# Search

First Try: Reflex Agent (only based on memory, prediction X)
↳ can be rational if needing quick decisions

Second Try: Planning Agent (decision based on possible consequences)
↳ completeness (gives an answer?), optimality (best answer?)
↳ a "replanning" agent solves the problem on-the-fly.

Search Problems consist of:
- state space
- successor function (actions & costs)
- start state & goal test
→ a solution is a sequence of actions that transforms
  the start state into an end state

World state vs Search state (abstraction)
→ things that don't change or don't matter for the solution
  don't need to be in the search state

State Space Graph: mathematical representation of search problem
↳ nodes are world state, arrows are successors.
Search Tree: encodes possible decisions as a chonological tree

Tree Search: expand on tree nodes, order matters!
↳ uses <u>fringe</u>, <u>expansion</u>, and <u>exploration strategy</u>
↳ main question: which fringe nodes to explore?

Search Strategies:
  Depth-First: expand the deepest node first.
    ↳ expands some left prefix, $O(b^m)$ time for finite tree,
       only stores siblings from path to root → space $O(bm)$
    ↳ not complete if infinite tree, not optimal and only
       finds the "left most" solution
  Breadth-First: expand the shallowest node first
    ↳ expands all nodes above the shallowest solution
    ↳ time $O(b^s)$, space $O(b^s)$
    ↳ Complete, optimal iff all costs are 1.

Iterative Deepening: Combine DFS & BFS
 ↳ Run DFS with depth limit increasing iteratively.
 ↳ ordering is BFS-like, but saves memory
 ↳ the "last layer" out costs the previous iterations,
    so asymptotics isn't that bad.
 Uniform Cost: explore "cheap" paths first

# Informed Search

 Heuristic: a function that estimates how close a state is to a goal
 ex) Manhattan Distance, Euclidean Distance

 Greedy Search: <u>Only</u> look at the lowest heuristic
 ↳ Similar to DFS, but only considers future costs.
 A* Search: Combines UCS and Greedy ($f(n) = g(n) + h(n)$)
  Is A* optimal? → can fail if too pessimistic (trapped)
 Admissible heuristics: always underestimates cost to goal
 ↳ heuristic h is admissible if $0 \leq h(n) \leq h^*(n)$ where
    $h^*$ is the true cost function.

"If A is an optimal goal, B is a suboptimal goal, h is admissible, A will exit the fringe before B."

→ Proof: Imagine B is on the fringe.

Some ancestor of A (n) is in the fringe, too.

Then, n will be expanded before B

↳ $f(n) \leq f(A)$ by admissibility

$f(A) \leq f(B) \rightarrow \underline{f(n) \leq f(B)}$

With the same argument, all ancestors of A are expanded before B!

How to design admissible heuristics?

↳ often relaxing constraints work (ex) Manhattan)

Heuristics should be informative, but not too costly to compute

* maximum of admissible heuristic is still admissible.

Graph Search: don't expand the same state twice.

↳ however, if the newly computed cost is better than the previously stored cost, expand it again.

→ still optimal

\* If a heuristic is <u>consistent</u>, the first expansion
  ensures optimality for that node. (not covered)

# CSP

"What is the best assignment to variables?"

Standard Search: state is a black box, goal & successors can be anything

CSP: State is defined by <u>variables</u> $X$ with values from <u>domain</u> $D$.
  ↳ the goal test is a set of <u>constraints</u> of allowable assignments

ex) map coloring with out adjacent states sharing colors
  ↳ variables: regions $\{R_1, R_2, ..., R_n\}$. domain: colors $\{red, green, blue\}$
  constraint: $(R_1 \neq R_2)$ (implicit), $(R_1, R_2) \in \{(red, green), (red, blue), ...\}$
                                                   (explicit)
  solution: assignment satisfying all constraints

Binary CSP: all constraints take at most two variables

ex) N-Queens: variables $\{X_{ij} \mid i \in \#rows, j \in \#columns\}$
                                         → column
domain $\{0, 1\}$. constraints? $\forall i,j,k \ (X_{ij}, X_{ik}) \in \{(0,0), (0,1), (1,0)\}$
           → row
$\forall i,j,k \ (X_{ij}, X_{kj}) \in \{(0,0), (1,0), (0,1)\}$, diagonal constraints... ,
also need to include $\sum_{i,j} X_{ij} = N$ to prevent trivial solution of
                                                                        all zeroes.

ex2) Different N-Queens formulation: Variable Q, domain $\{1 \cdots N\}$
 ↳ assign a queen in each row and assign column #s.

# Varieties of CSPs:
 Discrete/Continuous, Finite/Infinite domains, Unary/Binary/
 Higher Order Constraints, Soft Constraints (Preferences)

# How to Solve CSPs: Standard Search Formulation?
 Start with empty assignment, successor to assign a single variable
 ↳ BFS would be ineffective since the solution lives in the deepest layer!
 ↳ DFS works, but naively checking solutions doesn't check for early fails

Backtracking: <u>One variable at a time</u>, <u>check constraints on the fly</u>
                        (variable ordering)                      (fail-on-violation)
 ↳ strategies: ① filtering (detecting failures early) ② ordering (advantageous order?)
Filtering: Forward Checking - cross off violations when adding
 a variable to an existing assignment → exit on impossible variable
 ↳ however, it doesn't fail until the actual impossible assignment.
also, it only enforces constraints on the variable just assigned.
→ Constraint Propagation: reason from constraint to constraint

Arc consistency: An arc $X \rightarrow Y$ is consistent iff for every x in the tail, there is some y in the head which could be assigned without violating a constraint.

If an arc is inconsistent, remove an assignment from the tail such that the arc is now consistent.
If a tail is removed, check all arcs that had it as head need to be updated.
Detect early failure if a variable has no possible assignments
↳ Runtime: $O(n^2 d^3)$, can be reduced to $O(n^2 d^2)$.
However, detecting all future problems is NP-Hard.

Arc consistency only enforces constraints on pairs → needs backtracking
→ k=2 is arc consistency
K-consistency: for any k nodes, any assignments to (k-1) of the nodes can be "extended" to the last node
↳ "extended": there exists a valid assignment given other assignments
Strong n-consistency ensures a solution to a CSP.
↳ all of 1~(n-1) are consistent

Ordering: How to pick the variable/assignment to try next?

