# How to improve the performance of Neural Networks [15]

**Vanishing Gradients**
- Activation f$^n$
- Weight Initialization

**Normalization**
- Normalizing inputs
- Batch Normalization
- Normalizing Activations

**Optimisers**
- Momentum
- Adagrad
- RMS prop
- Adam

**Hyperparameter Tuning**
- No. of hidden layers
- Nodes/layer
- Batch size

**Overfitting**
- Reduce complexity, ↑ data
- Dropout layers
- Regularization ($l_1$ & $l_2$)
- Early stopping

**Gradient Checking & Clipping**

**Learning rate Scheduling**

## 1.) Vanishing Gradients

We improve the performance of Neural Networks by solving the vanishing gradient problem that occur during backprop. in deep neural networks using the below methods :-

## ☆ Activation functions

Activation f$^n$ adds non-linearity which help model learn complex non-linear patterns.
Without Activation f$^n$, a neural network would behave just like a simple linear eq$^n$ no matter how many layers you add.
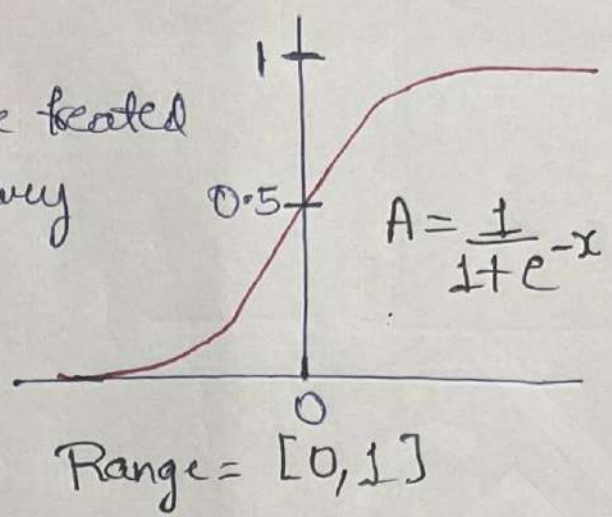
→ What is an Ideal / Good Activation f$^n$.

- The activation f$^n$ should be non-linear, eg sigmoid $= \frac{1}{1+e^{-z}}$

- It should be differentiable.

- It should be computationally inexpensive.

- It should be zero centered / Normalized.

- It should be non saturating :-
  like, Relu does not squeeze the input b/w a fixed range like sigmoid / tanh squeeze & Vanishing Gradient Occurs.

## → Sigmoid Activation $f^n$

**→ Advantages:**
- Output is in range [0,1], can be treated & used in output layer like binary classification as probability
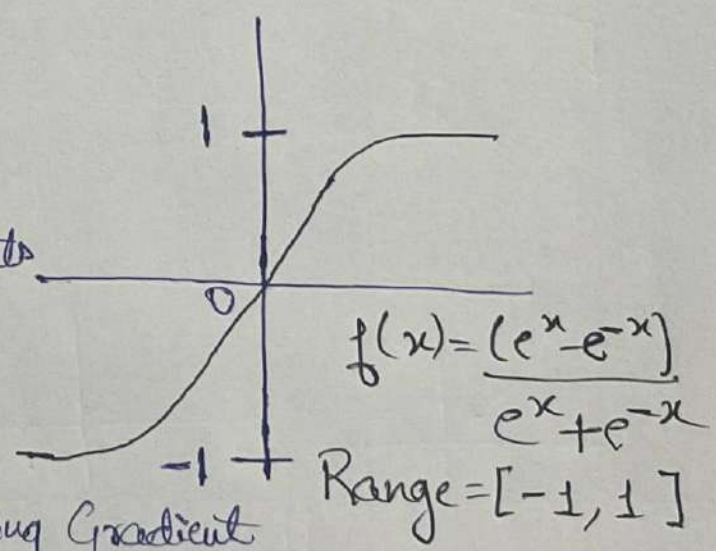
- Can capture Non-linear patterns
- Differentiable $f^n$.

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

Range = [0, 1]

**→ Disadvantages**

- Saturating $f^n$, squeezes the input in a range causes Vanishing gradient problem. [Rarely used in hidden layer now].

- Non-zero centered, mean is not = 0, training is slow because the gradient of one layer is either '+' or '-' for all nodes.

- Computationally expensive as we calculate exponents.

