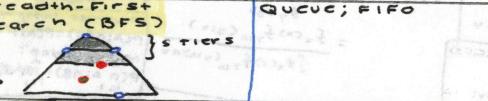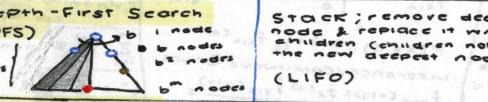## State Space side:

- **Fundamental counting principle:** if there are $n$ variable objects in a given world that can take on values $x_1, x_2, ..., x_n$, the total # of states is $\prod x_i$.
  - eg:
    Pacman can be in 120 distinct $(x,y)$ pos. PM can go NSEW. 2 Ghosts; can be in 12 diff pos. 30 Food pellets (can be eaten/not)
    size = $120 \cdot 4 \cdot 12^2 \cdot 2^{30}$

- **Completeness:** if a soln to a search prob exists, is the strategy guaranteed to return a soln in finite w/o resources?
- **Optimality:** if a strategy guaranteed to find the lowest cost path to a goal state?
- **Branching factor (b):** increase in # of nodes on the frontier each time a frontier node is dequeued & replaced w/ its children (b of b). At depth $k$ in the search tree, $\exists$ $O(b^k)$ nodes.
- max depth (m) · optimal path cost (C*)
- shallowest soln depth (s) · $\varepsilon$ (minimal cost btwn 2 nodes in ss graph)

## Uninformed Search Summary

| Search | Frontier | Completeness | Optimality | Time | Space |
|---|---|---|---|---|---|
| Depth-First Search (DFS) | Stack; remove deepest node & replace it w/its children (children now the new deepest nodes) (LIFO) | tree search - no; graph search: yes if finite (no cycles); no if infinite | No. Finds left-most solution. | $O(b^m)$ | $O(bm)$; be it maintains b nodes at each of the m depths of the frontier |
| Breadth-First Search (BFS) | Queue; FIFO 5 tiers | yes, if ∃soln, the depth of shallowest node must be finite. BFS will search this depth | No, be it doesn't take cost into consideration when determining node to replace on frontier. Only guarantees optimal if costs = 1, which reduces to UCS. | $O(b^s)$ | $O(b^s)$ |
| UCS (Uniform Cost Search) | Heap-based priority queue. Priority for node v is the backwards cost (path from start node to v) Note: despite the fact that UCS is complete & optimal, it's slow be it explores in all directions, (explores increasing cost contours) orders by path cost (forward cost) | yes. Same reason as BFS. | yes if all edge costs are non-neg. | $O(b^{C*/\varepsilon})$ | $O(b^{C*/\varepsilon})$ |
| Iterative Deepening gets DFS's space advantage w/BFS's time advantage | Stack | yes (same as BFS) | No. same as BFS. | $O(b^s)$ | $O(bs)$ |

## Informed Search Summary

| Search | Heuristic | Complete | Optimal |
|---|---|---|---|
| Greedy | estimate distance to nearest goal for each state | No | No, esp w/bad heuristic · Orders by goal proximity (backwards cost) |
| A* search | selects lowest total cost | yes | Graph: yes if heuristic is consistent; UCS is special case (h=0) · Tree: yes if heuristic is admissible; UCS optimal (h=0 consistent) |
|  | $f(n) = h(n) + g(n)$ |  |  |

### Local search

- **Hill Climbing / steepest ascent**
  - start wherever.
  - loop: move to best neighbor state
  - stop: if no neighbors better than curren
  - w/o random restarts? No
  - w/ random restarts? yes
  - ↑T to allow more "bad" moves
  - gradually reduce T

- **Simulated Annealing**
  - combines random walk w/ hill climbing

- **Local Beam Search**
  - create K iterations of local search algo, random initialized
  - for each iteration, gen all successors & choose best K to survive.