Variable ordering: Minimum Remaining Values (MRV)

↳ try variables with the fewest elements left in its domain

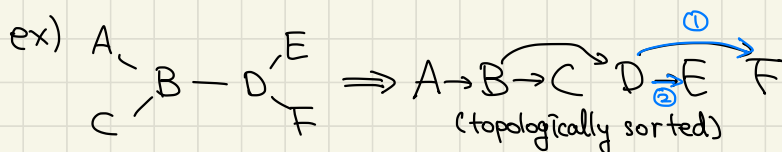"fail fast" ordering, tackle the hardest subproblems first

Value ordering: Least Constraining Value (LCV)

↳ given a choice of variable, choose the value that rules out the fewest values in remaining variables.

↳ being optimistic that the easiest path is correct

Reducing Structures: Disconnected graphs → independent subproblems

↳ If the constraint graph is a tree, CSP is solved in $O(cnd^2)$ time

ex)


$A \to B \to C \; D \rightleftarrows E \; F$
(topologically sorted)

Starting from the sink node, enforce arc consistency backwards  checking

→ edges with it as head are unchecked yet   nodes

↳ there are no multiple checks following a removal → $O(n \cdot d^2)$

then, just pick variables from the source. it is ensured

to give a solution (arc consistency ensures a valid assignment for all edges)

↳ no back tracking required! (infact, this is why NNs are DAGs)

If the CSP is not a tree, how about enforcing into one?

↳ Identify a cutset s.t. the remaining variables form a tree

⇒ improvement in runtime from exponential to (kind of) poly.

↳ try to cut out as little as possible when forming a tree

Iterative Improvement: start with some assignment, and improve inconsistent variables <u>locally</u> and greedily, suchthat the reassignment minimizes the # of remaining inconsistencies.

Difficulty of CSP: $R = \frac{\text{\# of constraints}}{\text{\# of variables}}$ → hard when not extreme

↳ R is big → almost trivial, R is small → large solution space

Local Search: no fringe, faster & effecient but incomplete & suboptimal.

# Adversarial Search

"How to choose actions in the presence of other agents?"

Types of games: Zero-sum (agents have opposite utilities), General games (independent utilities) → cooperation? indifference? Deterministic/Stochastic? # of players? Perfect information?

⇒ build a <u>strategy</u> to recommend an action based on current state

Adversarial Games: Deterministic, 2-player, zero sum, perfect information

- States: S (starts at $S_0$)
- Players: P = {MAX, MIN}
- Actions: A (depends on player/state)
- Transition Function: $S \times A \to S$
- Terminal Test: $S \to \{T, F\}$
- Terminal Utilities: $S \to R$ (reward = score)

A <u>value</u> of a state := best achievable outcome from that state →(=utility)

↳ for non-terminal states: $V(s') = \max\limits_{s \in \text{Successor}(s')} V(s)$

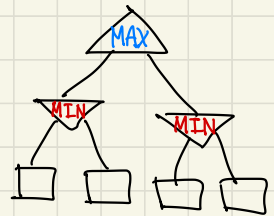A state is terminal when its value is (presumed to be) known

Minimax: when the adversary chooses, they try to <u>minimize</u>

↳ under opponent's control: $V(s') = \min\limits_{s \in \text{Successor}(s')} V(s)$

def <span style="color:blue">max</span><span style="color:red">min</span>-value(v):

   $V \leftarrow -(+) \infty$     → min and max are counterparts

   for each successor of v:

      $V = \frac{\text{max}}{\text{min}} (V, \frac{\text{min}}{\text{max}}\text{-value(successor)})$

return V

def value(v):

   call min-value or max-value depending whose turn it is

Minimax will be optimal against a perfect opponent. Otherwise?
↳ imperfect opponent → different modeling (Expectimax)

Efficiency: ≃ O(DFS) → time = $O(b^m)$, space = $O(bm)$
↳ not realistic in most game scenarios

Game Tree Pruning: Can we not traverse every single subtree?
↳ Intuition: once we see a value less than the current max
value, stop for that branch since the minimizer will return it
or something even worse (the maximizer never chooses that branch)
→ pass the current rolling "maximum min value" to min-value.
   $\underbrace{\text{maximum min value}}_{\alpha}$
   min-value stops exploring when $\underbrace{\text{its value}}_{n}$ drops below it ($n < \alpha$).
   (max-value version is symmetric) ⇒ Alpha-Beta Pruning

def $\frac{max}{min}$-value($V, \alpha, \beta$): ↳ no effect on minimax value for the root
   $V \leftarrow -(+) \infty$       however, intermediate values have
                                         different values
   for each successor of $v$:
        $V = \frac{max}{min} (V, \text{value}(\text{successor}, \alpha, \beta))$
        if $V \geq \beta (\leq \alpha)$, return $V$
        $\frac{\alpha}{\beta} = \frac{max}{min} (\frac{\alpha}{\beta}, V)$

Good child ordering improves pruning efficiency!
With "perfect ordering", time drops to $O(b^{m/2})$.
↪ this doubles solvable depth!

Depth-Limited Search: Just stop and appoximate after some depth
        ↦ similar to heuristics in A* search!
↪ an <u>evaluation function</u> guesses the utility of a state
Not guaranteed optimal play anymore, but use iterative
deepening for flexibility when computing
Eval(s) is usually a linear combination of game features
A bad evaluation function can cause an infinite loop...
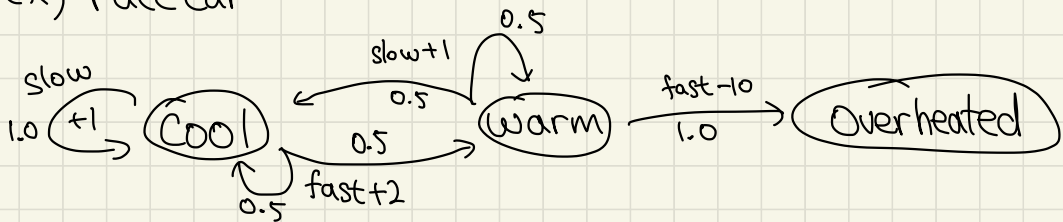
# Markov Decision Processes

