

CS188 MT1 CHEATSHEET

Note 1: Search

- admissible: nonnegative and underestimate
- consistency: $h(a) - h(c) \leq \text{cost}(A, C)$
- admissibility guarantees optimal tree search
- consistency guarantees optimal tree/graph search
- must give $h(c) \geq 0$ in goal state

Note 2: Game Trees

- AlphaBetaPruning

def max-value(state, α, β):

return V

def min-value(state, α, β):
initialize(s, α, β)

$V = \min(V, \text{value}(s, \alpha, \beta))$
if $V \leq \alpha$ return V
 $\beta = \max(\beta, V)$

Note 3: MDPs (non deterministic search) [know all info]

- Discounted utility guarantees convergence:

$$\begin{aligned} \text{- Discounted utility guarantees convergence:} \\ \cdot R(S_0, a_0, S_1) + rR(S_1, a_1, S_2) + r^2 \dots &= \sum_{t=0}^{\infty} r^t R(S_t, a_t, S_{t+1}) \\ &\leq \sum_{t=0}^{\infty} r^t R_{\max} = \frac{R_{\max}}{1-r} \end{aligned}$$

- Markovianness (memoryless) MDP:

$$\cdot T(s, a, s') = P(s' | s, a)$$

- Solving MDPs (optimal policy π yields max expected total utility/reward)

- $V^*(s)$: optimal value of state s ; optimal expected utility starting at s
- $Q^*(s, a)$: optimal value of q-state (s, a) ; optimal value of utility starting at s , taking action a , and acting optimally henceforth

Bellman Equation

$$\cdot V^*(s) = \max_{a, s'} \sum T(s, a, s') [R(s, a, s') + rV^*(s')]$$

• optimal q-value:

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

• simplified:

$$\cdot V^*(s) = \max_a Q^*(s, a)$$

★ $V^*(s)$ is simply the max expected utility over all possible actions from s .

Value Iteration: [returns the value of a state]

- max expected utility in state s given that the Markov Decision Process terminates in k timesteps

1. $\forall s \in S$, initialize $V_0(s) = 0$

2. Repeat until convergence:
 $V_{k+1}(s) \leftarrow \max_{a, s'} \sum T(s, a, s') [R(s, a, s') + rV_k(s')]$

• Runtime: $O(|S|^2 |A|)$

Policy Extraction [returns optimal action from each state]

- $\forall s \in S$, $\pi^*(s) = \arg \max_a Q^*(s, a) = \arg \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + rV^*(s')]$

• it's better to have optimal q-values of states, then just need argmax it.

$$(Q_{\text{new}}(s, a) = R(s, a, s') + r \max_a \dots)$$

- update each weight after each observation

Policy Iteration

- faster than value iteration

1. define any random policy

2. Repeat until convergence:

- assign $V^{\pi}(s) = 0$ for all terminal states

- evaluate current policy with policy evaluation.

Policy Evaluation: compute $V^{\pi}(s)$ for all states s , where $V^{\pi}(s)$ is expected utility starting at s , following π .

$$V^{\pi}(s) = \sum T(s, \pi(s), s') [R(s, \pi(s), s') + rV^{\pi}(s')]$$

- define π_{it} as policy iteration at i : each

$V^{\pi_i}(s)$ can be solved through systems of equations

3. use policy improvement to generate better policy:

- use policy extraction on the values generated by

policy evaluation to generate improved policy.

$$\pi_{\text{it+1}}(s) = \arg \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + rV^{\pi_i}(s')]$$

- UCT, A* guarantees optimal solution

- BFS optimizes all actions same cost

- Greedy search relies all on heuristic

• UCT - backward pass to origin

• A* - use $f(h) = g(h) + h(h)$

Note 4: RL

• definition:

- sample: (s, a, s', r) tuple

Model Based Learning

- approximate transition function $T(s, a, s')$ and $R(s, a, s')$

- generate π using value iteration/policy iteration

Model Free Learning

• Passive Reinforcement Learning

1. Direct Evaluation

1. follow a fixed policy, calculate multiple episodes, then calculate $V^{\pi}(s)$ as total reward / times visited

• slow to converge, learns at the end

2. Temporal Difference Learning

1. initialize $V_s, V^{\pi}(s) = 0$ ($\pi = \pi^*$)

2. Update as you go:

$$\cdot V^{\pi}(s) \leftarrow (1-\alpha) V^{\pi}(s) + \alpha \cdot \text{sample where}$$

