

# Lectures

## 1. Smoov & Curly's Bogus Journey

### MDP

state:  $S$

model:  $T(s, a, s') = \text{Prob}(s' | s, a)$

actions:  $A(s)$

reward:  $R(s), R(s,a), R(s,a,s')$

policy:  $\pi(s) \rightarrow a$

optimal policy:  $\text{argmax}(\text{long term reward})$

### Markovian properties

1. memoryless: model only depends on current state (no condition on anything past current state)
2. stationary: rules ( $T$ ) doesn't change over time

### Rewards

1. delayed reward (e.g. you can lose at first step in chess)  
credit assignment: which action caused delayed reward
2. minor changes matter

### Sequences of Rewards (assumptions of this course)

1. infinite horizons (stationary policies): if an agent doesn't live forever, it'll change its policy (max reward in limited steps) and it'll have different policy even if in the same state. without infinite horizon, we lose stationarity in policy
2. utility of sequences (stationary preferences): if  $U(s_0, s_1, s_2...) > U(s_0, s_1', s_2'...)$  then  $U(s_1, s_2...) > U(s_1', s_2'...)$ . this is why we do addition over rewards, otherwise this assumption may not hold

### Discounted Rewards

$$U(s_0, s_1, ...) = \sum (\gamma^t * R(s_t)) \quad [t=0 \dots \infty, 0 \leq \gamma < 1]$$

### Policies

optimal policy:

$\pi^* = \text{argmax}_{\pi} (E[\sum_t (\gamma^t * R(s_t)) | \pi])$  policy that maximizes long-term expected rewards

true utility of a state:

$U_{\pi}(s) = E[\sum_t (\gamma^t * R(s_t)) | \pi, s_0=s]$  expected rewards if we follow the optimal policy

$\pi^*(s) = \operatorname{argmax}_a (\sum_{s'} [T(s,a,s') * U(s')])$

reward: immediate feedback

utility: long-term reward

### Bellman equation!

$U(s) = R(s) + \gamma \max_a (\sum_{s'} [T(s,a,s') * U(s')])$  true value of a state

### Value Iteration

1. start with arbitrary U
2. update U based on neighbors
3. repeat until converge
4.  $U_{t+1}(s) = R(s) + \gamma \max_a (\sum_{s'} [T(s,a,s') * U_t(s')])$

### Policy iteration

1. start with arbitrary policy  $\pi_0$
2. evaluate: given  $\pi_t$ , calculate  $U_t = U(\pi_t)$
3. improve:  $\pi_{t+1} = \operatorname{argmax}_a (\sum_{s'} [T(s,a,s') * U_t(s')])$
4.  $U_{t+1}(s) = R(s) + \gamma \sum_{s'} [T(s, \pi_t(s), s') * U_t(s')]$

### Three forms of Bellman Equation

1. value:  $V(s) = \max_a [R(s,a) + \gamma \sum_{s'} (T(s,a,s') * V(s'))]$
2. quality:  $Q(s, a) = R(s,a) + \gamma \sum_{s'} [T(s,a,s') * \max Q(s', a)]$
3. continuation:  $C(s,a) = \gamma \sum_{s'} (T(s,a,s') * \max_a [R(s',a') + C(s, a)])$

when we know Q, we don't need T to get V

### Model-based vs Model-free

"model": transition matrix.

**model-based**: use model predictions of next state & reward to calculate optimal policy; dynamic programming (e.g. policy iteration, value iteration requires learning T)

**model-free**: purely sample from experience (MC-Control, SARSA, Q-learning)

### On-policy vs Off-policy

"policy": optimal actions

**on-policy**: attempts to evaluate / improve the policy & behave based on the policy (e.g. SARSA updates Q with Q associated with the policy it followed)

**off-policy**: doesn't behave on evaluated policy (e.g. Q-learning updates Q with maxQ instead of whichever action agent just took)

## 2. Reinforcement Learning Basics

agent --a --> env

agent <-- s, r-- env

## Behavior Structures

**plan:** fixed sequence of actions

**conditional plan:** if... then..

**stationary policy/universal plan:** map(state) -> actions (pro: always an optimal stationary policy; con: too many ifs, inefficient)

## Evaluate a policy

1. state immediate reward
2. truncate horizon
3. summarize sequence
4. summarize over sequences

## Evaluate a policy

1. value of returned policy
2. computational complexity (computation time)
3. experience complexity (sample size, data size needed) <important in rl>
4. space complexity

## 3. TD and Friends

### 3 types of models

1. model-based

$\langle s, a, r \rangle^* \rightarrow \text{model\_learner} \leftrightarrow T/R \rightarrow \text{MDP solver} \rightarrow Q^* \rightarrow \text{argmax} \rightarrow \pi$

2. value-function-based "model-free" (doesn't know T, learn Q from sar

$\langle s, a, r \rangle^* \rightarrow \text{value update} \leftrightarrow Q \rightarrow \text{argmax} \rightarrow \pi$

3. policy search

$\langle s, a, r \rangle^* \rightarrow \text{policy update} \leftrightarrow \pi$  (difficult, feedback from  $\pi$  is hard to use to update policy)

- more direct learning  $3 > 2 > 1$

- more supervised  $1 > 2 > 3$

### Computing estimates incrementally

$$V_T(S_i) = [(T-1) \cdot V_{T-1}(S_i) + R_T(S_i)] / T \rightarrow \text{(take T as 1/weight\_on\_reward)}$$
$$= V_{T-1}(S_i) + \alpha_T \cdot [R_T(S_i) - V_{T-1}(S_i)] \quad \text{where } \alpha_T = 1/T$$

### Properties of Learning Rate

for learning function  $V_T(S) = V_{T-1}(S) + \alpha_T \cdot [R_T(S) - V_{T-1}(S)]$

(notice:  $R_T$  is observed reward,  $V_{T-1}$  is estimation at S)

then  $V_T(s)$  **converges** to actual V

limit  $T \rightarrow \infty$   $[V_T(S)] = V(s)$  if:

1.  $\sum(\alpha_T) = \infty$  (sum alpha diverges)
2.  $\sum(\alpha_T^2) < \infty$  (sum squared alpha converges)

### Select learning rate

if  $p > 1$ ,  $\sum(1/x^p)$  converges!

### TD(1) Rule

for each episode (T)

for all S, initialize eligibility=0,  $V(S)$  as  $V(S)$  at T-1 (last episode),

after transit, add 1 to eligibility of the state just left ( $s_{t-1}$ )

for all S, update  $V(S)$  with TD, learning rate alpha

decay eligibility

### TD(1) Rule

Episode T

For all  $s$ ,  $e(s) = 0$  at start of episode,  $V_T(s) = V_{T-1}(s)$

After  $s_{t-1} \xrightarrow{f_t} s_t$  : (step t)

$e(s_{t-1}) = e(s_{t-1}) + 1$

For all  $s$ ,  
 $V_T(s) = V_{T-1}(s) + \alpha_T (r_t + \gamma V_{T-1}(s_t) - V_{T-1}(s_{t-1})) e(s)$   
 $e(s) = \gamma e(s)$

note: in  $r_t + \gamma V_{T-1}(s_t) - V_{T-1}(s_{t-1})$ ,  $s_t$ : new state,  $s_{t-1}$ : last state. this part is shared by all states in  $V_T(s)$  updating!

claim: TD(1) is the same as outcome-based updates (if no repeated states), wait till see all rewards and compare with est. value at the beginning of episode;

but if there were repeated states:

outcome based learning: ignores what happens during the episode

TD(1) : picks up what it learns the last time it saw the state <EXTRA LEARNING!>

why is TD(1) wrong?