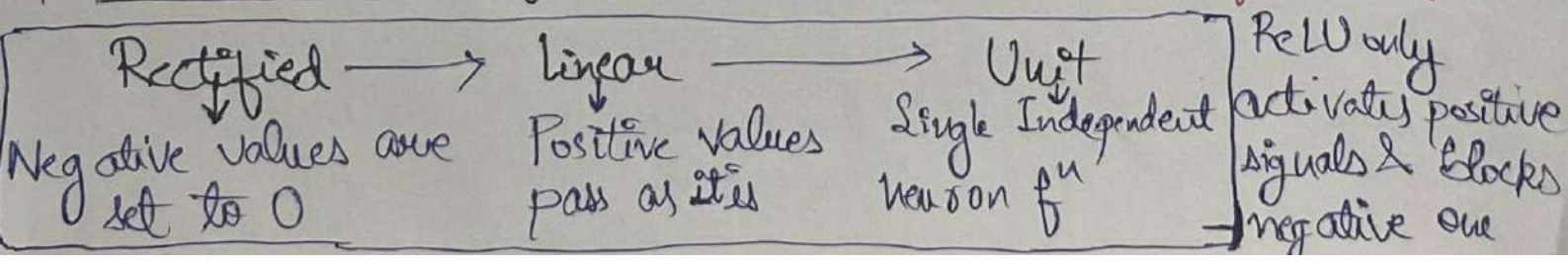## → Tanh Activation $f^n$

**→ Advantages:**

- Non-linear
- Differentiable
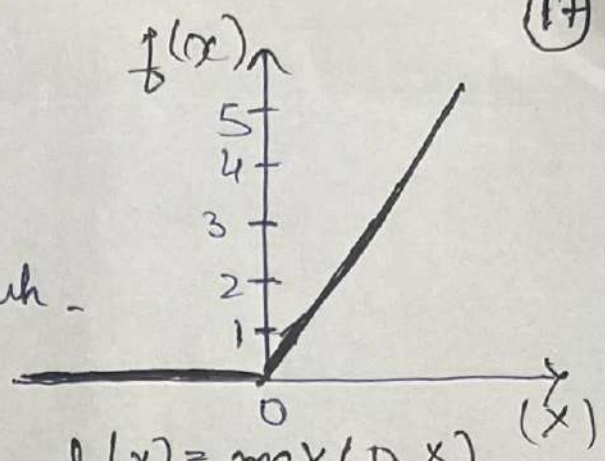- Zero centered, both +ve -ve gradients thus training is faster.

**→ Disadvantages:**

- Saturating $f^n$, prone to Vanishing Gradient

- This is also computationally expensive.

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

Range = [-1, 1]

## → ReLU Activation $f^n$ [One of the best Activation $f^n$ today]

| Rectified | → Linear | → Unit | ReLU only activates positive signals & blocks negative one |
|---|---|---|---|
| Negative values are set to 0 | Positive values pass as it is | Single Independent neuron $f^n$ | |

→ **Advantages:-**
- Non-linear
- Not-saturated in the positive region
- Computationally inexpensive
- Convergence is faster than sigmoid & tanh.



$$f(x) = max(0, x)$$
$$range = [0, \infty)$$

→ **Disadvantages:-**
- Not completely differentiable at 0.
- Not 0 centered {Batch Norm. helps solve this}

→ **Dying Relu Problem:-**

- Relu is = max(0,x), so for all negative value ReLU outputs 0 So the problem is for negative input gradient = 0, weight stops updating so it never learns again.

| Reasons | Because of these | Solutions |
|---|---|---|
| • Large Learning rate. | the value of $f^n$ become '−' so the output is 0 | • Set lower learning rate |
| • Bad weight Initialization | | • Use positive weights. |
| • Bias become '−'. | | • Use Bias +ve, 0.01 - scientific pro |
| • Negative Input data | | • Use variations of ReLU. |