• sample = $R(s, \pi(s), s') + rV^{\pi}(s')$

- learn at every timestep

- give exponentially less weight to older, potentially less accurate samples

- converge to learning true state values faster

• Active Reinforcement Learning

1. Q-Learning (converges to V^* when visit all states)

• Q Value Iteration update rule: - all Qs start as 0.

$$\cdot Q_{k+1}(s, a) \leftarrow \sum_{s'} T(s, a, s') [R(s, a, s') + r \max_a Q_k(s', a)]$$

$$\cdot Q(s, a) \leftarrow (1-\alpha) Q(s, a) + \alpha \cdot \text{sample where}$$

$$\cdot \text{sample} = R(s, a, s') + r \max_a Q(s', a)$$

2. Approximate Q-Learning

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

$$\text{update: } w_i \leftarrow w_i + \text{difference} \cdot f_i(s, a)$$

$$\cdot Q(s, a) \leftarrow Q(s, a) + \text{difference}$$

where $\vec{w} = [w_1, w_2, \dots]$ weight vector

$$\vec{f}(s) = [f_1(s), f_2(s), \dots]^T$$

$$f(s, a) = [f_1(s, a), \dots, f_n(s, a)]^T$$

are feature vectors for state s and q state (s, a)

if question asks "if $Q(\text{left}) = Q(\text{right})$, then it means

$$Q^* = Q(\text{left}) = Q(\text{right})$$
. Use system of eq.

Epsilon Greedy: takes random action & probability: $(1-\epsilon)$ optimal.

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

Note 5: CSP

• identification problem: is state s a goal state

• CSP is NP-hard

Solving CSP

1. Backtracking search:

optimization of DFS

- 1) Fix an ordering for variables, and select values for variables in this order. Assignments are commutative
- 2) select variables that don't conflict with previously assigned values. Backtrack otherwise

Improvements:

1. Filtering

- 1) Forward checking: whenever a value is assigned to a variable x_i , prune domains of unassigned variables that share a constraint with x_i that could violate constraint if assigned.

↓ generalize into Arc consistency Algorithm

interpret each undirected edge of constraint graph as two directed ones. Each directed edge is an arc

- 1) Begin by storing all arcs in the constraint graph for the CSP in a queue Q.

- 2) iteratively remove arcs from Q and enforce the condition that in each removed arc $x_i \rightarrow x_j$. For every remaining value v for the tail variable x_i , there is at least one remaining value w for the head variable x_j such that $x_i = v, x_j = w$ doesn't violate constraints. If some value v for x_i wouldn't work with any remaining x_j values, we remove v from the set of possible values for x_i .
- 3) if atleast one value is removed for x_i , when enforcing arc $x_i \rightarrow x_j$, add arcs $x_k \rightarrow x_i$ to Q, for all unassigned x_k .
- 4) continue until Q empty, or domain empty and backtrack.

Runtime: $O(ed^3)$ where $e = \# \text{arcs}$; $\max(n^d)$

2) Ordering

- more effective to compute next variable and corresponding value on the fly.

- 1) MRV [minimum remaining values] - when selecting which unassigned variable has the fewest valid remaining values (most constrained variable). This is intuitive in the sense that the most constrained variable is most likely to run out of possible values and result in backtracking if left unassigned, so its best to assign a value to it sooner than later.

- 2) LCV [least constraining value] - when selecting which value to assign next, a good policy to implement is to select the value that prunes the fewest values from the domains of the remaining unassigned values. Notably, this requires additional computation (e.g. rerunning arc consistency/forward checking or other filtering methods for each value to find the LCV), but still yield speed gains depend on usage.

Exploiting Structure of CSPs

- runtime from $O(d^n)$ to $O(nd^2)$

algorithm

1. First pick an arbitrary node in the constraint graph for the CSP to serve as the root of the tree

2. Convert all undirected edges in the tree to directed edges that point away from the root.

Then linearize (topologically sort) the directed acyclic graph.

$$\text{ex: } C \xrightarrow{A} B \xrightarrow{D} E \xrightarrow{F} A \xrightarrow{B} C \xrightarrow{P} E \xrightarrow{F}$$

- 3) Perform a backwards pass of arc consistency. Iterating from $i=n$ to $i=2$
- 4) Perform forward assignment, starting from x_1 to x_n , guarantees a correct sol.

