

* Artificial Neural network -

- inspired by biological neural networks.

- mimic network of neurons like human beings so computers can be made decision in human like manner.

Adv → • parallel processing capability

- sorting data on entire network

- capability to work with incomplete knowledge.

- Having memory distribution.

- having fault tolerance.

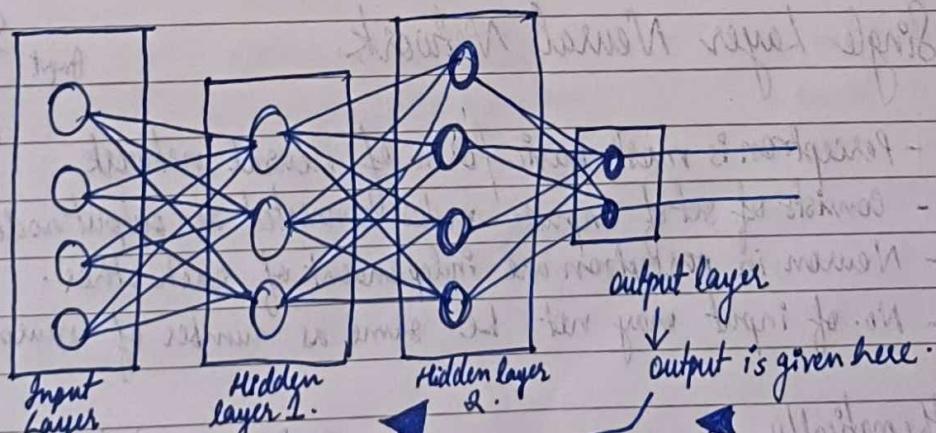
Disadv → • assurance of proper network structure

- unrecognizable behaviour of network.

- Hardware dependence

- difficult of showing issues to network

- duration of network is unknown.



accept input from
programmer in variety of diff
formats.

computation
necessary
to uncover
pattern buried
features.

communicate output after input
undergoes act of attention

in hidden layer.

* Types of ANN → 1) Feedforward ANN

→ atleast one layer of neurons

→ network intensity → output of input layer connected neurons

2) Feedback ANN

→ output loops back into network to achieve best internally → benefit → learns to assess & identify input patterns involved results

→ addressing optimization results

→ feed info into themselves

→ internal system error repairs

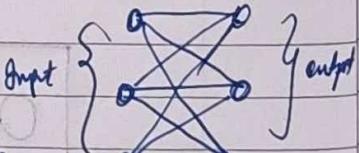
Applic'n of ANN

1) Social Media → Insta → people you may know

2) Sales & Marketing → e-commerce → suggest things acc. to hist. data

3) Health care → facial analysis

4) Personal Assistants → Siri, Alexa, etc



* Perceptions -

1) Single Layer Neural Network

- Perception is most basic form of neural network.

- consists of set of input nodes connected to output nodes using weighted connections

- Neuron in perception are independent of each other.

- No. of input may not be same as number of neurons (output)

* Mathematically -

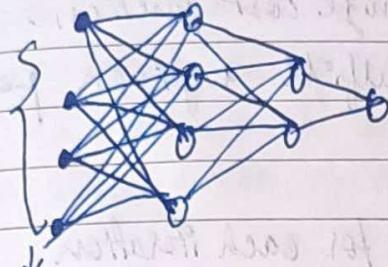
- 1) Perception → m inputs 'n' neurons
- 2) weights are labelled w_{ij} where $1 \leq i \leq m$
- 3) each neuron has its own weight denoted as w_{ij} , where i is input node no. & j is output node

* Limitations -

- uses just binary activation function
- only applicable to linear network
- provides an optimal sol'n due to supervised learning more training time to tackle linear inseparable problems

Multilayer Perceptron

If neural network requires complex decisions making we can create multiple layers of perceptron network.



layer 1.
(process inputs)

Its output is input for layer 2

along with output, weights are
also given

- atleast one layer b/w input & output layer.
- hidden layer between I/O layers.
- offers b/w soln to every categorization issue.

Input \rightarrow non-linear

multilayer network.

with linear discriminants use.

* Asked for multilayer networks -

- to solve complex problems
- segregate input output data
- intricate info is too complex for single layer to handle.
- multilayer go beyond limitation of single layers.

Linear model

- linear relⁿ b/w features & output.
- straight lines of hyperplane
- simpler, interpretable, linear relⁿ.
- can't capture non-linear relⁿ
- e.g. \rightarrow predicting house prices
- loan approval based on income.

Non-linearities

- non linear relⁿ:
- can be curves, spirals, or complex shapes
- more flexible, complex
- more complex & less interpretable.
- e.g. \rightarrow image recognition
- NLP.

