

Date: \_\_\_\_\_

Sun Mon Tue Wed Thu Fri Sat

**Perceptron:** Basic building block of an ann.

- used for binary classification.
- Supervised. learning.

**Components:**

1- **Inputs:** Perceptron takes binary inputs or real valued inputs.

2- **Weights:** parameters that the perceptron learns. It determine influence of each input on output

3- **Summation function:** perceptron computes the weighted sum

$$z = \sum_{i=0}^n w_i x_i$$

4- **Activation function:** The weighted sum is passed to activation function to get output in binary. Threshold or Sign function.

- **weight update:** During training, Based on difference in true & predicted

$$w_i = w_i + (\eta \cdot y_{act} - \hat{y}_i)$$

**Limitations:**

- Only learn linear separable data : becz it finds single linear decision boundary
- only binary outputs.
- Sensitive to noise
- Single layer so unable to learn complex relations
- Original weight update lacks adaptive mechanism.

**Bias:** - Shifts the activation function

- Adjust the decision boundary.
- handle asymmetry in data
- improve flexibility

**Learning rate:** Controls the size of step during weight update

Date:

Sun Mon Tue Wed Thu Fri Sat

### Multi Layer Perception:

- type of ANN made by multiple layers of neurons.
- Include input layer, 1 or more hidden layers and an output layer.
- Capable of learning complex data

### How Better than Single layer:

- layers learn progressively
- Network extract meaningful features from data.
- Shared parameters aid generalization.

### Gradient Learning:

- Gradient is a vector that points to steepest ascent
- we want the minimum loss so we descent to the <sup>minima</sup> ~~the~~
- Magnitude of gradient tells the rate of change of function
- In gradient descent the gradient guides the updates of parameters.
- As the gradient becomes smaller ~~sig~~ tells that model is reaching ~~in~~ a minimum.

Gradient Descent: it is an optimizing algorithm. used for finding minimum of a function.

- goal is to move towards the minimum of loss function by updating parameters.



Date:

Sun Mon Tue Wed Thu Fri Sat

- Steps:
- 1- Initialize parameters.
  - 2- Define the cost functions
  - ~~3- update Parameters.~~ 3- Compute Gradient
  - 4- update weights
  - 5- iterate until minimum

$$\text{Loss function} = \frac{1}{n} \sum_{i=0}^n (y^{(i)} - \hat{y}^{(i)})^2$$

$$\text{new\_parameter} = \text{old\_parameter} - \text{LR} \times \frac{\partial \text{Loss}}{\partial \text{parameter}}$$

Learning Rate: - hyperparameter

- Determine size of step taken toward descent

Too small LR: Converge slow, stuck at local minima.

Too big LR: Diverge, miss out minima.

Convergence: occurs when:

- 1- All point correctly classified
- 2- Error reaches global minimum.
- 3- optimal weights are set
- 4- LR is balanced.

Pattern Detection:

- Perceptron can identify black & white patterns.
- Struggles with grey pixels.

Date: .....

Sun Mon Tue Wed Thu Fri Sat

## Multi Layer Neural Network

- Forward pass: A learning example is fed to network and the output is compared with ground truth
- Backward pass: compute the derivative with respect to weights then weights are updated to make an optimized network.

### Activation Function:

- Sigmoid / Threshold function not differentiable
- Discontinuous at 0
- Sigmoid & tanh symmetric at 0
- ~~tanh~~ tanh in hidden & Sigmoid at output.

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

$$S'(x) = S(x) \cdot (1 - S(x))$$

$$\text{tanh} = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

$$\text{tanh}'(x) = 1 - \text{tanh}^2(x)$$

### Backpropagation:

- combined forward pass & backward pass to compute & propagate gradients.
- Iteratively adjusts weights to minimize error.



## Steps:

- 1- Initialization: initialize weights & biases with random values
- 2- forward pass:
  - 1- Input training set to network
  - 2- Apply activation function at each layer
  - 3- Compute predicted outputs.
- 3- Compute loss:
  - 1- Compare predicted output with actual output
  - 2- Use a loss function
- 4- Backward pass:
  - 1- Compute gradient of loss wrt. activation function
  - 2- compute gradient at each layer.
- 5- Optimization: 1- update the weights & biases
- 6- Iterative: Repeat 2, 3, 4 until minimum loss, or no. of epochs.