* this algorithm can be extended to CSPs that are close to a tree structure, use subset conditioning; remove that subset of variables that make it tree

2. Local Search Algorithm (Another CSP algorithm)

- start with a random assignment to values then iteratively select a random conflicted variable and reassign its value to one that violates fewest constraints til no more violations.

- $O(N)$ and high prob of success BUT

- 1) incomplete \Rightarrow won't necessarily converge and expensive on certain critical
- 2) suboptimal ratio: $R_c = \frac{\text{number of constraints}}{\text{number of variables}}$

Note 6: Propositional Logic

definitions:

valid: true in all models

satisfiable: atleast one model in which it is true

unsatisfiable: not true in all models

entailment: A entails B ($A \models B$) if in all models A is true, B is true:
models of A are subset of B

Theorems:

- 1) $(A \models B \text{ iff } A \Rightarrow B \text{ is valid})$
- 2) $(A \models B \text{ iff } A \wedge \neg B \text{ is unsatisfiable})$

Logic Tricks

- $V \wedge V \text{ and } V \wedge \Lambda \text{ are associative!}$
- $A \Rightarrow B = A \vee B$
- $A \wedge (B \vee C) = ((A \wedge B) \vee (A \wedge C))$
- $A \vee (B \wedge C) = ((A \vee B) \wedge (A \vee C))$

Algorithms for showing $KB \models q$,

- 1) simple model checking: enumerate all models $O(2^n)$
- 2) DPLL algorithm: backtracking with three improvements
 - 1) early termination
 - 2) pure symbols
 - 3) unit clauses

• CNF: $(P_1 \vee \neg P_2) \wedge \dots \wedge (P_i \vee \dots \vee P_n)$

• Local Search (WALK-SAT)

- randomly initialize all symbols and then flip some
Inferencing for showing $KB \models q$ (q is true in all models KB is true)

1. modus ponens: if KB contains $A \wedge B$, infer A, B

2. AND Elimination: $A \wedge B$, infer $A \wedge B$

3. Resolution: $A, B, \neg A \wedge B$, infer $\neg A \wedge B$

- sound & complete

Note 7: Probabilistic Inference

Bayes Rule: $P(A|B) = P(A \wedge B) / P(B) = P(B|A)P(A) / P(B)$

$P(C|B) = P(C \wedge A|B)P(A|B) + P(C \wedge \neg A|B)P(\neg A|B)$

Total Probability Rule: $P(B) = P(A \wedge B) + P(\neg A \wedge B) = P(B|A)P(A) + P(B|\neg A)(1 - P(A))$

Marginal Distribution:

$- P(X=a) = \sum_{b \in A} P(X=a, Y=b)$

Conditional Independence: $P(X|Y, Z) = P(X|Y)$
 $P(Z|Y, X) = P(Z|Y)$

Inference by enum:

$P(B_1, \dots, B_m) = \prod P(B_i, m) = \prod \sum_{e_i} P(B_i, e_i, j_i, m)$

$= \sum_e P(e) \prod_{i=1}^m P(B_i | e, j_i) \prod_{i=1}^m P(m | e)$

T/F
arc consistency not guaranteed to reduce/quash backtracking
in worst case still $O(d^n)$
set of values that remain doesn't depend on order in which
arcs are processed from q.

Midterm 2 CS188

Note 6: Probabilistic Inference

Joint Distribution: captures prob. of every possible outcome

- given joint PDF, we can trivially compute any desired $P(Q_1, \dots, Q_M | E_1, \dots, E_M)$ using slow, memory-heavy IBE

Inference By Enumeration: (IBE)

- given joint PDF collect all rows consistent with observed evidence variables; sum over all the hidden variables, normalize

Bayes Nets

- Directed Acyclic Graph; small local conditional prob-tables
- Each Node is conditionally independent of all its ancestors nodes in the graph given all of its parents $\xrightarrow{B \rightarrow A} \xrightarrow{A \rightarrow C}$ store $B, A|B, C|A$

Inference.

factors = unnormalized probability

1. Join (multiply together) all factors involving X : variable X elim
2. sum out X

Sampling: (faster, efficient for prob. reasoning) count samples

1. Prior Sampling

- just naive sample, go down the bayesnet and sample

2. Rejection Sampling

- reject samples that disagree with query

3. Likelihood Weighting

- declare all variables equal to the evidence

- weight each sample; multiply by the probability each time we get to evidence variable

4. Gibbs Sampling

- first set all variables to totally random value

- repeatedly pick one variable at a time, clear value, resample it given the values currently assigned to all other variables

