

U-1 (Data Science Honours)

DOMS

Page No.

Date

/ /

* Need for AI -

- handling massive data
- Automating tasks
- Improving accuracy
- Working 24/7
- Innovation across fields
- Enhanced decision making
- Personalized experiences
- Augmented human capabilities
- Addressing global challenges
- Exploration & discovery

* Applications -

- Image & speech recognition
- Recommendation system
- Fraud Detection
- Self driving cars
- Customer service
- Healthcare
- Manufacturing
- Agriculture

Logic programming (solving problems)

- Based on formal logic
- allows you to solve problems by declaring facts & rules.
- Facts → represent basic truths about world.
- Rules → define relationships b/w facts.

Benefits →

- Declarative (Focus on ~~how~~ what to solve)
- Readable → clear code resembling logic statements.
- Flexible → easy to add new facts or rules
- well suited for specific problems

Search

Heuristic Techniques -

- clearer way to find good solutions in complex problems
 - Weak Methods → (don't apply great deal of knowledge).
 - Not domain or problem specific
 - Heuristic Function → applied to a state in search space indicate a likelihood of success if state is selected.
 - given a search space, a current state & a goal state.
 - generate all successor states & evaluate each with our heuristic function.
 - select move yields best heuristic value.
 - faster than exhaustive search, especially for large problems.
 - finds good solutions even if optimal one is hard to find.
- eg → game playing (eg. chess), 8-puzzle, navigating robots or in games scheduling problems.



Constraint satisfaction problem -

→ consists of 3 components - V, D, C

Variable
 $\{N, NW, NE, M, MW, ME, S, SE\}$

Domains.

constraints
 adjacent regions.
 must have diff. colors.

→ finding soln that meets the constraints -

→ V is set of variables $\{V_1, V_2, V_3, \dots, V_n\}$

→ D is set of domains $\{D_1, D_2, D_3, \dots, D_n\}$ one for each variable.

→ C is set of constraints specify allowable combination of values.
 $C_p = (\text{scope}, \text{rel}) \quad \{C_1, C_2, C_3\}$

→ scope is set of variable that participate in constraint.

→ rel is relation that define values that variable to take.

→ Properties of CSPs -

- specific partial info.
- non-directional
- declarative
- additive
- independent.

→ example - Map Coloring, Sudoku, scheduling problems, configuration problems etc.

→ Variables & Constraints. find an assignment of values to all variables such as all constraints are satisfied.

* Local Search Techniques -

- Algorithms that iteratively explore solution space by making incremental changes to current solution, aiming to improve gradually.

Characteristics → Iterative Improvement → Adaptability
 → Greedy Approach → Scalability
 → May not guarantee optimality → Heuristic based.

Types → Hill Climbing Applications → Travel Salesman Problem
 → Simulated Annealing → Graph Coloring
 → Genetic Algorithms → Scheduling

Limitations → Local optima
 → sensitivity to initial solution
 → Difficulty in Escaping local optima
 → Limited Exploration
 → can get stuck in local optima
 → effectiveness depends on problem

* Greedy search → simple & intuitive algo used to solve optimization problems.
 → Hope of finding optimum

Characteristics → Local optimization Limitations → No Backtracking
 → Simple Implementations → May not find global optimum
 → efficiency → Depends on Heuristics
 → applicability → non optimal for some
 → simple & efficient → suboptimal solⁿ problems
 → good for approx. solⁿ → problem dependent

Big Data Learning (V2)

Honours Data Science

DOMS

Page No.

Date

/ /

* Characteristics of Big data -

Big data → large vol. of data available at various sources in varying degree of complexity, generated at different speed i.e. velocities & technologies, processing modules.

- volume → • big data characterized → enormous value → ~~reliability~~
 - large data needs specialized tools & technologies to store process & analyze
 - observe & tracks data from various sources.
- velocity → • data streams → high speed & dealt timing
 - processing of data → streamed data → real time results → fast
 - insights generated → relevant & actionable.
 - speed of generation of data.
- variety → • heterogeneous sources -
 - nature of data → structured & unstructured
 - social media data → sensors, audio, video, etc.
 - data → insights → manage & analyze.
- value → • business value → Big Data.
 - generate some sort of value → company doing analysis
 - insights
 - formal decisions
 - optimize operations
 - gain competitive edge in marketplace.
- veracity → • inconsistency / uncertainty in data
 - challenges in data quality & Reliability
 - Insights → accurate & trustworthy.

* Types of Data

	Structure	Unstructure	Semi-structured
	predefined schema (like tables)	No predefined format.	some internal structure
eg →	Databases, spreadsheets	Text documents, images, videos	Log files, JSON, XML
Analysis	• Easiest (SQL, statistics)	• Requires specialized technique	• Easier than unstructured
Storage	Relational databases, spreadsheets	File systems, cloud storage	• File system, cloud storage NO SQL databases
Uses	Transactions, reporting ML	Text analytics, image recognition	• Configuration, log analysis data exchange
Table?	Yes (directly)	No (needs processing)	Maybe (after processing)

Supervised ML

- Learns from labelled data with input-output pairs
- Labelled training data -
- Predicts output based on input data
- Receives feedback during training.
- Can test our model
- Desired output is given
- Also called classification.
- eg → classification, regression, object detection.