TD(1) not using all data points, only uses individual run. when it gets edge cases, it gets wrong estimate

### TD(0) Rule

if finite data repeated indefinitely, TD(0) finds max.likelihood estimate; without repeats, it doesn't perform well;

$$V_T(s_{t+1}) = V_T(s_t) + \alpha_T (r_t + \gamma V_T(s_t) - V_T(s_t))$$

$$V_T(s_{t+1}) = E[r_t + \gamma V_T(s_t)]$$

for TD(1) (outcome-based updates), repeat the sample over n over doesn't make difference because all data points are already used

Instead of updating all S, TD(0) only updates  $s_{t-1}$

### TD(lambda)

x lambda to eligibility each step

### k-step Estimators

TD(0) one step forward looking;

TD(1) all steps forward looking

### TD(lambda) Rule

Episode T

For all  $s$ ,  $e(s) = 0$  at start of episode,  $V_T(s) = V_{T-1}(s)$

After  $s_{t-1} \xrightarrow{f_t} s_t$  : (step t)

$e(s_{t-1}) = e(s_{t-1}) + 1$

For all  $s$ ,  
 $V_T(s) = V_{T-1}(s) + \alpha_T (r_t + \gamma V_{T-1}(s_t) - V_{T-1}(s_{t-1})) e(s)$

$e(s) = \gamma e(s)$

TD(1)

TD(0)

Episode T

For all  $s$ ,  $e(s) = 0$  at start of episode,  $V_T(s) = V_{T-1}(s)$

After  $s_{t-1} \xrightarrow{f_t} s_t$  : (step t)

$e(s_{t-1}) = e(s_{t-1}) + 1$

For all  $s$ ,  
 $V_T(s) = V_{T-1}(s) + \alpha_T (r_t + \gamma V_{T-1}(s_t) - V_{T-1}(s_{t-1})) e(s)$

$e(s) = \gamma e(s)$

Episode T

For all  $s$ ,  $e(s) = 0$  at start of episode,  $V_T(s) = V_{T-1}(s)$

After  $s_{t-1} \xrightarrow{f_t} s_t$  : (step t)

For all  $s$ ,  
 $V_T(s) = V_{T-1}(s) + \alpha_T (r_t + \gamma V_{T-1}(s_t) - V_{T-1}(s_{t-1}))$

$V_T(s) = V_{T-1}(s) + \alpha_T (r_t + \gamma V_{T-1}(s_t) - V_{T-1}(s_{t-1}))$

K-Step Estimators

$$\begin{aligned}
 E_1 \quad V(s_t) &= V(s_t) + \alpha_r (r_{t+1} + \gamma V(s_{t+1}) - V(s_t)) \quad \text{TD}(0) \\
 E_2 \quad V(s_t) &= V(s_t) + \alpha_r (r_{t+1} + \gamma r_{t+2} + \gamma^2 V(s_{t+2}) - V(s_t)) \\
 E_3 \quad V(s_t) &= V(s_t) + \alpha_r (r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \gamma^3 V(s_{t+3}) - V(s_t)) \\
 E_k \quad V(s_t) &= V(s_t) + \alpha_r (r_{t+1} + \dots + \gamma^{k-1} r_{t+k} + \gamma^k V(s_{t+k}) - V(s_t)) \\
 E_{\infty} \quad V(s_t) &= V(s_t) + \alpha_r (r_{t+1} + \dots + \gamma^{k-1} r_{t+k} + \dots \quad \text{TD}(0) \quad V(s_t))
 \end{aligned}$$

**TD(lambda) is weighted avg of k-step estimators**

$$\text{sum}_{k=1..inf} [\lambda^{k-1} * (1-\lambda)]$$

TD(lambda) empirical performance: sweet spot is in (0,1) interval

### [Reading] MC - sutton rl 5.1

Monte Carlo method: average the returns observed after visits to that state. As more returns are observed, the average should converge to the expected value

estimates each state independently

### [Reading] MC control - sutton rl 5.3

control: approximate optimal policy

## 4. Convergence

### Bellman Equation with Control

approximations in Q-learning update:

if we have the model, update with Bellman equation;

else, use TD(0) rule to update Q with Q\*

recall property of alpha: sum(alpha) diverge but sum(alpha^square) finite then Q converge

### Bellman Equations with Actions

actions:  $Q(s,a) = R(s,a) + \gamma \sum_{s'} T(s,a,s') \max_{a'} Q(s',a')$

$\langle s_{t+1}, a_{t+1}, r_t, s_t \rangle$

$$Q_t(s_{t+1}, a_{t+1}) = Q_{t-1}(s_{t+1}, a_{t+1}) + \alpha_t (r_t + \gamma \max_{a'} Q_{t-1}(s_{t+1}, a') - Q_{t-1}(s_{t+1}, a_{t+1}))$$

$Q_t(s,a) = Q_{t-1}(s,a)$ , otherwise TD(0)

Approximations

① If we know the model, synchronously update  $Q_t(s,a) = R(s,a) + \gamma \sum_{s'} T(s,a,s') \max_{a'} Q_{t-1}(s',a')$

② If we know Q\*, sampling asynchronously update  $Q_t(s_{t+1}, a_{t+1}) = Q_{t-1}(s_{t+1}, a_{t+1}) + \alpha_t (r_t + \gamma \max_{a'} Q^*(s_{t+1}, a') - Q_{t-1}(s_{t+1}, a_{t+1}))$

### Bellman Operator

mapping from Q to Q

$$[BQ](s,a) = r(s,a) + \gamma \sum_{s'} [T(s,a,s') * \max_{a'} Q(s',a')]$$

$BQ^* = Q^*$ : Bellman Equation

$Q_t = BQ_{t-1}$ : Value iteration

### Contraction Mapping

B is an operator,

if for all value function F, G and some  $0 \leq \gamma \leq 1$

$$\|BF - BG\|_{\infty} \leq \gamma * \|F - G\|_{\infty}$$

where  $\|Q\|_{\infty} = \max |Q(s,a)|$

then B is a contraction mapping (makes things closer)

non-expansion is a kind of contraction mapping

## Contraction Mapping Properties

if B is a contraction mapping

- 1)  $F^* = BF^*$  has a solution and it's unique
- 2)  $F_t = BF_{t-1} \dots F_t \rightarrow F^*$  (repeatedly apply contraction mapping, function converges; value iteration converges!)

unique because if there are two then they're not converging together which violates contraction mapping definition

since Bellman operator is contraction mapping, Bellman Equation only has one solution

## Bellman Operator Contracts

proof:---->

to remove max in  $|\cdot|$  on right work, we need to prove that we can remove max and the diff is still bounded.

*Bellman Operator Contracts*

Given  $Q_1, Q_2$

$$\begin{aligned} \|BQ_1 - BQ_2\|_\infty &= \max_{s,a} |[BQ_1](s,a) - [BQ_2](s,a)| = \\ &= \max_{s,a} \left| \left( R(s,a) + \gamma \sum_{s'} T(s,a,s') \max_{a'} Q_1(s',a') \right) - \left( R(s,a) + \gamma \sum_{s'} T(s,a,s') \max_{a'} Q_2(s',a') \right) \right| \\ &= \max_{s,a} \left| \gamma \sum_{s'} T(s,a,s') \left( \max_{a'} Q_1(s',a') - \max_{a'} Q_2(s',a') \right) \right| \\ &\leq \gamma \max_{s,a} \left| \sum_{s'} T(s,a,s') \right| \max_{s',a'} |Q_1(s',a') - Q_2(s',a')| \\ &\leq \gamma \max_{s',a'} |Q_1(s',a') - Q_2(s',a')| = \gamma \|Q_1 - Q_2\|_\infty \end{aligned}$$

## Max is a non-expansion

proof:

*Max is a non-expansion*

For all  $f, g$   $\left| \max_a f(a) - \max_a g(a) \right| \leq \max_a |f(a) - g(a)|$

Proof: WLOG  $\max_a f(a) \geq \max_a g(a)$

$$\begin{aligned} \left| \max_a f(a) - \max_a g(a) \right| &= \max_a f(a) - \max_a g(a) \\ &= f(a_1) - g(a_2) \quad \left\{ \begin{array}{l} a_1 = \arg \max_a f(a) \\ a_2 = \arg \max_a g(a) \end{array} \right. \\ &\leq f(a_1) - g(a_1) \\ &= |f(a_1) - g(a_1)| \\ &\leq \max_a |f(a) - g(a)| \end{aligned}$$

## Convergence

**Theorem** (proves Q-learning converges!)

let B be contraction mapping with fixed point  $Q^* = BQ^*$  (has an unique solution <--property of contraction mapping),

$Q_0$ : Q function

update  $Q_{t+1} = [B_t Q_t] Q_t$

then Q converges if:

1.  $[B_t Q_t]$  is non-expansion function (e.g. max|| in above proof)
2.  $B_t$  is a contraction mapping
3.  $\sum_t \alpha_t = \infty$ ,  $\sum_t \alpha_t^2 < \infty$

*Convergence* Define  $\alpha_t(s,a) = 0$  if  $s_t \neq s, a_t \neq a$

Theorem: Let B be a contraction mapping and  $Q^* = BQ^*$  be its fixed point. Let  $Q_0$  be a Q function and define  $Q_{t+1} = [B_t Q_t] Q_t$ .