* Gradient Descent algo (finds best fit line for giving training dataset).

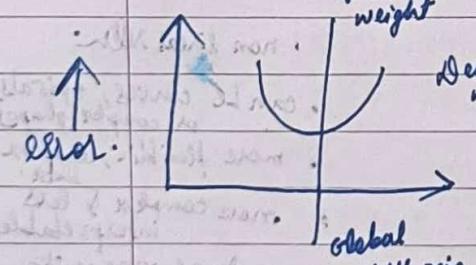
- optimization algo \rightarrow minimize cost function.
- locate least possible value \rightarrow fulfill a given fee function.

* Types -

- Batch \rightarrow processes all training ex. for each iteration
 - stochastic \rightarrow Process 1 training example per iteration. Faster than Batch
 - Mini Batch - More examples out of total training examples are processed per iteration.
- | | | |
|------------------|-------------------|--------------|
| cost function. | iterative updates | efficiency |
| steepest Descent | convergence | variants |
| | | local Minima |

* Back propagation -

- algo for supervised learning for ANN
- keeps adjusting weights of connected neurons with an α to reduce duration of output signal with target output.



reach global less mini using backpropagation

consists of multiple iterations known as epochs.

forward pass

Start

[Input Training Data]

[Provide weight.]

[Calculate target output.]

back propagation

No \rightarrow [Adjust weight]

Stop

Features → gradient descent method → case of simple perceptron network with diff. unit.

- weights are calculated in learning period of network:
- feed forward of input training pattern.
- calculation & back propagation of error.
- updation of weight

Adv → - simple, fast & easy

• only no. of inputs are tuned, not any other parameter.

• flexible & efficient.

• no need for user to learn any special function.

Disadv → • sensitive to noisy data & irregularities.

• performance → dependent → data.

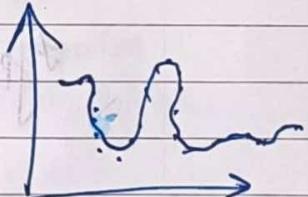
• too much time training.

matrix-based approach preferred over mini batch.



Overfitting -

- statistical model fits against its matching data.
- Algo/model can't perform well on unseen data
- low bias & high variance.
- match input data to target data non-existent.
- model → complex enough → match all datapoint & performs well.
- reasons → • noisy data • training data is too small • large no. of features.

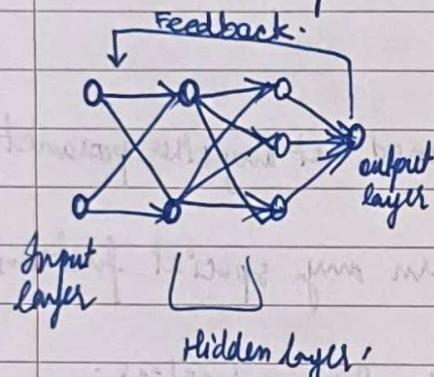


To avoid Overfitting -

- Cross Validation.
- Train with more data.
- Removing features
- Early stopping the training
- Regularization
- Ensembling

Recurrent Neural Network (RNN)

- Feed-Forward Neural Network
- Perceptrons are arranged to layers, hidden layers are not connected
- All nodes fully connected some layers are not connected ^{with outside world}
- need to access previous info in current iteration. ~~outside world~~
- commonly used in language translation, text to speech, etc.



- Adv →
 - remember each info through time
 - useful in time series prediction \rightarrow feature to number
 - extend effective pixel neighborhood ^{input LSTM}

Disadv → gradient vanishing & exploding problems -

- training is an difficult task -
- can't process very sequence if -
- using tanh or relu as activation function -

Types of RNN → one to one

one to many

many to one

many to many

- different views of it

- different ways

- this is nice view

- surface of pyramid

U-4 (DS Honors).

DOMS	Page No.
Date	/ /

* Hadoop Ecosystem -

- collection of open source s/w projects that facilitate storing, processing & managing big data.
- essentially framework provides tools for various tasks in big data analytics.

Core components -

- HDFS
- YARN
- MapReduce

Benefits, advantages -

- Scalability
- cost-effectiveness
- Flexibility
- Open source
- Fault tolerance.

Disadv -

- complexity
- Learning curve.
- Performance overhead
- Security concerns.
- Vendor lock-in

Appln → Log Analysis

- social Media Analytics
- Fraud Detection
- scientific computing
- recommendation systems

HDFS (Hadoop Distributed File System).

distributed file system designed to run on clustered & low cost hardware.

