

Constraint Programming: Powerful paradigm of declarative programming that lies at the intersection of AI, operations research (OR) and logic programming. It is used to solve complex combinatorial problems in various fields such as planning, resource optimization, robotics and intelligent transportation systems. The user specifies a problem in terms of constraints, and the machine solves it, making this paradigm attractive for addressing complex challenges.

Constraint Network: A constraint network $N = (X, D, C)$ is defined on a given vocabulary (X, D) , where we add a finite set of constraints C . Each constraint $c \in C$ is a pair $\langle \text{var}(c), \text{rel}(c) \rangle$, where $\text{var}(c) \rightarrow$ sequence of variables from X called the scope of the constraint and $\text{rel}(c) \rightarrow$ relation on $D^{|\text{var}(c)|}$ called the constraint relation which defines the allowed combinations of simultaneous values for the variables in $\text{var}(c)$.

A constraint language is a set $T = \{t_1, \dots, t_n\}$ of relations over a subset of \mathbb{Z}^n . CSP: A constraint satisfaction problem is the task of finding a solution that satisfies all constraints in a constraint network.

COP: A constraint optimization problem is defined by a constraint network (X, D, C) and an objective function $f: \Omega \rightarrow \mathbb{R}$ to be minimized or maximized. The goal is to find a solution to the constraint network that min/max the objective function. Ensures that each arc (x_i, x_j)

AC-3: The AC-3 algorithm meaning that for every value of x_i , there exists a compatible value of x_j . The input is a constraint network (X, D, C) . The output a network made arc-consistent, and the usage of the algo is for preprocessing before search to reduce domains.

B&B: Branch and Bound is a systematic search algorithm used in optimization problems (including constraint programming with an objective function like $C = f(X)$ where X represents the variables of the CSP, this is called the objective constraint). It combines backtracking search with pruning using bounds on the objective function.

Data Mining: Often referred to as Knowledge Discovery in Databases (KDD), is a multidisciplinary field that focuses on uncovering meaningful patterns, relationships and insights from large and complex datasets. The ultimate goal of DM is to transform raw data into actionable and valuable knowledge that can be used for decision-making, prediction and exploration.

Challenges in DM: Data Complexity, Scalability, Interpretability, Ethical concerns and User Centring and Human-in-the-loop.

Items: An item is one distinct object (e.g. $t_1, t_2, \dots, t_m = I$)

Itemset: An itemset P is any non-empty subset of those items I .

Frequency: The frequency of an itemset P in D , denoted as $\text{freq}(P)$ is said to be frequent if its frequency is above a given threshold α and an itemset P : P is closed iff there does not exist any superset $Q \supset P$ such that $\text{freq}(Q) = \text{freq}(P)$.

Minimality: P is generator iff there does not exist any subset $Q \subset P$ such that $\text{freq}(Q) = \text{freq}(P)$.

Given a dataset D , a min. Rule:

frequency α and a min. confidence β , a MNR $X \rightarrow Y$ is a valid rule such that there does not exist any rule $X' \rightarrow Y'$ with $X' \subseteq X$, $Y' \subseteq Y$, $\text{freq}(X') = \text{freq}(Y')$, $\text{conf}(X') = \text{conf}(Y')$ and $X' \neq X$.

Anti-Monotonicity Property: The anti-monotonicity property implies that if a set of items is infrequent, then any larger set that contains this set must also be infrequent. Let $P = \{a_1, \dots, a_m\}$

tail: The tail of P is the last item in P , i.e., $\text{tail}(P) = a_m$. i.e. I and a_m exists according to a fixed order on the items.

Up to the j th item is: prefix: The prefix of P is the prefix $P_{(j)} = \{a_1, \dots, a_j\}$ tree, is defined as a prefix $(P) = \{a_1, \dots, a_m\}$

SAT: Declarative Invariant of Data Mining leverages constraint programming (CP) and propositional satisfiability (SAT) to model and solve data mining tasks in a more flexible way. \exists independent constraint network: SAT is essentially a CP model for a data-mining problem such as itemset mining or association rule mining that does not directly encode any constraint in the space of the target patterns. Instead, it focuses on the rule relationships among the variables and the domain consistency properties needed to prune the search space. This approach still allows the use of global constraints but abstracts away from state-based filtering.

SAT-Based Approach: SAT-based approach, which can be used to translate the mining task into a purely Boolean. These methods have been successfully applied to mining closed and maximal itemsets, top-k frequent itemsets and association rules. Compared to CP, SAT-based approaches offer efficiency gains, parallelism, and modularity.

NP-Hard Problem: A problem is called NP-Hard (Non-deterministic Polynomial-time Hard) when its resolution is extremely challenging, even for powerful computers. This means that for large problem sizes, it is practically impossible to find an optimal solution in a reasonable time.

Efficiency (after): Although CP involves searching in a potentially very large solution space, modern solvers are often able to exploit constraint propagation techniques and heuristics search to efficiently find solutions, even for large problems.

CP Advantages: compactness: CP allows problems to be modeled concisely. constraints directly express complex relationships between variables, making models simpler and more readable than traditional linear programming equations.

Natural Modeling: Many problems, especially in planning, design or optimization, can be naturally formulated in terms of constraints. This leads to simpler and more readable models compared to other approaches like ILP.

Separation of Search and Modeling: In CP, problem modeling is separated from solution searching/constraint solving allowing designers to focus on defining requirements without worrying about the details of the solving algorithm.