Then,  $Q_t \rightarrow Q^*$  if:

1. For all  $u_1, u_2, s, a$ :  $\left| ([B_t u_1] Q^*)(s,a) - ([B_t u_2] Q^*)(s,a) \right| \leq (1 - \alpha_t(s,a)) |u_1(s,a) - u_2(s,a)|$
2. For all  $Q, u, s, a$ :  $\left| ([B_t u] Q^*)(s,a) - ([B_t u] Q)(s,a) \right| \leq \alpha_t(s,a) |Q^*(s,a) - Q(s,a)|$
3.  $\sum_t \alpha_t = \infty, \sum_t \alpha_t^2 < \infty$

(hidden condition: must visit (s,a) indefinitely often)

ps. in Q-learning, one Q is 1-step look-forward and the other Q is avg.

## Convergence Theorem Explained:

Convergence Theorem Explained

$$([B_t Q] w)(s,a) = Q(s,a) + d_t(s,a) (r_t + \gamma \max_{a'} w(s_t, a') - Q(s,a))$$

1. For all  $U_1, U_2, s, a$  :  $|([B_t U_1] Q^*)(s,a) - ([B_t U_2] Q^*)(s,a)| \leq (1 - \alpha_t(s,a)) |U_1(s,a) - U_2(s,a)|$

Convergence Theorem Explained

$$([B_t Q] w)(s,a) = Q(s,a) + d_t(s,a) (r_t + \gamma \max_{a'} w(s_t, a') - Q(s,a))$$

2. For all  $Q, U, s, a$  :  $|([B_t U] Q^*)(s,a) - ([B_t U] Q)(s,a)| \leq \gamma \alpha_t(s,a) |Q^*(s,a) - Q(s,a)|$

CONTRACTION

use Q function in two ways in convergence theorem:

1. w: acts as if we knew  $Q^*$ , plays the role of avg (sumT) → avg out stochasticity
2. Q: one-step look-back (in TD formula, prediction made at last state) → value iteration obj

## Generalized MDP

Generalized MDPs

$$Q^*(s,a) = R(s,a) + \gamma \bigoplus_{s'} \bigotimes_{a'} Q^*(s',a')$$

$\bigoplus_{s'} g(s') = \sum_{s'} T(s,a,s') g(s')$	$\bigotimes_{a'} f(s',a') = \max_{a'} f(s',a')$	regular MDP
$\bigoplus_{s'} g(s') = \min_{s' \in \{s'   T(s,a,s') > 0\}} g(s')$	$\bigotimes_{a'} f(s',a') = \max_{a'} f(s',a')$	pessimistic MDP risk aware
$\bigoplus_{s'} g(s') = \sum_{s'} T(s,a,s') g(s')$	$\bigotimes_{a'} f(s',a') = \sum_{a'} p(a a') f(s',a')$	exploration-sensitive max, min, mixture
$\bigoplus_{s'} g(s') = \sum_{s'} T(s,a,s') g(s')$	$\bigotimes_{a'} f(s',a') = \min_{a'} f(s',a')$	zero-sum game

## 5. AAA (Advanced Algorithm Analysis)

### More on VI

1. For some  $t^* < \infty$ , that's polynomial in  $|S|, |A|, R_{\max} = \max_{a,s} |R(s,a)|$ ,  $1/(1-\gamma)$ , number of bits of precision that are used to specify transition prob.

after  $t^*$  many steps,

$\pi(s) = \arg\max_a [Q_{t^*}(s,a)]$  is optimal

"some time before infinity, with value iteration, we'll get a Q that's close enough to optimal Q, such that when we take greedy policy respect to it, it's 100% optimal"

More on VI

1. For some  $t^* < \infty$  polynomial in  $|S|, |A|, R_{\max} = \max_{a,s} |R(s,a)|$ ,  $1/(1-\gamma)$ , bits of precision,  $\pi(s) = \arg\max_a Q_{t^*}(s,a)$  is optimal. Cramer's rule
2. If  $|V_k(s) - V_{k+1}(s)| < \epsilon \quad \forall s$   $H \sim \frac{1}{1-\gamma}$   
 $\max_s |V^k(s) - V^*(s)| < \frac{2\epsilon\gamma}{1-\gamma}$  swartighted
3.  $\|B^k Q_1 - B^k Q_2\|_{\infty} \leq \gamma^k \|Q_1 - Q_2\|_{\infty}$

\* but not polynomial bounded because of  $1/(1-\gamma)$  if  $\gamma = 1$

"greedy policy with respect to VI converges in reasonable time".

2. if two consecutive value estimation is small, value of greedy policy with respect to VI converges to optimal value

both 1&2 encourages us pick small gamma

- horizon =  $1/(1-\gamma)$  (small gamma  $\rightarrow$  shorter future  $\rightarrow$  easier problem  $\rightarrow$  shortsighted)

- iteration time

hence we still pick large gamma in practice

3. helps us decide how long to run VI

## Linear Programming

optimize linear object function under linear constraints

Bellman equation can't directly be solved by LP because of max

e.g.  $\max(-3, 7, 2, 5) = x \iff x \geq -3 \ \& \ x \leq 7 \ \& \ x \geq 2 \ \& \ x \leq 5$  (constraints) AND  $\min(x)$  (obj)

translate max in Bellman equation to linear obj & constraints ("PRIMAL")

obj:  $\min(\sum_s (V_s))$

constraints:

for any  $s, a$ ,

$$V_s \geq R(s, a) + \gamma \sum_{s'} T(s, a, s') V_{s'}$$

solve ^ with LP solver  $\Rightarrow$  get optimal V

choose greedy policy  $\Rightarrow$  get optimal policy

## Dual Linear Programming

add policy flow into  $\langle s, a \rangle$  pair and let the flow go thru so that we get max reward

concentrate more on policy

## Policy Iteration

initialize  $\rightarrow$  policy improvement  $\rightarrow$  policy evaluation

converges in finite time, at least as fast as VI (step-wise).

tradeoff? computation time

in policy evaluation step, we're doing an entire VI

converge time is

$\geq$  linear  $|S|$  AND  $\leq$  exponential  $|A|^{|S|}$

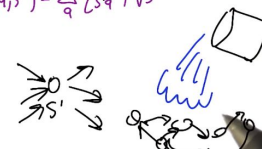
Linear Programming

Max  $\sum_s \sum_a q_{sa} R(s, a)$  "policy flow"

s.t.  $1 + \gamma \sum_s \sum_a q_{sa} T(s, a, s') = \sum_a q_{sa} V_{s'}$

$\forall s, a \quad q_{sa} \geq 0$

DUAL



Policy Iteration

initialize  $V_s, Q_0(s) = 0$

policy improvement  $V_s \pi_t(s) = \arg \max_a Q_t(s, a)$  (+20)

policy evaluation  $Q_{t+1} = Q^{\pi_t}$

open question convergence time

$\geq$  linear  $|S|$   $\leq$  exponential  $|A|^{|S|}$

1.  $Q_t \rightarrow Q^*$
2. Convergence is exact & complete in finite time.
3. Converges at least as fast as VI

Domination

$\pi_1 \geq \pi_2$  iff  $\forall s \in S \quad V^{\pi_1}(s) \geq V^{\pi_2}(s)$

strict domination

$\pi_1 > \pi_2$  iff  $\forall s \in S \quad V^{\pi_1}(s) > V^{\pi_2}(s)$



## Domination

$\pi_1 \succeq \pi_2$  iff for all  $s$   $V_{\pi_1}(s) \geq V_{\pi_2}(s)$

strict domination:

exists  $s$  that's  $V_{\pi_1}(s) > V_{\pi_2}(s)$

## Why Does Policy Iteration Work? (PI Property)

1. Bellman operator is **monotonic**. proof->
2. if we apply greedy policy one step to  $Q_t$ , it'll be equal to  $Q_t$  if not making it better (**greedy policy value improvement**)
3. PI doesn't get stuck in local optima  
 $\wedge$  because of value improvement property, greedy policy must strictly dominate at least one state in each iteration. If not, we must have found the optimal policy.  
 it's impossible to go back forth between 2 values (stuck) because of non-deprovement

