

Constraint satisfaction problems

Factor graph - (aka Markov random field) a set of **variables** $X = X_1, \dots, X_n$ where $X_i \in \text{Domain}_i$ and **factors** f_1, \dots, f_m , with each $f_j(X) \geq 0$. **Each factor is implemented as checking a solution rather than computing the solution.**

Domain - possible values to be assigned to a variable.

Scope of a factor f_j - the set of variables f_j depends on. **Arity** - the size of this set. “Unary factors” (arity 1); “Binary factors” (arity 2). “Constraints” (factors that return 0 or 1).

Assignment weight - each assignment $x = (x_1, \dots, x_n)$ yields a $\text{Weight}(x)$ defined as being the product of all factors f_j applied to that assignment.

$\text{Weight}(x) = \prod_{j=q}^m f_j(x)$ (x in its entirety is passed in to each f_j for simplicity of this notation, though in reality only a subset of x would be needed for f_j)

CSP - a factor graph where all factors are binary.

For $j = 1, \dots, f_j(x) \in \{0, 1\}$ (the constraint j with assignment x is said to be satisfied iff $f_j(x) = 1$.)

Consistent assignment x of a CSP - iff $\text{Weight}(x) = 1$ (i.e., all constraints are satisfied.)

Dependent factors $D(x, X_i)$ - a set of factors depending on X_i but not on unassigned variables.

Backtracking search - find maximum weight assignment of a factor graph.

- Backtrack**($x, w, \text{Domains}$)
- choose an unassigned **variable X_i (MCV)**
 - order **values** of X_i ’s Domain (**LCV**)
 - for each value v in the order:
 - $\delta \leftarrow \prod_{f_i \in D(x, X_i)} f_j(x \cup \{X_i : v\})$
 - if $\delta = 0$: continue
 - $\text{Domains}' \leftarrow \text{Domains}$ via **lookahead (forward checking)**
 - if $\text{Domains}'_i$ is empty: continue
 - Backtrack**($x \cup \{X_i : v\}, w\delta, \text{Domains}'$)

- Strategy: extends partial assignments
- Optimality: exact
- Time: exponential

Forward checking - one-step lookahead heuristic that preemptively removes inconsistent values from the domains of neighboring variables.

- After assigning a variable X_i , it eliminates inconsistent values from the domains of all its neighbors.
- If any of these domains become empty, stop the local backtracking search.
- if we unassign a variable X_i , have to restore the domain of its neighbors.

Most constrained variable - selects the next unassigned variable that has the fewest consistent values: fail early, prune early. **Useful when some factors are constraints.**

Least constrained value - assigns the next value that yields the highest number of consistent values of neighboring variables: prefers the value that is most likely to work. **Useful when all factors are constraints.**

Arc consistency of variable X_i - w.r.t. X_j is enforced when for each $x_i \in \text{Domain}_i$, there exists $x_j \in \text{Domain}_j$ such that any facts between X_i and X_j is non-zero.

AC-3 - a multi-step lookahead heuristic that applies forward checking to all relevant variables. After a given assignment, it performs forward checking and then successively enforces arc consistency w.r.t. the neighbors of variables for which the domain change during the process.

AC-3 only looks locally at the graph for nothing blatantly wrong; it can’t detect when there are no consistent assignments.

Beam search - extends partial assignments of n variables of branching factor $b = |\text{Domain}|$ by exploring the K top paths at each step. The beam size $1 \geq K \geq b^n$ controls the tradeoff between efficiency and accuracy. **Runtime is Linear to n : $O(n \underbrace{Kb \log(Kb)}_{\text{sorting top } K})$.** $K = 1$: greedy search ($O(nb)$

time); $K \rightarrow +\infty$: BFS ($O(b^n)$ time).

- Strategy: extends partial assignments
- Optimality: approximate
- Time: linear

Local search (iterated conditional modes) - modifies the assignment of a factor graph one variable at a time until convergence. At step i , assign to X_i the value v that maximizes the product of all factors connected to that variable.

- Initialize x to a random complete assignment (not partial)
- Loop through $i = 1, \dots, n$ until convergence:
 - Compute weight of $x_v = x \cup X_i : v$ for each v
 - $x \leftarrow x_v$ with highest weight

ICM may get stuck in local optima; adding randomness may help.

- Strategy: modify complete assignments
- Optimality: approximate
- Time: linear

Markov Networks

CSPs	Markov networks
variables	random variables
weights	probabilities
max weight assignment	marginal probabilities

Capture the uncertainty over assignment using the language of probability.

