

MDPs

Rewards w/ discounted utility

$$U(I_0, \dots, \infty) = \sum_{t=0}^{\infty} \gamma^t r_t \leq \frac{R_{\max}}{1-\gamma}$$

Value Iteration

$$V_0(s) = 0$$

$$V_{k+1}(s) = \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_k(s')]$$

repeat until convergence $O(S^2 A)$

Bellman Equations

$$V^*(s) = \max_a Q^*(s, a) \text{ optimal val of a state}$$

$$Q^*(s, a) = \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^*(s')]$$

optimal value of (s, a) is expected val of an agent utility after taking a firm s and acting optimally

Policy Evaluation

$$V_0^\pi(s) = 0 \quad O(S^2) \text{ per iteration}$$

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

Policy Extraction

$$\pi^*(s) = \arg\max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^*(s')]$$

Policy Iteration

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

$$\text{improvement: } \pi_{k+1}(s) = \arg\max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_k^{\pi_k}(s')]$$

Probability

$$P(A|B) = \frac{P(A) \cdot P(B|A)}{P(B)} = \frac{P(A \cap B)}{P(B)}$$

$$\text{Joint: } P(A, B) = P(B)P(A|B)$$

$$P(x_1, x_2, \dots, x_n) = \prod_i P(x_i | x_1, \dots, x_{i-1})$$

Bayes Net

can calculate any prob using chain rule

Inference

Variable elimination

1. Join (multiply) factors involving x
2. Sum out x

factors are unnormalized

$$P(A|B) = \frac{P(A) * P(B|A)}{P(A) * P(B|A) + P(\bar{A}) * P(B|\bar{A})}$$

Backprop

$$h_1 = f_1(x) \quad h_2 = f_2(x) \quad z = h_1, h_2$$

$$\frac{\partial z}{\partial x} = h_2 \frac{\partial f_1}{\partial x} + h_1 \frac{\partial f_2}{\partial x}$$

Model-Based Learning

$$\text{Transitions: } \frac{\#(s, a, s')}{\#(s, a)}$$

$$\text{Reward: } R(s, a, s')$$

$$n^p_k = \frac{n!}{n-k!}$$

$$n^c_k = \frac{n!}{(n-k)! k!}$$

Model Free Learning

Passive Reinforcement

Direct Eval

- fix policy π
- track utility & count of transitions
- compute est. val w/ utility/visits

TD Learning

$$\text{sample} = R(s, \pi(s), s') + \gamma V^\pi(s')$$

$$\text{update: } V^\pi(s) = (1-\alpha)V^\pi(s) + \alpha \text{sample}$$

$$V^\pi(s) = V^\pi(s) + \alpha (\text{sample} - V^\pi(s))$$

Q-Learning (off-policy learning)

$$Q_{k+1}(s, a) = \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma \max_{a'} Q_k(s', a')]$$

Learn $Q(s, a)$ values as you go

$$\text{sample} = R(s, a, s') + \gamma \max_{a'} Q(s', a')$$

$$Q(s, a) = (1-\alpha)Q(s, a) + \alpha [\text{sample}]$$

Q-val Iteration

$$Q_{k+1}(s, a) = \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma \max_{a'} Q_k(s', a')]$$

Approximate Q-Learning

- feature based representation
- use V & Q as linear value functions

$$Q(s, a) = w_1 f_1(s, a) + w_2 f_2(s, a) + \dots + w_n f_n(s, a)$$

$$\text{Transitions} = (s, a, r, s')$$

$$\text{difference} = [r + \gamma \max_{a'} Q(s', a')] - Q(s, a)$$

$$Q(s, a) = Q(s, a) + \alpha [\text{difference}] \text{ Exact } Q's$$

$$w_i = w_i + \alpha [\text{difference}] f_i(s, a) \text{ Approx. } Q's$$

ϵ -greedy

- prob $0 \leq \epsilon \leq 1$ act randomly and explore w/ prob ϵ
- act w/ $(1-\epsilon)$ w/ curr policy
- large ϵ is mostly random