## Policy iteration Proof

greedy policy converges to fixed point, optimal policy

**Policy Iteration Proof**

$\pi_1$  policy,  $B_1$  operator,  $Q_1$  value function wrt  $\pi_1$   
 $\pi_2$  greedy wrt  $Q_1$ ,  $B_2$  operator,  $Q_2$  value function wrt  $\pi_2$

non-deprovement

$B_2 Q_1 \geq Q_1$

$k \geq 0, B_2^{k+1} Q_1 \geq B_2^k Q_1$

$\lim_{k \rightarrow \infty} B_2^k Q_1 \geq Q_1$

$Q_2 \geq Q_1$

value improvement

monotonicity

transitivity

fixed point

$a \geq b, b \geq c \Rightarrow a \geq c$

$B_2$  is Monotonic

$$[B_1 V](s) = R(s, \pi_1(s)) + \gamma \sum_{s'} T(s, \pi_1(s), s') V(s')$$

$$[B_2 V](s) = R(s, \pi_2(s)) + \gamma \sum_{s'} T(s, \pi_2(s), s') V(s')$$

Thm: If  $V_1 \geq V_2$ , then  $B_2 V_1 \geq B_2 V_2$

Proof:  $(B_2 V_1 - B_2 V_2)(s) = \gamma \sum_{s'} T(s, \pi_2(s), s') (V_1(s') - V_2(s')) \geq 0$

Another Property in Policy Iteration

$\pi_1 \rightarrow B_1 \rightarrow Q_1 = B_1 Q_1$   
 $\pi_2$  greedy policy wrt  $Q_1$   
 $\downarrow B_2$   
 $Q_1 \leq B_2 Q_1$

$$Q_1(s, a) = R(s, a) + \gamma \sum_{s'} T(s, \pi_1(s), s') Q_1(s', \pi_1(s'))$$

$$\leq R(s, a) + \gamma \sum_{s'} T(s, \pi_2(s), s') Q_1(s', \pi_2(s'))$$

$Q_1 \leq B_2 Q_1$

## 6. Messing with Reward

### Why do we want to change reward function?

1. easier to solve/ repeat and similar to what it would have learned
  - faster - solvable
2. make agent behave the way we want it to through reward
3. efficiency:
  - speed (experience, computation)
  - solvability

How can we change MDP reward function without changing the optimal policy?

1. multiply by positive constant

if  $R'(s,a) = R(s,a)$  then  $Q'(s,a) = cQ(s,a)$

2. shift by constant (add)

if  $R'(s,a) = R(s,a)$  then  $Q'(s,a) = Q(s,a) + c/(1-\gamma)$

3. non-linear potential-based

potential-based shaping

## Reward Shaping

how to train a dolphin to jump through a hoop?

near hoop  $\rightarrow$  through hoop  $\rightarrow$  in air  $\rightarrow$  higher

soccer game reward function?

must give hints

score in goal + 100, near ball +1/distance, hitting ball +10, near goal +1/distance

what would happen? it causes spurious behavior (keep hitting the ball and moving away)

must avoid suboptimal positive loop!

## Potential-based Shaping

put bonus on certain states

by changing from state to state, how much more potential does an agent get?

## Change-in-State-based Bonuses

keep track of which state agent's in, when in desirable states, + bonus, otherwise - bonus

agent gets the difference in rewards between states, which can't accumulate

score in goal +100 near ball  $\Delta(1/\text{distance} - 1/\text{last\_distance})$

## Potential-based Shaping

if we subtract leaving-state's potential and add discounted arriving-state's potential when an agent moves from  $s$  to  $s'$ , new  $Q'(s,a) = Q(s,a) - \psi(s)$ .

this can be interpreted as, the value of a state is the expected total reward of the state minus potential loss of leaving the state

## Q-learning with Potentials

if set  $\psi = \text{optimal value } \max_a Q^*(s,a)$ , then  $Q(s,a) = 0$

if  $a$  is optimal else  $< 0$

### Potential-based Shaping

$$Q(s,a) = \sum_{s'} T(s,a,s') (R(s,a,s') + \gamma \max_{a'} Q(s',a'))$$

$$R'(s,a,s') = R(s,a) - \psi(s) + \gamma \psi(s')$$

$$Q'(s,a) = Q(s,a) - \psi(s)$$

infinite telescoping sum

$$Q = R + \gamma R + \gamma^2 R + \gamma^3 R + \dots$$
$$= R + \gamma R + \gamma^2 R + \gamma^3 R + \dots$$
$$= R + \gamma R + \gamma^2 R + \gamma^3 R + \dots$$
$$= R + \gamma R + \gamma^2 R + \gamma^3 R + \dots$$

### Q-learning with Potentials

$$\langle s, a, r, s' \rangle$$

$$Q(s,a) := Q(s,a) + \alpha (r - \psi(s) + \gamma \max_{a'} Q(s',a') - Q(s,a))$$

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

weird

$$Q(s,a) = Q^*(s,a) - \psi(s)$$

$$Q(s,a) = 0 \text{ if } a \text{ optimal}$$
$$< 0 \text{ otherwise}$$

strategy: explore and discover  $< 0$  then take optimal action

Q-learning with potential == Q-learning with Q-function initialized what potential would've been (action choices, same strategy at each step; but maybe different Q-values along the way)

injecting knowledge (initialize Q with potential instead of randomly), may improve performance (quicker convergence)

potential-based reward shaping preserves the original MDP's optimal policy

potential-based shaping learning time also depends on potential selection. it could take longer / stuck in suboptimal for a while before eventually finds optimal

## 7. Exploring Exploration

### Topics

bandits (no state transition, full stochasticity)

deterministic MDP (state transition, no stochasticity)

stochastic MDP (both)

### [Bandits] Confidence-based Exploration

two bandits (#payoff/#pull: 0/1, 1/2) what to do?

every time new data comes in, do following & pick:

maximize likelihood->always b (but never visit a)

maximize confidence->always b (pick high denominator, never a)

minimize confidence->alternate a&b (but never get to use the info collected)

better solution: balance b/w maximize likelihood & minimize confidence

"exploration-exploitation dilemma"

### Metrics for Bandits

1. identify optimal arm in the limit (can pull arms in order, but may take a long time to learn)

2. maximize (discounted) expected reward (collect reward along the way)

Gittins index:

success/pull-> map to an index-> choose greedily->optimal

^ opportunity cost: how high would the reward be (expected reward of bandit) if you were to stop trying

combines max likelihood & low-confidence

but difficult to generalize

3. maximize expected reward over finite horizon

computable but expensive

finite=1 == max likelihood

4. identify **near-optimal** arm (epsilon) with high probability  $(1-\delta)$  in time  $T(k, \epsilon, \delta)$ ,

# pulls before we identify near-optimal arm is **polynomial** in  $k, 1/\delta, 1/\epsilon$   
(figure out the optimal arm)

PAC: probably approximately correct

PAO: probably approx optimal

5. nearly maximize reward (epsilon) with high probability (1-delta) in time  $T'(k, \epsilon, \delta)$ ,

# pulls before near-optimal: polynomial in  $k, 1/\delta, 1/\epsilon$

(collect max reward on the way figuring it out, we can take longer if reward is high enough)

6. pull a non-near-optimal arm (epsilon) no more than  $T''(k, \epsilon, \delta)$  with high probability (1-delta)

**4&5&6 are equivalent! (proofs)**

Find Best  $\Rightarrow$  Few mistakes

Algorithm that finds  $\epsilon$  best with probability  $1-\delta$  in  $\tau(k, \epsilon, \delta)$  pulls. (A)

Want  $\epsilon'$  close arm on  $\tau'(k, \epsilon', \delta')$  pulls with probability  $1-\delta'$ .

$\epsilon = \epsilon', \delta = \delta'$ , run algorithm A and after  $\tau(k, \epsilon, \delta)$  pulls we know  $\epsilon'$  close arm.

Pull it forever.  $\tau'(k, \epsilon', \delta') \leq \tau(k, \epsilon, \delta)$

Few mistakes  $\Rightarrow$  Do well

We have an algorithm that pulls  $\epsilon$ -suboptimal arms at most  $\tau(k, \epsilon, \delta)$  times with probability  $1-\delta$ . B.

Want an algorithm  $\overset{C}{\text{that}}$ , with high probability  $1-\delta'$ , has a per-step reward  $\epsilon'$  close to optimal after  $\tau'(k, \epsilon', \delta')$  steps.