- **Genetic Algorithm**
  - resample K individuals at each step, selection weighted by fitness fcn
  - combine by pairwise crossover operators, plus mutation (for variety)

## Temporal Difference Learning

- uses exponential moving average
- sample = $R(s, \pi(s), s') + \gamma V^\pi(s')$
- update:
  $V^\pi(s) = (1-\alpha)V^\pi(s) + \alpha \cdot sample$

---

- **heuristics** - take in a state as an input & output a estimate. two main properties:
  1. **Admissibility (optimistic):**
     $0 \leq h(n) \leq h^*(n)$
     heuristic ⟵ ⟶ true cost
     → inadmissible heuristics break optimality by trapping good plans on the fringe
  2. **Consistency:** must under-estimate the cost/weight of each edge in a graph
     $\forall$ nodes A
     $h(A) - h(c) \leq cost(A, c)$

- $g(n)$ - backwards cost computed by UCS
- $h(n)$ - heuristic value fcn or estimated forward cost, used by greedy search
- $f(n)$ - fcn representing total cost, used by A*:
  $f(n) = h(n) + g(n)$
- consistency ⟹ admissability
- **dominance:** $h_1 \geq h_2$ if $\forall n$ $h_1(n) \geq h_2(n)$
- $h(n) = \max(h_1(n), h_2(n))$ of 2 admissable heuristics is admissable & dominates both

## Propositional Logic (PL)

- PL written in sentences composed of propositional symbols (PS)
- **model** - assignment of T/F to all PS
  - for N symbols, $2^N$ models
- **valid sentence** - true in all models
- **satisfiable** - true in ∃1 model
- **unsatisfiable** - false in all models
- **CNF (Conjunctive Normal Form)**
  $(P_1 \vee ... \vee P_i) \wedge ... \wedge (P_j \vee ... \vee P_N)$

- **Entailment** ⊨ ($\vDash$)
  → A entails B (A $\vDash$ B) if in all sentences where model A is true, model B is as well
  → if A $\vDash$ B, the models of A are a subset of the models of B
  → graph DFS guaranteed to return a path passing thru all possible grid locations if one exists, not tree DFS
  → same w/ BFS (both graph & tree) traverse whole kb
  · can prove entailment using 3 rules of inference
  1. **Modus ponens:** if our KB contains A & A⇒B, we can infer B
  2. **And-elimination:** if our KB contains (A∧B) we can infer A and can infer B
  3. **Resolution:** if our KB contains A & B we can infer (A∧B)

## Reinforcement Learning

- **Offline planning** - agents have full knowledge of transition & reward fcns (eg, solving an MDP).
- **Online planning** - no knowledge of T/R so have to try exploration (receiving feedback for actions it performs) & uses RL for exploitation (reward maximization)
- 2 types of RL:
  1. **Model-based learning** - attempts to estimate transition & reward fcns then uses those estimat to solve the MDP
  2. **Model-free learning:** attempts to estimate Q-values of states directly

## Model-Based Learning

$T(s, a, s') = \dfrac{\text{# times } s, a \to s'}{\text{# times we did } s, a}$

$R(s, a, s') = \dfrac{\sum \text{rewards from } s, a \to s'}{\text{# times we did } s, a \to s'}$

## First-Order Logic (FOL)

- FOL uses objects as its basic components
- each object represented as a constant symbol, each relationship by a predicate symbol, & each fcn by a fcn symbol.
- **Terms** - logical exps that refer to an obj.
  → simplest = constant symbols
- **atomic sentences:** relationships btwn objects
  Brother (John, Richard)
- **quantifiers** $\forall, \exists$:
  $\forall x \, \neg P \equiv \exists x \, \neg P$
  $\neg \forall x \, P \equiv \exists x \, \neg P$
  $\forall x \, P \equiv \neg \exists x \, \neg P$
  $\exists x \, \neg P \equiv \forall x \, \neg P$
  $\exists x \, P \equiv \neg \forall x \, \neg P$

