

DL Assignment

Date _____
Page _____

-Bishal Saini.

Unit 1

Q1
→

Discuss Various Norm form for Deep Learning ?

- In Deep Learning, norms are mathematical tools used to measure the size, magnitude, or length of vectors and matrices.
- They play a crucial role in optimization, regularization and defining loss function.
- Different types of norms provide different properties, which can influence the learning behaviour of Neural Networks.
- The most commonly used norm forms in Deep learning are as follow

1. L0 Norm :- - The L0 norm is not true norm in the mathematical sense, but it counts the number of non-zero elements in a vector

$$\text{Formula} :- \|x\|_p = \left(\sum |x_i|^p \right)^{\frac{1}{p}}$$

2. L2 Norm: - It is also known as Euclidean norm as it measures the straight lines or Euclidean distance from the origin point to point x

- It is easy for the optimization and penalize large values more strongly.

Formula :-

$$\|x\|_2 = \sqrt{\sum x_i^2}$$

3. L1 Norm :- It is also known as Manhattan Norm

- It measures distance by summing absolute values
- It is robust to outliers compared to L2

Formula :-

$$\|x\|_1 = \sum |x_i|$$

4. L ∞ Norm :- It is known as max norm

- It takes the largest absolute value among all elements
- It represents the least case deviation in any direction.

Formula :-

$$\|x\|_{\infty} = \max|x_i|$$

5. L p Norm :- The most general family of norms is the L p norm

- It measures the length or size of a vector in p -dimensional space

Formula :-

$$\|x\|_p = \left(\sum |x_i|^p \right)^{1/p}$$

2. Illustrate the use of Jacobian and Hessian matrix. What is gradient based matrix



i) Jacobian Matrix :- - The Jacobian Matrix is the collection of first order partial derivatives of vector-valued function

- It also generalized the gradient for vector-valued functions

- For a function $f: R^n \rightarrow R^m$:

$$J_{ij} = \frac{\partial f_i}{\partial x_j}$$

For E.g:- In backpropagation the chain rules relies on jacobian matrices to propagate gradient efficiently through layers.

ii) Hessian Matrix:- - The hessian captures second-order derivatives

- It is a squared matrix of second-order partial derivatives.

For a scalar function $f(x)$, the hessian is an $n \times n$ symmetric matrix:

$$H_{ij} = \frac{\partial^2 f}{\partial x_i \partial x_j}$$

- If hessian eigen values are large, optimization may suffer from sharp minima or saddle points.

Gradient Based Optimization :-

- The core principle of training Neural Network is to minimize a loss function $J(\theta)$ with respect to parameter.
- Gradient Descent :- Updates parameter in the opposite direction of the gradient.
$$\theta \leftarrow \theta - \eta \nabla J(\theta) \quad (\eta \text{ is learning rate})$$
- First order methods :- Use gradients
- 2nd Order methods :- Use curvature to adapt step sizes

Unit 21
→

What is Dropout

- Dropout is a regularization techniques used in Deep learning to prevent overfitting.
- The core idea is simple, during training randomly "dropout" (deactivate) a fraction of neurons in the network
- Dropout Working :- - At each training step, every neuron has a prob. probability p of being set to zero (dropped)
- Now it creates different Sub network for each iterations
- During testing all neurons are used, but their outputs are scaled ($1-p$) so that the expected activation remain constant.

Mathematical Formula :-

If h is a vector of hidden activations :-
 $h' = h \cdot r$, $r \sim \text{Bernoulli}(p)$
 where r is a random binary used

At test time

$$h_{\text{test}} = (1-p)h$$

E.g :- Suppose we have a hidden layer with 6 neurons

- Without Dropout :- All six are active every iteration
- With Dropout :- ($p = 0.5$). Only 3 neurons are active iteration and the chosen neuron changes randomly across steps.

2
→

Explain the Algorithm

a] Stochastic Gradient Descent :- Instead of using the entire dataset, SGD updates parameters using one sample or mini batch at a time

$$\theta_{t+1} = \theta_t - \eta \nabla L(\theta_t; x_i, y_i)$$

- Advantage :- Faster, good for large datasets and helps to escape local minima

- Disadvantage :- It can be noisy and unstable

b] Momentum Algorithm :- It adds velocity term to update, (which accumulates past gradients) to smooth out oscillations

$$v_i = \beta v_i - \eta \nabla L(\theta_t)$$

$$\theta_{t+1} = \theta_t - v_i$$

- Like a ball rolling downhill, momentum helps accelerate in the right direction

c] Nesterov Accelerated Gradient (NAG) :- It is an improvement over momentum that looks ahead before updating.