$\epsilon = \epsilon'/2$

$$\tau(k, \epsilon, \delta') = m \quad \epsilon' \geq \frac{1}{2} \cdot m + \frac{\epsilon'(n-m)}{2}$$

mistakes  $\epsilon'/2$

$$\tau'(k, \epsilon', \delta') = n$$

$$\begin{aligned} n\epsilon' &\geq m + \frac{\epsilon'}{2}n - \frac{\epsilon'm}{2} \\ \frac{n}{2}\epsilon' &\geq m - \frac{\epsilon'm}{2} \\ n &\geq \frac{2}{\epsilon'}(m - \frac{\epsilon'm}{2}) \end{aligned}$$

Do well  $\Rightarrow$  Find Best

We have an algorithm  $\overset{C}{\text{that}}$  gets within  $\epsilon$  (per step) of optimal after  $\tau(k, \epsilon, \delta)$  pulls

## Hoeffding Bound

- maximum likelihood confidence interval
- helps decide # pulls to get confident result
- used in PAC learning

### Strategy: combining arm Info

take  $c$  samples of each arm, then choose the arm with best estimate

to be  $1-\delta$  sure it's within epsilon,  
we need to learn each arm so that its estimate is within  
epsilon/2 with probability  $1-\delta/k$

interpretation:

$P(\text{all } k\text{-arm right}) = (1-\delta/k)^k$  <-assumes independence b/w arms

$P(\text{at least 1 arm wrong}) = k*\delta/k$  <- union bound, assumes overlap

$P(\text{all arms right}) = 1-\delta$

<bonferroni correction>

### How many samples?

solve  $n$  from confidence interval

How Many Samples?

want to be  $\epsilon/2$  accurate with probability  $1-\delta/k$

$$\sqrt{\frac{\frac{1}{2} \ln\left(\frac{2k}{\delta}\right)}{n}} \leq \epsilon/2$$

Quiz:  $C \geq \frac{2}{\epsilon^2} \ln\left(\frac{2k}{\delta}\right)$

PAC-style  
bounds for  
bandits



Hoeffding

Let  $X_1, \dots, X_n$  be iid with mean  $\mu$ . Let  $\hat{\mu}$  be our estimate of  $\mu$ :  $\sum_{i=1}^n X_i/n$ . The following is a  $100(1-\delta)\%$  confidence interval for  $\mu$ :

$$\left[ \hat{\mu} - \frac{z_\delta}{\sqrt{n}}, \hat{\mu} + \frac{z_\delta}{\sqrt{n}} \right] \quad z_\delta = \sqrt{\frac{1}{2} \ln \frac{2}{\delta}}$$

[Deterministic MDP] MDP Optimization Criteria

- no state, deterministic
- # mistake:  $O(K-1)$

## [Stochastic MDP]

general stochastic MDPs combine

1. stochastic (hoeffding bound until estimates are accurate)
2. sequential (unknown regions assumed optimistic)

## General Rmax

1. keep track of MDP (T, R)
2. **new param:**  $c$  -- # tries in an unknown state. if tries fewer than  $c$  times, use max likelihood estimate

(instead of assuming all  $\langle s, a \rangle = R_{\max}$ )

3. solve the MDP
4. take action from  $\pi^*$   
to make it work...

## Simulation Lemma

"if T and R errors are bounded, V is near-optimal"

if we have transitions and rewards off by  $\alpha$  or less  
(for fixed  $\pi$ , T diff and R diff are both bounded by  $\alpha$ )

$\alpha$  is bounded respect to  $\epsilon$

(each bandit must be  $< \epsilon$  so the overall is  $\epsilon$  optimal)

-- solve with Hoeffding bound & union bound

Simulation Lemma

Define  $G = \alpha + \gamma n \alpha R_{\max} / (1 - \gamma)$

$$\epsilon = |V_1(s) - V_2(s)| \leq \frac{G}{1 - \gamma}$$
$$\epsilon = \frac{G}{1 - \gamma} = \frac{\alpha + \gamma n \alpha R_{\max}}{1 - \gamma}$$

SA  $\rightarrow$   $\epsilon \leq \frac{\alpha}{1 + \gamma n \frac{R_{\max}}{1 - \gamma}}$

Hoeffding bound, union bound

## Explore-or-Exploit Lemma

if all transitions are either accurately estimated or unknown,  
optimal policy is

- either near optimal
- or an unknown state is reached quickly (in polynomial time)

$\wedge$  is true with high probability through  $R_{\max}$

$R_{\max}$  doesn't guarantee optimal policy, it obtains near-optimal result

## 8. Generalization

### Generalization Idea

need to generalize: zillions of states will take forever to solve  
methods: SL methods (e.g. linear function approximation)

what can we generalize? policy, value function, model

1. **policy:** mapping b/w states  $\rightarrow$  actions



generalization: for some similar states, use similar actions

2. **value function**: mapping states/action  $\rightarrow$  estimated return ( $Q(s,a)$ )  $\leftarrow$  most researches focus  
generalization: similar states have similar return

3. **model**: mapping between  $s \rightarrow s'$  &  $r$

$s,a \rightarrow s'$  : supervised examples

generalization: predict reward  $R$  & transition  $T$ .

challenge: needs model to be able to predict multiple states forward

### Basic Update Rule

- create gradient step = learning\_rate \* y\_diff \* gradient\_descent
- y\_diff = TD error =  $y^* - y_{\text{predict}}$   
= (current\_Qlearning\_pred -  $Q(s,a)$ )
- $y^*$  acts as "made-up label" here
- update current state  $Q(s,a)$  towards  $Q^*$  in Q-learning
- meanwhile update  $Q^*$  with bootstrapping

Basic Update Rule

General Form:

$$Q(s,a) = F(w^a, f(s))$$

$\langle s,a,r,s' \rangle$

$$\Delta w_i^a = \alpha \left( \underbrace{r + \gamma \max_{a'} Q(s',a')}_{\text{Bellman equation}} - \underbrace{Q(s,a)}_{\text{TD error}} \right) \frac{\partial Q(s,a)}{\partial w_i^a}$$

bootstrapping

gradient step

### Linear Value Function Approximation

$$Q(s,a) = \sum_{i=1}^n f_i(s) \cdot w_i^a$$

parameters provide generalization

- linear Q function:  $Q(s,a) = \sum_i (\text{feature\_function}_i(s) * \text{weight}(i, a))$
- each action gets a set of  $n$  weights to represent the  $Q$  values for that action.
- dot product weights & feature function
- weights give "importance" of all the features in contributing to an action's value

### Successes in research

#### linear approximation:

3-layer backprop: Tesauro (Backgammon, Jeopardy), Wang/Dietterich (shuttle scheduling, NASA), Charles (tic-tac-toe)

- experience: start to learn well, but bootstrapping aspect is dependent on previous predictions, once learning starts to go south, whole thing falls apart quickly

#### CMAC: Sutton

- like neural nets but first layer not learned, instead pre-decided how to breakup input space
- generalization of linear

#### Linear nets: Singh, Parr (cellphone, channel allocation)

- feature selection is the key!!! must be 1) computationally efficient, 2) value-informative
- example: taxi -- distance to passenger is a good feature

### aggressive generalization (not so linear):

#### Deep nets: Minh Etal (Atari)

- massive layers



- convolutions
- trained carefully (to avoid sampling bias, value oscillation, see more in paper notes)

TD gammon (Tesauro) -> sometimes don't work well (Boyan/Moore) -> but you can train it differently & it could work (Sutton)

### Failures could happen (Baird's Counter example)

- see paper session "Residual Algorithm"

Star problem:

- no rewards
- deterministic
- no choice
- linear value function approximation
- near tabular
- value function multiply representable (inf solutions)
  1. set all 0 ==> sum 0 ==>  $w = (0, 0, \dots, 0)$
  2.  $w_0 = -1$  & solve  $w_i$  ==>  $w = (-1, 1/2, \dots, 1/2, 7)$

- in 2nd case,  $w_0$  gets updated 6 times while other  $w$ 's update once respectively =>  $w_0$  gets pumped up quickly, so are other  $w$ 's except  $w_7$  => diverge!!
- in 1st case,  $w$  will all get stuck at 0 & never move

- if there are shared weights, can not converge

### Averagers function approximator

find convex combination of anchor points (e.g. linear interpolation)

for extrapolation, either or: 1. take the nearest anchor point value 2. between high&low points

requires: 1. pick basis/anchor, 2. pre-defined weights of anchor points, 3. convex comb

$$V(s) = \sum_{s_b \in B} \underbrace{w(s, s_b)}_{\text{basis/anchor}} \underbrace{V(s_b)}_{\text{parameters}}$$

s.t.  $w(s, s_b) \geq 0 \forall s, s_b$   
 $\sum_{s_b \in B} w(s, s_b) = 1, \forall s$

