

U-3 (DS Honors) clustering Methods.

cluster → a no. of similar things that occur together

→ Technique in which data points are arranged in similar groups dynamically without any pre-assign of groups

Appl'n → • Image Processing • Insurance • Recommendation engine

* properties of clustering algo -

- scalability
- ability to deal with different data types
- minimal requirements for domain knowledge for determining input parameters
- Interpretability & usability

* Typical Requirements → . clustering in datamining . • Discovery of clusters with arbitrary slope.

* Problems with Clustering →

- Not address all requir'.
- dealing with large no. of dimensions & data items → problematic
- effectiveness depends on distance
- distance measure shouldn't define it.
- result of clustering → different ways of interpretations

clustering Methods → 1) Partitioning Methods

K-Means
K-Medoids
K-Medians

2) Quality & choosing cluster

SSE
elbow
silhouette index

- 3) KNN
- 4) DBSCAN
- 5) K-Modes.

K-Means.

- Heuristic Method
- supervised learning algo.
- solves clustering problems
- group of unlabeled dataset in diff clusters
- $k \rightarrow$ defines no. of predefined clusters needs to be created.
 $K=2 \rightarrow 2$ clusters.
- each dataset belongs to one group has similar project.
- minimizes sum of b/w data points & clusters
- distance calculation \rightarrow Euclidean distance.

Properties \rightarrow

- Always k -clusters
- Atleast one item in cluster
- clusters are non-hierarchical & don't overlap
- every member of cluster is closer to cluster

Fadv \rightarrow

- efficient in computation
- easy to implement

Disadv \rightarrow

- only when mean is defined
- need to specify k , no. of clusters in advance
- trouble with noisy data & outliers
- not suitable to discover clusters with non complex shapes

Algo \rightarrow

- 1) define k
- 2) Select k random points as centroid.
- 3) compute distance from centroid $\rightarrow d = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$
- 4) Recompute centroids of clusters as $\rightarrow (x_c, y_c) = \left(\frac{\sum x}{m}, \frac{\sum y}{m} \right)$
- 5) Repeat 3 & 4 until one of follow is meet:-
 - a) Centroid don't change
 - b) Points remain in same cluster
 - c) Max no. of iterations are reached

Assessing quality & choosing no. of clusters

- 1) Internal Evaluation Measures → ~~SSE (sum of squared errors)~~
 - ↓ a) Sum of Squared Errors (SSE) -
 - Measures total distance from data points to assigned cluster center
 - Lower SSE → better clustering
 - b) Silhouette Coefficient -
 - measure how well data points fits its assigned cluster compared to other clusters
 - value closer to 1 → better clustering
- 2) External Evaluation Measures
 - a) Adjusted Rand Index (ARI)
 - b) Normalized Mutual Information (NMI)

Choosing no. of clusters

- 1) Elbow Method -
 - determine no. of clusters
 - Inertia is sum of all distances from each data point to centroid of that cluster
- 2) Silhouette Analysis
- 3) Gap Statistics
- 4) k-Means Centroid distance
- 5) Domain knowledge.

* KNN

- K-Nearest Neighbour
- Supervised learning approach
- assumes similarity b/w new case/data & available test case
- regression & classification (more preferred)
- lazy learning algo → not learning from training set immediately
 → stores dataset & time of classification → action

Fadv →

- simple to implement
- robust → noisy training data
- more effective → training data → large

Disadv →

- need to determine k value
- computation cost high
- sensitive → choice of distance

Why KNN? → discover category or class of selected dataset without difficulty

Working →

- select no. of k of neighbours
- calculate Euclidean distance of K no. of neighbours
- take K nearest neighbours as per calculated Euclidean distance
- k-neighbour, count no. of data point in each categories
- Assign new data point to category which no. neighbour is max.
- Our model is ready

Appn →

- Image recognition
- Recommender system
- Customer segmentation
- Anomaly detection

* 1-NN (1-Nearest Neighbor):

Classify a datapoint by assigning it class of its closest neighbor in training data.

