

PPT 2

* Feedforward Neural Networks :- (multilayered NN).

→ feedforward :- connections between nodes, or 'neurons' do not form cycle. This means that info moves only in one direction from input \rightarrow hidden layers \rightarrow output.

Key aspects :-

- I/p layer, hidden layer, o/p layer

- forward propagation

- Activation functn

- Training

- Input to the network K is an n-dimensional vector.

- contains L-1 hidden layers having n neurons each.

(here 2 hidden layers)

- 1 o/p layer containing K neurons (corresponding K classes)

- Each neuron in hidden layer + o/p layer can be

split into 2 parts:

(i) activation \circlearrowright
(ii) preactivation \circlearrowleft

↑ preactivation (a_i)

o/p layer \Rightarrow 0th layer. o/p layer \Rightarrow Lth layer.
 b_1, b_2, b_3 are biases.

Preactivation at layer i is given by :-

$$a_i(x) = b_i + w_i h_{i-1}(x)$$

($w_i x_i + b_i$).

activation at layer i:-

$$h_i(x) = g(a_i(x))$$

g is activation function. (Eg: tanh, linear, etc)

- The activation at o/p layer is given by $f(x) = h_L(x) = o(a_L(x))$.

where o is o/p activation function
(Eg: - Softmax, linear; etc).

Model description:-

Data:- You have dataset with N samples, each

sample consists of i/p x_i , o/p y_i .

model:- The model is neural network funtn.

$$\hat{y}_i = f(x_i)$$

$$\hat{y}_i = f(x_i) = o(w_3 g(w_2 g(w_1 x_i + b_1) + b_2) + b_3)$$

↑ 1st layer i/p then
activn. function → then provided to
→ then go on till next layer
o/p layers activation.

w_1, w_2, w_3 are weight matrices for each layer
 b_1, b_2, b_3 are bias vectors.

g is activation function applied to o/p of each layer.

o : final activation funtn or o/p layer transformation

* Parameters (Θ): $\Theta = \{w_1, w_2, w_3, b_1, b_2, b_3\}$

* Total no. of parameters :-

$$P = (n \times n + n) \times (L-1) + (n \times K + k)$$

↑ bias ↑
 l = no. of layers

→ o/p layer has K neurons.

* Algorithm used for this model :- Gradient Descent with Back Propagation.

• Objective / Loss function : → Goal of training the n/w is to find parameters Θ that minimize the error betn predicted o/p's (\hat{y}_i) and true o/p's (y_i) .

$$\min \left(\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^K (\hat{y}_{ij} - y_{ij})^2 \right)$$

→ mean square error . functn :

• minimize loss function . $L(\Theta)$.

* The loss function $L(\Theta)$ could take different forms depending on specific problem. (Eg: mean sq. error for regression or cross-entropy for classification)

* Learning Parameters of FeedForward Neural Networks :-

• Gradient Descent Algorithm :-

GD is an optimization algorithm used to minimize the loss function of a model by iteratively adjusting its parameters.

Steps :-

1) Initialization:-

Set initial Parameters : Initialize model parameters (weights + biases) randomly. This is Θ^0 .

$$\Theta^0 = [w_1^0, w_2^0, \dots, w_L^0, b_1^0, b_2^0, \dots, b_L^0]$$

$w_i^0 + b_i^0$ are initial values for each layer.

2) Set hyperparameters:-

- Learning Rate (η): Determines step size for each iteration.
- maximum Iterations: A stopping criteria for algo (Eg: 1000 iterations).

3) Iteration Loop:-

i. compute gradients:-

forward pass:- compute predictions \hat{y}_i using (Θ_t) .

Compute loss: calculate loss $L(\Theta_t)$ based on predicted & true values.

Backward pass:- Use backprop. to compute gradient of loss functn w.r.t each parameter.

For a parameter Θ_j^t in Θ_t :-

$$\nabla_{\Theta_j^t} L(\Theta_t) = \frac{d L(\Theta_t)}{d \Theta_j^t}$$

These gradients include \rightarrow gra. of weights w_i^t and biases b_i^t for each layer.

ii. Update Parameters:-

Adjust each parameter Θ_j^t by moving it in the direction of negative gradient

$$\Theta_{t+1} = \Theta_t - \eta \nabla_{\Theta_t} L(\Theta_t)$$

\hookrightarrow vector of all gradients at iteration t

3) Increment counter.

Repeat until max no. of iterations is reached or changes in parameters become very small (convergence).

$$\nabla_{\theta} L(\theta) = \left[\frac{\partial L(\theta)}{\partial \theta_1}, \frac{\partial L(\theta)}{\partial \theta_2}, \dots, \frac{\partial L(\theta)}{\partial \theta_n} \right]$$

* Algorithm :- gradient-descent()

$t \leftarrow 0$

max-iterations $\leftarrow 1000$

Initialize $\theta_0 = [w_1^0, \dots, w_L^0, b_1^0, b_2^0, \dots, b_L^0]$

While $t++ < \text{max_iterations}$ do

$$\theta_{t+1} \leftarrow \theta_t - \eta \nabla \theta_t ;$$

end.

where $\nabla \theta_t = \left[\frac{\partial L(\theta)}{\partial w_{1t}}, \dots, \frac{\partial L(\theta)}{\partial w_{Lt}}, \frac{\partial L(\theta)}{\partial b_{1t}}, \dots, \frac{\partial L(\theta)}{\partial b_{Lt}} \right]^T$

* Output functions and Loss Functions :-

Loss functions for different types of problems :-

1. Regression Problem (Predicting ratings) :-

Ex: you want to predict ratings, where output y_i is a real number.