SL method can be expressed as averagers? KNN

### Averagers' Connection to MDPs

- equivalent to a new transition function  $T'(s, a, s_b) = T(s, a, s') * w(s', s_b)$
- we can form a new MDP with function approximation & solve it with VI & Q-learning
- although in reality MDP doesn't transit to  $s_b$ , we can interpret the new MDP as "after transit to  $s'$ , add immediate second transition to  $s_b$ "

$$\begin{aligned} V(s) &= \max_a \left( R(s, a) + \gamma \sum_{s'} T(s, a, s') V(s') \right) \\ &= \max_a \left( R(s, a) + \gamma \sum_{s'} T(s, a, s') \sum_{s_b \in B} w(s', s_b) V(s_b) \right) \\ &= \max_a \left( R(s, a) + \gamma \sum_{s_b \in B} \left( \sum_{s'} T(s, a, s') w(s', s_b) \right) V(s_b) \right) \\ &= \max_a \left( R(s, a) + \gamma \sum_{s_b \in B} T'(s, a, s_b) V(s_b) \right) \end{aligned}$$

- not super practical because value functions are usually not formed in a convex way

## 9. POMDP

consider a "non-markov" environment, where the MDP is NOT observable to agent

MDP:  $A, S \mid Z$  (observables)

$T, R \mid O$  (observation function)

$z$ : observables

$o$ : observation function, probability of seeing  $z$  in state  $s$ ,  $O(s,z)$

(don't know where i am, but can observe something)

### POMDPs generalize MDPs

can represent MDP as POMDP:  $O(s,z)=1$  if  $s=z$

can also represent POMDP as MDP by expanding notion of state (by exploring & updating  $T, R$  & eventually make  $s$  "observable")

example: which color are we gonna see?

1. keep track of  $O(s,z)$ :  $\text{Prob}(s, z)$
2. normalize by  $\text{Prob}(z)$  (e.g. step 2, if we were in  $s_2$ , go right wouldn't see blue  $\rightarrow$  remove  $1/9$  by recalculate average)
3. get  $\text{Prob}(s|z) < \text{Baye's rule!}>$

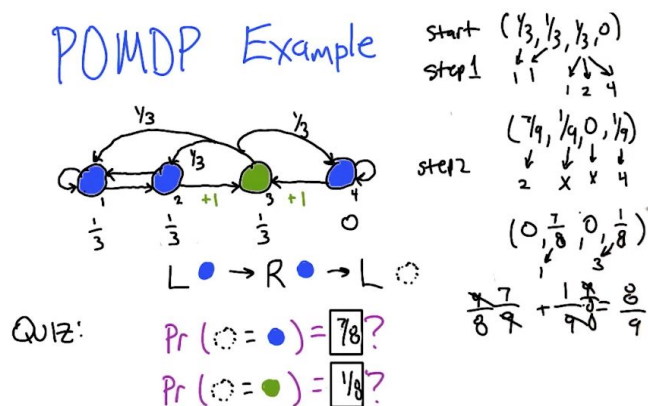
### State Estimation

we can expand POMDP & make it MDP

1. create belief state  $b$ :  $b(s)$ , probability we're in  $s$
2. turn POMDP to belief MDP where state in bMDP is probability distribution over states of the underlying MDP & POMDP  
belief MDP:  $b, a, z \rightarrow b'$  (parallel to MDP except  $b$  is prob. dist.)
3. actual reward is not observed, can only estimate from observable  $z$   
 $r = f(z)$
4. solve  $b'(s') = \text{Prob}(s' \mid b, a, z)$

$$\begin{aligned} \text{Pr}(s'|b,a,z) &= \text{Pr}(z|b,a,s') \text{Pr}(s'|b,a) / \text{Pr}(z|b,a) \\ &= \text{Pr}(z|b,a,s') \sum_s \text{Pr}(s|s',b,a) \text{Pr}(s|b,a) / \text{Pr}(z|b,a) \\ &= O(s',z) \sum_s T(s,a,s') b(s) / \text{Pr}(z|b,a) \end{aligned}$$

### Value Iteration in POMDPs



problem: we turned POMDP to MDP with inf states. if use VI, PI or LP in inf timesteps solve MDP -- must represent in finite way

$$V_0(b)=0$$

$$\text{for } t > 0, V_t = \max(R + \gamma \sum P \cdot V)$$

in POMDP scheme,  $V(b)$  is inf (continuous prob)

(from infinite to finite...)

**claim:**  $V(t)$  is the max of dot product between alpha & believe states where alpha space is finite

**interpretation:** exists finite number of linear function of believe states (these functions can be infinite) such that we can calculate value of each state by taking the maximum value of all the linear functions

### Piecewise-linear & Convex

- for each action, take upper surface over belief state (vector), then we can find upper surface over all actions over belief state (union bag of vectors for each actions) (this step convert infinite belief states into finite vector)

-  $\text{sum}(V)$ : sum over bag of vectors -- create bag of vectors for all  $\langle a, z \rangle$  pair & cross-sum

- still grow doubly-exponentially, but it's finite

- Q-function can be really big (sum(b), second step)

### Algorithmic approach

can get rid of all vectors that are not needed (only keep dominating vectors), shrinks vector space significantly (optimize second step)

### RL for POMDPs

VI doesn't necessarily converge!

### Value Iteration in POMDPs

$$V_0(b) = 0$$

$$V_t(b) = \max_a \left( R(b, a) + \gamma \sum_z \Pr(z|b, a) \cdot V_{t-1}(b') \right)$$

where  $b' = SE(b, a, z)$

CLAIM:  $V_t(b) = \max_{\alpha \in \Gamma_t} \alpha \cdot b = \max_{\alpha \in \Gamma_t} \sum_s b(s) \cdot \alpha(s)$

### Piecewise Linear and Convex II

Assume we have  $\Gamma_{t-1}$  s.t.  $V_{t-1}(b) = \max_{\alpha \in \Gamma_{t-1}} \alpha \cdot b$

Build  $\Gamma_t$  s.t.  $V_t(b) = \max_{\alpha \in \Gamma_t} \alpha \cdot b$

$$V_t(b) = \max_{a \in A} V_t^a(b)$$

$$V_t^a(b) = \sum_z V_t^{a,z}(b)$$

$$V_t^{a,z}(b) = \sum_s \left( \frac{R(s, a) \cdot b(s)}{|z|} \right) + \gamma \Pr(z|b, a) V_{t-1}(SE(b, a, z))$$

### Piecewise Linear and Convex III

$$\Gamma_t = \bigcup_a \Gamma_t^a$$

$$\Gamma_t^a = \bigoplus_z \Gamma_t^{a,z}$$

$$\Gamma_t^{a,z} = \left\{ \frac{1}{|z|} R(s, a) + \gamma \sum_{s'} \alpha(s') O(s', z) T(s, a, s') \mid \alpha \in \Gamma_{t-1} \right\}$$

$$V_{t-1}(SE(b, a, z)) = \max_{\alpha \in \Gamma_{t-1}} \sum_{s'} \alpha(s') O(s', z) \sum_s T(s, a, s') b(s)$$

$$|\Gamma_t| = |\Gamma_{t-1}|^{|A|} \cdot |A|$$

-- not gonna happen to POMDP because infinite # states, can only get near-optimal function approximation with near-optimal policy in finite steps

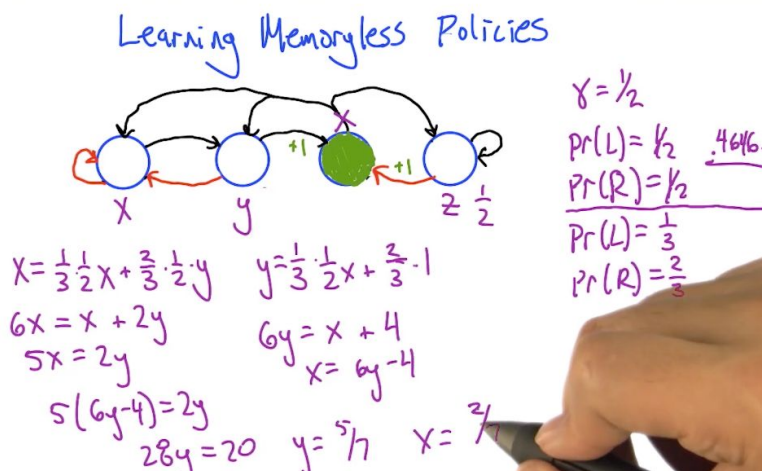
RL: don't know, need to interact with environment

1. model-based RL: learn model  $T$  (learn POMDP)
2. model-free RL: don't bother to learn model (map observations to actions)

	uncontrolled	controlled
observed	MC	MDP
partially observed	HMM	PomDP

... and repeat ...

learning memoryless policy: solve LP. may be random



optimal policy: figure out which MDP we're in, get high reward on the way

Bayesian RL blurs dilemma between exploration & exploitation. always exploits in belief space and maximize conditional total reward based on belief

"RL is actually planning in a kind of continuous POMDP, where the hidden states are parameters of MDP we're trying to learn. value function in the continuous POMDP is piecewise polynomial and convex"

algorithms work well: BEETLE (Bayesian exploration exploitation trade-off in learning)

Bayesian RL keeps learned **Bayesian posterior** and use it to **make decisions**  
but in "too expensive to use" category, Q-learning still wins

### Predictive State Representation (PSR)

motivation:

POMDP: track belief states (distribution of state) which are not observable. but do they even exist?