  $\neg(P \vee Q) \equiv \neg P \wedge \neg Q$
  $\neg(P \wedge Q) \equiv \neg P \vee \neg Q$
  $P \wedge Q \equiv \neg(\neg P \vee \neg Q)$
  $P \vee Q \equiv \neg(\neg P \wedge \neg Q)$

- formulating inference using FOL (i.e. KB⊨q)
  1. **Propositionalization:** translate the problem from FOL to PL so we can use our o/ne PL methods on it
  2. **Lifted inference:** directly do inference with FOL

## Approx. Q-learning
- allows us to generalize learning experiences to other states (repr. each state as a feature vector).

## Q-learning
Q-value samples:
sample = $R(s,a,s') + \gamma \max_{a'} Q(s',a')$
$Q(s,a) \leftarrow (1-\alpha)Q(s,a) + \alpha \cdot sample$
many similar situations by using feature-based representation of the states (repr. each state as a feature vector).

---

- **idea:** from each state, locally, move towards states w/ higher objective value until a max (hopefully global) is reached.
- **state space** - set of complete configurations
- **stochastic hill climbing** - selects action randomly from possible
  - uphill moves
  - **random sideways move** hill climbing - allows moves that don't strictly increase obj.

## Converting to CNF
$A \Leftrightarrow (B \vee C)$
1. Eliminate $\Leftrightarrow$
   $(A \Rightarrow (B \vee C)) \wedge ((B \vee C) \Rightarrow A)$
2. Eliminate $\Rightarrow$
   $(\neg A \vee B \vee C) \wedge (\neg(B \vee C) \vee A)$
3. $\neg$'s must only appear on literals
   $(\neg A \vee B \vee C) \wedge ((\neg B \wedge \neg C) \vee A)$
4. Distribute
   $(\neg A \vee B \vee C) \wedge (\neg B \vee A) \wedge (\neg C \vee A)$

## Reinforcement Learning (right column cont.)

### Direct Evaluation
$\hat{V}^\pi(s) = \sum_{\text{episodes}} \dfrac{\text{Utility starting from state } s}{\text{# episodes } s \text{ appeared in}}$

### Approximate Q-learning
diff = $[R(s,a,s') + \gamma \max_{a'} Q(s',a')] - Q(s,a)$
$w_i \leftarrow w_i + \alpha \cdot diff$

---

- **# of actions available** from a state (aka largest # of successors from any state)
  $b = 9$
  for
  throttle = {-1, 0, 1}
  steering = {-1, 0, 1}

- **a proof is a demonstration** of entailment btwn $\alpha$ & $\beta$
  → **sound algorithm:** everything it claims to prove is entailed
  → **complete algorithm:** everything that is entailed can be proved

2 ways to show entailment
1. $(A \vDash B \text{ iff } A \Rightarrow B$ is valid)
   → direct proof
2. $(A \vDash B$ iff $A \wedge \neg B$ is unsatisfiable)
   → proof by contradiction

### entailment algorithms
1. **Model Checking**
   N symbols → $2^N$ models to check
   ⟹ $O(2^N)$ in propositional logic
   ⟹ $O(\infty)$ in First order logic
2. **DPLL SAT Solver**
   - recursive depth 1st enumeration of models w/some extra features:
     1. early termination - stop if all clauses are satisfied or
        · $(A \vee B) \wedge (A \vee \neg C)$ is satisfied by A = true
        · any clause is falsified
        · $(A \vee B) \wedge (A \vee \neg C)$ is falsified by {A=F, B=F}
     2. **pure literals** - if all occurrences of a symbol in as-of-yet unsatisfied clauses have the same sign, give the symbol that value
        $(A \vee B) \wedge (A \vee \neg C) \wedge (C \vee \neg B)$
        → A is pure & (+), so set it to T.
     3. **unit clauses** - if a clause left w/a single literal, set symbol to satisfy clause.
        · if A=F in $(A \vee B) \wedge (A \vee \neg C)$ we have $B \wedge \neg C$, so set C=F, B=T
        → satisfying these can lead to further propagation/new unit clauses
3. **Forward Chaining**
   → useful for special case when our KB only has literals & implications, be then we can prove entailment in O(N) time (N=size of KB).
   → FC iterates through every statement where the premise (left side) is known to be true, adding it to the conclusion
4. **Backward Chaining**