- pick a new variable from that conditional distribution
* can calculate say $PCST(C, E)$ only using CPT's that connect S with neighbors

Note 8: Markov Models

Markov Models: time dependent, chain-like infinite Bayes Nets $\xrightarrow{W_0} \xrightarrow{W_1} \xrightarrow{W_2} \dots$

1. Initial Distribution: time = 0, probability of each state $PC(W_0)$

2. Transition Model: probability from one state to next: $P(W_{t+1}|W_t)$

* Markov/Memoryless Property: $t = t+1$ only dependent on $t = i$:

- $Pr(W_0, W_1, \dots, W_n) = Pr(W_0) \prod_{i=1}^n Pr(W_{i+1}|W_i)$ (Joint PDF)

Calculating Distribution of weather on t given day t .

1. Inefficient normal way

- compute joint distribution using Markov Prop.
- marginalize / sum out. basically TBE: $O(d^n)$ if d vars, n vals.

2. Mini Forward Algorithm

- By marginalization, we know $Pr(W_{t+1}) = \sum_{W_t} Pr(W_t, W_{t+1})$

- iteratively compute $Pr(W_0), Pr(W_1), \dots, Pr(W_{t+1}) = \sum_{W_t} Pr(W_{t+1}|W_t) Pr(W_t)$

stationary distribution: distribution that remains the same after the passage of time: $Pr(W_{t+1}) = Pr(W_t)$

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

- solve system of eq.

$$t = \sum_{W_t} Pr(W_t)$$

$$W_t$$

Hidden Markov Models (HMM) - allows us to observe evidence at each timestep which can potentially affect belief dist. at each state $\xrightarrow{W_0} \xrightarrow{W_1} \xrightarrow{W_2} \dots$

- $f_1 \perp\!\!\!\perp W_0, \forall i = 2, \dots, n, W_i \perp\!\!\!\perp f_2, \dots, f_{i-1}, f_i, \dots, f_{i-1} | W_{i-1}$

Representation

1. initial distribution/2. transition model/3. sensor model

def: belief distribution at time i with all evi. f_1, \dots, f_i observed:

$$B(W_i) = Pr(W_i | f_1, \dots, f_i)$$

$$B'(W_i) = Pr(W_i | f_1, \dots, f_{i-1})$$

Forward Algorithm (calculate Belief dist at given time i)

1. Time Elapse Update - determining $B'(W_{t+1})$ from $B(W_t)$:

$$B'(W_{t+1}) = \sum_{W_t} Pr(W_{t+1}|W_t) B(W_t)$$

2. Observation Update - determining $B(W_{t+1})$ from $B'(W_{t+1})$

$$B(W_{t+1}) = Pr(f_{t+1}|W_{t+1}) B'(W_{t+1})$$

Particle Filtering

- sampling but for HMM

1. randomly initialize particles from distribution

2. Time Elapse Update:

- for each particle in state t ; update its state from prob. dist. given by $Pr(T_{t+1}|t, t)$

3. Observation Update

- weight each particle with sensor model, $Pr(F_d|T_t)$ resampling bad if no evidence.

4) sum all weight for each state

3) normalize weights

4) resample particles from this new distribution \rightarrow gets you $B(T_{t+1})$ finalized

Note 8

$O(T \cdot \# \text{particles})$

Note 7: Decision Networks

def: Decision Network is a combination of Bayes Nets & Expectimax: gives effect of various actions on utility based on an overarching graphical probabilistic model

def: MEU(maximum expected utility)

- 1) instantiate all evidence that is known, run inference to calculate the posterior probabilities of all chance nodes parents of the utility node into which the action node feeds
- 2) go through each possible action and compute the expected utility of taking that action given the posterior probabilities computed in the previous step. The expected utility of taking an action a given evidence e and n chance nodes is computed with

$$EU(a|e) = \sum_{X_1, \dots, X_n} P(X_1, \dots, X_n|e) U(a, X_1, \dots, X_n)$$

where each X_i represents a value that the i^{th} chance node can take on. Weights correspond to the probabilities of each outcome.

3) finally select the action which yielded the highest utility to get the MEU:
ex) calculate optimal action, (should we leave or take umbrella) for weather:

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

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

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

Outcome Tree



def: Value of Perfect Information: answers is it worth it to invest in gaining new evidence. amount our max expected utility is expected to increase if we decide to observe new evidence.