**Based on linear & Non-linear transformations on ReLU**

**Linear**                            **Non-linear**

| Leaky Relu | Parametric Relu | ELU | SeLU |
|---|---|---|---|
| | | Exponential Linear Unit | Scaled Exponential Linear Unit. |
|  |  |  |  |
| $f(z) = max\left(\frac{z}{100}, z\right)$ | $f(z) = max(\alpha x, x)$ | | |
| $z \geq 0 \rightarrow z$ <br> $z < 0 \rightarrow \frac{z}{100}$ | $f(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha x & \text{otherwise} \end{cases}$ | $f(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha(e^x - 1) & \text{if } x \leq 0 \end{cases}$ <br> constant | $f(x) = \lambda \begin{cases} x & , x > 0 \\ \alpha(e^x - 1) & , x \leq 0 \end{cases}$ <br> $\alpha = 1.67, \lambda = 1.05$ fixed constants |
| • For − value output is not 0 it is 0.01z | • 'a' is a trainable parameter calculated during training | • This $f^n$ is always continuous & always differentiable | • Self-Normalizing without using batch Normalization |
| • Close to 0 centered | • More Flexible. | • Faster convergence | • More Research is being done on SeLU |
| • Solves Leaky ReLU Problem | • all same as Leaky Rest | • Computationally expensive <br> • Rest all same as ReLU. | |

# * Weight Initialization:

→ So step 1 is crucial to initialise the right value to 'w' & 'b' else we might encounter the below problems :-

- Vanishing Gradient Problem.
- Exploding Gradient Problem.
- Slow Convergence.

---

**Steps in Deep Learning**

1) Initialize the params

2) Chose an optimizing algo (eg GD)

3) Repeat these steps:
- Forward propogate input.
- Compute the cost $f^n$.
- Compute gradients of cost wrt params using backprop.
- Update each param using gradient

---

→ What shouldn't be done?
  (Wrong ways to initialise weights)

- Zero initialization = $\boxed{w=0, b=0}$, No training happens as weights & bias never update with this formula $\left[w = w - n\frac{dl}{dw}\right]$

- Non-zero constant = $\boxed{w=1.5, b=1.5}$, All nodes behave identical as same activation, behaves like linear model only.

- Random with small weights = $\boxed{np.random.randn(shape) \times 0.01}$ cause vanishing gradient problem & slow convergence.

- Random with large weights = $\boxed{np.random.randn(shape)}$, slow training & vanishing Grad., & saturation happens.

The learning is we have to initialise randomly but range matters.

→ What can be done?
  (right ways to initialise weights)
  Right way is to multiply Rand weights with $\frac{1}{n}$, n = no. of input to node

                        Methods

Xavier/Glorant (used with tanh)                    He init (used with ReLU)

Normal          Uniform                    Normal              Uniform

$= \sqrt{\frac{1}{n}}$     $= [-limit, limit]$         $= \sqrt{\frac{2}{n}}$      $= [-limit, limit]$ range

              $limit = \sqrt{\frac{6}{m+n}}$                          $limit = \sqrt{\frac{6}{n}}$

n = no. of inputs     m = input to node          n = no. of inputs      n = no. of inputs to
coming to the         n = no. of output to       to the node            node
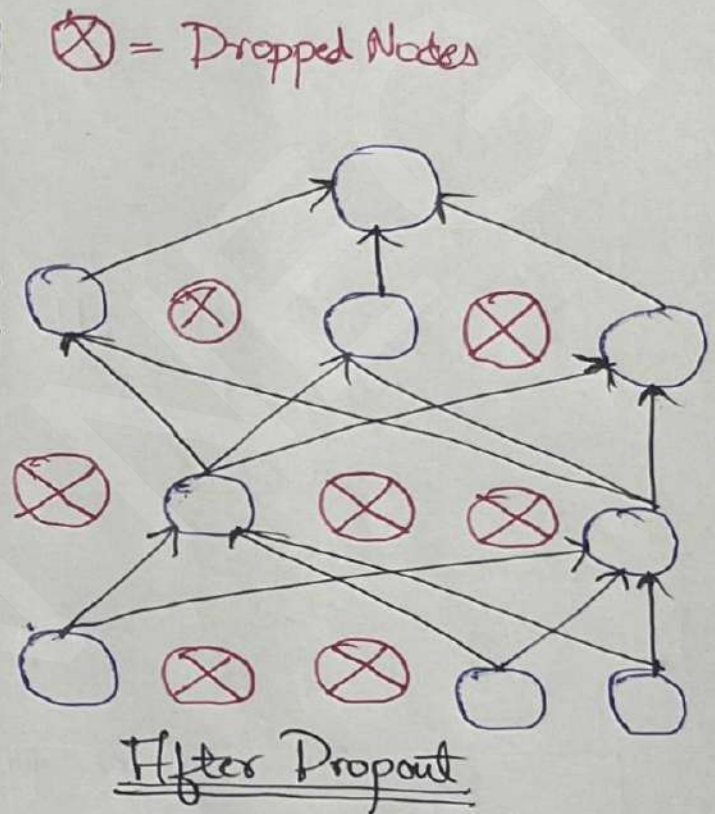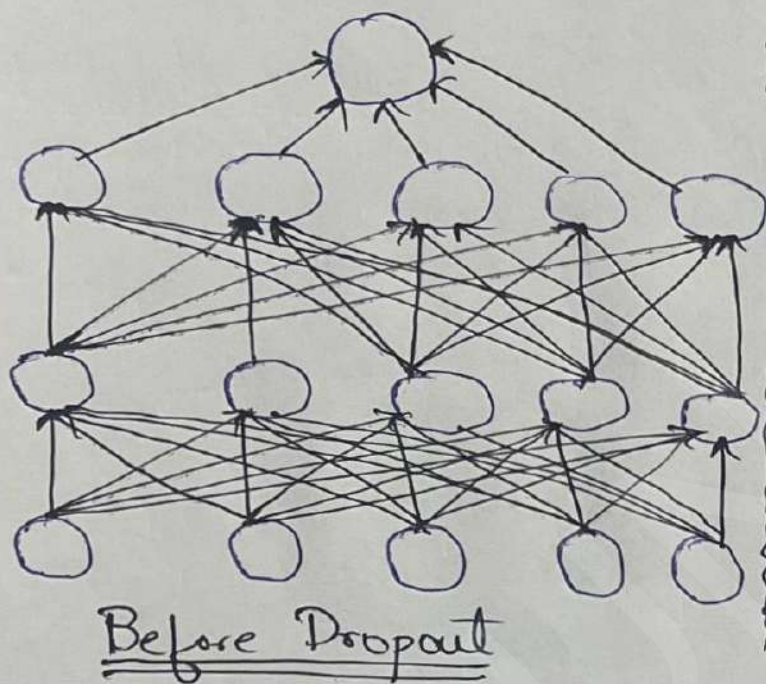node                  node