# (left column)

init all weights to 0 ($\bar{w}=0$)

For each training sample w/features $f(x)$ & true class label $y^* \in \{-1,+1\}$, do:

i) classify sample using current weights. Let $\hat{y}$ be the class predicted by curr $\bar{w}$:
$$y=\text{classify}(x)=\begin{cases}+1 & \text{if } h_w(x)=w^T f(x) > 0 \\ -1 & \text{o/w}\end{cases}$$

ii) Compare $y$ to $y^*$
   a) if $y=y^*$: do nothing
   b) if $y\neq y^*$: $\bar{w} \leftarrow \bar{w}+y^* f(x)$

> Note:
> Case 1: Misclassified as (−): $w \leftarrow w+f(x)$
> Case 2: Misclassified as (+): $w \leftarrow w-f(x)$

If you went through every training sample w/o having to update your weights, then terminate. o/w repeat ②.

## Multiclass Perceptron Algorithm
$\bar{w}=0$

For each training sample
i) Predict class w/current weights:
$y=\text{argmax } w_y f(x)$
ii) Compare $y$ to $y^*$
   a) if $y=y^*$: do nothing
   b) if $y\neq y^*$: lower score of wrong answer & raise score of correct answer:
   $w_y = w_y - f(x)$
   $w_{y^*} = w_{y^*} + f(x)$

## Perceptron Properties
- Separability: true if some parameters get the training set perfectly correct
- Convergence: if training is separable, perceptron will eventually converge (binary case)
- Mistake bound: max # of mistakes (binary) related to margin/degree of separability: mistakes $< k/\delta^2$
- If a perceptron misclassifies data, it might still misclassify after the weight update

## Perceptron Probs
- Noise: if data not separable, weights might thrash ⇒ avg weight vectors over time (avg'd perceptron) can help
- Mediocre generalization: finds "barely" separating soln
- Overtraining: test/held-out accuracy usually rises then falls

## Probabilistic Perceptron Decisions)
- instead of deterministic
- if $z=w\cdot f(x)$ very positive → prob of class +1 should approach 1 & vice versa

Sigmoid Fcn: $\phi(z)= \dfrac{1}{1+e^{-z}}$
probability that the class is +1, according to the classifier

softmax fcn:
$P(z)_i = (e^{z_i})/(\sum_{i=1}^K e^{z_i})$

ReLU: $\max(0,x)$
rectified linear units

## Optimization
Gradient Ascent:
- randomly initialize $\bar{w}$
- while $\bar{w}$ not converged:
  $\bar{w} = \bar{w} + \alpha \nabla_w f(w)$
  end
  $\alpha = $ learning rate

Gradient descent: $\bar{w}=\bar{w}-\alpha \nabla_w f(w)$

Stochastic gradient descent: only use 1 data point per iteration to compute $\nabla$. data point randomly sampled from data

Mini-batch gradient descent: uses a batch size of M data pts. at each iteration to compute $\nabla$

e.g. for least-sq. linear regression
$\text{Loss}(h_w) = \frac{1}{2}||y - X\bar{w}||^2$
$\nabla_w \text{Loss} = -X^Ty + X^TX \bar{w}$
$\bar{w} = \bar{w} - \alpha(-X^Ty+X^TX \bar{w})$ (for LS → descent)

# (second column)

allows us to turn linear combination of input features into a probability using logistic fcn:
$h_w(\bar{x})= \dfrac{1}{1+e^{-w^Tx}}$
- single layer NN w/softmax activation is the same as logistic regression