Pros →

- easy to implement understand
- Robust to outliers

Cons →

- sensitive → noise in data
- computationally expensive
 ↓ large datasets.

* K-Medians.

- clustering algo.
- similar to K-Means, instead of means to calculate
- uses median to calculate
- Datasets → outliers or non spherical clusters.

Adv → • Robust to Outliers

- suitable for non-spherical clusters
- Intuitive meaning of medoids

Disadv → • expensive

- sensitive to initial medoid selection
- not suitable for high dimensional data

Appn → • Customer segmentation • Image segmentation
• anomaly detection • Document clustering

* Density Based Spatial Clustering

* Hierarchical Clustering

Agglomerative

Divisive

* DBSCAN -

Applⁿ

- Density Based Spatial Clustering of ~~clustering~~ with Noise.
- groups together closely packed points ^
- unsupervised learning

$\epsilon \rightarrow$ radius of neighbourhood around point x .

minpt \rightarrow required to form a density cluster.

- Adv \rightarrow
- don't need to specify no. of clusters
 - flexible in shape & size of clusters
 - able to deal with noise & outliers
 - easy \rightarrow someone who knows dataset \rightarrow to set parameters.

- Disadv \rightarrow
- Input parameters \rightarrow difficult to determine
 - some situation very sensitive to input parameter
 - confused \rightarrow border points belong to two cluster
 - result \rightarrow distance metric
 - hard to guess parameters.



MACGREEN
educare series



Hierarchical Clustering

- method of cluster analysis in which data points are arranged in hierarchy of clusters.

Adv →

- simple to implement
- easy & results in hierarchy (more info)
- doesn't need prespecify no. of clusters.

Disadv →

- large clusters break
- difficult to handle diff sized clusters & shapes
- sensitive to noise & outliers
- can't be changed or deleted once done.

Hierarchical Clustering

Types

Agglomerative clustering

- bottoms up approach
- each item → own cluster
- identifies small cluster
- iteratively clustered are merged together
- also known as AGNES

Divisive clustering

- top down approach
- all item → one cluster
- large cluster
- large clusters are successively divided.
- also known as DIANA.

* Roles of Dendograms →

- diagram representing tree of hierarchy.
 - the similar two objects are less in height of that joins them.
 - Major info lost here
- 1) Visualization
2) Cluster analysis
3) Understanding Cluster Relationship
4) Anomaly Detection
5) Model Interpretation
6) communication & collaboration

* Choosing number clusters in hierarchical clustering

- ↳ 1) Elbow Method
- 2) Silhouette Analysis
- 3) Gap Statistic
- 4) Dendrogram Analysis
- 5) Domain Knowledge.

* Divisive clustering Techniques

- 1) Stopping Criteria.
- 2) Cluster quality metrics
- 3) Domain Knowledge

(specify min cluster size or max distance diameter).

U-4 (DS Honors)

