

CS188 Midterm Cheat Sheet

[simonxie2004.github.io](https://github.com/simonxie2004)

Lec1: Introduction

Lec2: Uninformed Search

- Reflex Agents V.S. Planning Agents:
 - Reflex Agents: Consider how the world IS
 - Planning Agents: Consider how the world WOULD BE
- Properties of Agents
 - Completeness: Guaranteed to find a solution if one exists.
 - Optimality: Guaranteed to find the least cost path.
- Definition of Search Problem:
 - "State Space", "Successor Function", "Start State" & "Goal Test"
- Definition of State Space: World State & Search State
- State Space Graph: Nodes = states, Arcs = successors (action results)
- Tree Search
 - Main Idea: Expand out potential nodes; Maintain a fringe of partial plans under consideration; Expand less nodes.
 - Key notions: Expansion, Expansion Strategy, Fringe
 - Common tree search patterns
 - (Suppose $b =$ branching factor, $m =$ tree depth.)
 - Nodes in search tree? $\sum_{i=0}^m b^i = O(b^{m+1})$
 - (For BFS, suppose $s =$ depth of shallowest solution)
 - (For Uniform Cost Search, suppose solution costs C^* , $\min(\text{arc_cost}) = \epsilon$)
 - Search Idea: Iterative Deepening
 - Run DFS(depth_limit=1), DFS(depth_limit=2), ...
 - Example Problem: Pancake flipping; Cost: Number of pancakes flipped
- Graph Search
 - Idea: never expand a state twice
 - Method: record set of expanded states where elements = (state, cost). If a node popped from queue is NOT visited, visit it. If a node popped from queue is visited, check its cost. If the cost is lower, expand it. Else skip it.

Strategy	Fringe	Time	Memory	Completeness	Optimality	
DFS	Expand deepest node first	LIFO Stack	$O(b^m)$	$O(bm)$	True (if no cycles)	False
BFS	Expand shallowest node first	FIFO Queue	$O(b^m)$	$O(b^m)$	True	True (if cost=1)
UCS	Expand cheapest node first	Priority Queue (g/cumulative cost)	$O(C^*/\epsilon)$	$O(b^{C^*/\epsilon})$	True	True

Lec3: Informed Search

- Definition of heuristic:
 - Function that estimates how close a state is to a goal; Problem specific!
- Example heuristics: (Relaxed-problem heuristic)
 - Pancake flipping: heuristic = the number of largest pancake that is still out of place
 - Dot-Eating Pacman: heuristic = the sum of all weights in a MST (of dots & current coordinate)
 - Classic 8 Puzzle: heuristic = number of tiles misplaced
 - Easy 8 Puzzle (allow tile to be piled intermediately): heuristic = total Manhattan distance
- Remark: Can't use actual cost as heuristic, since have to solve that first!
- Comparison of algorithms:
 - Greedy Search: expand closest node (to goal); Orders by forward cost $h(n)$; suboptimal
 - UCS: expand closest node (to start state); Orders by backward cost $g(n)$; suboptimal
 - A* Search: orders by sum $f(n) = g(n) + h(n)$
- A* Search
 - When to stop: Only if we dequeue a goal
 - Admissible (optimistic) heuristic: $\forall n, 0 \leq h(n) \leq h^*(n)$. A* Tree Search is optimal if heuristic is admissible. Proof: Suppose A is optimal, B is suboptimal, B is on fringe. Claim: n will be expanded first. Because $f(n) = g(n) + h(n) < f(B)$
 - Consistent heuristic: $\forall A, B, h(A) - h(B) \leq \text{cost}(A, B)$. A* Graph Search is optimal if heuristic is consistent.
- Semi-Lattice of Heuristics
 - Dominance: define $h_2 \geq h_1$ if $\forall n, h_2(n) \geq h_1(n)$
 - Heuristics form semi-lattice because: $\forall h_1, h_2, h_3, h_4(n) \in H$
 - Bottom of lattice is zero-heuristic. Top of lattice is exact-heuristic

```
function TREE-SEARCH(problem, fringe) return a solution, or failure
fringe ← INSERT(MAKE-NODE(INITIAL-STATE[problem]), fringe)
loop do
    if fringe is empty then return failure
    node ← REMOVE-FRONT(fringe)
    if GOAL-TEST(problem, STATE[node]) then return node
    for child-node in EXPAND(STATE[node], problem) do
        fringe ← INSERT(child-node, fringe)
    end
end
```

```
function A*-GRAPH-SEARCH(problem, frontier) return a solution or failure
reached ← an empty dict mapping nodes to the cost to each one
frontier ← INSERT(MAKE-NODE(INITIAL-STATE[problem]), 0, frontier)
while not IS-EMPTY(frontier) do
    node, node.CostToNode ← POP(frontier)
    if problem.IS-GOAL(STATE[node]) then return node
    end if
    if node.STATE is not in reached or reached[node.STATE] > node.CostToNode then
        reached[node.STATE] = node.CostToNode
        for each child-node in EXPAND(problem, node) do
            frontier ← INSERT(child-node, child-node.COST + CostToNode, frontier)
        end for
    end if
end while
return failure
```