BT: Backtracking is a systematic search method that constructs a partial instantiation by assigning variables one by one. When an instantiation violates a constraint, the method backtracks to try another value. It's simple to implement but inefficient if domains are large or constraints numerous.

FC: Constraint propagation method used during search, especially when assigning variables. At each step, when a variable is instantiated, FC checks the constraint between this variable and the uninstantiated ones, removing incompatible values from their domains.

Heuristics in CSP search: Heuristics significantly enhance search efficiency by guiding which variable to assign first and which value to choose for that variable.

Global constraint: A global constraint is a constraint defined over an arbitrary set of variables that captures a specific property or pattern, often recurring across multiple problem domains.

Advantages: completeness: A global constraint expresses a complex property in a single declarative statement. Searcher Propagation: By using advanced filtering algorithms they achieve higher levels of domain reduction thus decreasing the search space.

Encouraging reuse: Commonly used global constraints can be applied across various problems and domains, fostering modularity.

Steps in KDD process:

1. Data Selection: Identifying the relevant datasets for analysis, often from heterogeneous data.

2. Data Preprocessing: Cleaning, transforming and normalizing the data to ensure consistency and quality.

3. Data Transformation: Converting raw data into formats or features suitable for analysis.

4. Data Mining: Applying algorithms to extract patterns, associations or predictive models.

5. Knowledge Presentation: Visualizing and communicating results in a form that stakeholders can readily understand and act upon.

Frequent Pattern Mining: FPM has gained considerable attention due to its broad applicability and theoretical richness.

Originally introduced by Agrawal et al., involves identifying and analyzing recurring patterns within datasets.

Transactional dataset: collection of m itemsets, denoted as t_1, \dots, t_m which are called transactions.

Cover: The user of an itemset P , denoted as cover(P), refers to the set of transactions in D that contain P .

Maximal and Minimal Itemsets: Given a dataset D and a minimum frequency α , the itemset P is maximal iff P is frequent and there does not exist any superset $Q \supset P$ s.t. $\text{freq}(Q) > \alpha$.

Respectively, P is minimal iff P is frequent and there does not exist any subset $Q \subseteq P$ such that $\text{freq}(Q) > \alpha$.

Association Rule: An association rule captures information of the kind "if we have A and B, the chance to have C are high". It's an implication of form $X \rightarrow Y$, where X and Y are itemsets such that $X \cap Y = \emptyset$ and $X \neq \emptyset$. X represents the body of the rule and Y represents the head. The frequency of the rule $X \rightarrow Y$, i.e., $\text{freq}(X \rightarrow Y) = \text{freq}(XY)/\text{freq}(X)$. The confidence of a rule captures how often Y occurs in transactions containing X , that is, $\text{conf}(X \rightarrow Y) = \text{freq}(XY)/\text{freq}(X)$. A rule is known as a valid rule if its frequency and confidence are greater than or equal to user-specified thresholds (α and β). The problem of mining association rules consists in generating all valid rules.

Apriori Property: Direct consequence of anti-monotonicity which states that if an itemset is frequent, then all of its subsets must also be frequent. Let P be an itemset. If P is frequent, then all subsets of P , $\forall S \subseteq P$, must also be frequent, $\text{freq}(P) \geq \text{freq}(S)$.

Closure: A set S has a closure under an operation f if: $\forall x \in S \rightarrow f(x) \in S$. A closure operation is extensive, idempotent, monotone.

Closure: The closure of an itemset P is the intersection of all the transactions t in the cover of P :

Closure (P): $\text{closure}(P) = \bigcup_t \text{closure}(P \cap t)$. An itemset is said to be closed if it is equal to its closure: $P = \text{closure}(P)$.

CPIM: Model that relates 2 vectors of Boolean variables.

Decision Variable: x_1, \dots, x_n where $x_i: A \rightarrow$ item i is included in the transaction t .

Auxiliary Variable: t_1, \dots, t_m where $t_k: B \rightarrow$ the extracted itemset appears in transaction t_k of the dataset D .

Encoding CPIM Constraints:

Inequality: $\forall t \in D: t_k \rightarrow \bigwedge_{i \in I} \bigvee_{j \in J} x_i(t \wedge \text{closed}(x_i, t)) = 0$

MinFreq: $\exists t_k \forall x_i: \bigwedge_{j \in J} \text{closedness}(x_i, t_k) \rightarrow \bigvee_{i \in I} x_i(t_k)$

MaxFreq: $\exists t_k \forall x_i: \bigwedge_{j \in J} \text{closedness}(x_i, t_k) \rightarrow \bigvee_{i \in I} x_i(t_k)$

MaxFrequency constraint:

Voice I: $x_i = 1 \leftrightarrow \bigwedge_{t \in D} \text{cover}(x_i, t) \wedge \bigwedge_{t \in D} \bigvee_{j \in J} x_i(t \wedge \text{closed}(x_i, t)) = 0$

Voice II: $x_i = 1 \leftrightarrow \bigwedge_{t \in D} \text{cover}(x_i, t) \wedge \bigwedge_{t \in D} \bigvee_{j \in J} x_i(t \wedge \text{closed}(x_i, t)) = 0$

Monotonicity: closedness: The closure of a closed itemset is always closed.

Efficient Construction: The closure extension is always closed and checking for closedness is much faster than checking for every individual itemset.

Therefore by incrementally extending smaller closed itemsets, without having to check for every possible combination of items directly.