General Formula:

$$MEU(e) = \max_a \sum_s P(s|e) U(s, a)$$

$$MEU(e, E') = \sum_e P(e'|e) MEU(e, e')$$

$$MEU(p) = \max_a EU(a) = \max_a P(w) U(a, w)$$

$$MEU(e, E') = \sum_{e'} P(e'|e) MEU(e, e') = \sum_f P(F=f) MEU(F=f) = P(F=\text{good}) MEU(F=\text{good}) + P(F=\text{bad}) MEU(F=\text{bad})$$

VPI Properties

1) Nonnegative: $\forall E/e VPI(E'|e) \geq 0$

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

3) Order-Independence: $VPI(\sum_j E_j, E_k|e) = VPI(E_j|e) + VPI(E_k|e, E_j) = VPI(E_k|e) + VPI(E_j|e, E_k)$

Note 9: Machine Learning (Supervised Learning Algorithm (Bayes Structure))

def: classification problem: given datapoints, group them into one or more classes

were also given training data and their classes

Naive Bayes Classifier

- assume feature F_i is independent of all other features, given class label

prediction(F) = argmax $P(Y=y_i) \prod_j P(F_j=f_j|Y=y_i)$

remember $\theta = P(F_i=1|Y=\text{ham})$

How to get conditional probability y_i ? Parameter Estimation
tables used in our Bayes Net

- given sample, you're MLE for outcome x that can take on $|X|$ values from sample size N is,

$$\hat{P}_{MLE} = \frac{\text{count}(x)}{N}$$

Laplace Smoothing w/ strength k : $P_{AP, k}(x) = \frac{\text{count}(x) + k}{N + k|X|}$

similar with conditionals;
computing Laplace estimate for outcomes across diff. classes

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

$$P_{AP, co}(x) = \frac{1}{|X|}$$

example

$$1) B(W_1=\text{sun}) = \sum_{F_1} P(W_1=\text{sun}|F_1) B(W_1)$$

$$B'(W_1=\text{rain}) = \dots$$

$$2) B(W_1=\text{sun}) = \Pr(F_1=\text{good}|W_1=\text{sun}) B'(W_1=\text{sun})$$

normalize after

3) Iteratively compute up to $B(W_T)$ then normalize

Note 9: Machine Learning (continued)

Perception

- doesn't estimate distribution, instead use linear combination of features called activation
- activation $(x) = \sum w_i f_i(x) = w^T f(x) = w \cdot f(x)$
- classify $(x) = \begin{cases} + & \text{if } \text{activation}_w(x) > 0 \\ - & \text{if } \text{activation}_w(x) \leq 0 \end{cases}$

Algorithm for Getting Perfect W

1. initialize all weights to 0. $w = 0$
 2. For each training sample, with features $f(x)$ and true class label $y \in \{-1, 1\}$ do:
 - (a) classify the sample using current weights, let y^* be the class predicted by your current w :
 $y = \text{classify}(x) = \begin{cases} +1 & \text{if } \text{activation}_w(x) > 0 \\ -1 & \text{if } \text{activation}_w(x) \leq 0 \end{cases}$
 - (b) compare the predicted label y to the truth label y^* :
 - if $y = y^*$ do nothing
 - otherwise if $y \neq y^*$, update weights $w \leftarrow w + y^* f(x)$
 - if you went through every training sample w/out updating weights, then terminate. Repeat step 2 (else)
- Bias
- add feature to sample feature vector that is always 1.
 - add extra weight for this feature to your weight vector.
 $w^T f(x) + b = 0$.

Multiclass Perception

multiple weight vectors now. We predict it (sample) as the one with highest score.

- update: - add feature vector to the weight vector for true class y^*
- subtract feature vector from weight vector corresponding to predicted class y .

Note 10: Neural Networks

Non-linear Separators

Multilayer Perceptron:

Theorem (Universal Function Approximators): A two-layer neural network with a sufficient number of neurons can approximate any continuous function to any desired accuracy.

Measuring Accuracy: the accuracy of the binary perceptron after making m predictions can be expressed as:

$$\text{acc}(w) = \frac{1}{m} \sum_{i=1}^m \text{sgn}(w \cdot f(x^{(i)})) = y^{(i)}$$

where $x^{(i)}$ is datapoint i , w is our weight vector, f is our function that derives a feature vector from a raw datapoint, and $y^{(i)}$ is the actual class label of $x^{(i)}$. $\text{sgn}(x)$ is indicator function which is 0 if x is negative, 1 if x .