Lec4-5: CSP Problems

- Definition of CSP Problems: (A special subset of search problems)
 - State: Variables $\{X_i\}$, with values from domain D
 - Goal Test: set of constraints specifying allowable combinations of values
- Example of CSP Problems:
 - N-Queens
 - $\forall i, j, k, (X_{ij}, X_{jk}) \neq (1, 1), \dots, \text{and } \sum_{i,j} X_{ij} = N$
 - Formulation 1: Variables: X_{ij} , Domains: $\{0, 1\}$, Constraints:
 - Formulation 2: Variables Q_k , Domains: $\{1, \dots, N\}$, Constraints:
 - Cryptarithmic
 - $\forall (i, j), \text{non-threatening}(Q_i, Q_j)$
- Constraint Graph:
 - Circle nodes = Variables; Rectangular nodes = Constraints.
 - If there is a relation between some variables, They are connected to a constraint node.
- Simple Backtracking Search
 - One variable at a time
 - Check constraints as you go. (Only consider constraints not conflicting to previous assignments)
- Simple Backtracking Algorithm = DFS + variable-ordering + fail-on-violation

```
function BACKTRACKING-SEARCH(esp) returns solution/failure
return RECURSIVE-BACKTRACKING({}, esp)

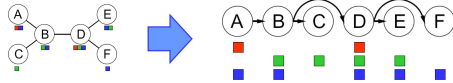
function RECURSIVE-BACKTRACKING(assignment, esp) returns soln/failure
if assignment is complete then return assignment
var ← SELECT-UNASSIGNED-VARIABLE(VARIABLES[esp], assignment, esp)
for each value in ORDER-DOMAIN-VALUES(var, assignment, esp) do
    if value is consistent with assignment given CONSTRAINTS[esp] then
        add {var = value} to assignment
        result ← RECURSIVE-BACKTRACKING(assignment, esp)
        if result ≠ failure then return result
        remove {var = value} from assignment
    end if
end for
return failure
```

- Filtering & Arc Consistency
 - Definition: Arc $X \rightarrow Y$ is consistent if $\forall x \in X, \exists y \in Y$ that could be assigned. (Basically X is enforcing constraints on Y)
 - Filtering: Forward Checking: Enforcing consistency of arcs pointing to each new assignment
 - Filtering: Constraint Propagation: If X loses a Value, neighbors of X need to be rechecked.
 - Usage: run arc consistency as a preprocessing or after each assignment
 - Algorithm with Runtime $O(n^2d^3)$

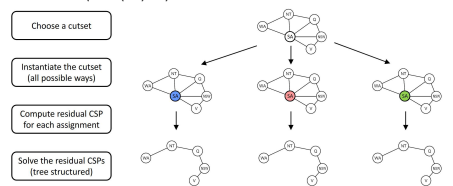
```
function ARC-3(esp) returns the CSP, possibly with reduced domains
inputs: esp, a binary CSP with variables  $\{X_1, X_2, \dots, X_n\}$ 
local variables: queue, a queue of arcs, initially all the arcs in esp
while queue is not empty do
     $(X_i, X_j) \leftarrow \text{REMOVE-FIRST}(queue)$ 
    if REMOVE-INCONSISTENT-VALUES( $X_i, X_j$ ) then
        for each  $X_k$  in NEIGHBORS[ $X_i$ ] do
            add  $(X_i, X_k)$  to queue
        end for
    end if
end while

function REMOVE-INCONSISTENT-VALUES( $X_i, X_j$ ) returns true if succeeds
removed ← false
for each  $x$  in DOMAIN[ $X_i$ ] do
    if no value  $y$  in DOMAIN[ $X_j$ ] allows  $(x, y)$  to satisfy the constraint  $X_i \leftrightarrow X_j$  then
        delete  $x$  from DOMAIN[ $X_i$ ]; removed ← true
    end if
end for
return removed
```