2) # Overfitting

→ Neural networks are prone to overfitting as they are complex.
→ Overfitting happens when NN perform well on training & poor on test.
Various ways to handle it are as follows:-

* ## Reduce Complexity & ↑ data - Basically by reducing no. of nodes or layers & by increasing the input data.

* ## Dropout Layers -

⊗ = Dropped Nodes



Before Dropout

After Dropout

→ Basically in each epoch we randomly drop nodes based on a factor p that we initialize.
$p = 0.5$ means randomly drop 50% nodes in each epochs.

→ Dropout matches the random forest analogy because in R-F also we introduce randomness & create multiple model variation & reduce overfitting through ensemble effect.

→ How does prediction work?
• For training we don't have many 'w' randomly — since dropout is ON
• During testing / inferencing dropout is OFF because we don't need randomness during prediction.
• So if more nodes active during testing, output will be larger so we do weight scaling. If $p = 0.3$, then $w = w(1-p)$ during testing.

→ Note :-

* Overfitting = Increase value of P.
Underfitting = Decrease the value of P.

* Start once by applying dropout only on last layer.
* General good range of dropout = 10% to 50%

→ Drawbacks of dropouts :-
* Convergence is slow.
* Its loss fⁿ value changes a lot, difficult to interpret gradients/model.

# ☆ Regularization -

* Regularization adds a penalty term to the loss fⁿ, the penalty term punishes the complex models [Large weights should be penalised]
* Large so that the model does not overfit.

## Types



| L₂ Regularization (most common) | L₁ Regularization |
|---|---|
| • If weight is small → penalty is small | • Same weight & penalty ratio as L₁ |
| If weight is large → penalty is large | • But penalty ↑ or ↓ linearly |
| Penalty form: $\lambda \Sigma w^2$  [λ controls the strength] | Penalty form: $\lambda \Sigma |w|$ |
| $Loss = L(y, \hat{y}) + \lambda \Sigma w^2$ | $Loss = Loss + \lambda \cdot Penalty\ term$ |
| Weights are pulled toward zero But never zero. | Weights directly pushed to zero Many weights become exactly zero. |

# ☆ Early Stopping.

* Early stopping helps us stop training before overfitting starts.
* We split the data into train & validate data & train epoch by epoch. After each step we check validation loss. If validation loss is improved — continue training else stop training at that epoch if validation loss does not improve for some time.
* Early stopping stops the model where it performs best on unseen data, not where training accuracy is maximum.