Artificial Neural Network

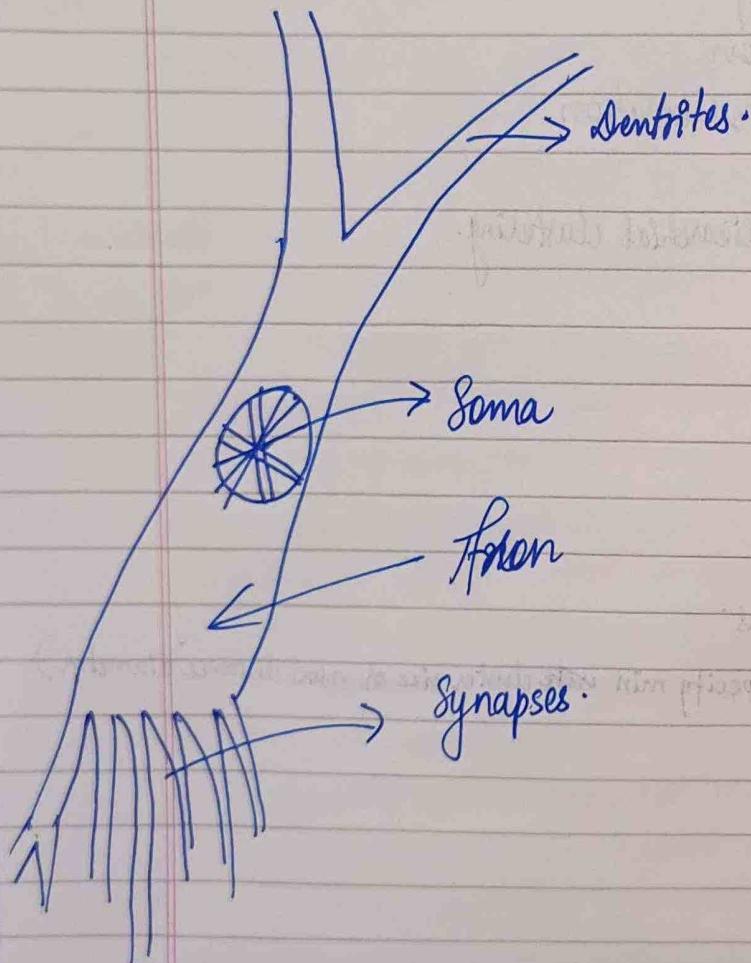
classmate

Date _____

Page _____

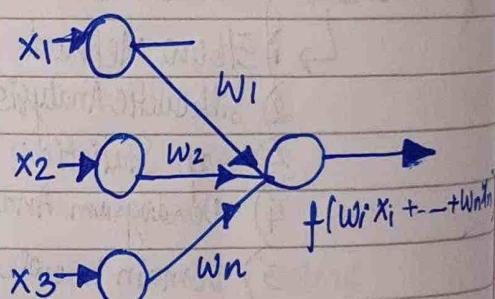
Biological Neuron Model

- Dendrites → accept inputs.
- soma → process input.
- Axon → turns processed input into output.
- Synapses → electrochemical contact b/w neurons



Artificial Neuron Model

- ANN
- Edge or connection or link
- Weight or connect in s



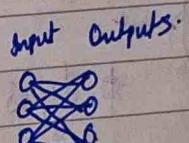
* Models of a Neuron -

- Input Layer → similar to dendrites → input layer receives various inputs or features
- Weights → determine significance of each input.
- Activation Function → determine output of ANN based on weighted sum of inputs.
- Outputs → calculates using activation function & passed to next layer

* Network Architectures:

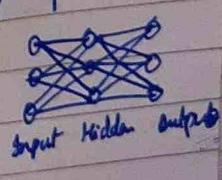
i) ^{Neural} feedforward Network (FNN) -

- Adv →
- simple & easy to implement & train.
 - efficient for learning static patterns & mappings.
 - widely used for inductive learning & basic tasks.
 - unidirectional



ii) Multilayered feed-forward / Multilayer Perceptron -

- Adv →
- complex decision making → create multiple layers of perception.
 - intricate info is too complex for layer to handle
 - multilayer go beyond limitation of single layer.
 - atleast one layer b/w input & output layers.



iii) Radial Basis Function Network (RBFN) -

- Adv →
- Fast learning • High efficiency
 - Adaptability • Parametric control
 - Good generalization
 - Robustness • Overfitting

iv) Functional Link Neural Network (FLNN) -

- Adv →
- Reduced complexity • Improved interpretability • Flexibility
 - Data efficiency • Computational efficiency

2) Feedback Network -

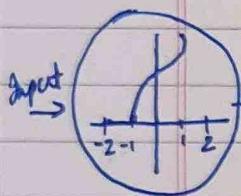
suggests feedback network has feedback paths -

- Recurrent networks - feedback networks with closed loops
- Fully recurrent network - all nodes connected to all other nodes of each node.
- Jordan Network - closed loop network works as both input & output. output will go to input - \hookrightarrow feedback -



Activation Function -

decides whether artificial neural network fire for a given set of inputs.



It's crucial in determining accuracy & computationally efficient of model.