Regression Analysis -

- statistical technique used to estimate relⁿ b/w a dependent variable & one or more independent variable.
- identify general trend b/w variables
- use model to predict dependent variable based on value of its independent variables.

Unsupervised ML

- Learns patterns & structures from unlabelled data.
- doesn't require labelled
- discovers hidden patterns or structure in data.
- Doesn't receives feedback.
- desired output is not given
- clustering
- K-Mean clustering, PCA

* Process \rightarrow • Data Collection • Parameter Estimation
 • Model selection • Evaluation.
 (Linear regression, multiple etc).

* Types of Regression.

1) Linear Regression \rightarrow ML algorithm \rightarrow supervised learning
 \rightarrow predict dependent variable based on independent.
 \rightarrow regression line is best fit line for model.
 $x \rightarrow$ Independent variable $y \rightarrow$ output.

2) Simple linear Regression \rightarrow one independent (or input) variable.

3) Multiple linear Regression \rightarrow more than one independent variables.
 \rightarrow multivariate regression.

4) Logistic Regression \rightarrow used to model & estimate probability of an event
 \rightarrow estimate probability of an event occurring
 \rightarrow is a curve
 \rightarrow sigmoid curve or S-curve in short
 $\rightarrow 0 \rightarrow$ totally ^{un}certain & $1 \rightarrow$ certain



Clustering -

cluster \rightarrow a no. of similar things that occur together.

\rightarrow technique in which data points are arranged in similar groups dynamically without any pre-assignment of groups.

partitioning a set of data in a set of meaningful subclasses.

properties of a cluster -

- all data points in a cluster should be similar to each other.
- data points from diff clusters should be as different as possible.

Types \rightarrow Hard clustering \rightarrow each data point \rightarrow only one cluster

\rightarrow soft clustering \rightarrow each data point \rightarrow separate cluster probability or likelihood of data.

Applⁿ \rightarrow

- Customer Segmentation
- Marketing
- Seismology
- Image Processing
- Insurance
- Land Use
- Recommendation Engine
- Healthcare
- Urban Planning



K-Means

- heuristic method.
- unsupervised learning algo.
- solve clustering problems.
- groups are unlabelled dataset in diff clusters.
- $k \rightarrow$ define no. of predefined clusters.
- adv \rightarrow
 - efficient in computation.
 - easy to implement
- disadv -
 - applicable only when mean is defined
 - need to specify k .
 - trouble with noisy data.
 - not suitable to discover clusters.

Hierarchical Clustering -

- hierarchical cluster analysis or HCA.
- method of cluster analysis
- data points are arranged in hierarchy

Agglomerative

- bottom up approach
- each item \rightarrow cluster
- iteratively merged together

Divisive

- top down approach
- all item \rightarrow one cluster
- large clusters successively divided

* Data dimensionality -
no. of attributes or features used to represent each data point in a dataset.

Type → High → many features relative to no. of observations.
→ Low → Features relative to no. of observations.

Impact → .curse of dimensionality
• overfitting
• Interpretability
• Computational complexity

Dimensionality Reduction
→ Feature Selection
→ Feature Extraction

Techniques → PCA (Principal Component Analysis)
→ + distributed Stochastic Neighbor Embedding
→ Linear discriminant analysis
→ Autoencoders

Applⁿ → Image & Video Processing
→ Text Mining & NLP
→ Bioinformatics
→ Financial Modeling

* Spark Programming Model → Built for distributed computing on big data
→ Lazy evaluation of transformations & actions.
→ Fault tolerance through lineage info.
→ In memory processing for high performance.

Benefits	Adv	Disadv
→ Scalability	→ Flexibility	→ Complexity
→ High Performance	→ Interactivity	→ Resource Management Overhead
→ Fault Tolerance	→ Integration	→ Overhead in memory Processing
→ Ease of Use	→ more approachable	→ Debugging & Monitoring
→ Unified Framework	for programmers	

* MLlib Library -

- algorithm for Big Data - offers various ML algo to handle large datasets
- provides easy to use tools for ML workflows high level APIs.
- ML library for Apache Spark.

Adv → Scalability
 Variety of algo
 ease of use
 Integration with Spark
 Active community
 flexibility
 simplified development

Disadv → Learning Curve
 Complexity for Beginners
 limited out of box features
 dependency on Spark ecosystem

* Content based recommendation System -

- suggest items based on features of items & user preferences.
- don't rely on user-item interaction but analyze item attributes.

Adv → • Interpretability
 • Cold start Mitigation
 • User Personalization
 • Transparency
 • User Independence

Disadv → Data Sparsity
 Echo chambers
 limited context consideration
 lack of diversity in recomm

Appln → E-commerce
 media streaming
 News Aggregation