# 3) Normalization

→ Problem without Normalization :-

- eg Age→ 0 to 1000, Salary→ 10,000 to 10,00,000, Rating→ 1 to 5
- Large-value features dominated & small value features got ignored
- Loss surface becomes zig-zag & learning became slow & unstable.

## ＊ Normalising inputs:

- Normalising makes all inputs comparable & in same scale, centered.
  So model treats all input features fairly.

<div align="center">Two methods</div>

| Standardisation | Normalisation |
|---|---|
| ▪ Centeres data around mean & scales to standard deviation | • Scales data to a fixed range (usually 0 to 1) |

$$= \frac{X_i - \mu}{\sigma} \quad \begin{array}{l}\mu = mean \\ \sigma = std. dev\end{array}$$

$$= \frac{X_i - X_{min}}{X_{max} - X_{min}}$$

## ＊ Batch Normalisation : [Normalization done using batch of data]

- It is used to normalize the output of each layer during training
- It makes training faster, more stable & less sensitive to 'w' initialization
- Output of layer 1 ──→ input to layer 2, so even if input is standardized internal values keep changing & we apply Batch Normalisation.
  ↳ Internal Covariate Shift - Distribution of inputs to each layer keep changing as weights update.

- So, standardization or Normalization only happen once on input data.
  But Batch Normalization is applied at every layer [every hidden]

> Steps - Neuron give output = $Z = WX + b$ → Normalize Activation $(Z - \mu/\sigma)$
>   Scale = how wide the values should be          ← — Scale & shift $(= \gamma \times Norm. value + \beta)$
>   Shift = Where center should be                     This is done as at this step also
>   ↳ And then Activation fⁿ is applied          all neuron will have same value out
>                                                                      so we introduce γ-scale, β-shift.

- During Training - Batch norm. normalizes Activations using mini batch $\mu$ & $\sigma$ ＊＊＊
- During Testing — To normalize use running $\mu$ & $\sigma$ 'learned during training
  Each node is normalized independently but using values from all samples in the batch ＊＊

# 4) Gradient Checking & Clipping

These are also techniques used to improve performance & stability of neural networks during training.

→ **Gradient Checking** = It is a debugging technique used to verify whether backpropogation gradients are computed correctly or not.

> If backprop grad ≠ Manually calculated numeric gradient → means bug in code

→ **Gradient Clipping** = It is used to limit the magnitude of gradients to prevent exploding gradients during training.

> If gradient Becomes very large → clip it to a fixed limit.

---

# 5) Optimisers

→ We use optimizers to speed up the training process.

→ We studied gradient descent & its variation optimizer algos that updates model weights to reduce loss during training but there are few challenges:-

- Single learning rate for every weight / both directions.
- Stuck in local minima possibility.
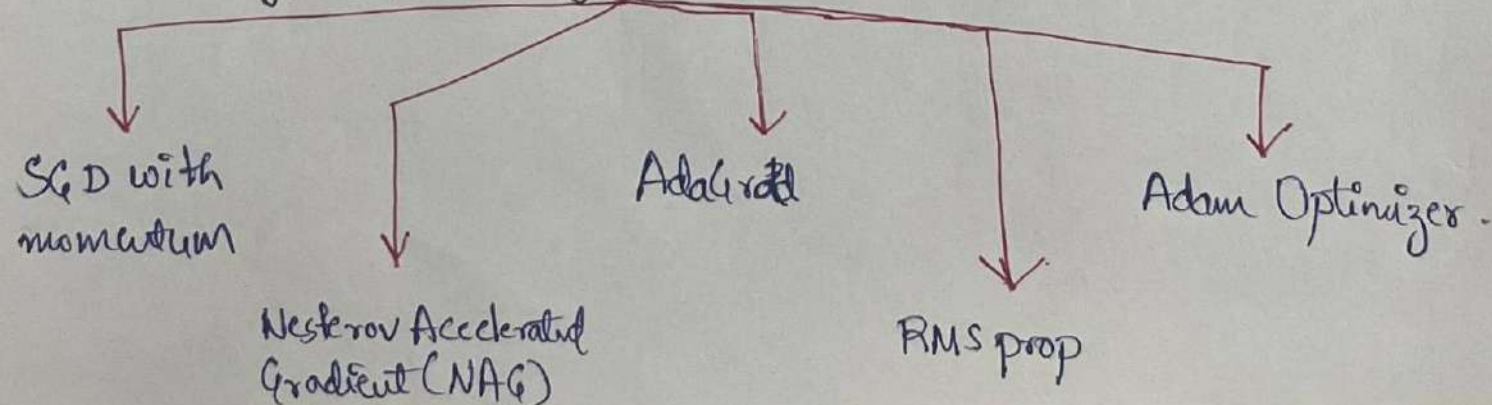- Saddle point handling- where slope is zero, so no weight update

→ Exponential weighted Moving Average (EWMA):-   **★★**
A technique that computes a weighted average where recent values get higher weightage & older value decay exponentially.

$$V_t = \beta V_{t-1} + (1-\beta) O_t$$

→ In normal average all past values get equal weight in real life recent loss/gradient is more important.