Types of activation function -

$$1) \text{ Identity func}^n = g(x) = x$$

$$2) \text{ Binary step func}^n = g(x) = 1 \text{ when } x > 0 \text{ otherwise } 0.$$

$$3) \text{ Logistic/Sigmoid} = (s(x)) = \frac{1}{1+e^{-x}}$$

$$4) \text{ Tan H} = \text{Tan H}(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

$$5) \text{ Rectified Linear Unit (ReLU)} = R(x) = \max(0, x)$$

* Perceptron

- simple artificial neural network.
- single layer of processing units \rightarrow neurons.
- fundamental building block of more complex neural networks.

key components \rightarrow

- 1) Inputs
- 2) Weights
- 3) Bias
- 4) Activation Function
- 5) Output.

Training \rightarrow

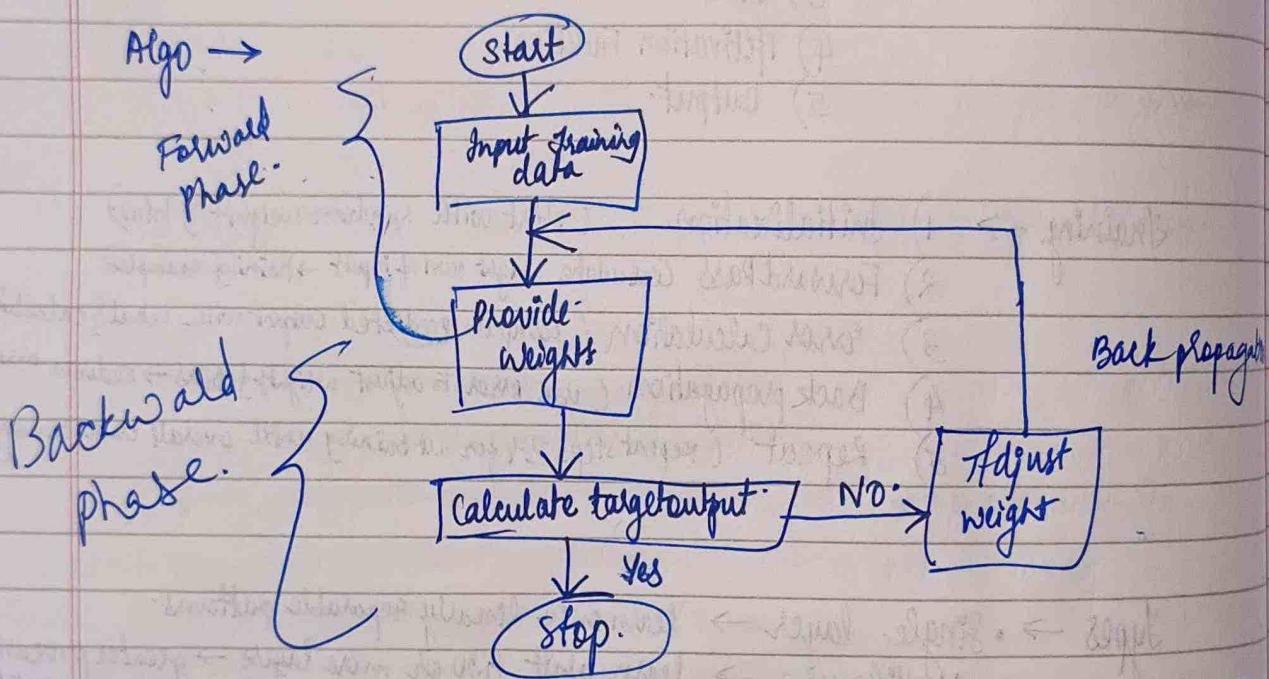
- 1) Initialization (start with random weights & bias)
- 2) Forward Pass (calculate weight sum of inputs \rightarrow training example)
- 3) Error Calculation (compare predicted output with actual & calculate error)
- 4) Back propagation (use error to adjust weights & bias \rightarrow reduce error)
- 5) Repeat (repeat steps 2-4 for all training until overall error converges)

Types \rightarrow

- Single layer \rightarrow learn only linearly separable patterns.
- Multi-layer \rightarrow learn about two or more layers \rightarrow greater processing power

* Back Propagation -

- algo for supervised learning for ANN
- keeps adjusting weights of connecting neurons with an object to reduce deviation of output signal with target output.
- reach global loss min using backpropagation.
- consists of multiple iterations known as epochs.