## Neural Networks
- many probs require non-linear solns
- NN's: multi-layer perceptrons that are more expressive
- Universal Fcn Approx. Thm: 2-layer NN w/sufficient # of neurons can approximate any continuous fcn

## Accuracy
- accuracy of binary perceptron after making n predictions:
$$l^{acc}(\bar{w}) = \frac{1}{n}\sum_{i=1}^n (\text{sgn}(w\cdot f(x)) == y_i)$$
indicator fcn that derives feature / weight vector / data point / actual class label of $x_i$

- softmax fcn defines probability of classifying $x^{(i)}$ to class j as:
$$\sigma(x_i)_j = \frac{e^{f(x_i)^T w_j}}{\sum_{l=1}^K e^{f(x_i)^T w_l}} = P[y_i=j|f(x_i);w]$$

- likelihood of a set of weights:
$$l(\bar{w}) = \prod_{i=1}^n P(y_i|f(x_i);w)$$

- log-likelihood (want to maximize this quantity by finding $\bar{w}^*$):
$$\log l(\bar{w}) = \log \prod_{i=1}^n P(y_i|f(x_i);w)$$
$$= \sum_{i=1}^n \log P(y_i|f(x_i);w)$$
↳ maximize by computing the $\nabla$
$$\nabla_w l(w) = \left[\frac{\partial ll(w)}{\partial w_1} \cdots \frac{\partial ll(w)}{\partial w_n}\right]$$

## Backpropagation
- allows you to efficiently compute gradients for each param. in the NN.
- represents NN as a dependency graph of operators & operands (computational graph)

## Backpropagation Algorithm
1) Forward pass: compute output of node
2) Backwards pass: partial derivative of last node wrt variable at current node

Hyperbolic tangent:
$g(t) = \dfrac{e^z - e^{-z}}{e^z + e^{-z}}$



$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial g}\cdot\frac{\partial g}{\partial y}=4$
$\frac{\partial f}{\partial g} = \frac{\partial}{\partial g}(z\cdot g)=z=4$
$\frac{\partial f}{\partial z} = \frac{\partial}{\partial z}(z\cdot g)=g=5$
$\frac{df}{df}=1$

# (third column — DBN)



- M×N grid w/o walls/obstacles
- Pacman trying to figure out location of ghosts at timestep T ($G_T$)
- 3 sensor variables: $D_t, E_t, F_t$
- guess Y = is pacman's guess of $G_T$
- Utility $U(Y, G_T)$ a predetermined fcn that's larger when Y closer to $G_T$

(a) Pacman doesn't observe sensor values & doesn't know any of the (PT's in the DBN. $U_g(Y,g_T)=$ 100 $Y=g_T$, −1 o/w uniform distribution & ghost moves randomly for a long time. what's MEU?

$$MEU(\emptyset) = P[Y=G_T]U(Y=G_T) + P[Y\neq G_T]U(Y\neq G_T)$$
$$= \frac{1}{MN}\cdot 100 + \frac{MN-1}{MN}\cdot(-1)$$

(b) Assume Pacman knows everything all conditional distributions on the DBN. What maximizes his expected utility when no evidence is available?
$$Y = \text{argmax}_{g_T} \sum_{g_{T-1}} P(g_T|g_{T-1}) P(g_{T-1})$$

(c) Pacman can choose to uncover a sensor value for all T to help max utility.
(i) Which of the following is the strongest valid statement about VPI's for sensor variables?
$VPI(E)$ ?, $VPI(F)$ ? $VPI(D)$
(ii) Pacman chooses F. What maximizes E[utility]?
$$Y = \text{argmax}_{g_T} \sum_{c\in l_t} P(f_T|c_T)P(d_T)P(c_T|d_T,g_T)P(g_T|f_{0:T-1})$$
(iii) Pacman chooses E. Using Fcns for the belief distribution & general utility Fcn, what's the expression for $EU(Y|e_{0:T})$?
$$EU(Y|e_{0:T}) = \sum_{g_T} B(G_T=g_T)U(Y,g_T)$$
(iv) Minimum set of variables to maximize expected utility: $G_0, E_{1:T}, D_{1:T}$ (Markov blanket of all ghost variables → would contribute a positive VPI)