## Types of Optimizers

```
                  Types of Optimizers
        ┌──────────────┼──────────────┬──────────────┐
        ↓              ↓              ↓              ↓
    SGD with       AdaGrad                      Adam Optimizer.
    momentum           ↓                 ↓
        ↓          Nesterov Accelerated  RMS prop
                   Gradient (NAG)
```

SGD with momentum

AdaGrad

Adam Optimizer.

Nesterov Accelerated Gradient (NAG)

RMS prop
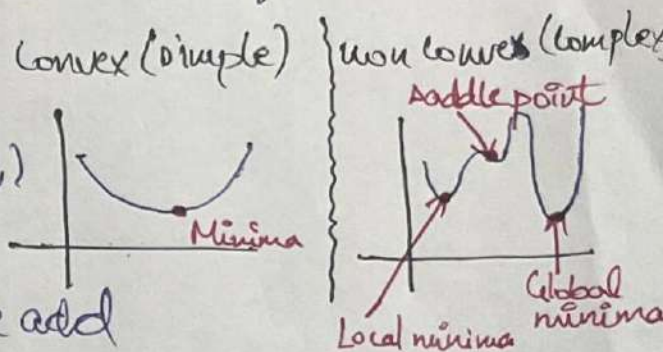
# ✱ SGD with momentum

→ In complex deep learning the loss f^n is also complex so to reach minima is tough because:-

1) Local minima
2) Saddle point (one side gradient ↑ & other direc ↓)
3) High Curvature - very sharp bend/curve.

→ So SGD deals with above problems so we add momentum to it & <u>convergence speed increases</u>.

Convex (simple) | non convex (complex)
saddle point
Minima | Local minima | Global minima

Physics - If a ball is moving downhill, it doesn't stop immediately, it keeps moving because of momentum.

So instead of using current gradient it also remembers past gradients -

| Update = Current gradient + past direction (memory) |, this

- reduces zig zag and move faster in correct direction & smooth learning

SGD

SGD with momentum

# ✱ Nesterov Accelerated Gradient (NAG)

→ Problem with SGD with momentum was it could overshoot oscillation
→ NAG solves this by [ let me first see where my speed will take me then I'll correct before actually moving

NAG calculates the gradient after looking ahead (peeking) in the momentum direction, not at the current position.

# ✱ AdaGrad (Adaptive Gradient)

→ Main idea = Different parameters should have diff learning rates.
→ This help in handling sparsity by giving rare features higher learning rate, so sparse features learn quickly & converge fast.

→ Disadvantage :- { Now, we don't use AdaGrad in Deep learning }
- learning rates keep decreasing continuously because it accumulates all past gradients, which can make learning rate extremely small & cause training to stop early.

# ⭐ RMS Prop (Root Mean Square Prop)

→ RMSProp improves AdaGrad by using an exponential moving average of squared gradients instead of accumulating all past gradients, thus prevents learning rate from becoming too small.

AdaGrad remembers everything, RMS prop remembers only recent history.

# ⭐ Adam Optimizer (Adaptive Moment Estimation) ⭐⭐⭐

→ Currently most powerful optimizer, used in ANN, CNN, RNN.
→ Mix of Momentum & RMS prop.

| Update = Direction from momentum + learning rate control from RMSProp |

→ Default parameter in Deep Learning.

# 6) Learning Rate Scheduling

→ Learning rate decides how big step your model will take while learning/training.

Large LR → fast but unstable

Small LR → stable but very slow.

→ LR scheduling idea = Start with higher learning rate & gradually reduce it during training.