- A set of states $s \in S$, A set of actions $a \in A$
- Transition function $T(s,a,s') := P(s'|s,a)$
- Reward function $R(s,a,s')$ (sometimes just $R(s)$ or $R(s')$)
- A start state, maybe a terminal state

"Markov"ness: action outcomes only depend on <u>current</u> state

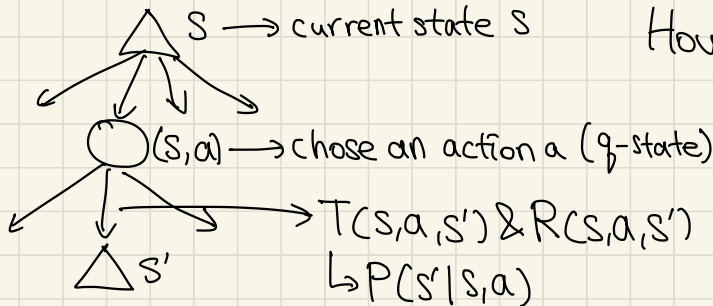For an MDP, we want a policy $\pi^* : S \rightarrow A$

↳ the optimal policy maximizes the expected utility

ex) race car



optimal policy: $\pi^*(\text{cool}) = \text{fast}$, $\pi^*(\text{warm}) = \text{slow}$, $\pi^*(\text{over}) = \text{end}$

MDPs can be formulated as a search tree (expectimax)



$S \rightarrow$ current state $s$

$(s,a) \rightarrow$ chose an action $a$ ($q$-state)

$T(s,a,s')$ & $R(s,a,s')$
↳ $P(s'|s,a)$

How to model rewards?
    now or later?
↳ decay rewards
    exponentially
$(1, \gamma, \gamma^2, \cdots | \gamma \in [0,1])$

Each round, the reward will be multiplied by discount factor $\gamma$.
↳ Sooner rewards have higher rewards than later ones
↳ It also helps rewards converge rather than approach infinity
$$U([r_0,...,r_\infty]) = \sum_{t=0}^{\infty} \gamma^t r_t \leq R_{max}/(1-\gamma) \text{ (bounded)}$$

How to solve MDPs? → think like expectimax, kind of
↳ states are repeated in subtrees → cache them!
↳ do depth-limited computation until changes are small
$V^*(s) :=$ expected utility of starting in s & acting optimally
$Q^*(s,a) :=$ expected utility of the q-state (s,a) & acting optimally
$\pi^*(s) :=$ the optimal action from state s
Bellman Equations (similar to expectimax): ∘ immediate rewards · discounted expected utility
$$V^*(s) = \max_a Q^*(s,a), \quad Q^*(s,a) = \sum_{s'} T(s,a,s')\left[\underbrace{R(s,a,s')}+\underbrace{\gamma V^*(s')}\right]$$
$$\Rightarrow \underline{V^*(s)} = \max_a \sum_{s'} T(s,a,s')\left[R(s,a,s')+\gamma \underwave{V^*(s')}\right] \rightarrow \text{how to solve this?}$$
Time Limited Values: $V_k(s) :=$ optimal value of s if the game
ends in k more steps (depth-k expectimax for s)
                                                      ┌ given all $V_k$ values
                                                      │    for all s'
$$V_0(s) \leftarrow 0, \quad V_{k+1}(s) \leftarrow \max_a \sum_{s'} T(s,a,s')\left[R(s,a,s') + \gamma \underline{V_k(s')}\right]$$
↳ repeat until convergence, which yields $V^*$ ($O(s^2A)$ each step)

Bellman Equation for $Q^*$? $Q^*_{(s,a)} = \sum_{s'} T(s,a,s')\left[R(s,a,s') + \gamma\left(\max_{a'} Q^*(s',a')\right)\right]$
↳ Leads to Q-value iteration algorithm for RL

But how do we get information about actions (policies)?
↳ Imagine we have the optimal values $V^*(s)$. How should we act?
Do a mini-expectimax: $\pi^*(s) = \arg\max_a \sum_{s'} T(s,a,s')\left[R(s,a,s') + \gamma V^*(s')\right]$
(argmax returns the "key value" of the largest value in a dict)
⟹ Policy extraction, since it gets optimal policies by values.
If we have optimal Q-values, $\pi^*(s) = \arg\max_a Q^*(s,a)$
↳ extracting policies are a lot easier with q-values!

Issues with value iteration: ① slow, ② "max" rarely changes
③ policies converge much faster than values
⟹ Policy-based methods can be more efficient!

Policy Evaluation: what are the consequences of a policy?
↳ rather than computing maximizer nodes, just do what policy tells
⟹ S will take $\pi(s)$ and land in q-state $(s, \pi(s))$
                                      → no more max, always choose $\pi(s)$
↳ $V^\pi(s) = \sum_{s'} T(s, \pi(s), s')\left[R(s, \pi(s), s') + \gamma V^\pi(s')\right]$

Turn $V^\pi(s)$ into iterations: $V_0^\pi(s) = 0$

$$V_{k+1}^\pi(s) \leftarrow \sum_{s'} T(s, \pi(s), s')[R(s, \pi(s), s') + \gamma V_k^\pi(s')]$$

↳ efficiency is $O(s^2)$, no more factor of $a$ when maxing

↳ without $\max_a$, this is just a set of linear equations!

Policy Iteration: Alternate between Policy $\overbrace{\text{evaluation}}^{\pi \rightarrow V^\pi}$ & $\overbrace{\text{extraction}}^{V \rightarrow \pi^V}$

① calculate utilities for some fixed policy until convergence

② update policy using one-step lookahead with calculated utilities

⇒ still optimal, could converge faster than value iteration

# Reinforcement Learning

Still assume MDP, looking for a policy $\pi(s)$

↳ What if we don't know T or R? (no measure of "good")

⇒ Must try out actions to learn from them!