Exploration fn

$$\text{Modified } Q\text{-update: } Q(s, a) = \alpha R(s, a, s') + \gamma \max_{a'} f(Q(s', a'))$$

$$f(u, n) = u + \frac{k}{n} \text{ encouraged to explore states less visited}$$

$$N(s', a')$$

each softmax output is between 0 and 1
total sums to 1 (prob distribution)

Search Problems

- state space = all possible states
- successor fn: takes in state, action and computes cost and successor state
- start state: initial state
- goal test: test if state is goal

Manhattan distance = $|x_1 - x_2| + |y_1 - y_2|$

Admissibility: $\forall n, 0 \leq h(n) \leq h^*(n)$

$g(n)$: total backward cost comp. by UCS (one cost)

$h(n)$: estimated forward cost

$g(n) + h(n)$: total cost used by arc

consistency: $\forall A, C, h(A) - h(C) \leq \text{cost}(A, C)$

$H(n) = 0$ to be admissible & consistent

consistency \Rightarrow admissibility $h(A) \leq \text{cost}(A, C) + h(C)$

Dominance: $h_a \geq h_c$ if $\forall n, h_a(n) \geq h_c(n)$

Search:

DFS: explores deepest node, LIFO stack

complete: no, can get caught in cycles Time: worst case $O(b^m)$

optimal: no

BFS: selects shallowest node, FIFO queue

complete: yes, if sol'n exists

optimal: no, doesn't account for cost.

optimal if all edges = 1

UCS: selects lowest cost. Heap based

complete: yes, if sol'n exists

optimal: yes, as long as positive edges

Informed Search

Greedy: selects lowest heuristic node. same as UCS, but forward cost instead of backward

complete: no

optimal: no, bad heuristic

A*: selects lowest estimated total cost to goal. PQ using back & forward

ast for priority

complete: yes

Graph Search

add nodes to set, never expand a state twice

Games

deterministic: actions determined

stochastic: actions random

minimax: $O(b^m)$

Alpha-Beta Pruning

Time Comp: $O(b^{m/2})$

α - min β - max

CSP

consists of variables domains constraints

N var. domain size $d \rightarrow O(d^N)$ assignment

unary: only on 1 variable

binary: constraint b/w two var.

Backtracking search:

- fix var order & assign vals

- don't assign previously assigned values

- if no value, backtrack

Forward checking

- after assigning a value to a var, prune the domains of unassigned variables that share constraint w/ it.

along a path
if value never decreases

arc consistency

$X \rightarrow Y$ is consistent iff for every x in the tail there's some y which could be assigned w/o violating a constraint.

If x loses a value, neighbors of x need to be rechecked

- detects failure earlier than forward checking

$O(n^2 d^3)$ can be reduced to $O(n^2 d^2)$

MRV: choose variable w/ fewest valid remaining values

LCV: select value that prunes the least values from domains of unassigned vars

k-consistency

- guarantees for any set of k nodes in csp, an assignment to any $k-1$ nodes guarantees k th node will have one consistent value

$k=2 \equiv$ arc consistency

higher $k \rightarrow$ harder to compute

Tree-structured

- can be solved in $O(nd^2)$ from $O(d^n)$

- backward pass arc consistency

- forward pass assignment w/o backtrack

cutset conditioning

find smallest subset of vars in constraint graph w/ removal makes a tree.