# (far right column)

are linear classifiers (can represent any fcn, incl. non-linear ones)

## Loss
- want to predict Q-values as close to $q((x,y),a)$ as possible
  1) MSE loss — works bc it defines squared difference btwn the 2 values we want to be similar
  2) Indicator — (zero-one loss) won't work bc it doesn't give reasonable metric for how close our estimate is & doesn't yield a gradient
  3) Cross-entropy loss isn't suitable since we aren't doing classification

# (bottom middle — decision network)



- can be a (event, subway) decision

- say we eliminated R & C to find $P[E=true]$, find $P[E=true|H=true]$
$$P[E=true|H=true] = \frac{P(H=true|E=true)P(E=true)}{P(H=true)}$$
$$= \frac{P(H=true|E=true)P(E=true)}{P(H=true|E=true)P(E=true) + P(H=true|E=false)P(E=false)}$$
- $EU[\text{event, subway}|H=true] = P[E=true|H=true]U_1(\text{event}, E=true) + P(E=false|H=true)U_1(\text{event}, E=false) + U_2(\text{subway}, H=true)$
- $VPI(H) = P(H) MEU(H) + P(\neg H) MEU(\neg H) - MEU(\emptyset)$
- $MEU(R=true) = \max_{a_1,a_2} \sum_r P(e|r)[U_1(e,a_1) + \sum_h P(h|e)U_2(h,a_2)]$

- What are all possible individual variables s.t. we can decompose MEU(X=x) = MEU_1(X=x) + MEU_2(X=x)? R,C,E,H (by linearity of E)
- Let X be any boolean variable in any decision network. The following conditions are jointly sufficient to show that $VPI(X)>0$: $P(x)$ is not...

## CSPs (Constraint Satisfaction)

- type of identification problem where we have to see if a state is a goal state or not
- defined by 3 factors:
  1) Variables: $(X_1, ..., X_N)$ that can take on a single value from a defined set of values
  2) Domain: set $\{x_1, ..., x_d\}$ representing all possible values a variable can take on
  3) Constraints: define restrictions on values of variables
- CSPs are NP-Hard
- given problem w/ N vars of domain size $O(d)$, for each var, there are $O(d^N)$ possible assignments
- types of constraints
  1) Unary - involve a single variable
     ↳ not repr. in constraint graphs, just used to trim domain
  2) Binary - involve 2 variables
  3) Higher-order - involve ≥2 variables
- **every discrete, finite CSP can be represented as a SAT problem & vice versa, & a correct repr. of a discrete, finite CSP has exactly same # satisfying assign.**

### Solving CSPs
- Backtracking search
  ↳ optimization on DFS, where it also
  1) Fixes an ordering for variables & selects values for variables in that order
  2) Only selects values that don't conflict w/ any prev values. If no values exist, backtrack to prev variable & change its value.
- Filtering: can we detect inevitable failure early? ie, keep track of domains for unassigned variables and cross off bad options
  - Forward checking: whenever a value is assigned to $X_i$, prune the domains of the unassigned variables that share a constraint w/ $X_i$ that would violate that constraint
  - arc consistency: an arc $X \to Y$ is consistent iff $\forall x$ in the tail, $\exists y$ in the head that could be assigned w/o violating a constraint
    ↳ limitations: after enforcing $AC'$, can have 0, 1, or >1 solns left
    → still runs in backtracking search

### Orderings for CSPs
1) MRV: Minimum Remaining Values
   ↳ choose most constrained variable next
2) LCV: Least constraining Value
   ↳ choose the value that prunes the fewest domains of remaining unassigned variables