→ black boxed initially, often acts probabilistically

Agent —— $a$ ——→ Environment (MDP)

Agent ←—— $s', R$ —— Environment (MDP)

Offline (MDP) vs. Online (RL)

Passive RL          vs. Active RL
  ↳ Model-Based RL      ↳ Exploration vs. Exploitation
  ↳ Model-Free RL

Model-Based Idea: Learn an approximate model, and
solve for values assuming it is correct
① Learn the distribution $\hat{T}(s,a,s')$ & $\hat{R}(s,a,s')$
② Solve the model with iteration
③ Run the learned policy, repeat if unsatisfactory

Model-Free: Don't know T and R, first learn $V(s)$.
Direct Evaluation: Just average all experiences afterwards
when that state was visited (no state dependency)
  ↳ Bellman updates don't work b/c they depend on T & R.
⇒ How do we take the weighted average without knowing them?
$$V_{k+1}^{\pi}(s) \leftarrow \sum_{s'} T(s, \pi(s), s')\left[R(s, \pi(s), s') + \gamma V_{k+1}^{\pi}(s')\right]$$
  ↳ take samples of outcomes s' and average them.
$$V_{k+1}^{\pi}(s) \leftarrow \frac{1}{n}\sum_{i}\left[R(s, \pi(s), s_i') + \gamma V_k^{\pi}(s_i')\right]$$
              ↳ samples will already be weighted by frequency

Temporal Difference Learning: learn from _every_ experience!
 keep a running average of $V(s)$ until $s$ is visited again
 $\rightarrow V^{\pi}(s) \leftarrow (1-\alpha)V^{\pi}(s) + \alpha \cdot \text{sample}$
 $\equiv V^{\pi}(s) \leftarrow V^{\pi}(s) + \alpha(\text{sample} - V^{\pi}(s))$
Exponential Moving Average: $\overline{X}_n = (1-\alpha)\overline{X}_{n-1} + \alpha \cdot X_n$
 ↳ recent samples are emphasized, past estimates are "forgotten"
 ⇒ still only does evaluation, we want new, better policies

Q-Learning: sample-based Q-value iteration
 $$Q(s,a) \leftarrow (1-\alpha)Q(s,a) + \alpha \cdot \underbrace{\text{sample}}_{\longrightarrow R(s,a,s') + \gamma \max_{a'} Q(s',a')}$$
 ⇒ converges to optimal policy (off-policy learning)
 ↳ as long as Q-value can converge (# of trials, lr decay, etc)
 how we choose to collect samples does not matter!

Active RL: how to act to collect data?
 ↳ The learner can choose what it wants to explore!
 Simplest scheme: $\varepsilon$-greedy (act randomly with probability $\varepsilon$)
 ↳ not really deliberate in exploring other states
 ⇒ somehow represent "novelty" to promote exploration!

Exploration Function: $f(u,n) = u + k/n$ ($n$: visit count, $u$: utility)

$\hookrightarrow$ high when $N$ is low!

$$Q(s,a) \xleftarrow[\alpha]{\text{(weighted update)}} R(s,a,s') + \gamma \max_{a'} \underline{f(Q(s',a'), N(s',a))}$$

Regret: how effectively did we learn? (optimally learn the optimal)

$\hookrightarrow$ less regret means faster learning

Feature Representation Formulas:

$$Q(s,a) = \vec{w} \cdot \vec{f}(s,a)$$

$$\text{diff} = \left[ R(s,a,s') + \gamma \max_{a'} Q(s',a') \right]$$
$$- Q(s,a)$$

$$w_i \leftarrow w_i + \alpha \cdot \text{diff} \cdot f_i(s,a)$$

# Probability

Observed variables (evidence): what the agent knows
Unobserved variables: agent needs to reason about these
Model: agent knows how to relate observed to unobserved

Random Variables: Aspect of the world we might have <u>uncertainty</u>
↳ each RV has a domain, discrete, boolean, continuous, tuples, etc.
Probability Distribution: Assigns each value of a RV a <u>probability</u>
↳ $P(X=v)$ denotes the probability $X$ takes on value $v$.
⇒ $\forall x, P(X=x) \geq 0, \sum_x P(X=x) = 1$ (basic rules for PD)
Joint Distribution: Probability of set of RVs, $P(x_1, x_2, \ldots, x_n)$
↳ the size of JD grows exponentially as variables increase
Events: Set of possible outcomes, $P(E) = \sum_{(x_1,\ldots,x_n) \in E} P(x_1, \ldots, x_n)$
↳ E acts like a filter for which JD we are interested in
Marginal Distribution: Collapsed rows by <u>eliminating</u> RVs in JD
↳ Acts as if we have no knowledge of the eliminated RV
↳ $P(X_1=x_1) = \sum_{x_2} P(X_1=x_1, X_2=x_2)$ (sum up all possible $X_2$ over $X_1$)

Conditionals: $P(a|b) = \frac{P(a,b)}{P(b)}$ ($P(a)$ given that b already holds)

↳ simple relation between <u>joint</u> and <u>conditional</u> probability

⇒ $P(b)$ can generally be found by marginalization over b

Conditional Distribution: PD over some variables when others are fixed

↳ Acts like taking a subset of the JD then renormalizing probabilities

Probabilistic Inference: compute a desired probability based on others

↳ generally compute conditionals, new evidence cause beliefs to be <u>updated</u>

Inference by Enumeration: $\overbrace{E_1,...,E_k}^{\text{evidence}} = e_1,...,e_k, \overbrace{Q}^{\text{query}}, \overbrace{H_1,...,H_r}^{\text{hidden}} \in X_{(1-n)}$

⇒ $P(Q|e_1,...,e_k)$ (observed $e_1,...,e_k$, then what is $P(Q)$?)

1) Select entries consistent with evidence

2) Sum out H to get JD of E and Q,   3) Normalize

↳ Leads to runtime & space complexity $O(d^n)$ (Ineffecient!)

Product Rule: <u>$P(x,y) = P(y)P(x|y)$</u> (derived from $P(x|y) = \frac{P(x,y)}{P(y)}$)

↓

Chain Rule: <u>$P(x_1,x_2,...,x_n) = \prod_i^n P(x_i|x_1,...,x_{i-1})$</u>

↳ $P(x_1,x_2,x_3) = P(x_1)P(x_2|x_1)P(x_3|x_2,x_1) = \cancel{P(x_1)} \frac{P(x_2,x_1)}{\cancel{P(x_1)}} \cdot \frac{P(x_3,x_2,x_1)}{\cancel{P(x_2,x_1)}}$

Bayes Rule: $P(x,y) = P(x)P(y|x) = P(y)(x|y) \rightarrow P(x|y) = \frac{P(x)P(y|x)}{P(y)}$

↳ useful for "flipping" probabilities when finding one is easier than other

↓ $P(\text{effect}|\text{cause})$     ↓ $P(\text{cause}|\text{effect})$

ex) M: meningitis, S: stiff neck. $P(+m) = 0.0001$, $P(+s|+m) = 0.8$,
$P(+s|-m) = 0.01$. What is $P(+m|+s)$?

