

Unit-5

challenges in Neural Network

Optimization

Neural Networks face challenges in terms of data requirements, overfitting, explainability, security, and bias etc. They require large amounts of labeled data for effective training, can suffer from overfitting where they perform well on training data but fail to generalize etc in their predictions.

① Training Data Requirements

require a large amount of labeled data to learn effectively, and also annotation such datasets take lot of time.

costly even impractical in some cases.

② Overfitting

occurs when a neural network learns to perform exceptionally well on the training data but fails to generalize to new, unseen data.

③ computational complexity

Neural networks can be computationally demanding, particularly as the model size and complexity increase. It needs high performance GPU's or specialized hardware like TPU (Tensor Processing units).

④ hyperparameter tuning

Networks have several hyperparameters such as (learning rate, number of layers, activation function). Selecting them is most important.

⑤ Data Privacy

→ Neural Networks often require large amount of data to train effectively, this raises concerns about data privacy and security, especially when dealing with sensitive information.

Optimizer	complexity	Computation expense
Gradient descent	low	high
Stochastic gradient descent	low	low
Stochastic gradient descent + Momentum	medium	low

Adasrad	High	medium
Adam	High	High

⑥ Learnings Rate

- is too high, the model overshoots the best solution.
- is too low, training becomes very slow.

⑦ Batch size problem

- small batches can be noisy, while large batches need more memory.
- so we should try each to find which works best for us.

Applications of large-scale

Deep learning

is to use of very large neural networks trained on massive datasets, often uses powerful computational resources like GPU or TPU.

① Computer vision

Used to analyze and understand visual data like images and videos.

Applications

- ① Image Recognition
- ② Object Detection
- ③ medical imaging
- ④ Video Analysis

② Natural language processing

Deep learning models process and understand human language.

Application

- ① Language Translation.
- ② Chatbots and virtual assistants
- ③ Sentiment Analysis
- ④ Text Summarization.

③ Speech Recognition

Converts spoken language into text or commands.

Application

- ① Voice Assistants.
- ② Transcription Services
- ③ Voice-controlled systems.

④ Recommendation System

↳ Predict user preferences and recommend products or content.