Characteristics & features of HDFS -

- Fault tolerance
- streaming data access
- Large Data Sets
- write once Read many models
- Highly portable

Architecture →

Job Tracker
Name Node.

Task Tracker
Data Node 1
Task Tracker
Data Node 2.
Task Tracker
Data Node 3

HDFS → master slave architecture.

- Name Node → single Name Node → master slave.
- Data Nodes → multiple Data Nodes → storage attached to it

Adv → Scalability
→ Cost-efficient.
→ Reliability
→ Fault Tolerance
→ Integration.

Disadv → Complexity
→ Learning Curve
→ Overhead
→ Security Concerns
→ Not ideal for small data

Map Reduce -

- programming model & software framework,
- developed by Google

Map() Reduce()

Characteristics

- Very Large Scale data:
- write once & read many data.
- Map & reduce main operations → simple code.
- all maps → reduced operation starts.
- map & reduce → physical processor
- No. of map tasks & reduced tasks.

Adv →

- Scalability
- cost-efficient
- Fault Tolerance
- Ease of programming
- security & authentication

Disadv →

- Limited functionality
- Verbosity
- Shuffle Bottleneck.

* Python with Hadoop Streaming

- allows to write MapReduce Programs
- leverage power of distributed processing in Hadoop.

Benefits/Adv

- Readability & Maintainability
- Rich Libraries
- Familiarity for Python Developers.

Disadv -

- Limited functionality
- Performance overhead
- Debugging challenges.

* Spark -

- open source framework for large scale data processing.
- addresses some of limitations of traditional MapReduce & significant improvements in terms of speed, performance & flexibility.

Core Concepts -

- In Memory Processing
- Unified Platform

Spark Components

- Spark Core
- spark SQL
- spark Streaming
- MLlib
- GraphX

Fav →

- speed
- flexibility
- Ease of Use
- Rich Ecosystem

Disadv →

- Hardware Requirement
- Learning curve
- complexity

* Pyspark -

- Is Python API for Apache Spark .
- spark's functionalities accessible & usable from Python Programming Language

Benefits →

- Readability & Maintainability

- Familiar for Python development.
- Rich Ecosystem of Python library
- Apache Spark functionalities.

Disadv →

- Performance overhead

- Hardware requirement
- Learning Curve.

U-5 (DST Honors).

* Data Warehousing -

- process of constructing & using a data warehouse.
- integrating data from multiple heterogeneous sources → analytical reporting.

Data warehouse → store historical info from multiple sources to allow you to analyze & report on related data.

- Goals →
- reporting / analysis → maintain organization's historical info
 - foundation → decision making.

organizations → info → datawarehouse.

- increasing customer focus
- repositioning products & managing products.
- analysing operations & looking for profit
- managing customer relationships, making environmental corrections & managing rest of corporate assets.

Characteristics of DW

- Subject Oriented
- Integrated
- Non-Volatile
- Time Variant
- access & highspeed query.
- end user → time sensitive.
- large amt → historical data
- queries → large amt of data.

* Data Mining:-

process of extracting knowledge & insights from large datasets.

Techniques → • statistic Analysis • Machine Learning Algo
• Data Visualization Tools

Data Mining process -

- 1) Data collection
- 2) Data Preparation
- 3) Data Selection
- 4) Model Selection
- 5) Model Building
- 6) Evaluation
- 7) Deployment

Benefits

- Uncover hidden meaning
- Improved Decision making
- Enhanced customer understanding
- Fraud Detection

Challenges

- Data Quality
- Complexity of Technologies
- Privacy Concerns

Characteristics -

- Focus on Pattern & Relationship
- Predictive Capabilities
- Statistical Foundation
- Machine Learning Integration
- Data Driven Insight
- Iterative Process
- Focus on Specific Q's
- Subject to Data Quality

A Data Analysis Using Hive -

Hive - powerful tool for data warehousing & data analysis on large dataset stored in HDFS.

Key Functionalities

- SQL like interface
- Data Warehousing Concepts
- Scalability
- Integration with Hadoop Ecosystems

Process

- Data Loading
- Schema Definition
- Data Querying
- Data Analysis
- Data Exploration

~~Advantages of Hive~~

- Ease of Use
- Scalability
- Flexibility
- cost effectiveness
- Integration with Hadoop Systems

Disadv-

- Limited functionality
- performance overhead
- Learning Curve

Characteristics of Hive for Data Analysis

- SQL like Interface
- Data Warehousing Capabilities
- Batch Processing
- Schema on Read
- Distributed Processing

* Data Ingestion -

process of acquiring, importing & preparing data for storage analysis in data warehouse, data lake or other target systems

Process

- Data Extraction
- Data Transformation
- Data Validation
- Data Loading

Benefits →

- Improved Data Quality
- Enhanced data accessibility
- Faster time to insights
- Streamlined Analytics workflow

Challenges -

- Data Volume & Variety
- Data Quality issues
- Real-time Data integration.
- Security & Compliance.