↳ $P(+m|+s) = \dfrac{P(+s|+m) \, P(+m)}{P(+s) \rightarrow ?} = \dfrac{P(+s|+m) \, P(+m)}{P(+s,+m) + P(+s,-m)}$ (marginals)

$= \dfrac{P(+s|+m) \, P(+m)}{P(+s|+m) \, P(+m) + P(+s|-m) \, P(-m)}$ (product rule) $= \dfrac{0.8 \cdot 0.0001}{0.8 \cdot 0.0001 + 0.01 \cdot 0.9999} \approx \underline{0.008}$

# Bayes Nets

Independence: X,Y are independent if $\forall x,y \in (X,Y), \underline{P(x,y) = P(x)P(y)}$

↳ also implies $\forall x,y$, $\underline{P(x|y) = P(x)}$ (y reveals nothing about x)

⇒ Independence is a modeling assumption! Empirical JD are "close".

Conditional Independence: Independent when a third variable is observed
X is cond. ind. of Y given Z iff:
$\forall x,y,z$: $P(x,y|z) = P(x|z) \, P(y|z)$, $P(x|z,y) = P(x|z)$

Decomposition of Chain Rule: $P(x_1, x_2, x_3) = P(x_1) P(x_2|x_1) P(x_3|x_2 x_1)$
can be reduced to simpler structures → $P(x_3|x_2,x_1) = P(x_3|x_1)$

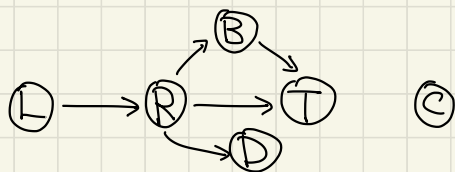Bayes' Nets: describing complex JD using local conditionals

Graphical Models: nodes → variables, arcs → interactions

ex) N independent coin flips: $\bigcirc\!\!\!\!x_1$  $\bigcirc\!\!\!\!x_2$ $\cdots$ $\bigcirc\!\!\!\!x_N$

ex) Traffic: R(raining), T(traffic)  $\bigcirc\!\!\!\!R \longrightarrow \bigcirc\!\!\!\!T$

ex) Traffic II: R,T, L(low pressure),

D(roof drips), B (ballgame), C (cavity)



Semantics: DAG topology + $\overset{\text{→this is a table (CPT)}}{\text{conditional } P(x|a_1,\ldots,a_n)}$ where

$a_i \to X$ is an edge in the DAG. $P(x_1,\ldots,x_n) = \prod_{i=1}^{n} P(x_i \mid \text{Parents}(x_i))$

ex) (Cavity) $\_P(cavity)$              $P(+cav, +catch, -tooth)$

$\searrow P(toothache \mid cavity)$    $= P(+cav)\cdot P(+catch \mid +cav)$

(Toothache)   (catch) $P(catch \mid cavity)$    $\cdot P(-tooth \mid +cav)$

Why is this true? $P(x_1,\ldots,x_n) = \prod_{i=1}^{n} \overline{P(x_i \mid x_1,\ldots,x_{i-1})} \; P(x_i \mid \text{Parents}(x_i))$

* Not all JD can be represented from BN!   $\searrow$ this is a core assumption
                                              of the world modeling

ex) N independent coin flips: $P(h,h,t,h) = \prod_{i=1}^{n} P(x_i \mid \cancel{\text{Parents}(x_i)}) = 0.5^4$

ex) Traffic: $\bigcirc\!\!\!\!R \longrightarrow \bigcirc\!\!\!\!T$  $P(+r,-t) = P(+r \mid \emptyset) \cdot P(-t \mid +r) = P(+r) P(-t \mid r)$

Reverse Causality? $\bigcirc\!\!\!\!T \to \bigcirc\!\!\!\!R$ still possible to reconstruct the JD!

$\hookrightarrow$ Direction of the edges do not mean direction of causality!

* Topology really encodes <u>conditional probabilities</u>

Size of the BN: N boolean variables → $2^N$ entries of JD
↳ N node BN with ≤ k parents → $N \cdot 2^{(k+1)}$ entries of JD
⟹ if k≪N, $N \cdot 2^{(k+1)} < 2^N$ → faster local CPTs & queries!

# Bayes' Nets: Independence

$$F \rightarrow \textcircled{S} \rightarrow A$$

ex) Alarm ⊥⊥ Fire | Smoke → Alarm doesn't care about source of smoke
BN often give rise to additional conditional independence.
ex) $\textcircled{X} \rightarrow \textcircled{Y} \rightarrow \textcircled{Z} \rightarrow \textcircled{W}$ : z⊥⊥x|y, w⊥⊥x,y|z → w⊥⊥x|y ?? how

D-Seperation: Algorithm for determining conditional independence from graphs
↳ study properties of triples, then compose them into complex paths

1) Causual Chains: $\textcircled{X} \rightarrow \textcircled{Y} \rightarrow \textcircled{Z}$  $P(x,y,z) = P(x)P(y|x)P(z|y)$
↳ Z⊥⊥X is false, however Z⊥⊥X|Y is true! $(P(z|x,y) = P(z|y))$
   (not guaranteed)

2) Common Cause: $\textcircled{X} \xleftarrow{\quad} \textcircled{Y} \xrightarrow{\quad} \textcircled{Z}$  $P(x,y,z) = P(y)P(x|y)P(z|y)$
↳ Z⊥⊥X is false, however Z⊥⊥X|Y is true! $(P(z|x,y) = P(z|y))$
   (not guaranteed)

3) Common Effect: $\textcircled{X} \searrow \textcircled{Z} \swarrow \textcircled{Y}$  $P(x,y,z) = P(x)P(y)P(z|x,y)$
↳ X⊥⊥Y is true, however X⊥⊥Y|Z is false!! (it is likely that
   (guaranteed)
one is actually contributing to Z, which decreases likelihood of the other)

|  | Active | Inactive |
|---|---|---|
| causal | $X \to Z \to Y$ | $X \to \boxed{Z} \to Y$   $X \perp\!\!\!\perp Y \mid Z$ |
| common cause | $X \leftarrow Z \to Y$ | $X \leftarrow \boxed{Z} \to Y$   $X \perp\!\!\!\perp Y \mid Z$ |
| common effect (v-structure) | $X \to \boxed{Z} \leftarrow Y$    $X \to Z \leftarrow Y$, $Z \to \boxed{W}$ | $X \to Z \leftarrow Y$   $X \perp\!\!\!\perp Y$ |

General Case: entire graph is just repetition of the three canonical cases!