Loss function: Squared error loss is suitable for this type of problem. It measures the squared deviation betn true & pred value

$$L(\theta) = \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^3 (\hat{y}_{ij} - y_{ij})^2$$

2. Classification Problem:-

Ex: you want to classify image into one of K classes, so the true output y is a probability

distribution. Softmax functn is used here:-
Output functn:-

$$O(a_i)_j = \frac{e^{q_{Lj}}}{\sum_{i=1}^K e^{q_{Li}}} \quad \text{This converts raw scores into probabilities.}$$

Loss functn:- cross entropy loss is suitable for classification problems. It measures how well predicted prob distribution matches true distribution.

$$L(O) = - \sum_{c=1}^K y_c \log \hat{y}_c$$

y_c is 1 for true class, 0 otherwise.
 $L(O) = -\log \hat{y}_{\text{true}}$.

where \hat{y}_{true} is predicted probability for true class.

* Choice of o/p function:-

- for regression:- o/p functn should be linear bco you want pred values \hat{y}_i to be real numbers.
- for classifn:- o/p functn should be softmax to ensure that o/p's are probabilities (they sum to 1 and are all b/w 0 & 1).

• Log-Likelihood: for classification problem, the term $\log \hat{y}_{\text{true}}$ represents the log-likelihood of data given the model parameters.

Maximizing this log-likelihood helps in finding parameters Θ that best fit the training data o/p's.

	Real values	Prob.
o/p Activn	linear	Softmax
loss fun'	Squared error	Cross Entropy

* Back Propagation Algorithm:-

- Gradient descent :- Iteratively updates the network parameters to minimize loss by using gradients computed through forward + backward propagation.
- Steps seen earlier.

• Forward Prop. Algo:-

Purpose: To compute network's predictions + intermediate values during forward pass.

1. Layer by layer computation:

for each layer k from 1 to $L-1$:

Calculate a_k by adding bias to weighted sum of activations from previous layer.

$$a_k = b_k + w_k h_{k-1}$$

Compute h_k by activating applying activation functn g to a_k :

$$h_k = g(a_k)$$

2. O/p layer :- Compute final layer's pre-activation

$$a_L = b_L + w_L h_{L-1}$$

$$\hat{y} = o(a_L) \quad o \rightarrow \text{o/p activation functn.}$$

• Back Propagation :- computes gradients of loss w.r.t each parameter by propagating errors backward from o/p to i/p layer.

Calculate gradient of loss functn

$$\nabla_{a_L} L(\theta) = -(y - \hat{y})$$

Propagate gradient backwards :-

- Gradient w.r.t weights

$$\nabla_{w_k} L(\theta) = \nabla_{a_k} L(\theta) h_{k-1}^T$$

(gradients of pre activation \times activations from prev layer)

- Gradient w.r.t bias :-

$$\nabla_{b_k} L(\theta) = \nabla_{a_k} L(\theta) \text{ gradient of pre-activations.}$$

- Gradient w.r.t prev layer's Activations:

$$\nabla_{h_{k-1}} L(\theta) = w_k^T \nabla_{a_k} L(\theta).$$

" Pre layer's Actvtn :-

$$\nabla_{a_{k-1}} L(\theta) = \nabla_{h_{k-1}} L(\theta) \odot g'(a_{k-1}).$$

* How to compute g' :-

1. Logistic functn :-

$$g(z) = \sigma(z)$$

$$= \frac{1}{1+e^{-z}}.$$

$$\left(\frac{d}{dx} \frac{1}{x} = -\frac{1}{x^2} \right).$$

$$g'(z) = \frac{-1}{(1+e^{-z})^2} \cdot \frac{d}{dz} (1+e^{-z})$$

$$= \frac{-1}{(1+e^{-z})^2} \cdot (-e^{-z})$$

$$= \frac{1}{1+e^{-z}} \times \frac{e^{-z}}{1+e^{-z}} = \frac{1}{1+e^{-z}} \times \frac{(1+e^{-z}-1)}{1+e^{-z}}$$

$$= g(z) \cdot (1-g(z)).$$

2. tanh

$$g(z) = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}.$$

$$g'(z) = \frac{(e^z + e^{-z}) \frac{d}{dz} (e^z - e^{-z}) - (e^z - e^{-z}) \frac{d}{dz} (e^z + e^{-z})}{(e^z + e^{-z})^2}$$

$$= \frac{(e^z + e^{-z})(e^z + e^{-z}) - (e^z - e^{-z})(e^z - e^{-z})}{(e^z + e^{-z})^2}$$

$$= \frac{(e^z + e^{-z})^2 - (e^z - e^{-z})^2}{(e^z + e^{-z})^2}$$

$$= 1 - \frac{(e^z - e^{-z})^2}{(e^z + e^{-z})^2}$$

$$= \underline{1 - (g(z))^2}$$