Application

- ① E-commerce
- ② Streaming services
- ③ social media

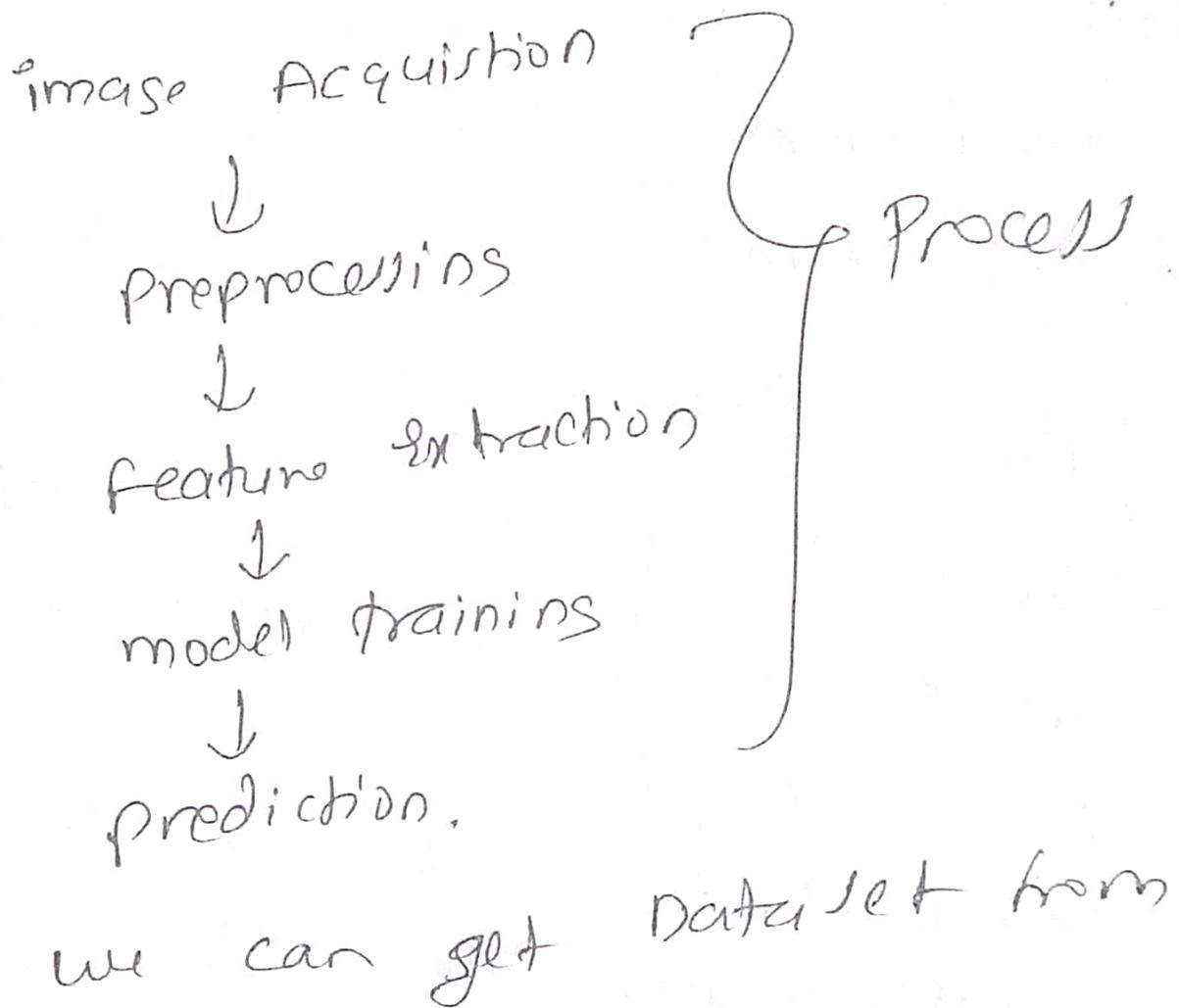
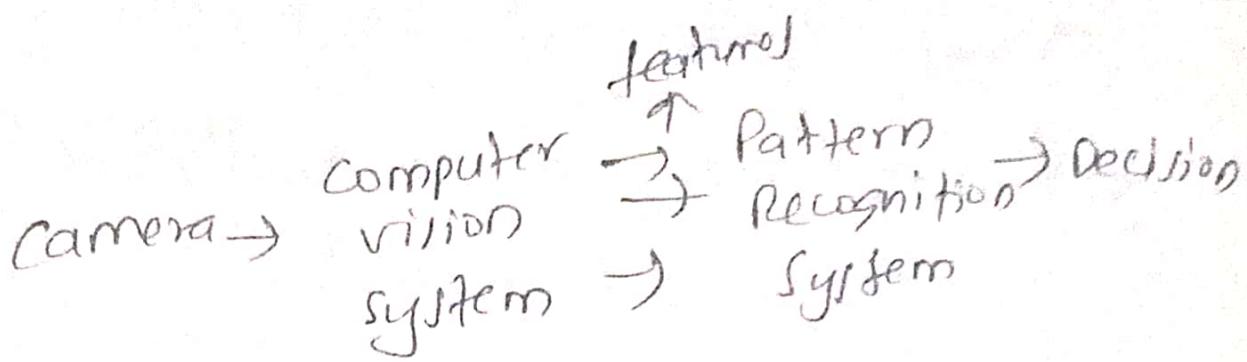
⑤ Autonomous vehicles

↳ Enables self-driving cars to perceive and navigate their environment.

Application

- ① object detection
- ② path planning
- ③ control systems
- ④ Healthcare
- ⑤ Gamins
- ⑥ Finance
- ⑦ Robotics
- ⑧ Generative models

Computer Vision



→ MNIST

→ ImageNET

→ COCO (Common objects in context)

Computer vision model

- ① YOLO (You only look once)
- ② SSD (single shot multibox detector)
- ③ Transformer in vision.
- ④ U-Net.
- ⑤ GAN.

Speech Recognition

→ It is the process of converting spoken language into text using AI and machine learning techniques. It is widely used in applications like voice assistants (Google Assistant, Siri, Alexa), hands-free control.

How it works

- ① Audio input
- ② preprocessing
- ③ Acoustic modeling
 - ↳ maps audio features
- ④ Language modeling
 - ↳ predicts the most likely word sequences.
- ⑤ Decoding
 - ↳ converts phonemes into readable text.

popular models

- ① DeepSpeech.
- ② Google Speech-to-Text API
- ③ kaldi (for research & enterprise use)
- ④ vosk (lightweight & offline support)

Automatic
Speech
Recognition

Natural
language
processing

→ text to
speech

Acoustic
modeling

↓

Voice → speech
enhancement → feature
extraction → phonetic
unit
recognition

↓
Hello
(rest)

Natural Language processing

it is a branch of artificial
intelligence that enables computers
to understand, interpret, generate,
and respond to human language

key components

- ① text processing.
- ② syntactic analysis. (syntax)
- ③ semantic analysis (meaning)

④ Speech & conversational NLP.