- Advanced Definition: K-Consistency
 - K-Consistency: For each k nodes, Any consistent assignment to k-1 nodes can be extended to kth node.
 - Strong K-Consistency: also k-1, k-2, ..., 2-Consistent; Can be solved immediately without searching / backtracking
 - Problems of Arc-consistency: only considers 2-consistency
 - Example of being NOT 3-consistent:
- Advanced Arc-Consistency: Ordering
 - Variable Ordering: MRV (Minimum Remaining Value): Choose the variable with fewest legal left values in domain
 - Value Ordering: LCV (Least Constraining Value): Choose the value that rules out fewest values in remaining variables. (May require re-running filtering.)
- Advanced Arc-Consistency: Observing Problem Structure
 - Suppose graph of n variables can be broken into subproblems with c variables: Can solve in $O(n^c d^c)$
 - Suppose graph is a tree: Can solve in $O(nd^2)$. Method as follows
 - Remove backward: For $i = n-2$, apply RemoveInconsistent(Parent(X_i), X_i)
 - Assign forward: For $i = 1, \dots, n$, assign X_i consistently with Parent(X_i)
 - *Remark: After backward pass, all root-to-leaf are consistent. Forward assignment will not backtrack.

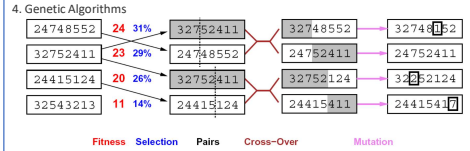


- Advanced Arc-Consistency: Improving Problem Structure
 - Idea: Initiate a variable and prune its neighbors' domains.
 - Method: instantiate a set of vars such that remaining constraint graph is a tree (cutset conditioning)
 - Runtime: $O(d^c * n \cdot d^2)$ to solve CSP.



- Iterative Methods:
 - Local Search
 - Algorithm: While not solved, randomly select any conflicted variable. Assign value by min-conflicts heuristic.
 - Performance: can solve n-queens in almost constant time for arbitrary n with high probability, except a few of them.
 - Hill Climbing
 - function HILL-CLIMBING(problem) returns a state
 - current ← make-node(problem, initial-state)
 - loop do
 neighbor ← a highest-valued successor of current
 if neighbor value \leq current value then
 return current state
 end if
 current ← neighbor
 end loop
 - Hill Climbing
 - Remark: Stationary distribution: $p(x) \propto e^{E(x)/kT}$

```
function SIMULATED-ANNEALING(problem, schedule) returns a solution state
inputs: problem, a problem
schedule, a mapping from time to "temperature"
local variables: current, a node
next, a node
T, a "temperature" controlling prob. of downward steps
current ← MAKE-NODE(INITIAL-STATE[problem])
for t ← 1 to ∞ do
    T ← schedule[t]
    if T = 0 then return current
    next ← a randomly selected successor of current
    ΔE ← VALUE[next] - VALUE[current]
    if ΔE > 0 then current ← next
    else current ← next only with probability  $e^{-\Delta E/T}$ 
end for
```



Lec6: Game Trees (MiniMax)

- Zero-Sum Games V.S. General Games: Opposite utilities v.s. Independent utilities
 - Examples of Zero-Sum Games: Tic-tac-toe, chess, checkers, ...
 - Value of State: Best achievable outcome (utility) from that state.
 - For MAX players $\max_{s' \in \text{children}(s)} V(s')$
 - For MIN players, min...
 - Search Strategy: Minimax
 - def value(state):
 if the state is a terminal state: return the state's utility
 if the next agent is MAX: return max-value(state)
 if the next agent is MIN: return min-value(state)
 end def

```
def max-value(state):
    initialize v = -∞
    for each successor of state:
        v = max(v, value(successor))
    return v

def min-value(state):
    initialize v = ∞
    for each successor of state:
        v = min(v, value(successor))
    return v
```

- Minimax properties:
 - Optimal against perfect player. Sub-optimal otherwise.
 - Time: $O(b^m)$, Space: $O(bm)$
- Alpha-Beta Pruning
 - Algorithm:
 - α : MAX's best option on path to root
 - β : MIN's best option on path to root

```
def max-value(state, α, β):
    initialize v = -∞
    for each successor of state:
        v = max(v, value(successor, α, β))
        if v ≥ β return v
        α = max(α, v)
    return v

def min-value(state, α, β):
    initialize v = ∞
    for each successor of state:
        v = min(v, value(successor, α, β))
        if v ≤ α return v
        β = min(β, v)
    return v
```