## Function Composition:

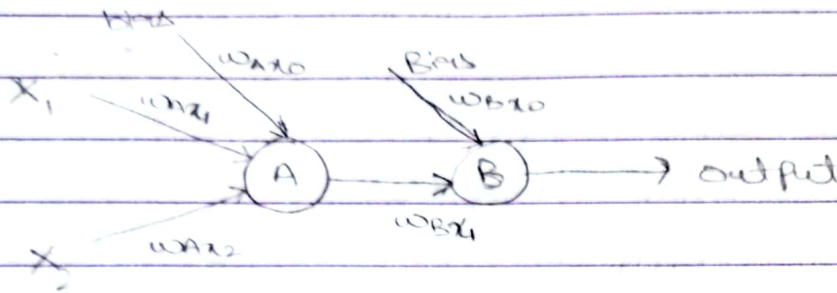
- 1- NN are compositions of functions with each acting as a function applied to output of previous layer
- 2- differentiating error wrt parameter helps minimize loss
- 3- Chain rule break computation of gradient layer by layer
- 4- ~~Chain rule~~ By applying chain rule errors can be efficiently, guiding parameter update.

3

5-

Date: \_\_\_\_\_

Sun Mon Tue Wed Thu Fri Sat



Neuron A = tanh

Neuron B = Sigmoid

error = MSE

Weights: Neu-A:  $w_{Ax0}, w_{Ax1}, w_{Ax2}$

Neu-B:  $w_{Bx0}, w_{Bx1}$

$$MSE = \frac{1}{n} \sum (y_{pred} - y_i)^2$$

$$Neu-A = \tanh \left( \overset{z_A}{w_{Ax0} + x_1 \cdot w_{Ax1} + x_2 \cdot w_{Ax2}} \right)$$

$$Neu-B = \sigma \left( \underset{z_B}{w_{Bx0} + Neu-A \cdot w_{Bx1}} \right) \rightarrow y_{pred} \text{ or } f$$

$$Error = \left( y_{act} - \sigma(w_{Bx0} + w_{Bx1} \cdot \tanh(w_{Ax0} + x_1 w_{Ax1} + x_2 w_{Ax2})) \right)$$

Error =  $\frac{(y - f)^2}{2}$  idk 2 kaha se aya.

Errors

Date: \_\_\_\_\_

Sun Mon Tue Wed Thu Fri Sat

$$\rightarrow \frac{\partial \text{error}}{\partial w_{Bx0}} = \frac{\partial \text{error}}{\partial y_{\text{pred}}} \cdot \frac{y_{\text{pred}}}{z_B} \cdot z_B$$

$$\rightarrow \frac{\partial \text{error}}{\partial w_{Bx1}} = \frac{\partial \text{error}}{\partial y_{\text{pred}}} \cdot \frac{y_{\text{pred}}}{z_B} \cdot z_B$$

$$\rightarrow \frac{\partial \text{error}}{\partial w_{Ax0}} = \frac{\partial \text{error}}{\partial y_{\text{pred}}} \cdot \frac{y_{\text{pred}}}{z_B} \cdot \frac{z_A \text{New}_A}{w_{Ax1} z_A} \cdot z_A$$

$$\rightarrow \frac{\partial \text{error}}{\partial w_{Ax1}} = \frac{\partial \text{error}}{\partial y_{\text{pred}}} \cdot \frac{y_{\text{pred}}}{z_B} \cdot \frac{\text{New}_A}{z_A} \cdot z_A$$

$$\rightarrow \frac{\partial \text{error}}{\partial w_{Ax2}} = \frac{\partial \text{error}}{\partial y_{\text{pred}}} \cdot \frac{y_{\text{pred}}}{z_B} \cdot \frac{\text{New}_A}{z_A} \cdot z_A$$

$$\frac{\partial \text{error}}{\partial y_{\text{pred}}} \cdot \frac{y_{\text{pred}}}{z_B} \quad \text{common in all}$$

$\frac{\partial y_{\text{pred}}}{\partial z_B}$  same chess

$$\ast \frac{\partial \text{error}}{\partial y_{\text{pred}}} = \frac{\partial \frac{(y-f)^2}{2}}{\partial y_{\text{pred}}} = \frac{-\cancel{2}(y-f)}{\cancel{2}} = -(y-f)$$

$$\ast \frac{\partial y_{\text{pred}}}{\partial z_B} = \frac{\partial S'(z_B)}{\partial z_B} = S'(z_B)$$

$$\frac{\partial z_B}{\partial \text{New}_A}, \frac{\partial \text{New}_A}{\partial z_A}, \text{common in last 3}$$



$$S' = S \cdot (1 - S)$$

$$\tanh' = (1 - \tanh^2)$$

Date:

Sun Mon Tue Wed Thu Fri Sat

$$\frac{\partial Z_B}{\partial \text{Neu}_A} = \frac{\partial (w_{Bx_0} + \text{Neu}_A \cdot w_{Bx_1})}{\partial \text{Neu}_A} = w_{Bx_1}$$

$$\frac{\partial \text{Neu}_A}{\partial z_A} = \frac{\partial (\tanh(z_A))}{\partial z_A} = \tanh'(z_A)$$

Calculating:

$$\rightarrow \frac{\partial z_A}{\partial w_{Bx_0}} = \frac{\partial (w_{Bx_0} + \text{Neu}_A \cdot w_{Bx_1})}{\partial w_{Bx_0}} = 1$$

$$\rightarrow \frac{\partial z_A}{\partial w_{Bx_1}} = \frac{\partial (w_{Bx_0} + \text{Neu}_A \cdot w_{Bx_1})}{\partial w_{Bx_1}} = \text{Neu}_A$$

$$\rightarrow \frac{\partial z_B}{\partial w_{Ax_0}} = \frac{\partial (w_{Ax_0} + x_1 w_{Ax_1} + x_2 w_{Ax_2})}{\partial w_{Ax_0}} = 1$$

$$\rightarrow \frac{\partial z_B}{\partial w_{Ax_1}} = x_1 \quad \rightarrow \frac{\partial z_B}{\partial w_{Ax_2}} = x_2$$

Weight updates:

$$\rightarrow \frac{\partial e}{\partial w_{Bx_0}} = -(y - f) \cdot S'(z_B) \cdot 1$$

$$\rightarrow \frac{\partial e}{\partial w_{Bx_1}} = -(y - f) \cdot S'(z_B) \cdot \text{Neu}_A$$

$$\rightarrow \frac{\partial e}{\partial w_{Ax_0}} = -(y - f) \cdot S'(z_B) \cdot w_{Bx_1} \cdot \tanh'(z_A) \cdot 1$$

$$\rightarrow \frac{\partial e}{\partial w_{Ax_1}} = -(y - f) \cdot S'(z_B) \cdot w_{Bx_1} \cdot \tanh'(z_A) \cdot x_1$$



Date:

Sun Mon Tue Wed Thu Fri Sat

## Data Sets:

- good data set must contain overrepresentation (diversity) of the real world

- 1- Handle missingy data: 1- Imputation 2- Deletion
- 2- Handle categorical data: 1- One hot encoding 2- Label encoding  
3- Ordinal encoding
- 3- Scaling Normalization: 1- MinMax Scaling 2- Standardization.
- 4- Handle outliers: 1- Identification 2- transformation (log, sqrt)  
3- Removal
- 5- Feature Engineering: 1- Create New feature 2- Reduce feature
- 6- Handle Imbalance data: 1- Resampling 2- Synthetic Data Generation

## Over-fitting:

Memorizing  $\rightarrow$  Overfitting

- 1- Increase training Set
- 2- Early Stopping

## Hyperparameter techniques:

- 1- Grid Search: Grid of parameters
- 2- Random Search: Choose random combinations
- 3- Bayesian Optimization: probabilistic model for parameters.
- 4- Gradient-based: Use Gradient based algos
- 5- Ensemble: Combine different models into one.
- 6- Automated: in

Date:

Sun Mon Tue Wed Thu Fri Sat

Training: Coming up with weights for the network is done before deploying into production

Inference: using network without updating weights

Optimal Model.

- Make data not biased
- Preprocess the data
- Split train, test, validation
- Proper loss functions
- Suitable optimizer
- Proper hyperparameter tuning.

Batch Size :- faster convergence

- Efficient memory usage



## Bias:

- Simple model complex data
- Don't capture patterns
- leads to under-fitting
- \* Use complex model
- \* Add features
- \* Add polynomial featur
- \* Decrease  $\lambda$

## Variance:

- Model's sensitivity to noise
- overly complex model captures noise
- leads to over fitting
- \* Simplify model
- \* Reduce feat
- \* Reduce poly feat
- \* Increase  $\lambda$

Regularization: ~~it~~ introduces penalty terms to reduce influence of inputs on outputs

Lasso L1: - adds sum of absolute values to loss function  
- Encourages sparsity, some coefficient may become 0

Ridge L2: - adds sum of squared values to the loss function  
- Spreads out the impact of each feature across all features, reducing impact of any single feature

$\lambda \rightarrow$  Regularization Strength: controls tradeoff between underfit overfit

$\rightarrow$  Prevent of overfit, Encourage Simplicity, Improve Generalization