↳ All it takes to block a path is <u>a single active segment</u>

ex) $A-B-C-D-E \to A-B-C, B-C-D, C-D-E$

If there are <u>any</u> path $X \rightsquigarrow Y$ that is active, <u>not D-seperated</u>.

\* $X \perp\!\!\!\perp Y \mid \{Z\}$ is guaranteed iff. $X$ and $Y$ are <u>D-separated given $\{Z\}$</u>.

↳ This does not imply anything about $X \perp\!\!\!\perp Y \mid Z$ when $X$ and $Y$ aren't D-seperated, only that it's <u>not</u> guaranteed!

# Bayes' Nets: Inference

Inference: Calculating some useful quantities from a JD

  ex) Posterior: $P(Q|E_1=e_1,\ldots,E_k=e_k)$

    Most likely: $\underset{q}{\text{argmax}}\ P(Q=q|E_1=e_1,\ldots,E_k=e_k)$

Inference by Enumeration is slow b/c it expands to a full JD

Variable Elimination can take short cuts when marginalizing.

Factors: ① $P(X,Y) \rightarrow$ sums to 1 ② $P(x,Y) \rightarrow$ sums to $P(x)$

③ $P(Y|x) \rightarrow$ sums to 1 ④ $P(Y|X) \rightarrow$ sums to $|X|$

⑤ $P(y|X) \rightarrow$ sums to... unknown! In general, $P(Y_1\ldots Y_N|X_1\ldots X_M)$

  has a dimension equal to # of unassigned variables

Enumeration: $\sum_t \sum_r P(L|t)P(r)P(t|r)$ vs Elimination: $\sum_t P(L|t) \sum_r P(r)P(t|r)$

If we have evidence, start with consistent entries only.

General VE procedure: $P(Q|E_1=e_1,\ldots,E_k=e_k)$

  While $\exists\ H_i$: Join all factors mentioning $H_i$, then eliminate $H_i$.

  Finally, normalize the JD to match the original query.

  ↳ basically reordering to lessen redundant multiplications, worst case

    exponential runtime w.r.t. size of the BN.

# Bayes' Nets: Sampling

Sampling is like repeated simulation.

Basic Idea: Draw $N$ samples from sampling distribution $S$.

Compute an approximate posterior probability. Show that this converges to the true probability $P$ as $N$ grows.

Step 1) $u \leftarrow$ uniform$(0,1)$ (kind of given)

Step 2) Convert $u$ into an outcome based on subintervals in $[0,1)$

Prior Sampling: Naively repeat sampling from start to finish

for $i = 1 \ldots n$: Sample $X_i$ from $P(X_i | \text{Parents}(X_i))$

Rejection Sampling: Only sample those that are absolutely needed

for $i = 1 \ldots n$: Sample $X_i$ from $P(X_i | \text{Parents}(X_i))$

if $X_i$ not consistent with evidence: reject & return early

↳ Rejects a LOT of samples, and evidence is not utilized.

Likelihood Weighting: what if we just force the evidence?

↳ just do it, but keep track of the likelihood that it ACTUALLY happens with a weight factor

$w \leftarrow 1.0$, for $i = 1 \ldots n$:

    if $X_i$ is an evidence variable:

        $X_i \leftarrow$ observation $x_i$ for $X_i$

        Set $w \leftarrow w \times P(X_i | \text{Parents}(X_i))$

    else:       $\longrightarrow$ basically means "this is equivalent to $w$ # of samples, where $w \in [0,1)$"

        Sample $X_i$ from $P(X_i | \text{Parents}(X_i))$

$\hookrightarrow$ Pretty good, just that it ignores evidence that comes later

Gibbs Sampling: Kind of like local search, perturb one observation

1) Fix evidence 2) Initialize all other variables

3) Repeat: Choose a non-evidence variable X.

    Resample X from $P(X | \text{all other variables})$

$\Rightarrow P(X | \text{all other variables})$ is very efficient due to cancellation

   with BN assumptions

# Decision Network

Bayes' Nets, but with additional types of nodes!
- Action Node (some domain, agent's choice) ▭
- Utility Node (based on its parents' outcomes) ◇

Goal: Maximize expected utility, given the evidence!
Action Selection: 1) Instantiate all evidence.
2) Set action in every way  3) Calculate posteriors
4) Calculate expected utility  5) Choose maximizing action

Almost looks like expectimax/MDP, but with BN distribution
* MEU <u>can</u> decrease with additional information, but it
doesn't mean that we are less happy, it just means that
the initial assumptions were inaccurate descriptions of reality.
$$\underline{MEU(E=e) = \max_{a} \sum_{s} P(s|e) U(s,a)} \text{ . (same for multiple evid.)}$$
Value of Information: compute the value of acquiring evidence
↳ Value := expected gain in MEU with new evidence
✳ $\underline{MEU(E') = \sum_{e'} P(E'=e') MEU(E'=e')}, \underline{VPI(E') = MEU(E') - MEU(\emptyset)}$

VPI Properties:

1) Nonnegativity: $\forall E', e$, $VPI(E'|e) \geq \emptyset$

2) Nonadditivity: $VPI(E_j, E_k|e) \neq VPI(E_j|e) + VPI(E_k|e)$

3) Order-independent: $VPI(E_j, E_k|e) = VPI(E_j|e) + VPI(E_k|E_j, e)$
$$= VPI(E_k|e) + VPI(E_j|E_k, e)$$

\* If Parents $(U) \perp\!\!\!\perp Z |$ Current Evidence, then $VPI(Z| Curr.Evi) = \emptyset$

POMDP: MDP, but states update their probabilities over time

↳ Solve using truncated expectimax to approximate utilities

# Hidden Markov Models

"What if the state of the world evolves over time?"