PSR: predictive state is probabilities of future predictions

calculate prediction (probability of outcome) based on belief state

can we estimate belief based on prediction? -- need more experiments

### PSR Theorem

any n-state POMDP can be represented by a PSR with no more than n tests, each of which is no longer than n steps long

nice thing about the test model: there are pred. prob = 1/0 in some circumstances, can think tests as "what do i need to do in order to figure out where i am" -- find set of tests to find exactly where i am (w. prob  $\rightarrow 1 \Rightarrow$  turn POMDP to MDP)

PSR work: Holmes & Isbelle

also recent work uses least square instead of probabilistic, can learn PSR efficiently

## 10. Options

options: abstract actions which are "small goals" for a set of states

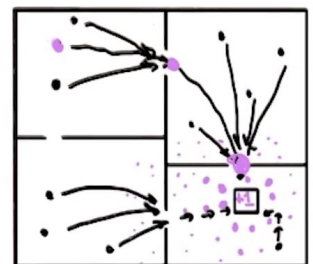
### Temporal Abstraction Options (generalizations)

options =  $\langle I, \pi, \beta \rangle$

$I$  = initiation set of states (options / "super actions" are legal to be executed)

$\pi$  = policy  $(s, a) \rightarrow [0, 1]$  but in probability form:  $P(\text{in state } s \text{ and take } a)$

$\beta$  = set of states where option terminates, but in probability form:  $P(s \text{ will end in } b)$ , mapping  $s \rightarrow [0, 1]$



## Temporal Abstraction Option Function

example: pac man game, options: eat dot, eat big dot, eat ghost, avoid ghost

make option that's just an action:

I: the state where agent can take this action

PI: for the state,  $p(a|s)=1$  and all other  $p(a'|s) = 0$

Beta: 1 everywhere (once action is executed, option is terminated)

R: discounted expected reward of a set of rewards

F: expectation of ending up in  $s'$  after  $k$  steps

create Semi-MDP, a generalization of MDP that inherits bellman related properties (semi because it depends on timestep  $k$  and hence not really markov)

Handwritten notes on a piece of paper:

$$V_t(s) = \max_o \left( R(s, o) + \sum_{s'} F(s, o, s') V_t(s') \right)$$
$$R = E [ r_1 + \gamma r_2 + \gamma^2 r_3 + \dots + \gamma^{k-1} r_k \mid s, k \text{ steps} ]$$
$$F = E [ \gamma^k \Delta_{s_k, s'} \mid s, k \text{ steps} ]$$

MDP:  $\overset{+}{\underset{1}{\bullet}} \overset{+}{\underset{1}{\bullet}} \overset{+}{\underset{1}{\bullet}} \overset{+}{\underset{1}{\bullet}} \overset{+}{\underset{1}{\bullet}} \dots$  atomic actions

SMDP:  $\bullet \xrightarrow{\quad} \bullet \xrightarrow{\quad} \bullet \xrightarrow{\quad} \bullet \dots$  variable time actions

once create SMDP and bellman operator, can pretend it as MDP

but need to be careful because timestep  $k$  of options is variable

it's ok to think SMDP as MDP because:

1. can turn non-markovian to markovian by keeping track of enough history
2. that is what option does. it collects  $r$  &  $s$  along the way. within one option, it's not markovian because don't know  $k$  yet; but once transit out of the option, this "big step" is markovian when look back

## Notes on Options

1. inherit optimality, convergence and stability

given the set of options we have, we can converge to the optimal policy with respect to options (hierarchical optimality)

BUT, learning with respect to options won't guarantee optimal policy (with respect to actual actions)

it's ok because actions are relative. we are limited by options of our choice, but were also limited by actions too

2. avoids "boring part" of MDP, skip irrelevant state changes and focus on places where we can make decisions

## Goal Abstraction

options can be interpreted as accomplishing goals

1. sequencing actions
2. parallel goals

things can happen in parallel. temporal abstraction is really about managing competing goals



e.g. pac-man, eat small dot, eat big dot, eat ghost, avoid ghost.. all happen at the same time  
predator-prey scenario: go around to find prey, but predators are out there to look out for

Beta (termination distribution) can be interpreted as

1. i've accomplished a goal or
2. there's a more important goal i need to focus on and interrupt what i was paying attention to

## Goal Arbitration

"modular RL"

different modules try to accomplish some goals, how to arbitrate?

1. greatest mass Q-learning: pick action that has the most mass (chosen by more Qs) (vote on actions)
  - but could end up picking least favorable actions because no other actions get enough votes
2. top Q-learnings: pick action that contribute the most value  $\max Q(s,a)$  (negotiation)
  - but one action may benefit a bunch of agents
3. negotiated W-learning: agent with most to lose gets to pick which action to take (min loss)

voting's problem: Arrow's impossibility theorem

it's impossible to construct a fair voting system due to individual different utility units which is impossible to compare

## Monte Carlo Tree Search (MCTS) - policy search algorithm

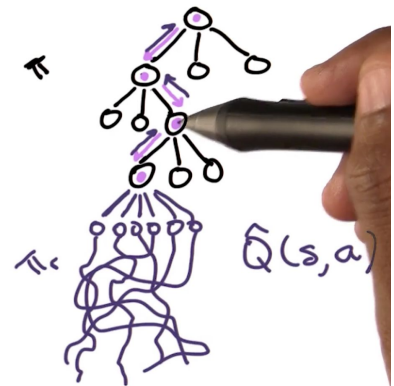
**SELECT:** decide what action to take

**EXPAND:** when you don't know enough to make choices, explore available actions that we could take. collect concrete states we could end up in

**SIMULATE:** simulate different scenarios (behave randomly for a while) (random policy, rollout policy) how far to go? depends on horizon & discount factor. and get an estimation of Q function

**BACKUP:** back up information we've learned all the way up top with bellman equation

... and repeat....



initial  $\pi$ : from estimation of Q -> updating Q each loop -> repeat many times -> confident to expand tree -> find optimal policy ->...

after figuring out policy, act greedily & make a move

how to be smarter about it?

- options instead of actions
- instead of random, apply constraints (goals on avoiding damage)

## MCTS Properties

1. useful for large state space

2. need lots of samples to get good estimate
3. planning time is independent of  $|S|$  ("voting survey")
4. running time is exponential in the horizon  $O((|A|^X)^H)$ 
  - H: how far in the future to look, depends on gamma
  - X: num steps to run

## 11. Game Theory

### 2-player zero-sum finite deterministic game of perfect information

zero-sum: means constant sum in all scenarios. not necessarily 0

strategy: minimax

can be solved with LP

### Fundamental Result

in a 2-player, zero-sum deterministic game of perfect information (e.g. game tree):

1. minimax == maximin
2. always exists an optimal **pure strategy** for each player

### Von Neumann Theorem

in a 2-player, zero-sum **non-deterministic** game of perfect information (e.g. game tree):

1 & 2 above still hold

in a 2-player, zero-sum **non-deterministic** game of **hidden** information (e.g. minipoker):

1. minimax != maximin anymore!!
2. pure strategy doesn't work any more!!