Features → gradient descent method → case of single perceptron network with diff unit
 → weights are calculate in learning period of network.
 → feedforward of input training pattern.
 calculation & backpropagation of error
 updatation of weight.

Adv → • simple, fast & easy
 • only no. of input architecture are tuned, not any parameters
 • flexible & efficient
 • No need for user to learn any special function.

Disadv → • sensitive to noisy data & irregularities
 • performance → dependent → data
 • too much time taking
 • matrix based approach preferred over min batch

* Generalized Delta Learning Rule -

- most common method for training backpropagation network.
- used to train ANN & MLPs.
- network to learn & adjust its weights & biases based on error it makes.
- extends concept of delta rule used for training single layer perceptron.

Steps → 1) Forward Pass 2) Cost Calculation 3) Propagation
4) Weight Update 5) Repeat.

Benefits →

- efficiently learns complex relationships
- flexible & adaptable
- widely used & well studied.

Limitations →

- gets stuck in local minima
- computational cost
- sensitive to initialization.

Powerful training tool for ANNs. but limitations require careful consideration & exploration of alternative methods.

* Limitations of MLP -

- limited representation power
- Local Minima
- sensitivity to initialization
- overfitting
- computational cost
- Interpretability
- Black Box Nature

U-5 (DS Honors) Convolutional Neural Network

classmate

Date _____

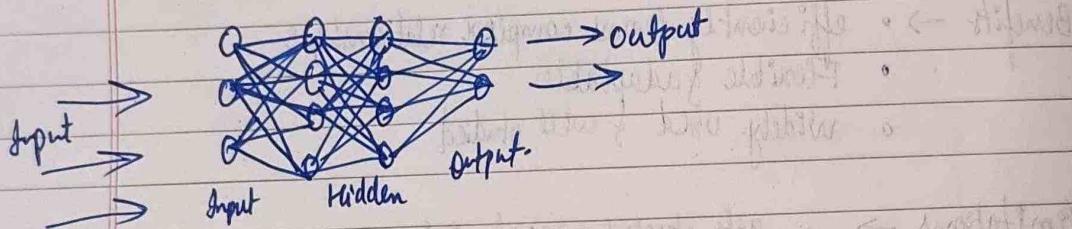
Page _____



convolutional Neural Network - (CNN)

- deep learning neural network
- doesn't require preparation → can operate on raw data
- feed forward NN with up to 20 to 30 layers

- Key features → Convolutional Layers Activation functions
• Pooling Layers • Fully connected Layers



- Adv →
- detecting patterns & features in img, audio & video
 - robust to translation, rotation & scaling invariance
 - end to end training, no need for manual feature extraction
 - handles large amt of data & high accuracy

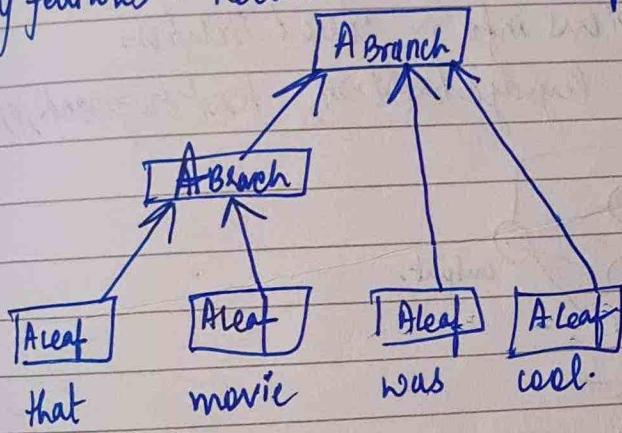
- Disadv →
- computationally expensive
 - lot of money needed
 - prone to overfitting
 - large amt of labelled data for training
 - Interpretability is limited, hard to understand what network has learned.

