agent function maps from precept history to actions $f: P^* \to A$ — not every function can be implemented, some program

agent program $I$ runs on some machine $m$ to implement $f$: $f = Agent(I, m)$

**environments agent measures**
- partially observ.
- stochastic
- multi-agent
- static — has time to calculate
- continuous time — continuously w. controller

Memory
- prepare for contingency
- behave randomly

**Search problem consists of:**
- state space
- set of actions for each state
- transition model
- step cost function
- start state + goal test

**DFS**
- $b$ = branching factor
- $m$ = depth
- time: $O(b^m)$
- space $O(bm)$
- not optimal

**BFS**
- time: $O(b^s)$
- space: $O(b^s)$
- optimal if costs = 1

**UCS** — optimal
- solution cost = $C^*$
- arcs cost = $\epsilon$
- $O(b^{C^*/\epsilon})$ time
- $O(b^{C^*/\epsilon})$ space

Tree Search — exponential
graph search — quadratic
- required mem prop to its runtime

Max of admissible heuristics is admissible + better
- admissible: $heuristic(x) \leq actual\ cost(x)$
- consistent: $h(A) - h(C) \leq cost(A\ to\ C)$
- $h(A) \geq cost(A\ to\ C) + h(C)$

**Hill climbing** — just climb up

**local search** — path is irrelevant — just find the solution
- improve your current state
- constant space, works for online + offline search

**non-determinism**
- solutions are contingency plans
- use AND-OR search to find then

outcomes don't depend deterministically on some hidden state

**Local beam search**
- $k$ copies of local search w/ random starts
- for each iteration
  - generate all successors for k states
  - choose best k to be new states
- like evolution
- come over here — the grass is greener

**Simulated annealing**
- allow "bad" move depending on temp
- high temp = more bad moves allowed, gets out of local min
- gradually reduce temp
- finds global optimum w/ prob 1 if slow enough cooling sched

**Partial observability:**
- belief state — set of all environments the agent could be in,
- $n$ physical states: $2^n$ belief states
- Transition model
  - determ: $Result(b,a) = $ union of $Result(s,a)$ for each $s$ in $b$
  - nondeterm: $Results(b,a) = $ union of $Results(s,a)$ for each $s$ in $b$
  - step cost$(b,a,b') = $ stepcost$(s,a,s')$ for all $s$ in $b$

Sensorless blank state
- everything work the same

**Minimax**
- efficiency just like DFS
- we can't go all the way down so we have eval fns

**CSPs**
- constraint graphs
- each constraint relates 2 variables
- types of CSPs
  - Discrete variables
    - $n$ variables, domain size $d$
    - $O(d^n)$ complete assignments
  - continuous variables
  - LP

**Alpha-beta pruning**
- have eval-time available at node $n'$
- $\alpha$ is best value max can get it $n$ becomes worse than $\alpha$ we can prune $n's$ other children
- doubles solvable depth

types of constraints
- unary constraints: one variable
- binary constraints: two variables

**Backtracking search**
- one variable at a time
- check constraints as you go
- DFS + these 2 improvements

ACS constraint arc — arc $X \to Y$ is consistent iff for every $X$ in tail there is some $Y$ in head which could legally be assigned

variable ordering for backtracking
- use minimum remaining value — variable w/ fewest values left in domain
- then do LCV (least constraining value)

$a \models b$ iff in every world where $a$ is true, $b$ is also true

Speedups for backtracking: ordering, filtering, structure

Successor state axiom: $X_t \iff [X_{t-1} \land \neg (Some\ Action_{t-1}\ made\ it\ false)] \lor [\neg X_{t-1} \land (Some\ Action_{t-1}\ made\ it\ true)]$

$P(a|b) = \frac{P(a,b)}{P(b)}$     $P(y)P(x|y) = P(x,y)$

chain rule: $P(x_1, x_2, x_3) = P(x_1)P(x_2|x_1)P(x_3|x_1,x_2)$

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

Independence: $X$ and $Y$ are ind. iff $\forall x,y\ P(x|y) = P(x)$ or $P(y|x) = P(y)$ $P(x,y) = P(x)P(y)$

$X$ is ind of $y$ given $z$ iff $\forall x,y,z\ P(x|y,z) = P(x|z)$    $\forall x,y,z\ P(x,y|z) = P(x|z)P(y|z)$