Popular NLP libraries

Python libraries

- NLTK
- spacy
- Transformers
- Gensim

Deep Learning NLP models

- BERT (Google)
- GPT (openAI)
- T5 & BART
- XLNet

Applications

- ① chatbot & virtual assistants
- ② machine translation
- ③ text summarization
- ④ search engines
- ⑤ spam detection
- ⑥ customer support automation
- ⑦ healthcare

input in Natural language



Speech Recognition



Natural Language
Understanding



Natural Language
Generation



Output in
Natural language

Optimization for Train Deep Model

Basic Algorithms

① Gradient Descent:-

The foundation of most optimization algorithms. It minimizes the loss function by iteratively updating the model parameters in the direction of the negative gradient.

Update rule

$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta} L(\theta_t)$$

Gradient of the loss
function w.r.t θ

\downarrow

learning rate

↳ model para
-meter at
step t

variants

① Batch Gradient Descent

↳ uses the entire dataset to compute the gradient.

② Stochastic Gradient Descent (SGD)

↳ uses a single data point to compute the gradient.

③ mini-Batch Gradient Descent

↳ uses a small subset of the data to compute the gradient.

2. Momentum

→ Accelerates SGD by adding a fraction of the previous update vector to the current update. This helps in smoothing the optimization path and escaping local minima.

Update rule

$$v_{t+1} = \gamma v_t + \eta \nabla_{\theta} L(\theta_t)$$

momentum term. $\theta_{t+1} = \theta_t - v_{t+1}$

(velocity vector)

3. Nesterov Accelerated Gradient

→ A variant of momentum that "looks ahead" by computing the gradient at the estimated future position.

update Rule

$$v_{t+1} = \gamma v_t + \eta \nabla_{\theta} L(\theta_t - \gamma v_t)$$

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

4. Adaptive Learning Rate Method

These methods adjust the learning rate for each parameter based on the history of gradients.

a. AdaGrad

↳ Adapt the learning rate for each parameter based on the sum of squared gradients.

Update rule

$$G_t = G_{t-1} + (\nabla_{\theta} L(\theta_t))^2$$

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{G_t + \epsilon}} \nabla_{\theta} L(\theta_t)$$

↳ Accumulated
 Squared
 gradients

small
 constant
 for stability

b. Adam (Adaptive moment estimation)

Combines the benefits of momentum and RMSprop, maintains both a moving average of gradients and a moving average of squared

gradient

update rule

$\beta \rightarrow$ Exponential decay rate.

$$\hat{m}_t = \frac{m_t}{1 - \beta^t}, \quad \hat{v}_t = \frac{v_t}{1 - \beta^t}$$

first moment

$$\theta_{2+1} = \theta_t - \frac{m}{\sqrt{\hat{v}_t + \epsilon}} \hat{m}_t$$

second moment

GD

uses the entire dataset to compute the gradient

updates parameter once per epoch

slower

high memory usage

best for small datasets

SGD

uses one data point to compute the gradient.

updates parameter after every data point.

faster.

low memory usage.

best for large datasets.

Parameter initialization strategies

Parameter initialization is a critical step in training deep learning models. Poor initialization can lead to issue like vanishing gradients, exploding gradients or slow convergence.

① Zero initialization

- initializes all weights to zero.
- rarely used in practice.

② Random initialization

- initializes weights with small random values.
- it allows neurons to learn different features.
- if the weights are too large or too small, it can lead to exploding or vanishing gradients.

3. Glorot initialization

→ Scales the random initialization based on the number of input and output neurons to maintain consistent variance across layers.

4. He initialization

A variant of Glorot initialization designed for ReLU activation functions. It scales the weights by the number of input neurons only.

5. lecun initialization

Similar to Glorot initialization but scaled specifically for tanh activation functions:

6. orthogonal initialization

initialize weights as an orthogonal matrix, ensuring that the input and output have the same norm.
commonly used in RNN and LSTM.

7. pre trained initialization (transfer learning)

use weights from a pretrained model as the starting point.

8. Biases initialization

Biases are typically initialized to zero or a small constant value (e.g 0.01).

Algorithms with Adaptive learning rates

gt is used to dynamically adjust the learning rate during the training to improve convergence, prevent divergence and enhance overall model performance.