- Appn →
- Image Recognition
 - Object detection
 - Image Segmentation
 - Video analysis
 - medical image analysis

Recursive Neural Network (RvNNs)

- Type of ANN that handles hierarchical, structured data like NLP, parse trees etc.
- excel at tasks that involve understanding relationship & dependencies

Key features → • Recursive structure • Compositionality • Bated activation function



- Applic →
- NLP → (Text summarization, machine transltn).
 - Code analysis & generation → (functionality of code, identifying bugs)
 - Biological sequence analysis → (analyzing biological sequences like DNA)

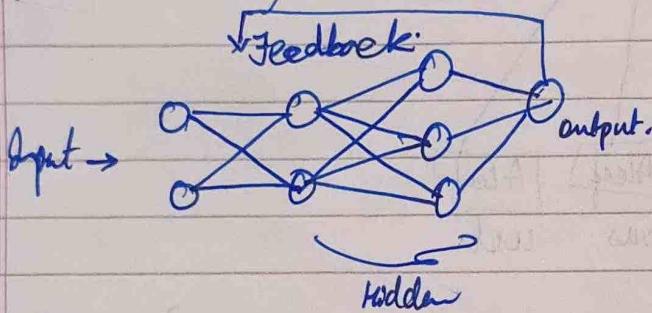
- Disadv →
- Training Complexity
 - Internal decision making
 - Data sparsity (large amt of training data)

- Applic →
- Sentiment analysis
 - Machine translation
 - Text Summarization
 - Code generation
 - Protein structure prediction



Recurrent Neural Network (RNN)

- Feed forward Neural Network
- Preceptors are arranged in layers, hidden layers are not connect with outside world.
- all nodes fully connected, some layers aren't
- need to access previous info in current iteration
- commonly used in language translation, text to speech, etc.



- Adv \(\rightarrow\)
- remember each info throughout time
 - useful in time series prediction \rightarrow feature to symbol input
 - extend effective pixel neighborhood
 - handling variable length sequence
 - generative Modeling

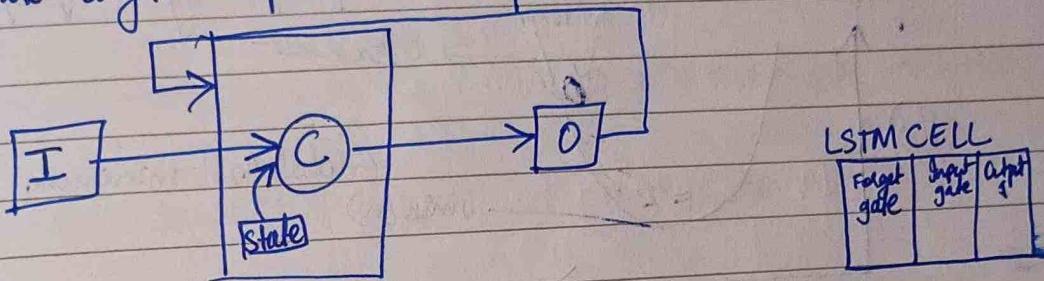
- Disadv \(\rightarrow\)
- gradient vanishing & exploding problems
 - ~~gradient~~ is difficult task to train
 - can't process very long sequence if using tanh or relu as activation func.

