

Unit - 2

Challenges in Neural Network Optimization:-

→ Optimization tries to minimize the error or loss and tries to reduce them by adjusting weights & biases.

→ Challenges:-

a) Local Minima:-

- The grandfather of all optimization problems, local minima is a permanent challenge in the optimization of any deep learning alg.
- The local minima problem arises when the gradient encounters many local minimums that are ~~different~~ like valleys, where the error is small, but it's not the lowest.

b) Saddle Points:-

Saddle Points are that regions where the gradient (rate of change) is zero, making it diff to move in right direction.

c) Vanishing Gradients:-

They occur during backpropagation where the gradient (changes in weights) become very small, preventing the network from learning effectively.

d) Overfitting:-

It happens when the model learns the training data too well, including noise & irrelevant details, making it perform poorly on new, unseen data.

c) Under fitting:-

It occurs when model is too simple or not trained enough and it fails to capture important patterns in data.

④ Gradient-Descent Alg:-

→ It is an optimization alg used to min loss function by updating model parameters (weights ~~& bias~~) in the direction of negative gradient.

→ It is a first order optimization alg.

→ working:-

1) Start with Initial Parameters:-

• Initially, the parameters (weights) of model are set to random values.

2) Compute the Gradient:-

The gradient (^{rate of change} ~~change in weight~~) of loss function is calculate w.r.to each parameter (weight).

• The gradient represents rate of change of loss function in direction each parameter.

3) Update the Parameters:-

• The parameters are updated in opposite direction of gradient to min loss function. Using formula

$$\theta = \theta - \alpha \cdot \nabla_{\theta} J(\theta)$$

θ → model parameters (weights)

α → learning rate

$\nabla_{\theta} J(\theta)$ → gradient of loss function $J(\theta)$

4) Repeat process for several epochs until parameters converge to values that min loss function.

→ Disadv

- Can be slow for large datasets
- Can get stuck in local minima

② Stochastic Gradient Descent:-

→ SGD is a variant of basic GD alg. Unlike traditional one, which uses entire dataset to compute the gradient & update model parameters, SGD updates the parameters using only one data point at a time.

→ This makes the alg. faster but more noisy.

→ working:-

1) Initialization:-

Start with random values for model parameters.

2) Select a single data point:-

Instead of using whole dataset, SGD selects one data point randomly from the dataset.

3) Compute the Gradient:-

The gradient of loss function is computed with r. to the model parameters, using selected data point.

4) Update parameters:-

The parameters are updated based on the gradient, using formula:

$$\theta = \theta - \alpha \cdot \nabla_{\theta} J(\theta)$$

$\nabla_{\theta} J(\theta)$ → gradient of loss function $J(\theta)$ w.r.t to parameters, based on selected data point.

5) Repeat process for each data point in the dataset. After one full pass through the data, this is called one epoch. The parameters are updated in each epoch.

6) Iterate the process for multiple epochs until model parameters converge to values that min loss function.

→ Adv:

- Faster Updates
- Escaping local minima

→ Disadv:

- Noise can make it less stable
- Can take longer for convergence.

④ Mini-Batch Gradient Descent:

→ It is an optimization alg that strikes a balance b/w Batch GD & SGD.

→ Instead of using the entire dataset or just a datapoint, MBGD updates the model parameters using a small subset of data, called mini-batch.

→ Working:

1) Divide the dataset into Mini-batches:

- The entire dataset is divided into small batches.
- Each mini-batch contains 32 - 256 data points.
- Size of mini-batch can act as hyperparameter. If chosen very less it can become noisy like SGD, if chosen more it can be slow like GD.

2) Initialize Parameters.

3) Compute the Gradient:-

For each mini-batch, the gradient of loss function is computed based on average of gradients of data points in that mini-batch.

4) Update Parameters:-

$$\theta = \theta - \alpha \cdot \nabla_{\theta} J(\theta)$$

$\nabla_{\theta} J(\theta) \rightarrow$ gradient of $J(\theta)$ based on mini-batch.

5) Repeat for each mini-batch till an epoch.

6) Iterate for multiple epochs until convergence.

Adv:-

- Faster Updates & req. less computation
- Memory efficient
- More stable

Disadv:-

- Choosing mini-batch size

⑩ AdaGrad (Adaptive Gradient Alg):-

\rightarrow It is an optimization alg. that adapts the learning rate for each parameter individually based on its historical gradient.

\rightarrow It aims to ^{im}prove learning rate by adjusting it for each parameter, allowing more frequent updates for parameters with smaller gradients and fewer updates ~~with~~ ^{for} parameters with large gradients.

→ Working:-

1) Initialization of Parameters ~~(Initialization)~~

2) Compute Gradients:-

For each parameter, AdaGrad computes gradient of the loss function wrt to that parameter (weights)

3) Accumulate Squared Gradients:-

• AdaGrad keeps track of the sum of squared gradients for each parameter. This helps in adjusting learning rate for each parameter.

$$G_{t,i} = G_{t-1,i} + g_{t,i}^2$$

where,

$G_{t,i}$ → accumulated squared gradient for parameter i at t

$g_{t,i}$ → gradient of loss wrt i at t

4) Update Parameters:-

• The parameters are updated using

$$\Theta_i = \Theta_i - \frac{\alpha}{\sqrt{G_{t,i}} + \epsilon} \cdot g_{t,i}$$

Θ_i → parameter

α → global learning rate

ϵ → small value added to prevent div by '0'.