CPT = conditional prob table

Bayes nets: $n$ variables, max domain size $\leq d$, max parents = $k$, BN has size $O(n \cdot d^k)$

full joint distribution has size $O(d^n)$
$$P(x_1 \dots x_n) = \prod_i P(x_i | parents(x_i))$$

**Variable elimination:**
- have summations inward
- start w/ initial factors
- while still hidden variables
  - pick hidden var $H$
  - join all facts mentioning $H$
  - eliminate (sum out) $H$
- join all remaining facts + normalize

time is linear in network size if you eliminate in form leaves to root on poly tree

$B$ is ind of $E$
$B$ not ind of $E$ given $A$
$J$ ind of $M$ given $A$

**Markov models**
$P(X_t | X_{t-1})$ tells you how the state evolves over time
$P(X_0 \dots X_t) = P(X_0) \prod_t (P(X_t | X_{t-1}))$ — joint distribution

**Forward algorithm for markov models:**
$$P(X_t) = \sum_{x_{t-1}} P(X_{t-1} = x_{t-1}) P(X_t | X_{t-1} = x_{t-1})$$

Transition: $X_1 \to X_2 \to X_3$ — initial dist: $P(X_0)$
- initial model: $P(X_t | X_{t-1})$
- sensor model: $P(E_t | X_t)$

joint dist: $P(X_0, X_1 \dots X_T, E_T) = P(X_0) \prod_{t=1}^{T} P(X_t | X_{t-1}) P(E_t | X_t)$

more on back

**Stationary distribution:**
$P_\infty = P_{\infty+1} = T^T P_\infty$

example: $X_{t-1}$

| | $P(X_t \| X_{t-1})$ | |
|---|---|---|
| | sun | rain |
| sun | .9 | .1 |
| rain | .3 | .7 |

Solving for $P_\infty$
$$\begin{bmatrix} .9 & .3 \\ .1 & .7 \end{bmatrix}\begin{bmatrix} P \\ 1-P \end{bmatrix} = \begin{bmatrix} P \\ 1-P \end{bmatrix}$$
$.9P + .3(1-P) = P$

more Hmms:

Filtering: $P(X_t | e_{1:t})$ — input to decision process of rational agent

Prediction: $P(X_{t+k} | e_{1:t})$ $k > 0$ — evaluation of possible action sequences

Smoothing: $P(X_k | e_{1:t})$ $0 \le k < t$ — better estimate of past states

most likely explanation: $\arg\max_{X_{1:t}} P(X_{1:t} | e_{1:t})$

filtering alg: $P(X_{t+1} | e_{1:t+1}) = \alpha P(e_{t+1} | X_{t+1}) \sum_t P(X_t | e_{1:t}) P(X_{t+1} | X_t)$

DBN — repeat fixed Bayes net structure at each time
- every Hmm is a single variable DBN
- every discrete DBN is an Hmm
- DBNs have way fewer parameters than Hmms

Value of Information
- change in expected utility by looking at unavailable
- non negative
- non additive $VPI_e(E_i, E_j | e) \ne VPI(E_i | e) + VPI(E_j | e)$
- order independent $VPI_e(E_i, E_j | e) = VPI(E_j, E_i | e)$

MDPs are fully observable
but probabilistic search probs

$V^\pi(s_0) = P(S_1 | S_0, \pi(S_0)) \times P(S_2 | S_1, \pi(S_1))$ ← reward sequence for states + policies probabilities

- expected utility of $\pi$ in $S_0 = V^\pi(S_0)$
  Sum over all possible state sequences of
  (discounted sum of rewards)(prob of state sequences)

$Q^*(s,a)$ = expected utility of taking action $a$ in state $s$ and then acting optimally

$V^*(s) = \max_a Q^*(s,a)$ = expected utility of starting in $s$ and acting optimally