- Appln \(\rightarrow\)
- NLP - speech recognition \rightarrow time series forecasting
 - Music generation \rightarrow image captioning

- Types of RNN \(\rightarrow\)
- one to one
 - one to many
 - many to one
 - many to many

* Long Short Term Memory (LSTM)

- advanced RNN
- sequential network.
- capture long term dependencies in sequential data



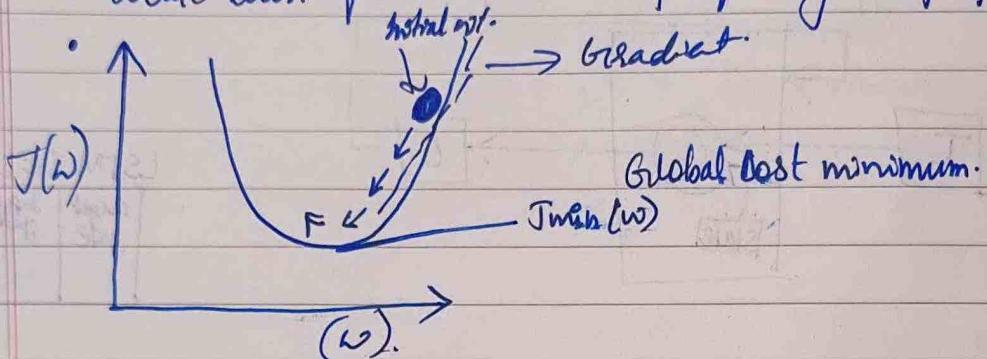
- Fav →
- Capture long term dependencies
 - Overcomes vanishing gradient problems.
 - Wide range of applications
 - Flexible architecture.
 - state of art performance (leading result in various tasks).

- Disadv →
- computationally expensive
 - less interpretable
 - hyperparameter tuning
 - Large data requirements
 - Susceptible to overfitting

- Appn →
- Speech Recognition
 - Recommender System
 - Video Analysis
 - Time Series
 - Language Modelling

Gradient Descent -

- finds best fit line for giving training dataset.
optimization algo \rightarrow min. cost funcn
 - locate least possible value \rightarrow fulfil a given function.



Working → 1) Initialization 2) calculate gradient
3) Update parameters 4) Repeat (until loss function converges)

Benefit → Simple & easy to implement
- Flexible & adaptable

- Flexible & adaptable
 - widely used & studied

limitations →

- Might get stuck at local minima
- sensitive to initialization
- expensive cost
- Tuning Parameters.

Optimization

Types →

- 1) Batch → process all training examples for each iteration
- 2) Stochastic → process 1 training example for each iteration - Faster than BATCH
- 3) MiniBatch → More examples out of total
in training examples are processed per iteration.

V-6 (DS Honors). Appin Perspective.

classmate

Date _____
Page _____

* Text Preprocessing → First step of NLP → preparing data for further processing.

1) Tokenization

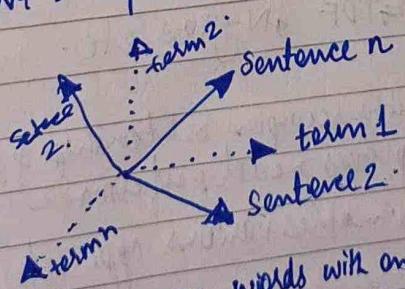
- process of dividing text into meaningful pieces.
- pieces are known as tokens
- whitespace/unicode tokenization → entire text is split into words by splitting them at whitespaces.
- regular expression tokenization → regular expression pattern to get tokens

example → "the sky is blue." becomes → ["the", "sky", "is", "blue".]

- 2) Normalization → text to consistent format for easier processing
- 3) Stop word removal → eliminating the, a, is, etc
- 4) Part of speech tagging → identifies grammatical role of each word
- 5) Feature engineering → Bag of words, n-grams, TF-IDF, etc

* Document Representation.

- In Vector space model → each term/word is an axis/dimension.
- Text/document is represented as a vector in multi-dimensional space



- Ways →
- 1) Hot encoding → words with one hot vector i.e. associate each word with index
 - 2) n-grams → n-words sequence, we take words from sentence & consider it as estimate following words probability given previous words.
 - 3) Bag of words → uses term frequency & used in sentiment analysis.



MACGREEN
educare series

* Feature Selection / Extraction -

subset of relevant features from original set based on their importance

objectives → , Reduce Dimensionality

- Improve Performance
- Avoid overfitting

benefits → Reducing Training time & computational cost -

- Improves model interpretable → focus feature
- prevent overfitting

Limitation → may discard potentially useful info

↳ require careful selection of feature selection technique

* Feature Extraction -

creating new features from original set that captures more info aspects of text

objective → enhances representation of text data.