Markov Network - a factor graph which defines a joint distribution over random variables

$X = (X_1, \dots, X_n)$: $\mathbb{P}(X = x) = \frac{\text{Weight}(x)}{Z}$

$Z = \sum_{x'} \text{Weight}(x')$ - sum all the possible assignments’ weights (normalization constant)

Marginal probability - the probability of when one particular variable X_i is assigned with a

particular value v : $\sum \mathbb{P}$ when $X_i = v$

$\mathbb{P}(X_i = v) = \sum_{x: x_i=v} \mathbb{P}(X = x)$

Gibbs sampling - approximately computes marginal probabilities. **Follows the template of local search: change one variable at a time. But unlike ICM, Gibbs is a randomized algo .**

- Initialize x to a random complete assignment.
- Loop through $i = 1, \dots, n$ until convergence:
 - Set $x_i = v$ with probability $\mathbb{P}(X_i = v | \underbrace{X_{-i}}_{\text{all vars except } X_i} = x_{-i})$
 - Increment $\text{count}_i(x_i)$ (how often this assignment is encountered. **can just track particular vars we’re interested in.**)
- Estimate $\hat{\mathbb{P}}(X_i = x_i) = \frac{\text{count}_i(x_i)}{\sum_v \text{count}_i(v)}$

ICM	Gibbs sampling
max weight assignment	marginal probabilities
choose best value	sample a value
converges to local optimum	marginals converge to correct answer

$P(A|B, C) = \frac{P(B|A, C)P(A|C)}{P(B|C)}$

Bayesian Networks

Explaining away - suppose two causes positively influence an effect. Conditioned on the effect, further conditioning on one causes reduces the probability of the other cause.

Bayesian network - a directed acyclic graph that specifies a joint distribution over random variables $X = (X_1, \dots, X_n)$ as a product of local conditional distributions, one for each node:

$\mathbb{P}(X_1 = x_1, \dots, X_n = x_n) = \prod_{i=1}^n p(x_i | x_{\text{Parents}(i)})$

Probabilistic program - randomizes variable assignment such that we can write down complex Bayesian networks that generates assignments without having to explicitly specify assignments probabilities. **Unlike normal classification (e.g., neural nets), Bayesian networks provide a different paradigm where we think about going from output to the input.**

Probabilistic inference strategy - to compute the probability $P(Q|E = e)$ of query Q given evidence $E = e$:

- Remove vars that aren’t ancestors of the query Q or the evidence E by marginalization
- Convert Bayesian network to factor graph
- Condition on the evidence $E = e$
- Remove nodes disconnected from the query Q by marginalization
- Run probabilistic inference algorithm

Filtering question - asks for the distribution of some hidden variable H_i conditioned on only the evidence up until that point. **Useful for real-time object tracking as the future can’t be seen.**

Smoothing question - asks for the distribution

of some hidden variable H_i conditioned on on the evidence including the future. **Useful when all the data have been collected and we want to retrospectively go and figure out what the hidden state H_i was.**

Forward-backward algorithm - computes the exact value of $P(H = h_k | E = e)$ a smoothing query) for any $k \in \{1, \dots, L\}$ in the case of an HHM of size L .

- for $i \in \{1, \dots, L\}$, compute $F_i(h_i) = \sum_{h_{i-1}} F_{i-1}(h_{i-1})p(h_i | h_{i-1})p(e_i | h_i)$
- for $i \in \{L, \dots, 1\}$, compute $B_i(h_i) = \sum_{h_{i+1}} B_{i+1}(h_{i+1})p(h_{i+1} | h_i)p(e_{i+1} | h_{i+1})$
- for $i \in \{1, \dots, L\}$, compute $S_i(h_i) = \frac{F_i(h_i)B_i(h_i)}{\sum_{h_i} F_i(h_i)B_i(h_i)}$

with the convention $F_0 = B_{L+1}$. We get

$P(H = h_k | E = e) = S_k(h_k)$

Particle (partial assignment) filtering -

approximates the posterior density of state variables given the evidence of observation variables by keeping track of K particles at a time. Initialize $C \leftarrow [\{\}]$. For $t = 1, \dots, n$:

- proposal**: for each old particle $x_{t-1} \in C$, sample x from the transition probability distribution $p(x | x_{t-1})$ and add x to a set C' .
- weighting**: weight each x of the set C' by $w(x) = p(e_t | x)$ (e_t is the evidence observed at time t).
- resampling**: sample K elements from the set C' using the probability distribution induced by w and store them in C : these are current particles x_t .

Maximum likelihood - if we don’t know the local conditional distributions, we can learn them using max likelihood.

$\max_{\theta} \prod_{x \in D_{\text{train}}} p(X = x; \theta)$ **HMM parameters:**

$\theta = (p_{\text{start}}, p_{\text{trans}}, p_{\text{emit}})$.

Laplace smoothing - for each distribution d and partial assignment $(x_{\text{Parents}(i)}, x_i)$, add λ (a constant) to $\text{count}_d(x_{\text{Parents}(i)}, x_i)$, then normalize to get probability estimates. **Hallucinate λ occurrences of each local assignment.** Larger $\lambda \Rightarrow$ more smoothing \Rightarrow probabilities closer to uniform.

Expectation-Maximization (EM) - estimates the parameter θ through maximum likelihood estimation by repeatedly constructing a lower-bound on the likelihood (E-step) and optimizing that lower bound (M-step):

- E-step: evaluate the posterior probability $q(h)$ that each data point e came from a particular cluster h : $q(h) = P(H = h | E = e; \theta)$
- M-step: use the posterior probabilities $q(h)$ as cluster specific weights on data points e to determine θ through maximum likelihood.

Like K-means, EM converges to a local optima. EM can handle partial data.