- Properties:
 - Meaning of Alpha: maximum reward for MAX players, best option so far for MAX player
 - Meaning of Beta: minimum loss for MIN players, best option so far for MIN player
 - Have no effect on root value; intermediate values might be wrong.
 - With perfect ordering, time complexity drops to $O(b^m/m^2)$
- Depth-Limited Minimax: replace terminal utilities with an evaluation function for non-terminate positions
 - Evaluation Functions: weighted sum of features observed
 - Iterative Deepening: run minimax with depth_limit = 1, 2, 3, ... until timeout

Lec7: Game Trees (Expectimax, Utilities)

- Expectimax Algorithm:
 - def value(state):
 if the state is a terminal state: return the state's utility
 if the next agent is MAX: return max-value(state)
 if the next agent is EXP: return exp-value(state)
 end def

```
def max-value(state):
    initialize v = -∞
    for each successor of state:
        v = max(v, value(successor))
    return v

def exp-value(state):
    initialize v = 0
    for each successor of state:
        p = probability(successor)
        v += p * value(successor)
    return v
```

Assumptions vs Reality:	Rational & Irrational Agents	Adversarial Ghost	Random Ghost
Minimax Pacman	Won 5/5	Won 5/5	Won 5/5
Expectimax Pacman	Avg. Score: 483	Avg. Score: 493	Avg. Score: 503

- Axioms of Rationality
 - Orderability: $(A > B) \vee (B > A) \vee (A \sim B)$
 - Transitivity: $(A > B) \wedge (B > C) \Rightarrow (A > C)$
 - Continuity: $A > B > C \Rightarrow \exists p \in [0, 1-p, C] \sim B$
 - Substitutability: $A \sim B \Rightarrow [p, A; 1-p, C] \sim [p, B; 1-p, C]$
 - Monotonicity: $A > B \Rightarrow (p \geq q \Rightarrow [p, A; 1-p, B] \geq [q, A; 1-p, B])$
- MEU Principle
 - Given any preferences satisfying these constraints, there exists a real valued function U s.t.:
 - $U(A) \geq U(B) \Leftrightarrow A \succeq B$
 - $U([p_1, S_1; \dots; p_n, S_n]) = \sum p_i U(S_i)$
- Risk-averse v.s. Risk-prone
 - Def. $L = [p, X, 1-p, Y]$
 - If $U(L) < U(EMV(L))$, risk-averse
 - Where $U(L) = pU(X) + (1-p)U(Y)$, $U(EMV(L)) = U(pX + (1-p)Y)$
 - i.e. if U is concave, like $y = \log_2 x$, then risk-averse
 - Otherwise, risk-prone. i.e. if U is convex, like $y = x^2$, then risk-prone

Lec8-9: Markov Decision Process

- MDP World: Noisy movement, maze-like problem, receives rewards.
 - "Markov": Successor only depends on current state (not the history)
- MDP World Definition:
 - States, Actions
 - Transition Function $T(s, a, s')$ or $Pr(s' | s, a)$, Reward Function $R(s, a, s')$
 - Start State, (Probably) Terminal State
 - MDP Target: optimal policy $\pi^*: S \rightarrow A$

4. Discounting: Earlier is Better! No infinity rewards!

$$U([r_0, \dots, r_{inf}]) = \sum_{t=0}^{\infty} \gamma^t r_t = \leq R_{max} / (1 - \gamma)$$

5. MDP Search Trees:

1. Value of State: expected utility starting in s and acting optimally.

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

2. Value of Q-State: expected utility starting out having taken action a from state s and (thereafter) acting optimally.

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

3. Optimal Policy: $\pi^*(s)$

6. Solving MDP Equations: Value Iteration

1. Bellman Equation: $V^*(s) = \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^*(s')]$

2. Value Calculation: $V_{k+1}(s) \leftarrow \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_k(s')]$

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

4. Complexity (of each iteration): $O(S^2A)$

5. Must converge to optimal values. Policy may converge much earlier.

7. Solving MDP Equations: Q-Value Iteration

1. Bellman Equation: $Q^*(s, a) = \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma \max_{a'} Q^*(s', a')]$

2. Policy Extraction: $\pi^*(s) = \arg \max_a Q^*(s, a)$

8. MDP Policy Evaluation: Evaluating V for fixed policy

1. Idea 1: remove the max'es from Bellman, iterating

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

2. Idea 2: is a linear system. Use a linear system solver.

9. Solving MDP Equations: Policy Iteration

1. Idea: Update Policy & Value meanwhile, much much faster!

2. Algorithm:

▪ Initialize $\pi_0(s) = \text{some default action for all } s$

▪ for i of policy iteration:

Policy evaluation:

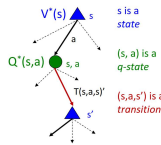
▪ Initialize $V_0^{\pi_i}(s) = 0$ for all s

▪ for k of policy evaluation:

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

Policy improvement:

$$\pi_{i+1}(s) = \arg \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^{\pi_i}(s')]$$



7. Scaling up RL: Approximate Q Learning

1. State space too large & sparse?

Use linear functions to approximately learn $Q(s, a)$ or $V(s)$

2. Definition: $Q(s, a) = w_1 f_1(s, a) + w_2 f_2(s, a) + \dots$

3. Q-learning with linear Q-functions:

Transition := (s, a, r, s')

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

Approx. Update weight: $w_i \leftarrow w_i + a \cdot \text{Difference} \cdot f_i(s, a)$

Lec 10-11: Reinforcement Learning

1. Intuition: Suppose we know nothing about the world. Don't know T or R.



2. Passive RL I: Model-Based RL

1. Count outcomes s' for each s, a ; Record R ;

2. Calculate MDP through any iteration

3. Run policy. If not satisfied, add data and goto step 1

3. Passive RL II: Model-Free RL (Direct Evaluation, Sample-Based Bellman Updates)

1. Intuition: Direct evaluation from samples.

Improve our estimate of V by computing averages of samples.

2. Input: fixed policy $\pi(s)$

3. Act according to π . Each time we visit a state,

write down what the sum of discounted rewards turned out to be.

4. Average the samples, we get estimate of $V(s)$

$$\text{sample}_1 = R(s, \pi(s), s'_1) + \gamma V_k^{\pi}(s'_1)$$

$$\text{sample}_2 = R(s, \pi(s), s'_2) + \gamma V_k^{\pi}(s'_2)$$

...

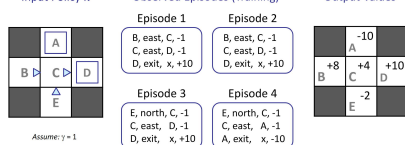
$$\text{sample}_n = R(s, \pi(s), s'_n) + \gamma V_k^{\pi}(s'_n)$$

$$V_{k+1}^{\pi}(s) \leftarrow \frac{1}{n} \sum \text{sample}_i$$

5. Problem: wastes information about state connections.

Each state learned separately. Takes long time.

Input Policy π Observed Episodes (Training) Output Values



4. Passive RL II: Model-Free RL (Temporal Difference Learning)

1. Intuition: learn from everywhere / every experience.

Recent samples are more important. $\hat{x}_n = (1 - \alpha)x_{n-1} + \alpha x_n$

$$\hat{x}_n = x_n + (1 - \alpha) \cdot x_{n-1} + (1 - \alpha)^2 \cdot x_{n-2} + \dots$$

2. Update:

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

$$\text{Update to } V(s): V^{\pi}(s) \leftarrow (1 - \alpha)V^{\pi}(s) + (\alpha)\text{sample}$$

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

3. Decreasing learning rate (alpha) converges

4. Problem: Can't do policy extraction, can't calculate $Q(s, a)$ without T or R

5. Idea: learn Q-values directly! Make action selection model-free too!

5. Passive RL III: Model-Free RL (Q-Learning + Time Difference Learning)

1. Intuition: Learn as you go.

2. Update:

1. Receive a sample (s, a, s', r)

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

3. Incorporate new estimate into a running average:

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

4. Another representation: $Q(s, a) \leftarrow Q(s, a) + a \cdot \text{Difference}$

where $\text{diff} = \text{sample} - \text{orig}$

3. This is off-policy learning!

6. Active RL: How to act to collect data

1. Exploration schemes: eps-greedy

1. With probability eps , act randomly from all options

2. With probability $1 - \text{eps}$, act on current policy

2. Exploration functions: use an optimistic utility instead of real utility

1. Def. optimistic utility $f(u, n) = u + k/n$

suppose $u = \text{value estimate}$, $n = \text{visit count}$

2. Modified Q-Update: $Q(s, a) \leftarrow \alpha R(s, a, s') + \gamma \max_{a'} f(Q(s', a'), N(s', a'))$

3. Measuring total mistake cost: sum of difference between expected rewards and suboptimality rewards.