Techniques → BOW, TF-IDF, n-grams, etc.

Benefits → . capture complex relationship & pattern in text data

- improves model performance

- enables various apps of ML algorithms.

Limitation → . Cost is High

- require careful selection of feature extraction technique

- May introduce irrelevant or redundant information.

J. Techniques -

- 1) Bag of Words → Keep frequency count of all unique words & consider it feature -
Steps → i) Identify unique word from document.
ii) Find frequency of ~~unique~~ unique words from a single sentence.
- 2) Term Frequency - Inverse Document Frequency (TF-IDF).

TF → Frequency of each word in document

IDF → Assign lower weight of words that appear more frequently,
basically depicts rarity of word in document.

TF → Term: frequency in document
Total no. of words in document

$$\text{IDF} = \log_e \left[\frac{\text{Total documents}}{\text{documents with term p.}} \right]$$

* Topic Modelling Algorithms - Latent Dirichlet Allocation (LDA) -

- LDA → helps extract topics from a given corpus.
- classifies data into documents by words per topic, these are modelled based on dirichlet distributions & processes
- Assumptions by LDA → i) Documents are a mixture of topics
ii) Topics are a mixture of tokens.

Working → 1) Preprocessing
2) Latent Variable
3) Generative Process
4) Inference

1) Preprocessing → remove noise, tokenization & text normalization
2) Latent Variable → each topic → probability distribution over words
3) Generative Process → each document → probability distribution over topics
4) Inference → infer latent variable given observed data -
(topic distribution) (document & word)



(LDA)

- * Adv →
 - Unsupervised
 - Probabilistic
 - Scalable
 - Interpretable

- * Disadv →
 - Topic coherence
 - No. of topics
 - Overfitting.

- * Appn →
 - Document clustering • Topic summarization
 - Recommendation system • Trend analysis

* Stemming & Lemmatization

- technique used in NLP to reduce words to their basic roots
- remove words of reded inflectional or derivational variations

• Stemming

- process of removing prefixes & suffixes from word to obtain word's base form
- resulting stem may not always be a valid word

• Way faster than Lemmatization

- ex → original word → fishing, fishes, fished
stemmed words → fish, fish, fish

Lemmatization

- transform words to their base form or lemma, conserving word's meaning & part of speech
- uses linguistic rules & morphological analysis to achieve this
- resulting lemma is a valid word representing base form

- example original words
walking, walks, walk
lemmatized → walk, walk, walk

* Text Similarity Measures -

- Used to find similar text.
- Google, Quora, Stack overflow use it to find similar questions.
- Important in NLP for tasks like information retrieval, document clustering & machine translation.

feature to consider when choosing text similarity measure

- Task at hand
- Nature of text data
- Computational resource
- Interpretability

- Types →
- 1) Lexical - Jaccard Similarity, Levenshtein Distance (Focus on words by characters, not on meaning)
 - 2) Semantic - Cosine Similarity (Capture meaning, intent, relationship between words)
 - 3) Featurebased - kernels, Tree kernels (Utilize extracted features like structure, topics)
 - 4) Neural-Network based - sentence BERT, Siamese Networks (Learn complex semantic relationship).
capture non-local relationships.

Appn → Info Retrieval, document clustering, machine translation, plagiarism detection, chatbots etc

* Jaccard Similarity → Ratio of common words to total words. (No imp to duplication of words)

$$JS = \frac{|A \cap B|}{|A \cup B|} = \frac{|A \cap B|}{|A| + |B| - |A \cap B|}$$

* Cosine Similarity →

$$\text{Similarity}(A, B) = \frac{A \cdot B}{|A| \times |B|} = \frac{\sum (A_i \times B_i)}{\sqrt{\sum A_i^2} \times \sqrt{\sum B_i^2}}$$

- suitable → words → repeated & important
- ratio of dot product of two vectors of words to their product of magnitude.