- # correct predictions / total # predictions

Loss Functions and Multivariate Optimization

gradient vector: $\nabla_w \text{ll}(w) = \left[\frac{\partial \text{ll}(w)}{\partial w_1}, \dots, \frac{\partial \text{ll}(w)}{\partial w_n} \right]$

initialize weights w

for $i = 0, 1, 2, \dots$

$$w \leftarrow w + \alpha \nabla_w \text{ll}(w)$$

↑ changing our w each time with its derivative. want w to go to 0 with loss.

- want stepsize α to be correct.

$$\begin{aligned} \hat{y} &= h(x) \\ h(x) &= \frac{1}{1 + e^{-w^T x}} = \frac{1}{1 + e^{-r}} \\ r &= w^T x \end{aligned}$$

Test Stuff

HW8: Combining Factors Notes

Joining over a variable: (Multiply)

ex) given $P(A|C)$ and $P(B|A, C)$, join over C

$$P(A|C) \cdot P(B|A, C) = P(A, B|C)$$

$$P(A|B) \cdot P(C|B) = P(A, C|B)$$

$$\sum_c P(B, c) = P(B)$$

ex) given $P(A|C)$ and $P(B|A, C)$ join over A , sum over C

$$\sum_a P(A|C) \cdot P(B|A, C) = P(A, B|C) = \sum_a P(A, B|C) = P(B|C)$$

Logistic Regression Example

Given: $\hat{y} = h(x) = s(\sum_i w_i x_i) = \frac{1}{1 + \exp(-\sum_i w_i x_i)}$ $s(r) = \frac{1}{1 + \exp(-r)}$

Loss func $L = -[y \ln h(x) + (1-y) \ln(1-h(x))]$

Problem: find dL/dw_i hint: $s'(r) = s(r)(1-s(r))$

Solution: Chain Rule/Backpropagation

given: $L = -[y \ln h(x) + (1-y) \ln(1-h(x))]$
 $\hat{y} = h(x) = s(w^T x) = s(r) = \frac{1}{1 + e^{-r}} = \frac{1}{1 + e^{-w^T x}}$

steps: $L = -[y \ln \hat{y} + (1-y) \ln(1-\hat{y})] \Rightarrow \frac{dL}{d\hat{y}} = -\frac{y}{\hat{y}} + \frac{1-y}{1-\hat{y}}$
 $\hat{y} = s(r) = \frac{1}{1 + e^{-r}} \Rightarrow \frac{d\hat{y}}{dr} = \frac{ds(r)}{dr} = \hat{y}(1-\hat{y})$
 $y = w^T x \Rightarrow \frac{dy}{dw_i} = \frac{d}{dw_i} \sum_j w_j x_j = x_i$

multiply:

$$\frac{dL}{dw_i} = \frac{dL}{d\hat{y}} \cdot \frac{d\hat{y}}{dr} \cdot \frac{dr}{dw_i} = (-y + \hat{y}) x_i$$

stochastic gradient descent update for w_i , w step size α

$$w_i = w_i - \alpha \frac{dL}{dw_i} = w_i - \alpha(-y + \hat{y}) x_i = w_i - \alpha(n(-y + \hat{y}) x_i)$$

Test Stuff

- In IBE, the whole joint table is established: calculate the order from that

- CPT with 2 evidence, independent variables:
 $P(X_{t+1}|D_1, \dots, D_t, E_1, \dots, E_t) = P(X_{t+1}, D_{t+1}, E_{t+1}|D_1, \dots, D_t, E_1, \dots, E_t)$
 $= P(D_{t+1}, E_{t+1}|X_{t+1}, D_1, \dots, D_t, E_1, \dots, E_t) \cdot P(X_{t+1}|D_1, \dots, D_t, E_1, \dots, E_t)$
 $= P(D_{t+1}, E_{t+1}|X_{t+1}) P(E_{t+1}|X_{t+1}) \cdot P(X_{t+1}|X_{t+1}, D_1, \dots, D_t, E_1, \dots, E_t) P(X_{t+1}|D_1, \dots, D_t, E_1, \dots, E_t)$
 $= P(D_{t+1}|X_{t+1}) P(E_{t+1}|X_{t+1}) \cdot P(X_{t+1}|X_{t+1} = x) P(X_{t+1}|D_1, \dots, D_t, E_1, \dots, E_t)$
 $= P(D_{t+1}, E_{t+1}|D_1, \dots, D_t, E_1, \dots, E_t)$

MDPs

- at convergence, provide enough info to obtain optimal policy:

• model based learning of T and R

• Q-Learning to estimate $Q(s, a)$

NOT TD learning, or Direct Eval.