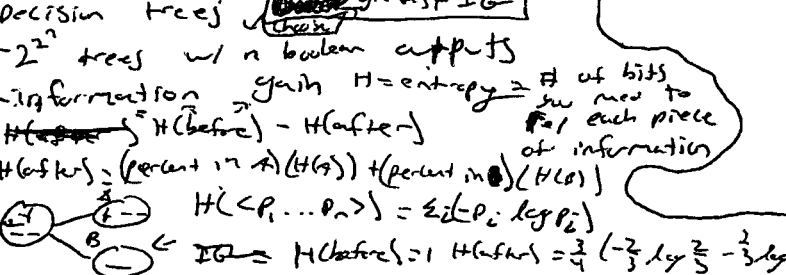So... $V^*(s) = \max_a \sum_s P(s'|a,s)[R(s,a,s') + \gamma V^*(s')]$
- Bellman equation

Value Iteration
Do Bellman update until you converge
(all updates smaller than $\frac{\epsilon(1-\gamma)}{\gamma}$)

$V_{k+1}(s) = \max_a \sum_s P(s'|a,s)[R(s,a,s') + \gamma V_k(s')]$

ie $V_{k+1} = B V_k$

we know it will converge
error is reduced by a factor of
at least $\gamma$ every iteration
— exponentially fast
— for policy iteration
  — do V.I
  but keep track of best policy

Q learning
sample = $R(s,a,s') + \gamma \max_{a'} Q(s',a')$
$Q(s,a) \leftarrow (1-\alpha) Q(s,a) + \alpha [\text{sample}]$
— converges to optimal policy
even if you are acting suboptimally

Decision trees ← choose greatest IG
- $2^{2^n}$ trees w/ n boolean inputs
- information gain $H$ = entropy = # of bits you need to for each piece of information
$H = H(before) - H(after)$
$H(after) = (\text{percent in } A)(H(A)) + (\text{percent in } B)(H(B))$
$H(\langle P_1 ... P_n \rangle) = \sum_i -P_i \log_2 P_i$
← $IG = H(before) = 1$ $H(after) = \frac{3}{4}(-\frac{2}{3} \log \frac{2}{3} - \frac{2}{3} \log \frac{2}{3}) + \frac{1}{4}(-1 \log(1))$

Viterbi alg (max)
- for each state at $t$, keep track of max prob of any path to it
$P(e_{t+1} | X_{t+1}) \max_{X_t} P(X_{t+1} | X_t) \cdots$

| $O(|X|^2 T)$ time | $O(|X| T)$ space | $O(|X|^T)$ paths |
|---|---|---|

$O(|states|^L)$ for time step

- Decision network = Bayes nets + action nodes (can not have parents, given) + utility node (depends on action + chance nodes)
$EU = \sum_w P(w|e,a) U(a,w)$ for all $a$
- return action w/ highest EU

Utility = $\sum_i P_i U(S_i)$
MEU = max utility

- Neural networks
- inputs $a_i$ come to neuron $j$ for neuron $i$
- each $a_i$ w/ a weight $w_{i,j}$
- total input $in_j = \sum_i w_{i,j} a_i$
- output $a_j = g(in_j) = g(w \cdot x)$   $g$ = activation function

- Perceptron learning (neural networks)
- if $y \neq thresh(x)$ (there is an error)
- if $w \cdot x < 0$ but output should be $y = 1$
  - false neg
  - increase weights on pos inputs, decrease on neg inputs
- if $w \cdot x > 0$ but output should be $y = 0$
  - false pos
  - decrease weights on pos, increase weights on neg
- $w = w + \alpha(y - h_w(x)) x$   $\alpha$ = learning rate

- data is lin sep if you can draw a line
- if it is — perceptron learning will converge to perfect
- if not perceptron learning will converge to minimal error solution as long as $\alpha$ is decreased correctly

MLP = multi layer perceptron
- feed forward net
  - at least one hidden layer
  - can rep any function
MLP's minimize squared error (loss function)

Bayesian learning
- predictors use likelihood weighted average over hypotheses
$P(X_{n+1} | x) = \sum_k P(X_{n+1} | x, h_k) P(h_k | x) = \sum_k P(X_{n+1} | h_k) P(h_k | x)$

We can express $V$ as a $Q$ as weighted linear functions
$V_w(s) = w_1 f_1(s) + ... + w_n f_n(s)$   (approx)
$Q_w(s,a) = w_1 f_1(s,a) + ... + w_n f_n(s,a)$
update: $w_i = w_i + \alpha [R(s,a,s') + \gamma \max_{a'} Q(s',a') - Q(s,a)] f_i(s,a)$

— this approximation may diverge!!

(· = dot product)