Markov Model: $P(x_1)$, $P(X_t | X_{t-1})$ ← Same for all (stationary)

↳ past & future independent of present, only dependent on previous

$$P(x_t) = \sum_{X_{t-1}} P(x_{t-1}, x_t) = \sum_{X_{t-1}} P(x_t|x_{t-1}) P(x_{t-1}) \to \text{converges as } t\to\infty !$$
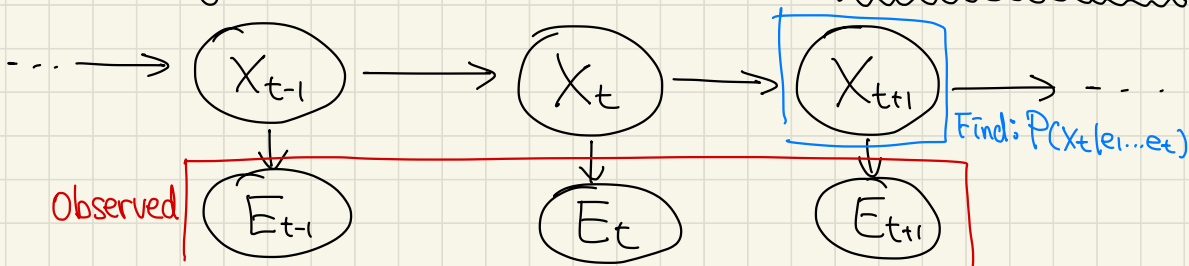
Stationary Distribution: $\underline{P_\infty(X) = P_{\infty+1}(X) = \sum_x P(X|x) P_\infty(x)}$

↳ This can be solved as a system of linear equations!

However, Markov models are generally not good modeling of reality

Hidden Markov Models (HMM): observe outputs at every time step!
↳ defined by: Initial $P(X_1)$, Transitions $P(X_t | X_{t-1})$, Emissions $P(E_t | X_t)$



Observed

Find: $P(X_t | e_1 \ldots e_t)$

Independence Properties: 1) $X_{t+1}$ is only dependent on $X_t$.
  2) Current observation is independent of all else given the current state.
  * It is not the case that evidences are always independent!


Filtering: Tracking and updating $B_t(X) = P_t(X_t | e_1 \ldots e_t)$ over time
↳ idea: start at $P(X_1)$ and derive $B_t(X)$ using $B_{t-1}(X)$
Two steps: Passage of Time & Observation (Incomplete & Complete)
↳ Passage of Time: $B_t(X) = P_t(X | e_{1:t}) \Rightarrow P_t(X_{t+1} | e_{1:t}) = \sum_{X_t} P(X_{t+1} | \underbrace{X_t, e_{1:t}}_{\text{actually irrelevant!}}) P(X_t | e_{1:t})$

$= \sum_{X_t} P(X_{t+1} | X_t) P(X_t | e_{1:t}) \Rightarrow \underline{B'(X_{t+1}) = \sum_{X_t} P(X' | X_t) B(X_t)}$

↳ Observation: $B'(X_{t+1}) = P(X_{t+1} | e_{1:t+1}) \Rightarrow P(X_{t+1} | e_{1:t+1}) = P(X_{t+1}, e_{t+1} | e_{1:t}) / \underset{\text{constant}}{P(e_{t+1} | e_{1:t})}$

$\propto_{X_{t+1}} P(X_{t+1}, e_{t+1} | e_{1:t}) = P(e_{t+1} | \underset{\text{irrelevant!}}{e_{1:t}}, X_{t+1}) P(X_{t+1} | e_{1:t}) = P(e_{t+1} | X_{t+1}) P(X_{t+1} | e_{1:t})$

$\Rightarrow \underline{B(X_{t+1}) \propto_{X_{t+1}} P(e_{t+1} | X_{t+1}) B'(X_{t+1})}$ ("reweighting" beliefs after observing)
↳ need renormalization after derivation!

Forward Algorithm: $P(x_t | e_{1:t}) \propto_{x_t} P(e_t | x_t) \sum_{x_{t-1}} P(x_t | x_{t-1}) P(x_{t-1}, e_{1:t-1})$
→ How do we deal with large state spaces?

Particle Filtering: Approximate Inference for Markov models
↳ Representation of $P(x)$ is a list of $N$ samples (particles)

Passage of Time: $x' = \text{sample}(P(x'|x))$ (generate the next step)
Observation: $w(x) = P(e|x)$, $B(x) \propto \overbrace{P(e|x)}^{w(x)} B'(x)$ (downweight w.r.t. likelihood)
Resample: Choose new samples based on $B(x)$'s distribution ($\approx$ renormalizing)

Dynamic Bayes' Nets*: Multiple Markov / Observation nodes in BN!

## Machine Learning: Naïve Bayes

"How to acquire a model from data/experience"

Types of Problems: Supervised, Reinforcement, Unsupervised
  ↗ labels          ↗ reward func.        ↗ no labels, just features

Supervised Learning $\Longrightarrow$ Classification: Discrete domains
                      $\longrightarrow$ Regression: Real-valued domains

Classification: Dataset $(x, y) \xrightarrow{\text{extraction}}$ Features $\xrightarrow{\text{ML}}$ Predict $y$
↳ ML learns patterns between features and labels from data!

ex) Spam Filter: Dataset (Email, {spam, ham}) → predict spams!
↳ What features do we want to look at? words (FREE), symbols ($),...

ex) Digit Recognition: Dataset (Pixel grid, {0,...,9}) → predict digits!
↳ Features: Pixel $(x,y) = $ On/Off, shape patterns (components, loops,...)

Model-Based Classification: Build a BN where both label and features
are RVs. Instantiate any observed variables, and find distribution of y.
Naïve Bayes: All features ($F_i$) are independent effects of the label ($Y$).
⇒ $P(Y)$: Prior. $P(F_i|Y)$: Probability of feature, given the label.