$$v_i = \beta v_i - \eta \nabla L(\theta_t - \beta v_{t-1})$$

$$\theta_{t+1} = \theta_t - v_t$$

- It provides more accurate update by considering the future positions.

d) AdaGrad :- It adjusts learning rate individually for each parameter, scaling by past squared gradients.

$$g_t = \nabla(\theta_t)$$

$$\hat{G}_t = G_t - 1 + g_t^2$$

$$\theta_{t+1} = \theta_t - \eta \frac{g_t}{\sqrt{G_t + \epsilon}}$$

- It works well for sparse data
- The learning rates keeps shrinking - may stop learning too early.

e) RMS Propagation (Root Mean Square) - It's a refinement of AdaGrad that uses an exponentially decaying average of past square gradients instead of summing them forever.

- It prevents learning rate from shrinking too much

$$\text{Formula} := E[g^2]_t = \beta E[g^2]_{t-1} + (1-\beta)g_t^2$$

$$\theta_{t+1} = \theta_t = \frac{\eta}{\sqrt{E[g^2]_{t+1}}}$$

f) Adam (Adaptive Moment Estimation) - It combines Momentum + RMS (Root Mean Square) propagation

- It tracks both:-

1^{st} moment (means of gradient)
 2^{nd} moment (Variance of gradient)

Rule:-

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t$$

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) g^2 t$$

g) RMS Prop with Momentum :- It combines RMS Prop's adaptive learning rate with Momentum acceleration

Rule:-

$$v_t = \beta v_{t-1} + \frac{n}{\sqrt{E[g^2]t + \epsilon}} g_t$$

$$\theta_{t+1} = \theta_t - v_t$$

Unit 3

1. What are RNN. Explain the various representation.



- A RNN (Recurrent Neural Network) is a class of neural network designed to process sequential data by maintaining an internal hidden state or memory.
- Unlike traditional FFN network which assume inputs and outputs are independent, an RNN uses information from previous steps in the sequence to influence the processing of the current steps.
- Mathematical Representation

$$h_t = f(w_h h_{t-1} + w_x x_t + b)$$

$$y_t = g(w_y h_t + c)$$

x_t = input at time step t

h_t = hidden state at time t

y_t = output at time t

w = weight matrices

f, g = activation functions

- Various Representation of RNNs

a] One to One : A single i/p produces a single o/p. This is the standard configuration of RNN

Use Case : Image Generation Classification

$$Tx = 1, Ty = 1$$

b) One to Many :- A single i/p is processed to generate a sequence of o/p. Th. hidden state is updated and passed along to generate subsequent o/p.

- Use Case :- Music Generation

- $T_x = 1, T_y > 1$

c) Many to One :- A sequence of i/p processed and a single o/p is generated at very end, which summarized the entire sequence.

- Use Case :- Sentiment Analysis

- $T_x > 1, T_y = 1$

d) Many to Many (with different lengths) :-

- I/p sequence length \neq o/p sequence length

- Use Case :- Video Captioning

- $T_x = T_y$ (Synchronised or)

- $T_x \neq T_y$ (Asynchronous)

2
→

Describe LSTM model

- A Long Short term Memory is a Model is a type of RNN designed to overcome the limitations of traditional RNN, particularly the problem of learning long-term dependencies in sequential data, which is often hindered by the vanishing gradient problem.
- Structure of LSTM:-

- a] Cell State (c_t) - carries long-term memory
- b] Hidden State (h_t) - short term information for output
- c] Gates :- Control the flow of information.
 - Forget Gate
 - I/P gate
 - o/p gate

Working of LSTM

- a] Forget Gate :- Decides what information to discard from Cell State.

$$f_t = \sigma(w_f \cdot [h_{t-1}, x_t] + b_f)$$

- b] Input Gate :- Decides what new information to store

$$i_t = \sigma(w_i \cdot [h_{t-1}, x_t] + b_i)$$

$$C_t = \tanh(w_c \cdot [h_{t-1}, x_t] + b_c)$$

c) update cell state - Combines old memory with new information

$$C_t = f_t + C_{t-1} * i_t$$

d) o/p Gate :- Decides what part of the cell state to o/p

$$o_t = \sigma(w_o \cdot [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh(h_t)$$

- The cell state acts like a conveyor belt that carries long-term memory with minor modifications.
- Gates mechanism allows sensitive remembering or forgetting.
- It's application are speech Recognition, Time Series Forecasting, video processing.

Dev
4/10/25

$$(i_t * f_t)(1 - o_t) = ?$$

$$(i_t * f_t)(1 - o_t) = 1$$

$$(i_t * f_t)(1 - o_t) \cdot o_t = 1$$