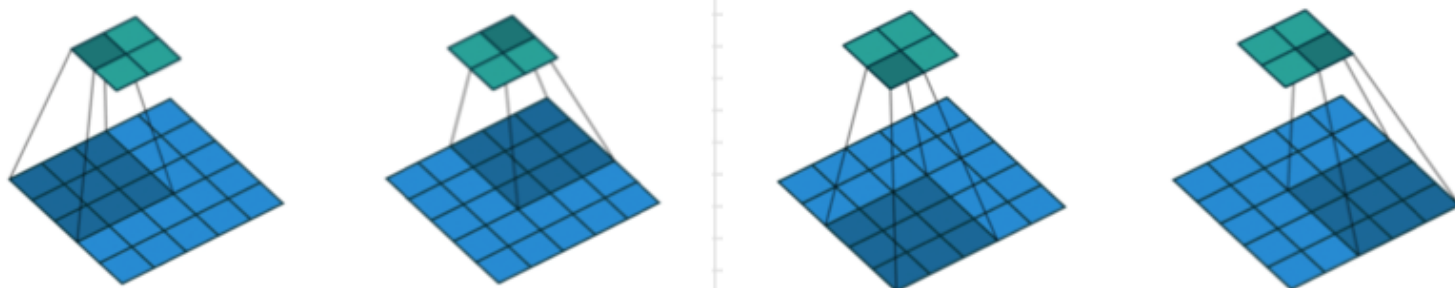
^ We START WITH THE CONVOLUTION WINDOW IN UPPER LEFT of the input tensor

^ then SLIDE IT OVER ALL LOCATIONS BOTH DOWN & TO THE RIGHT.

^ **STRIDE** LENGTH :
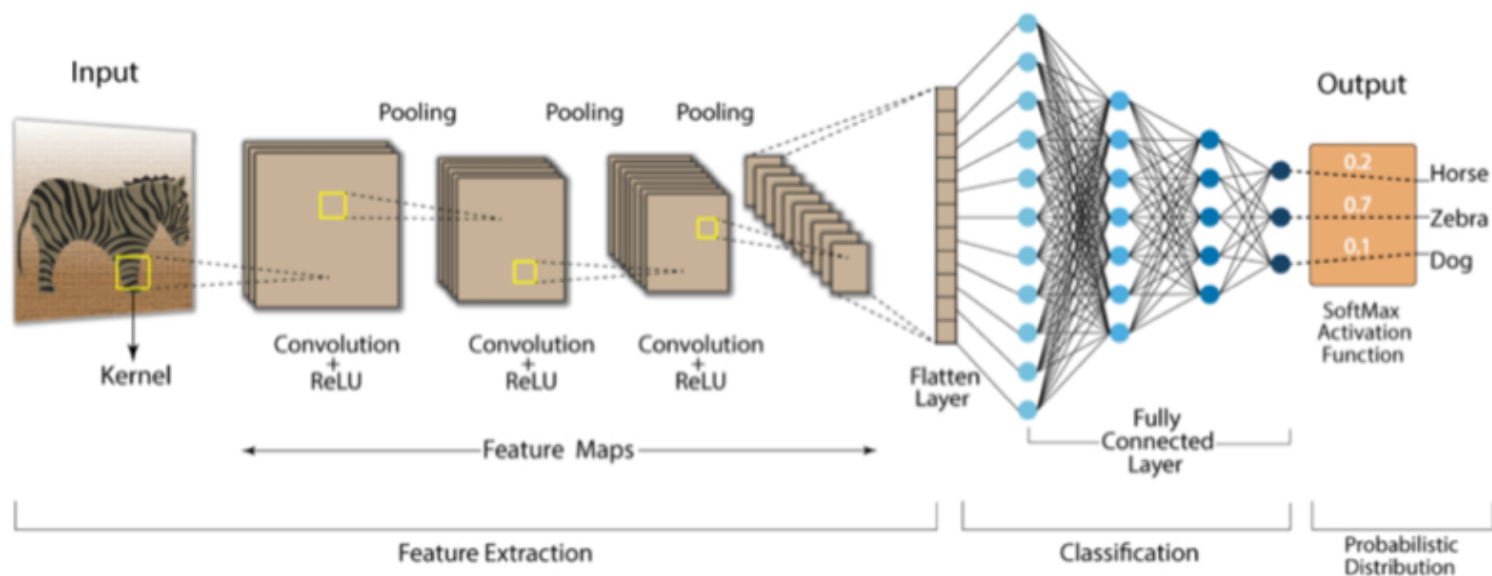
- # OF STEPS WE TAKE WHEN SLIDING OUR FILTER across an IMAGE.

^ It is a DOWNSAMPLING METHOD that SKIPS PIXELS WHEN MOVING THE FILTER across IMAGE
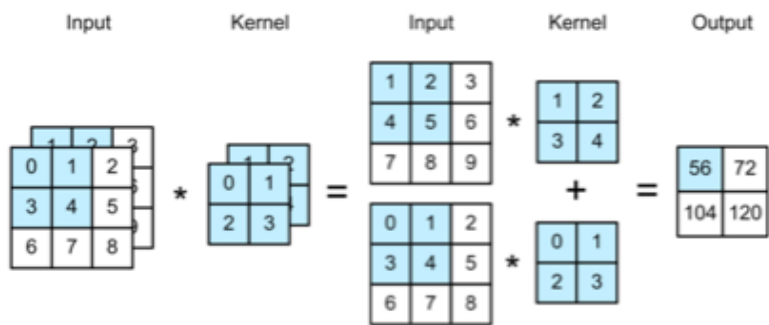


# Convolution Neural Network (CNN)



Input — Kernel

Pooling  Pooling  Pooling

Convolution + ReLU  Convolution + ReLU  Convolution + ReLU

Flatten Layer

Fully Connected Layer

Output

SoftMax Activation Function

0.2 — Horse
0.7 — Zebra
0.1 — Dog

← Feature Maps →

Feature Extraction          Classification          Probabilistic Distribution

## MULTIPLE INPUT CHANNELS

^ We OFTEN HAVE MULTIPLE INPUT CHANNELS:

- Images may come with 3 CHANNELS FOR RGB

^ We Implement MULTI-CHANNEL FILTERS to DEAL WITH MULTI-CHANNEL INPUTS.

| | Input | | | Kernel | | Input | | | Kernel | | Output | |

^ WOULD LIKE TO APPLY MULTIPLE DIFFERENT FILTERS OVER THE IMAGE

^ DIFFERENT FILTERS DETECT DIFFERENT IMPORTANT INFORMATION

## POOLING

^ DOWNSAMPLING A FEATURE MAP & REDUCING ITS SIZE → CAN EMPHASIS DOMINANT FEATURES.

^ MAX-POOLING NORMALLY BETTER THAN AVERAGE-POOLING

^ NORMAL TO INSERT POOLING LAYER BETWEEN SUCCESSIVE CONVOLUTIONAL

### Extract features

^ EXTRACTION FROM LOW-LEVEL TO HIGH-LEVEL

   - FLATTEN LAST LAYER'S OUTPUT.

   - APPLY FULLY CONNECTED LAYERS

^ BY ADDING FULLY CONNECTED LAYER, NETWORKS LEARNS NON-LINEAR COMB.

^ IN CNN:

   - CONVOLUTION CAN BE SEEN AS A NEURAL NETWORK LAYER WHERE WEIGHTS ARE SHARED.