### CSP Structure
- has $O(n \cdot d^n)$ leaves in search tree
- general runtime is $O(d^N)$ but can be reduced to $O(d^2 n)$ (linear in # of variables) by doing the following:
  1) Pick arbitrary node in the constraint graph to serve as a root node
  2) Convert the undirected graph edges to edges that point away from the root. Linearize/topologically sort the graph s.t. all edges point rightwards
  3) Perform a backwards pass of arc consistency from $i=N$ to $i=2$ for all arcs $Parent(X_i) \to X_i$
  4) Perform a forwards assignment, assigning each $X_i$ a value consistent

### Cutset Conditioning
- find the smallest subset of variables in a graph s.t. their removal results in a tree (a cutset for the graph)
  ↳ leaves us w/ tree w (n-c) variables
    ↳ solvable in $O((n-c)d^2)$
    ↳ runtime of cutset conditioning on a general CSP is $O(d^c(n-c)d^2)$
      ↳ may still need to backtrack $d^c$ times.

### K-Consistency
- 1-consistency (node consistency): each single node's domain has a value that meets its unary constraints
- 2-consistency (arc consistency): for each pair of nodes, any consistent assignment can be extended to another
- K-consistency: For each K nodes, any consistent assignment to k-1 can be extended to the $k^{th}$ node

## MDPs (Markov Decision Processes)
- models used to solve nondeterministic search problems
- Some properties:
  - $\gamma$: discount factor
  - $T(s, a, s')$: Transition function - probability that taking action "a" from state "s" results in "s'"
    ↳ $= P(s' | s, a)$
  - $R(s, a, s')$: Reward function -
  - $Q(s, a)$: Action states
- Finite horizons - defines lifetime for agents before they get terminated
- discount factors - model exponential decay of rewards over time, so instead of maximizing additive utility:
  $$U([s_0, a_0, s_1, a_1, ...]) = R(s_0, a_0, s_1) + ... + R(s_n, a_n, a_{n+1})$$
  we maximize discounted utility:
  $$U([s_0, a_0, s_1, a_1, ...]) = R(s_0, a_0, s_1) + \gamma R(s_1, a_1, s_2) + \gamma^2 R + ...$$
  $$= \sum_{t=0}^{\infty} \gamma^t R(s_t, a_t, s_{t+1})$$
  <span style="color:red">guaranteed to be finite valued as long as $|\gamma| < 1$:</span>
  $$\sum_{t=0}^{\infty} \gamma^t R_{max} = \frac{R_{max}}{1-\gamma}$$
- Markovian/memorylessness encoded in transition function: $T(s, a, s') = P(s' | a, s)$

### Solving MDPs
- want to find optimal policy $\pi^*: S \to A$, a function mapping each state $s \in S$ to action $a \in A$. An explicit policy $\pi(s)$ defines a reflex agent; given a state s, an agent implementing $\pi$ will select $a = \pi(s)$.
- $U^*(s)$ or $V^*(s)$: optimal value of a state s; expected value of the utility of an optimally-behaving agent that starts in s will receive.
- $Q^*(s, a)$: optimal value of a Q-state

### Bellman Equations:
$$U^*(s) = \max_a \sum_{s'} T(s, a, s')[R(s, a, s') + \gamma U^*(s')]$$
$$Q^*(s, a) = \sum_{s'} T(s, a, s')[R(s, a, s') + \gamma U^*(s')]$$
$$U^*(s) = \max_a Q^*(s, a)$$

## Value Iteration
- time-limited value for a state s with a time-limit of K timesteps ($U_K(s)$) represents the maximum expected utility attainable from s given that the MDP terminates in K timesteps. (eg, expectimax depth-k)
- idea: VI is a DP algo that uses an iteratively longer time to compute time-limited values until convergence
- algorithm:
  1) Initialize. $\forall s \in S \quad U_0(s) = 0$
  2) Repeat update rule until convergence: $\forall s \in S$
  $$U_{K+1}(s) \leftarrow \max_a \sum_{s'} T(s, a, s')[R(s, a, s') + \gamma U_K(s')]$$