Q-Learning: converges to optimal Q-values for all states if:

• fixed policy taking actions uniformly at random

• an ϵ greedy policy

• NOT greedy or fixed optimal policy

Probability

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

given

$P(A, B, C, D)$ joint table, can get

$P(A, D)$ by $\sum_B \sum_C P(A, B, C, D)$

use SUM and product formulas/methods to play around.

$$\begin{array}{c} A \rightarrow B \\ \downarrow \\ C \\ \downarrow \\ P(B, C) \\ = \sum_a P(B|a) P(c|a) P(a) \end{array}$$

DFS

- not complete - complete

BFS

- complete

RISK STUFF

$$U(\text{EMV}(L)) \text{ vs. } U(L) \text{ so... lets say } U = x^2$$

and $L = \frac{1}{3}(3) \text{ and } \frac{2}{3}(6)$
then $U(L) = \frac{1}{3}(3^2) + \frac{2}{3}(6^2) = 3 + 24 = 27$
 $U(\text{EMV}(L)) = U(\frac{1}{3}(3) + \frac{2}{3}(6)) = U(5) = 25$

RISK PRONE: if second derivative is non-zero (+) then $\frac{\partial^2 U}{\partial L^2} > 0$
(-) then $\frac{\partial^2 U}{\partial L^2} < 0$

Gradient Ascent Path

- orthogonal to contours, goes toward max
- Gradient Descent goes toward minimum!

CS188 Beginnings

Search Problems - Note 1

Completeness - if there exists a solution to the search problem, are we guaranteed to find it?

Optimality - is the strategy optimal?

Branching factor b - the increase in number of nodes on the fringe each time a fringe is dequeued and replaced with its children

* At depth k, (in the search tree), there are $O(b^k)$ nodes

Search Strategies

Depth-First-Search

completeness: not complete, can be infinite loop if cycles exist

optimality: not optimal

time complexity: $O(b^m)$ where m = max depth

space complexity: $O(bm)$

Breadth-First-Search

completeness: complete

optimality: not optimal unless all edges are same cost.

time complexity: $1+b+b^2+\dots+b^s = O(b^s)$ where s is our goal depth

space complexity: $O(b^s)$

Uniform Cost Search

- selects lowest cost fringe node from start for exploration

completeness: complete

optimality: optimal if all edge costs are non-negative

time-complexity: $O(b^{C^*/\epsilon})$

Greedy Search

- always selects fringe node with lowest heuristic value

completeness: not guaranteed

optimality: not guaranteed

A* Search

- selects lowest backward cost

completeness: complete

optimality: optimal

* IN Graph Search, we never expand a state twice. If we visit a node, we put it in a closed set to not visit again.

↳ ruins optimality, so we need consistency.

* IN TREE search, we can visit/expand same nodes repeatedly, but this leads to not always complete (inf. recursion), so we add the used set condition which makes it a graph search.

MDP Questions

- When asking for transition function and reward function, give $T(s, a, s')$ and $R(s, a, s')$ for as many possibilities

- MDP cannot be solved with expectimax search because self-loops = infinite tree

- Q-Learning can converge to optimal policy even if actions are suboptimal explorations: caveat is that you need to explore enough and learning rate small enough and not decrease too fast.

- TD learning learns values for a particular policy (on-policy)

Gradient Loss Func for Approximate Q-Learning

- suppose $L_1 = (y - \sum_k w_k f_k(x))$ then $\frac{dL_1}{dw_m} = \begin{cases} -f_m(x) & \text{if } y - \sum_k w_k f_k(x) > 0 \\ f_m(x) & \text{if } y - \sum_k w_k f_k(x) < 0 \end{cases}$

$$w_m \leftarrow w_m - \alpha \frac{\partial L_1}{\partial w_m}$$

- default update: we have $Q(s, a) = \sum_a Q(s, a)$

$$L_2 = \frac{1}{2} (y - \sum_k w_k f_k(x))^2; \frac{dL_2}{dw_m} = -(y - \sum_k w_k f_k(x)) f_m(x) \text{ where}$$

$$y = r + \gamma \max_a Q(s', a')$$

$$\text{so } w_m \leftarrow w_m + \alpha (r + \gamma \max_a Q(s', a') - Q(s, a)) f_m(x)$$