$O((n-c)d^2)$ (remove c arcs)

$O(d^c (n-c)d^2)$ total time

Local search (iterative improvement)

- random assign vars

- selects var that violates most constraints and reset to value that violates the least constraints

- not optimal or complete

Logic

$\alpha \models \beta$ iff in every world where α is true, β is also true



To prove: model checking

Theorem proving

- search for a sequence of proof steps leading from α to β

$\alpha \Rightarrow \beta \equiv \neg \beta \Rightarrow \neg \alpha$

$\alpha \Rightarrow \beta \equiv \neg \alpha \vee \beta$

$\alpha \Leftrightarrow \beta \equiv \alpha \Rightarrow \beta \wedge \beta \Rightarrow \alpha$

$\neg(\alpha \wedge \beta) \equiv \neg \alpha \vee \neg \beta$

$\neg(\alpha \vee \beta) \equiv \neg \alpha \wedge \neg \beta$

$\alpha \wedge (\beta \vee \gamma) \equiv (\alpha \wedge \beta) \vee (\alpha \wedge \gamma)$

$\alpha \vee (\beta \wedge \gamma) \equiv (\alpha \vee \beta) \wedge (\alpha \vee \gamma)$

DPLL

Early Termination: return as soon as satisfied

1. Pure symbols: only shows up positive or neg.

can be assigned immediately

2. Unit clause: clause w/ just one literal or a disjunction w/ one literal and many

literals.

If KB contains A and $A \Rightarrow B$, we can infer B

If KB contains $A \wedge B$ we can infer A or B

If KB contains A and B we can infer $A \wedge B$

Probability

Chain: $P(A|B) = \frac{P(A)P(B|A)}{P(B)}$

Joint: $P(A, B) = P(A)P(B|A)$

Query var Evidence
 $P(Q_1 \dots Q_n | e_1 \dots e_n)$ Hidden vars

Inference by enumeration (compute everything)

Bayes Net

- acyclic graph of nodes node Parents
 - cond. dist. for each node $P(X_i | A_i, \dots, A_n)$
 - can calculate any prob using chain rule
- $$P(X_1, X_2, \dots, X_n) = \prod_{i=1}^n P(X_i | \text{parents}(X_i))$$

Inference

- Variable elimination
1. Join factors involving X
 2. Sum out X
- factors are unnormalized

Sampling

Rejection Sampling

reject samples that don't match evidence

Likelihood weighting

set evidence var and weight = 1
 multiply weight by sampled variables that affect evidence variables

- $P(T+C, +e)$, set $w_3=1$
- sample t
- $w_3 = w_3 \cdot P(+c|t)$...

Gibbs

- set evidence, then random assign rest
 - pick one, clear and resample w/ assigned vars
 - resample $P(-a|+c, +b)$
- X from $P(X | \text{markov-blanket}(x))$

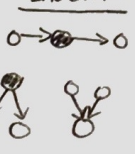
D-Separation

1. shade observed
2. Enumerate path X-Y
3. If all triples on path active, active
4. If no active path, X, Y are cond. indep.

Active



Inactive



Laplace Smoothing

$$P_{MLE}(x) = \frac{\text{count}(x)}{N}$$

$$P_{LAP, k}(x) = \frac{\text{count}(x) + k}{N + k|x|}$$

$$P_{LAP, k}(x|y) = \frac{\text{count}(x, y) + k}{\text{count}(y) + k|x|}$$

$$P_{LAP, \infty}(x) = \frac{1}{|x|}$$

Decision Networks

$$EU(a|e) = \sum_{x_1, \dots, x_n} P(x_1, \dots, x_n | e) U(a, x_1, \dots, x_n)$$

highest utility is MEU

$$Ex: EU(\text{leave} | \text{bad}) = \sum_w P(w | \text{bad}) U(\text{leave}, w)$$

$$EU(\text{take} | \text{bad}) = \sum_w P(w | \text{bad}) U(\text{take}, w)$$

$$MEU(F = \text{bad}) = \max_a EU(a | \text{bad})$$

VPI \rightarrow how much max expected utility gained from info

$$VPI(E|e) = MEU(e, E') - MEU(e)$$

$$= \sum_{e'} P(e'|e) MEU(e, e')$$

$$VPI(F) = MEU(F) - MEU(\emptyset)$$

Properties of VPI

nonnegativity: $\forall E', e \quad VPI(E'|e) \geq 0$

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

If we know both gives more info than just one

$$VPI(B, C) > VPI(B) + VPI(C)$$

If $B=C$ then $VPI(B, C) < VPI(B) + VPI(C)$
 one gives just as much info as other