→ Common types (step decay, Reduce on Plateau, Cosine Annealing)
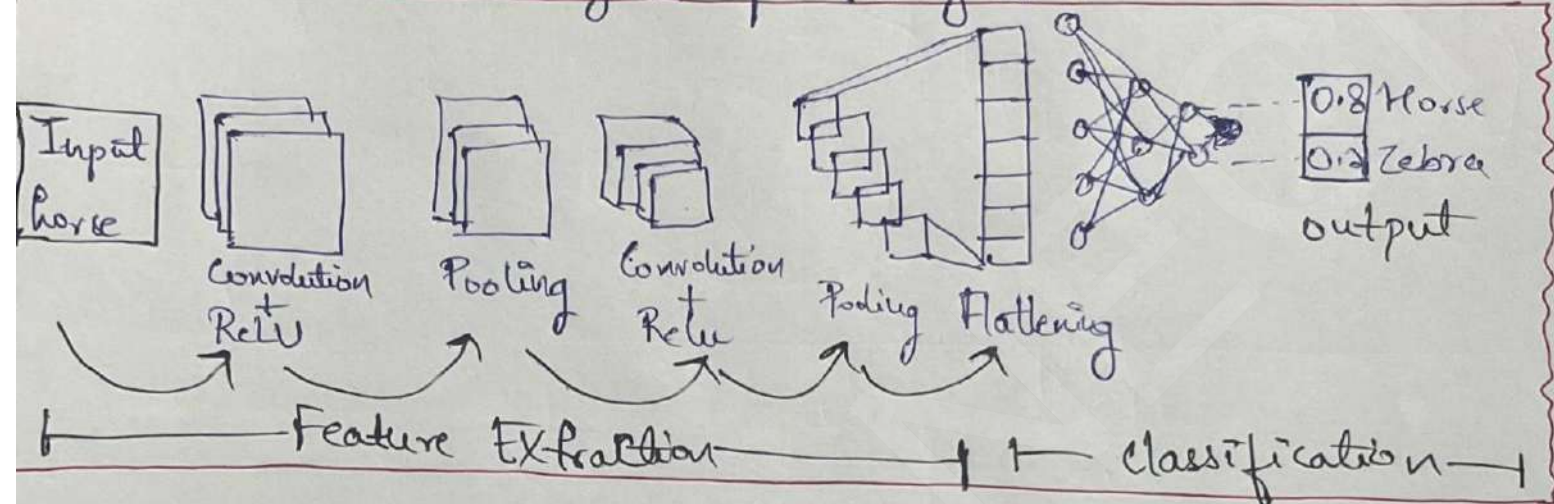
# IT'S NEVER TOO LATE

-TARUN NEGI

# BONUS

## Convolutional Neural Network (CNN)

→ CNN is a deep learning model mainly used for images

→ Images have many pixels & too much data, normal N·N treat images like normal numbers & don't understand shape but CNN understands edges, shapes & weights.



| Input horse | Convolution + ReLU | Pooling | Convolution ReLU | Pooling | Flattening | | 0.8 Horse / 0.2 Zebra output |

Feature Extraction ————————— Classification

★ <u>Filters/Kernel</u> - Small matrices that slide over input image matrix

★ <u>Convolution operation</u> - Each filter convules (slides) over the input image, computing the dot product b/w the filter & the input patch. This generates a feature Map [matrix]
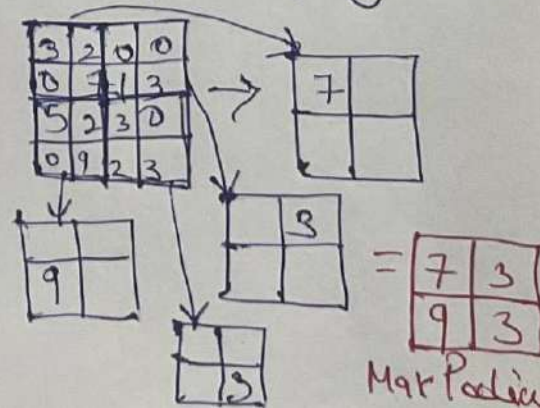
★ <u>Stride</u> = The step size by which the filter moves. Larger strides result in smaller feature maps.