Some Neural Network functions

$$f(x_1, x_2) = \max(x_1, -x_1) - 1$$

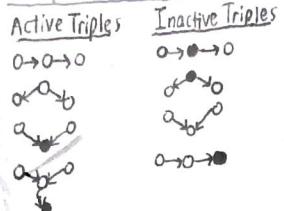
$$f(x_1, x_2) = x_2 - \max(x_1, -x_1) \text{ or } x_2 - \max(x_1, x_2, -x_2)$$

uninformed

informed

CS188 Test Stuff

D-Separation for Bayes Nets



Paths: All it takes is a subset of a path to be inactive to be inactive
 ex) $D \rightarrow E \rightarrow B \rightarrow A$ but $D \rightarrow A$ is inactive because of just one inactive

Conditional Independence: all paths (undirected) connecting the two

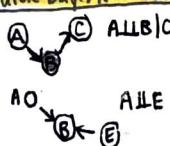
Variables need to be inactive

D-Separation Algorithm

given query $X_i \perp\!\!\!\perp X_j | \{X_{k_1}, \dots, X_{k_n}\}$

- 1) shade all evidence nodes
- 2) check all undirected paths from X_i and X_j
 - exhaust all paths to check that all paths are inactive

Quick Bayes Nets Independence



HMM Updates

standard HMM Elapse Time Update

$$P(X_t | e_{1:t-1}) = \sum_{x_{t-1}} P(X_t | x_{t-1}) P(x_{t-1} | e_{1:t-1})$$

Observation update

$$P(x_t | e_{1:t}) \propto P(x_t | e_{1:t-1}) P(e_t | x_t)$$

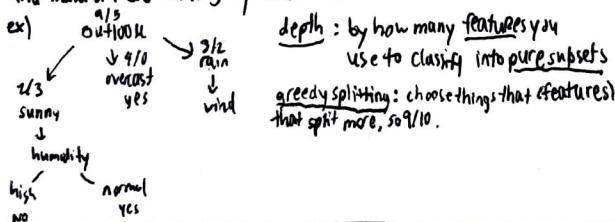
More Probability!

$$P(+s|-o, -w, +e, +h) = \frac{P(+s, -o, -w, +e, +h)}{P(+s, -o, -w, +e, +h) + P(-s, -o, -w, +e, +h)}$$

- * Power of Bayes Nets: always know Joint Probability
- * can always convert conditional probability with Joint Prob. stuff

Decision Trees

- Suppose given training set data, divide the data by each feature and make a tree sorting by the label.



Bayes Nets Inference

- once you have final factors, say $f_2(td, ta)$ and $P(A)$, and you want to calculate $P(td|ta)$, you just multiply and normalize.

$$\text{ex) } P(td|ta) = \frac{f_2(td, ta) \cdot P(td)}{\sum_a P(a) f_2(a, ta)}$$

MDP Questions

When they ask to define a game as an MDP, give state space

- 1) state space
- 2) Transition model $T(\text{state}, \text{action}, \text{state}')$
- 3) Reward Function

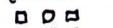
Classification Problems

• Naive Bayes can classify every dataset that perceptron can

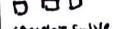
Naive Bayes

- linear

- circular?



- CANNOT solve



- Perceptrons

- guaranteed to learn separating decision boundary for separable dataset

- not guaranteed to find max-margin separating hyperplane

Loss functions

- must be differentiable (continuous range and well defined gradients)

ex) negative log-probability of correct digit

VPI Questions

- if $U(S, A)$, check to see if any VPI query have/relate to $S \cap A$. If they don't, $VPI = 0$.
- The more directly the query variables are related to S, A , the larger their VPI.
- check independence.

Value Iteration vs. Policy Iteration

V.i : $O(|S|^2 |A|)$

P.i : $O(|S|^3)$

Neural Network Representation / Graph Questions

- 1) if the graph doesn't go through the origin, add, or the function must have the bias term b .
- 2) if a graph is strictly one line (straight), can't be a function with one ReLU. Must be a 2-dimensional ReLU layer.
- 3) if graph is too hirs connected, there must have 2-D ReLU layer

4) A 2D ReLU can do 2 slope changes

Epsilon Greedy

- if there are 3 available actions and one optimal one, there are these probabilities:

optimal one: $(1-\epsilon) + \epsilon/3$

two: $\epsilon/3$

three: $\epsilon/3$