Naïve Bayes for Digits: One feature $F_{ij}$ for every pixel grid position $(i,j)$
→ $P(Y)$ (likelihood of every digit), $P(F_{ij}|Y)$ (on/off when the label is y)

Naïve Bayes for Text: $W_i$ is the word at position i ($W_i \in$ {Dictionary})
Moreover, each $P(W_i|Y)$ is assumed to be the same. → identically distributed
↳ This assumption reduces the # of parameters, also generalizes better!
↳ However, it will be insensitive to word ordering! (design choice)
→ $P(Y)$ (spam/ham), $P(W|$ {spam, ham}) (likelihood of word given the
    type of email)

In general, the joint probability will be $P(Y, F_1, ..., F_n) = P(Y) \prod_{i=1}^{n} P(F_i | Y)$.

↪ total # of parameters is linear w.r.t. $n$!

⇒ Computing $P(Y | F_1, ..., F_n)$ is just inference in BN.

↪ Inference by Enumeration: $P(Y | F_1, ..., F_n) \propto P(Y, F_1, ..., F_n) = P(Y) \prod_{i=1}^{n} P(F_i | Y)$.

⇒ $P(y_j) \prod_{i=1}^{n} P(F_i | y_j)$, then normalize to get $P(y_j | f_1, ..., f_n)$.

We also need to estimate the CPTs → Let $\theta$ denote all parameters! (trainable)

Parameter Estimation: Empirically learn using training data    $P(Data | \theta)$

↪ Maximum Likelihood: choose $\theta$ that <u>maximizes the probability of data</u>!

→ solve $\underset{\theta}{argmax}(f(\theta))$ where $f(\theta)$ is the probability of data happening

Useful fact: $\underset{\theta}{argmax} f(\theta) = \underset{\theta}{argmax} \ln(f(\theta))$ (easier analytic solution)
with differentiation

↪ For Naive Bayes, $P(y) = \frac{\# \text{ of } y}{total}$, $P(f|y) = \frac{\# \text{ of } f \text{ AND } y}{\# \text{ of } y}$

Empirical Risk Minimization: we want models to perform well on unseen data

↪ More training data, or regularize model complexity

However, training data could misrepresent the true distribution!

In general, we don't want to assign $\emptyset$ probabilities for uncertain $\theta$s.

↪ need smoothing or regularization

Laplace Smoothing: $P_{LAP_k}(X) = \frac{C(x) + k}{\sum_x [C(x) + k]} = \frac{C(x) + k}{\underbrace{N + k|X|}}$

↳ intuitively, act as if we observed $k$ more events of each outcome

Tuning: find the optimal smoothing value $k$ via held-out dataset
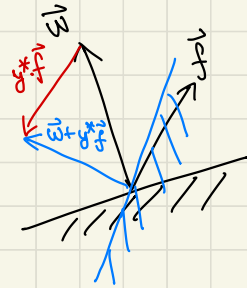
# Perceptrons

Binary Classifier: $activation_w(x) = sgn\left(\sum_i w_i \cdot f_i(x)\right) = sgn(\vec{w} \cdot \vec{f}(x))$

↳ dot product signifies the correlation between weight & feature

In the feature vector space, data are points, and weight vectors are hyperplanes.

⇒ we need to learn the weight vector from data!

Weight Updates: $y = \begin{cases} +1 & \text{if } \vec{w} \cdot \vec{f}(x) \geq 0 \\ -1 & \text{if } \vec{w} \cdot \vec{f}(x) < 0. \end{cases}$

update if $y$ is wrong$(y \neq y^*)$, $w \leftarrow w + \overbrace{y^*}^{\text{correct label}} \cdot f$

↳ Intuitively, we are shifting the hyperplane to reflect observed data

Multiclass Decision: $\vec{w}_y$ for each class, $y = \underset{y}{argmax} \; \vec{w}_y \cdot \vec{f}(x)$

↳ update $w_y = w_y - f(x)$ for wrong answer, $w_{y^*} = w_{y^*} + f(x)$ for correct answer
   (if $y \neq y^*$)

If the data are perfectly seperable, the perceptron will converge

↳ However, it might have problems if not. (thrashing, suboptimal)
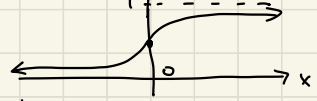
# Logistic Regression

Non-Seperable Data: Any linear boundary will make at least one mistake

↳ interpret the line as a probabilistic decision (50:50)

Perceptron scoring: $z = \vec{w} \cdot \vec{f}(x) \to$ want $Pr \to 1$ if positive, $\to 0$ if negative

↳ Sigmoid: $\emptyset(z) = \frac{1}{1+e^{-z}}$ achieves this behavior!

$\Rightarrow Pr(y=1|x,w) = \frac{1}{1+e^{-(w \cdot f(x))}}, Pr(y=-1|x,w) = 1 - \frac{1}{1+e^{-(w \cdot f(x))}}$

↳ increasing $w$ will make the boundary sharper (best $w$?)

MLE of Log. Reg.: Log likelihood $= \sum_i \log Pr(y^{(i)}|x^{(i)}, w)$

$\Rightarrow$ a probilistic intepretation can also improve the seperable case!

Multiclass Log. Reg: $\forall z_i \in \{z_1, ..., z_n\}, softmax(z_i) := \frac{e^{z_i}}{\sum_j e^{z_j}}$

↳ transforming original activations into "softmax" activations

$\Rightarrow P(y|x,w) = \frac{e^{w_y \cdot f(x)}}{\sum_{y'} e^{w_{y'} \cdot f(x)}}$ for perceptron interpretation

Deep Neural Network: Cascading logistic regression of multiple layers

↳ a hidden layer $h_i^{(\ell)} := \emptyset(\vec{w}_i^{(\ell)} \cdot \vec{h}^{(\ell-1)}) = \emptyset(\sum_j w_{ji}^{(\ell)} \cdot h_j^{(\ell-1)})$

$\to$ in matrix form, $\vec{h}^{(\ell)} = \emptyset(W^{(\ell)} \times \vec{h}^{(\ell-1)})$ where $W^{(\ell)}$ is the matrix $\begin{bmatrix} - \vec{w}_1^{(\ell)} - \\ \vdots \\ - \vec{w}_n^{(\ell)} - \end{bmatrix}$

↳ still uses MLE, but now it is iterative