Data Ingestion Tools

- ETL (Extract, Transform, Load)
- Data Integration Platforms
- Cloud-based Data Ingestion Services

"Traditional" ML struggles

- slow training time
- large datasets
- data silos hinder.

• Adv.

- Distributed Processing
- In-Memory Computing
- Unified Platforms.

Feature

- spark MLlib
- spark DataFrames
- spark SQL
- MLflow Integration

Benefits

- Faster Training
- efficient Resource Utilization
- Improved Model Accuracy
- Simplified Workflow

Use Case -

- Recommendation System
- Fraud Detection
- Customer Segmentation
- Image & Text Recognition
- Social Network Analysis.

U-6 (DS Honors)

* NLP - (unlocking meaning from Text)

empowers machines to understand & process human language.

Steps

- 1) Text Preprocessing → Tokenization
Normalization (removing stop words)
- 2) Feature Extraction → Bag of words (BOW) (text as a collection of word counts)
TF-IDF (weights to words based on their frequency)
Word Embeddings (words as numerical vectors)
- 3) Applying NLP Techniques → Sentiment Analysis
Named Entity Recognition (NER)
Text Summarization
Machine Translation

Features -

- Automatic Text Analysis
- Improves Machine Understanding
- Enables Powerful Appn
- Uncovers Hidden Insights
- Improves Human Computer Interaction

Disadv -

- Data Dependency
- Limited Context Understanding
- Computational Requirements
- Explanability challenges
- Potential for Bias

Appn of NLP

→ Comm'g & Interaction → Machine Translation, chatbots, smart replies, -

Content Creation & Summarization → sentiment analysis

Text Summarization, machine generated content

Information Extraction & Retrieval →

Question answering systems, name entity recognition & text classification

Add'l appn

→ Voice Search, Speech Recognition, paraphrasing text, legal / financial text processing

* Sentiment Analysis (NLP) -

- 1) Preprocessing Text (Clean & prepare ^{text} data)
- 2) Feature Extraction (Bow, TF-IDF or word embeddings)
- 3) Model Training → Text samples with predefined sentiment labels.
- 4) Sentiment Prediction → Predict sentiment of new, unseen test data.

* Computer Vision

Empowers machine to interpret & understand content of digital images & videos.

Strengths

- Automate Visual Analysis
- Improve Machine Perception
- Unlock Valuable Insights
- Enhanced Human-Computer Interaction
- Wide Range of Applications

Disadvantages

- Reliance on Quality Data
- Limited Generalizability
- Computational Demands
- Privacy Concerns
- Vulnerability to Adversarial Attacks

~~Appn -~~

- Transportation
 - Self Driving Cars
 - Advanced Driver Assistance Sys.
 - Traffic Monitoring & Management
- Security & Surveillance
 - Facial Recognition
 - Video Analytics
 - Object Tracking
- Retail & E-commerce
 - Self Checkouts
 - Product Recognition
 - Visual Search
- Manufacturing & Quality Control
 - Automated Visual Inspection
 - Inventory Management
 - Robot Vision
- Healthcare & Medical Imaging
 - Medical Image Analysis
 - Surgical Robotics
 - Patient Monitoring



2484523

Inclusive of all taxes

Official Use

15 SIZ

* Steps in Computer Vision -

1) Image Preprocessing -

- Resizing / Cropping - adjusting img size for uniformity or focusing on specific regions
- Noise Reduction - removing unwanted artifacts or noise from image to ^{an extent}
- Color Normalization - Adjusting color variations for consistency. improve clarity

2) Feature Extraction -

- Edge Detection - Identifying boundaries & edge within images.
- Color Histograms - Analyzing distribution of colors in image.
- Local Features - Extracting keypoints or regions of interest with image

3) Applying Machine Learning Algorithms -

- Image Classification - Categorize image content
- Object detection - Identifying & locating objects within image
- Image Segmentation - Grouping pixels into meaningful segments corresponding to objects or regions

Object Detection - Spotting Objects in Image -

- crucial (but appn) aims to identify & localize objects within an img or video -
- 1) Object proposal - system generates regions in img that contain objects.
 - 2) Feature Extraction - Features are extracted from regions.
 - 3) Classification & Bounding Box Prediction - ML classifies each region, predicting whether it contains an object & its bounding box.