5) Repeat for each data point or mini-batch (depending on variant of GD used) as the model parameters are updated after each iteration.

→ Disadvantages

- Learning rate will diminish over time
- Not ideal for non-convex functions.

- RMS Prop (Root Mean Square Propagation):
① Adam (Adaptive Moment Estimation):-

→ It is an ~~advanced~~ optimization alg. that modifies AdaGrad to address the issue of rapidly diminishing learning rates.

→ It uses moving average of squared gradients to normalise the gradient, allowing for more stable & effective updates, especially in non-stationary setting (ex:- RNNs).

→ Working:-

1) Initialization

2) Compute Gradients

3) Calculate moving Average of Squared Gradients:-

Instead of accumulating all squared gradients, RMS uses decaying average of squared alg.

$$v_t = \beta v_{t-1} + (1-\beta) g_t^2$$

v_t → moving avg of squared gradients at t .

g_t → gradient of loss function

β → decay factor (typically close to 1)

which controls how much weight is given to past gradients.

④ Update Parameters:-

$$\theta = \theta - \frac{\alpha}{\sqrt{v_t} + \epsilon} \cdot g_t$$

5) Repeat for each data point or mini-batch and the model parameters are updated iteratively.

→ Adv:-

- Avoids diminishing learning rates.
- Better performance in non-convex problems
- Faster convergence

→ Disadv:-

- Sensitive to hyperparameters (α, β_1, β_2)
- Choice of decay factor

⑤ Adam (Adaptive Moment Estimation):-

→ It is an ~~optim~~ advanced optimization alg. that combines the ideas of Momentum & RMS prop.

→ It computes adaptive learning rates ~~for~~ ^{for} each parameter, using both first moment (mean) & second moment (variance) of gradients.

→ Adam is one of the most widely used optimizers in DL due to its efficiency.

→ working:-

1) Initialization:-

- Adam starts by initializing the model parameters randomly

- Additionally, two moment estimates are initialized to zero:

$m_0 = 0$ (first moment), gradient _{mean}

$v_0 = 0$ (second moment ~~estimate~~ estimate)

2) Compute Gradients:-

3) Update Moment Estimates:-

Adam computes 2 moment estimates.

• First moment estimate: The moving average of gradients (m_t)

• Second moment estimate: Moving avg of squared gradients (v_t)

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

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

β_1 & $\beta_2 \rightarrow$ hyperparameters that control decay rates for moment estimates

4) Bias Correction:

• Since m_t & v_t are initialized to 0, their estimates are biased toward '0'.

To correct this bias,

$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t}$$

$$\hat{v}_t = \frac{v_t}{1 - \beta_2^t}$$

\hat{m}_t & $\hat{v}_t \rightarrow$ bias corrected estimates

5) Update parameters:-

$$\theta = \theta_{t-1} - \frac{\alpha}{\sqrt{\hat{v}_t} + \epsilon} \cdot \hat{m}_t$$

6) Repeat for data point & minibatch and the model parameters are updated iteratively.

Disadvantages

Memory Usage, Hyperparameter Sensitivity, Risk of Overfitting.

④ Large Scale Deep Learning:-

- It involves using DNNs to process & learn from vast amounts of data, which requires ~~has~~ powerful hardware, like GPUs or specialized H/Ws to handle large scale computations.
- NNs are trained on massive datasets that could include millions or billions of examples.
- As the n/w processes more data, it improves its ability to make accurate predictions.
- The larger the dataset, the better model can generalize & perform complex task.
- Big Data - Large scale datasets that contain a lot of info like images, videos, text & sensor data.
- Using GPUs instead of traditional CPUs speeds up training.
- On large scale DL, the task is often distributed across multiple machines or even data centers to speed up training process.

→ App:-

Self Driving Cars, Recommendation Systems, Healthcare

→ Challenges:-

Data Quality
Computational Power (large)
Overfitting

⑫ Computer Vision :-

- It is a field in AI that enables machines to interpret & understand visual info from the world, like images or videos.
- DL, particularly CNNs, play crucial role in computer vision.
- CNNs are used to process image data. They work by applying filters to detect edges, textures, other features in the image at diff levels.
- CNNs apply multiple layers like convolutional layers to extract features, pooling to extract features and in reducing computation ~~cost~~ by reducing spatial size of image.
- Image Classification
- Object detection (classifying multiple objects in images)
- Semantic segmentation aims to label each pixel in an image with its corresponding object class.
- GANs are used to generate realistic images that resemble specific dataset.

→ Apps :-

Healthcare
Autonomous Driving
Surveillance & Security

④ NLP

- NLP in deep learning refers to a process, understand & generate human language.
- NLP enables computers to interact with & understand human lang. in a way that is meaningful & contextually relevant.
- Deep & particular NNs have revolutionized NLP tasks.
- Text Classification (spam detection, topic categorization)
- Machine Translation (sentence-to-sentence RNNs or transformers)
- Named Entity Recognition (Names of people, org. locations)
- Post Summarization
- Sentiment Analysis
- Question Answering
- Speech Recognition & Synthesis
 - ↓
 - (converting spoken lang to text) (generating human-like speech from text)
- Language Generation (Chatbots & Lang. Models)
- RNNs, LSTMs, Transformers (GPT, BERT)
 - ↓
 - Generative Pre Trained Transformer