## Policy Extraction
- if you're in a state s, you should take the action a which yields the max expected utility
  $$\forall s \in S \quad \pi^*(s) = \arg\max_a Q^*(s, a)$$
  $$\pi^*(s) = \arg\max_a \sum_{s'} T(s, a, s')[R(s, a, s') + \gamma U^*(s')]$$

### Q-Value Iteration
$$Q_{K+1}(s, a) \leftarrow \sum_{s'} T(s, a, s')[R(s, a, s') + \gamma \max_{a'} Q_K(s', a')]$$

## Policy Iteration
- algorithm:
  1) Define an initial policy. Can be arbitrary but converges faster the closer the initial policy is to optimal policy.
  2) Repeat the following until convergence:
     i) Policy Evaluation: compute $U^\pi(s) \ \forall s \in S$, until convergence
     $$U^\pi(s) = \sum_{s'} T(s, \pi(s), s')[R(s, \pi(s), s') + \gamma U^\pi(s')]$$
     $$U_{K+1}^{\pi(s)} \leftarrow \sum_{s'} T(s, \pi_i(s), s')[R(s, \pi_i(s), s') + \gamma U_K^{\pi_i(s')}]$$
     2) Policy Improvement: For fixed values, get a better policy using policy extraction:
     $$\pi_{i+1}(s) = \arg\max_a \sum_{s'} T(s, a, s')[R(s, a, s') + \gamma U^{\pi_i}(s')]$$
     ↳ if $\pi_{i+1} = \pi_i$, we have converged ⟹ $\pi_i = \pi^*$

## MDP Recap
- Value iteration: Used to compute optimal values of states by iterative updates until convergence
- Policy evaluation: compute values for a particular policy
- Policy extraction: turn your values into a policy
- Policy iteration: compute optimal values for a particular policy
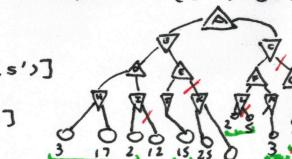
## Alpha-Beta Pruning
- Minimax → $O(b \times b \times \cdots \times b) = O(b^d)$ [branching factor b, depth d]
  ↳ also worst case $\alpha$-$\beta$
- best case:
  → odd depth: $O(b \times 1 \times b \times 1 \times \cdots \times b)$
  → even depth: $O(b \times 1 \times b \times 1 \times \cdots \times b \times 1) = O(b^{d/2}) = O(\sqrt{b^d})$

- $O(|S||A|)$

### Approx. Q-learning
Linear value fcns:
$$V(s) = w_1 f_1(s) + w_2 f_2(s) + \cdots + f_n(s) w_n$$
$$= \vec{w} \cdot \vec{f}(s)$$
$$Q(s, a) = w_1 f_1(s, a) + \cdots + w_n f_n(s, a)$$
$$= \vec{w} \cdot \vec{f}(s, a)$$
$$\text{difference} = [R(s, a, s') + \gamma \max_{a'} Q(s', a')] - Q(s, a)$$
$$w_i \leftarrow w_i + \alpha \cdot \text{difference} \cdot f_i(s, a)$$

exact Q-learning:
$$Q(s, a) = Q(s, a) + \alpha \cdot \text{difference}$$

### $\epsilon$-greedy
- explore randomly w.p. $\epsilon$
- exploit w.p. $(1-\epsilon)$
- exploration fcns:
$$Q(s, a)$$
$$\leftarrow (1-\alpha)Q(s,a) + \alpha[R(s,a,s') + \gamma \max_{a'} f(s', a')]$$
where f is an exploration fcn.
- common choice for $f(s,a) = Q(s,a) + \frac{K}{N(s,a)}$
  K := predetermined value
  $N(s,a) = $ # times Q-state $(s,a)$ was visited