Order-Independence

$$VPI(E_j, E_k | e) = VPI(E_j | e) + VPI(E_k | E_j) = VPI(E_k | e) + VPI(E_j | E_k)$$

HMM

Mini Forward algorithm

$$Pr(W_{t+1}) = \sum_{W_t} Pr(W_t, W_{t+1})$$

$$= \sum_{W_t} Pr(W_{t+1} | W_t) Pr(W_t)$$

Stationary: (same as mini forward)

$$P(W_{t+1}) = P(W_t)$$

$$= \sum_{W_t} P(W_{t+1} | W_t) P(W_t)$$

$$P(W_{\infty+1}) = P(W_{\infty})$$

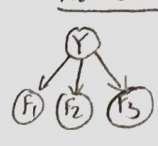
Vertex Alg

$$m_t[x_t] = \max_{x_{t-1}} P(x_{t-1}, x_t, e_{1:t})$$

$$= \max_{x_{t-1}} P(e_t | x_t) P(x_t | x_{t-1}) \max_{x_{t-2}} P(x_{t-1}, e_{1:t-1})$$

$$= P(e_t | x_t) \max_{x_{t-1}} P(x_t | x_{t-1}) m_{t-1}[x_{t-1}]$$

Naïve Bayes



$$\text{prediction}(f_1, \dots, f_n) = \underset{y}{\text{argmax}} P(y | F_1=f_1, \dots, F_n=f_n)$$

$$= \underset{y}{\text{argmax}} P(y=y) \prod_{i=1}^n P(F_i=f_i | y)$$

$$P(Y, F_1=f_1, \dots, F_n=f_n) = \begin{cases} P(Y=y, F_1=f_1, \dots, F_n=f_n) \\ P(Y=y_n, F_1=f_1, \dots, F_n=f_n) \end{cases}$$

cond. indep. evidence

Likelihood

$$\mathcal{L}(\theta) = P_\theta(x_1, \dots, x_n)$$

$$\mathcal{L}(\theta) = \prod_{i=1}^n P_\theta(x_i)$$

$$\frac{d}{d\theta} \mathcal{L}(\theta) = 0 \text{ (gradient)}$$

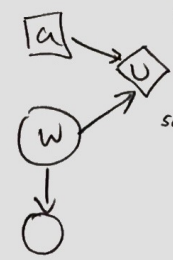
$\log(\mathcal{L}(\theta)) = \log \left(\prod_{i=1}^n P_\theta(x_i) \right)$
 if no theta, prob just a factor so can be taken out of derivative

Prior Sampling

$$SpS(x_1, \dots, x_n) = \prod_i P(x_i | \text{parents}(x_i))$$

$$= P(x_1, \dots, x_n)$$

sample procedure consistent



similar to Bayes Net

Markov Mule

- Transition model $P(x_t | x_{t-1})$

- Stationary Assumption = transition probs are same

- x_{t+1} is indep. x_0, \dots, x_{t-1} given x_t

$$P(x_0, \dots, x_T) = P(x_0) \prod_t P(x_t | x_{t-1})$$

Markov Blanket: consist of a variable's parents, children, children other parents

Belief dist.

$$B(w_t) = P(w_t | f_1, \dots, f_t)$$

belief \equiv prob w_t given evidence

$$= P(w_t | f_{1:t})$$

$$B'(w_t) = P(w_t | f_1, \dots, f_{t-1}) \text{ evidence up to } t-1$$

$$= P(w_t | f_{1:t-1})$$

Forward Alg

$$B'(w_{t+1}) = \sum_{w_t} P(w_{t+1} | w_t) B(w_t)$$

$$B(w_{t+1}) \propto Pr(f_{t+1} | w_{t+1}) B'(w_{t+1})$$

$$B(w_{t+1}) \propto Pr(f_{t+1} | w_{t+1}) \sum_{w_t} Pr(w_{t+1} | w_t) B(w_t)$$

Time elapse:

determine $B'(w_{t+1})$ from $B(w_t)$

Obs update:

$$B(w_{t+1}) \text{ from } B'(w_{t+1})$$

HMM

Initial distribution $P(x_0)$

Transition model: $P(x_t | x_{t-1})$

Sensor model: $P(e_t | x_t)$

Joint =

$$P(x_0, x_1, \dots, x_T, E_T)$$

$$= P(x_0) \prod_{t=1}^T P(x_t | x_{t-1}) P(e_t | x_t)$$

$$F_1 \perp\!\!\!\perp W_0, W_1$$

$$\forall i = 2, \dots, n \quad W_i \perp\!\!\!\perp \{W_0, \dots, W_{i-2}, F_1, \dots, F_{i-1}\} | W_{i-1}$$

$$\forall i = 2, \dots, n; F_i \perp\!\!\!\perp \{W_0, \dots, W_{i-1}, F_1, \dots, F_{i-1}\} | W_i$$

F_i are observations

Max likelihood geom dist

$P(F_i=1|Y=\text{ham})$ corresponds to counting the # of ham emails in which word i appears and dividing it by the total # of ham emails

Particle simulation

observe weights

resample: sample weights by using weights to create distributions for samples.

Assign particle per new distribution

Time-elapsed update: based off current particle

assignment, use prob of next timestep as distribution and assign based off a new sample

Observation update: calculate weight of each particle given new evidence

2. calc total weight for each state (if needed)

3. if sums to 0, reinitialize

4. else normalize and resample particles as above

Each particle is moved by sampling its next position from the transition model

$x' = \text{sample}(P(x'|x))$

Fix observation, similar to likelihood weighting

$w(x) = P(e|x)$

$f(x) \propto P(e|x)f'(x)$

Resample: rather than tracking weighted samples, we resample N times, choose from our weighted sample distribution equivalent to renormalizing the distribution

Perceptron

$$\text{activation } w(x) = \sum w_i f_i(x) = w^T \cdot f(x)$$

$$\text{classify}(x) : \begin{cases} + & \text{if activation}(x) > 0 \\ - & \text{if activation}(x) < 0 \end{cases}$$

Binary Perceptron Alg

1. Init all weights to 0

2. For each train sample, $f(x)$ and true class label

$y^* \in \{-1, +1\}$ do:

$$y = \text{classify}(x) = \begin{cases} +1 & \text{if act.} = w^T f(x) > 0 \\ -1 & \text{if act.} = w^T f(x) < 0 \end{cases}$$

compare predicted label to true label

if $y = y^*$, do nothing

else update

$$w \leftarrow w + y^* \cdot f(x)$$

Multi-class Perceptron

$$w_0 = [-2 \ 2 \ 1] \quad f(x) = [-2 \ 3 \ 1]$$

$$w_1 = [0 \ 3 \ 4]$$

$$w_2 = [1 \ 4 \ -2] \quad s_0 = 11 \quad s_1 = 13 \quad s_2 = 8$$

would predict s_1

- if s_1 not correct, subtract feature from s_1

- add feature to correct weight

POMDPs

State S , Actions A , Transition func $P(S'|S, a)$, Rewards R ,

Observations O , Obs func $P(o|S)$ (or $O(S, o)$)

Gradient Ascent

init w

α = learning rate

for iter = 1, 2, ...

$$w \leftarrow w + \alpha * \nabla g(w)$$

Sigmoid func

$$g(z) = \frac{1}{1 + e^{-z}}$$

↑ for descent

Hyperbolic Tangent

$$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

ReLU

$$g(z) = \max(0, z)$$

$$g'(z) = g(z)(1 - g(z))$$

$$g'(z) = 1 - g(z)^2$$

$$g'(z) = \begin{cases} 1, & z > 0 \\ 0, & \text{otherwise} \end{cases}$$

$$f(x) = g(h(x)) \quad f'(x) = g'(h(x)) h'(x)$$

Decision Tree

Split on node w must info gained

↳ all positive or all negative