★ <u>Padding</u> = Adding zero around input matrix to preserve the spatial dimension. The edges of the matrix also get equal chance/weightage while the convolution slides over it.
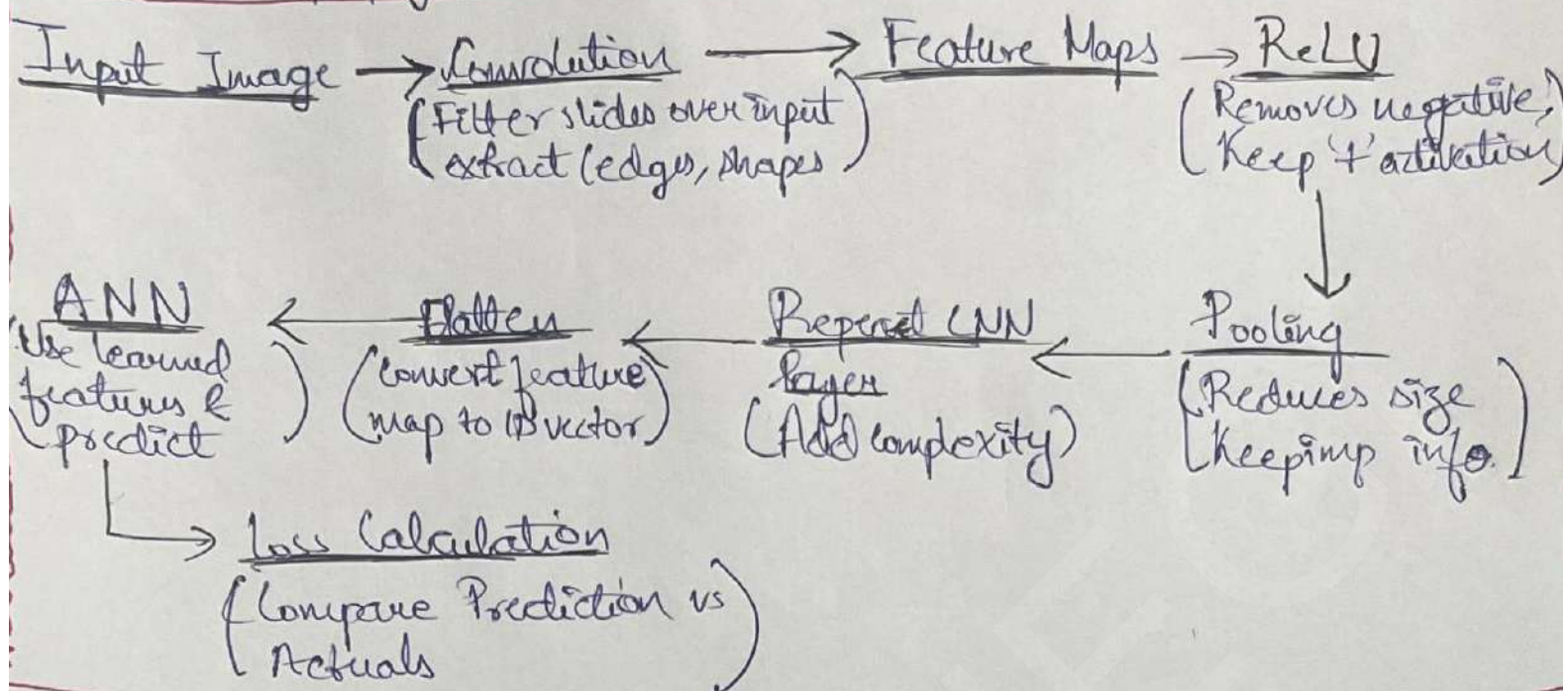
★ <u>Pooling</u> = Reduces the dimensions of feature map, keeping the most important information.

• Max pooling - Takes max value from each patch of feature map.

• Avg pooling - Takes the avg value from each patch of the feature map.



Max Pooling

→ **Forward propogation in CNN:- (BONUS)**

Input Image → Convolution ——→ Feature Maps → ReLU
                (Filter slides over input    (Removes negative,
                 extract (edges, shapes)      Keep +' activation)

                                                        ↓

ANN ⟵ Flatten ⟵ Repeat CNN ⟵ Pooling
Use learned   (convert feature)   layer        (Reduces size
features &     (map to 1D vector)  (Add complexity)  keeping info.)
predict

    └→ Loss Calculation
        (Compare Prediction vs)
        (Actuals)

→ **Backpropogation in CNN:** ✱✱✱

- Error flows backwards, CNN updates its weights.
- Back Calculate loss/gradient & adjust ANN weights using GD.
- Backprop through the flatten layer. (Just reshaping, no learning here)
- Backprop through pooling & Relu.
- And update filter weights through CNN.

→ **Transfer learning**

- Transfer learning means using a pretrained model for a new task.
- Model is already trained on large dataset.
- CNN layers are reused for feature learning.
- New fully connected layers are added for classification.
- Saves time, data & computation, works best when dataset is small.