-ache

① Adagrad (Adaptive Gradient Algorithm)

→ Adjust learning rate individually for each parameter based on past gradients.

Update Rule

$$g_t = g_{t-1} + \nabla \theta_t^2$$

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

2. RMSprop (Root Mean Square propagation)

- improves Adagrad by introducing a moving average of squared gradients.
- well-suited for non-stationary problems like deep networks.

Update rule

$$g_t = r g_{t-1} + (1-r) \nabla \theta_t^2$$

$$\theta_{t+1} = \theta_t - \frac{m}{\sqrt{g_t} + \epsilon} \nabla \theta_t$$

typical choice $r = 0.9$.

3. Adadelta

- Extension of Adagrad, addressing its learning rate decay problem.
- no need for manual learning rate tuning.

4. Adam (Adaptive moment Estimation)

- combines the benefits of Adagrad and RMS prop.
- maintains both first-moment (mean) and second-moment (uncentered variance) estimates of gradients.

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

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

$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t}, \quad \hat{v}_t = \frac{v_t}{1 - \beta_2^t}$$

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{v_t} + \epsilon} \hat{m}_t$$

$\beta_1 = 0.9, \beta_2 = 0.999$ default values

5. AdamW

- A variation of Adam with decoupled weight decay.
- Helps improve generalization in large models.

Adagrad

Best for
sparse data
(NLP, recommender
systems)

Weakness

Learning rate
keeps decreasing

Needs tuning
of decay
factor.

RMSprop

RNNs, non-
stationary
problems.

Can still be
sensitive
to hyperparameters

Adadelta

adaptive
learning rate
without manual
tunings.

Can generalize
poorly on
some tasks.

Adam

most deep learning
tasks (CNNs, RNNs,
transform).

Requires
careful
weight
decay
tunings.

AdamW

large-scale
models,
regularization

Second-order Methods

4.1 leverage curvature information (Hessian matrix) of the loss function to improve optimization.

unlike first-order methods, which rely only on gradients, second-order methods use the Hessian to guide updates.

offer

- ① Better direction
- ② Better Step-size

~~Newton~~ Newton's method

it is mostly used second-order method. it uses a second-order Taylor expansion to approximate the loss function $J(\theta)$ near a point θ_0 .

Newton update Rule

Solving for the critical point of the quadratic approximation gives the ~~newton~~ Newton update.

$$\theta^* = \theta_0 - H^{-1} \nabla_{\theta} J(\theta_0)$$

for a quadratic function with positive definite H , ~~newton~~ Newton's method jumps directly to the minimum.

For non-quadratic functions, the update is iterated.

Training Algorithm

- ① compute the gradient.
- ② compute the Hessian.
- ③ update parameters.
- ④ conjugate gradients
is an iterative method to avoid computing H^{-1} explicitly. It addresses the inefficiency of steepest descent, which often follows a zig-zag path.

Alg

1. initialize α_0 and d_0
2. For each iteration i :
 - compute step size ϵ_i via line search.
 - update parameters.
 - compute new gradients.
 - update conjugate direction.

4. BFGS and L-BFGS

BFGS is a quasi-Newton method that approximates the inverse Hessian H^{-1} iteratively without explicit computation.

- Limited-memory BFGS reduces memory waste by storing only a few vectors instead of the full matrix.
- suitable for large scale optimization problems.

use

Better direction,
Better step-size

Challenges

computational cost

scalability.

optimization strategies and
Meta Algorithms

- ① Gradient-Based optimization.
- ② Adaptive Learning rate methods.
- ③ Second-order methods.

→ Newton method

→ Quasi-Newton method
(L-BFGS)

- ④ Learning rate scheduling.

meta Algorithms

- ① Ensemble method

→ Bagging → Trains multiple models

independently on different subsets
of the data and averages their
predictions.

Boosting → sequentially trains models with each new model focusing on the errors of the previous one.

Stacking → combines predictions from multiple models using a meta-model.

Transfer Learning

Reuses pre-trained models and fine tunes them for a new task. Reduces training time and data requirements.

Neural Architecture Search (NAS)

Automates the design of neural network architectures using reinforcement learning, evolutionary algorithms, or gradient based methods.

Curriculum learning

- ↳ Trains the model on easier examples first and gradually introduces more complex examples.

self-supervised learning

- ↳ uses unlabeled data to pre-train models by defining pretext tasks

meta-learning (learning to learn)

- ↳ trains models to quickly adapt to new tasks with minimal data.