Von Neumann theorem doesn't hold any more

mixed strategy: prob. distribution over actions

pure strategy:  $\text{prob}(\text{action\_pure}) = 1$

**non-zero-sum, non-deterministic** game of **hidden** information (e.g. prisoners' dilemma)

### Nash Equilibrium

n players with strategies  $s_1, \dots, s_n$

$s_1^*, s_2^* \dots s_n^*$  are a NE iff exists  $s_i^* = \text{argmax}_{s_i} [\text{Utility}(s_1^*, s_2^* \dots s_n^*)]$

"no reason for anyone to change their strategy"

can be mixed or pure

### Nash Equilibrium Theorems

1. in the n-player **pure strategy** game, if elimination of strictly dominated strategies eliminates all but one combination, that combination is the **unique N.E.**

2. any N.E. will survive elimination of **strictly dominated** strategies



3. if  $n$  is finite and exists an  $i$  such that  $s_i$  is finite, then (mixed) N.E. exists -- **there exists a N.E. for any finite games**

### n-repeated game

when we play game prisoners' dilemma  $n$ -times, we'll get  $n$ -repeated N.E. (repeating game won't change NE)

## 12. Game Theory II

### Iterated Prisoners' Dilemma (IPT)

what happens if # of round is unknown?

uncertain end with discounted factor  $\gamma$  -> expected # rounds =  $1/(1-\gamma)$

### Tit for Tat (TfT)

famous IPD strategy

1. cooperate in first round
2. copy opponent's previous move thereafter

what's the best response to TfT?

for high  $\gamma$ , always cooperate; for low  $\gamma$ , always defect (can be calculated with geometric formula)

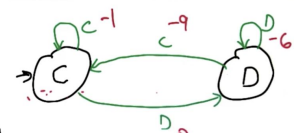
NOT subgame perfect

### Finite State Strategy

our choice impact our payoff and opponents' future decision

- "the matrix is all you need" is not applied any more (only true for 1-step games)
- it becomes a MDP to us. solve MDP -> VI

Best Response To A  
Finite-state Strategy



states labeled with opponent's choice  
edges labeled with our choice  
edges annotated with our payoff for that choice  
our choice impacts our payoff and future decisions of the opponent

mutual best response: pair of strategies where each is the best response to others. also a N.E.!  
(e.g. TfT vs TfT)

### Folk Theorem Idea

in **repeated** games, the possibility of **retaliation** opens the door for **cooperation**

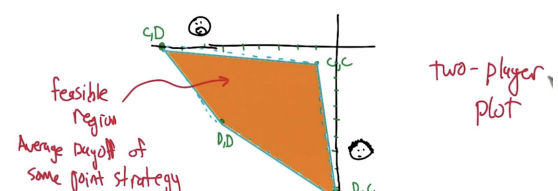
(different from definition in math) in game theory, Folk Theorem describes the set of payoffs that can result from **Nash strategies in repeated games**

### Feasible Region

the convex region created by N.E.s == average payoff of some point strategy == Nash strategies

Repeated Games and the Folk Theorem

General idea: In repeated games, the possibility of retaliation opens the door for cooperation.



## Minimax Profile

a pair of payoffs, one for each player, that represent the payoffs that can be achieved by a player defending himself from a **malicious adversary**

- assuming opponent's goal is to minimize my payoff
- find minimize payoff in respect to each of my action, then pick the strategy which maximizes it

minimax profile: pure strategy

security level: mixed strategy

## Security Profile

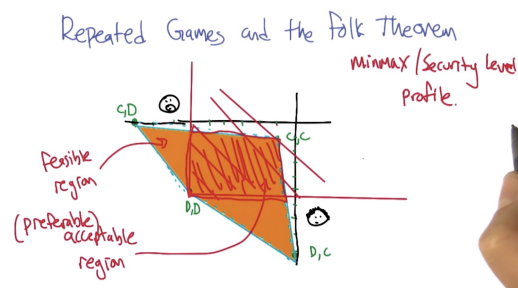
overlap between feasible region & acceptable region

acceptable region: within acceptable region, the area that has higher payoff for both players

## Folk Theorem

Any feasible payoff profile that **strictly dominates the minimax / security level profile** can be realized as a **Nash Equilibrium payoff profile**, with sufficiently large discount factor.

- security profile is N.E. with sufficiently large gamma
- reason: if it strictly dominates the minimax profile, opponent can use it as a **threat**



## Grim Trigger

"once opponent defects, I defect forever"  
has N.E.

but it's an implausible threat

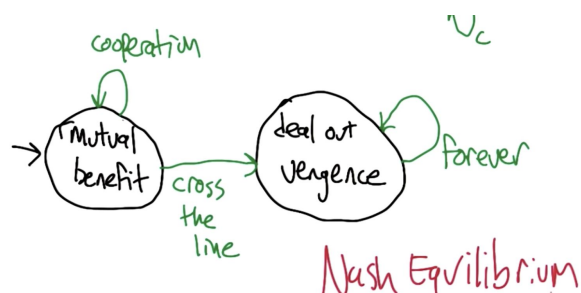
not subgame perfect

(Grim Trigger is supposed to do CCCCC...

but if we alter history & insert D

player will do DDDDD...

changed behavior => not subgame perfect



## Implausible Threats

after giving threat, payoff(threat happens) << payoff(don't do it)

e.g. robber with dynamite: "give me \$ or i blow both of us up" -> payoff\_robber(blow up) <<

payoff\_robber(leave without money)

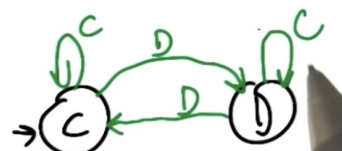
## Subgame Perfect

- always best response independent of history
- if alter history at some point, do players have better action? if not, subgame perfect

Grim Trigger: NOT subgame perfect

TfT: NOT subgame perfect

## Pavlov



- take advantage of opponent until they pull a trigger
- cooperate if agree, defect if disagree

Pavlov vs Pavlov is Nash! and subgame perfect! < -- plausible threat

### Computational Folk Theorem

For any 2-player bi-matrix average-reward repeated game, can build a Pavlov-like machine for any game, and use it to construct subgame perfect N.E. for any game in **polynomial time**

how does it work?

1. if mutual beneficial, Pavlov is possible -> N.E.
2. if not, zero-sum like -> solve an LP -> N.E.
3. if neither works, at most one player improves his behavior (zero-sum like) -> N.E.

### Stochastic Games and Multi-agent RL

generalization of repeating game

MDP: RL :: Stochastic game: Multi-agent RL

### Zero-sum Stochastic Games

replace max in Q-learning with minimax

with additional agent,

✓ value iteration works (but unclear if it can be solved in polynomial time)

✓ minimax-Q converges

✓ unique solution to Q

✓ policy can be computed independently

✓ update efficiently: updating Q step can be computed in polynomial time (with LP)

✓ Q-function sufficient to specify policy

Zero-sum Stochastic Games

$$Q_i^*(s_i(a,b)) = R_i(s_i(a,b)) + \gamma \sum_{s'} T(s_i(a,b), s') \min_{a', b'} Q_i^*(s', (a', b'))$$

$$\langle s_i(a,b), (r_i, f_i), g \rangle: Q_i(s_i(a,b)) \leftarrow r_i + \gamma \min_{a', b'} Q_i(s', (a', b')) \quad \text{minimax-Q}$$

### General-sum Stochastic Games

replace minimax above with Nash

everything above breaks

General ~~Zero-sum~~ Stochastic Games

$$Q_i^*(s_i(a,b)) = R_i(s_i(a,b)) + \gamma \sum_{s'} T(s_i(a,b), s') \min_{a', b'} Q_i^*(s', (a', b'))$$

$$\langle s_i(a,b), (r_i, f_i), g \rangle: Q_i(s_i(a,b)) \leftarrow r_i + \gamma \min_{a', b'} Q_i(s', (a', b')) \quad \text{Nash minimax-Q}$$

### Ideas

1. repeated stochastic game: folk theorem
2. cheap talk -> seek cooperation -> correlated equilibrium
3. cognitive hierarchy (assume other players have more limited resources)-> best responses
4. side payment (coco value)

- value iteration ~~works~~ doesn't work

- ~~minimax-Q converges~~ doesn't converge

- No unique solution to Q\*

- Policies can ~~not~~ be computed independently

- update ~~not~~ efficient P = PPA

- Q functions ~~not~~ sufficient to specify policy

incompatible

insufficient

## 13. Game Theory III

### Correlated Equilibria (CE)

rely on shared source of information (correlated distribution)

#### C.E. Facts

1. C.E. can be found in **polynomial** time
2. all mixed Nash are correlated (so C.E. exists)
3. all convex combination of mixed Nash are correlated

### CoCo Values (Cooperative-competitive values)

"split extra payoff by creating side-payment"

coco value of a game:

U: payoff to player, U': payoff to other players

U, U': general sum game

$$\text{CoCo}(U, U') = \text{maxmax}((U+U')/2) + \text{minimax}((U-U')/2)$$

1. maxmax: cooperative minmax: competition
2. two pair <U,U'> may be different -> computed separately

### CoCo Properties

1. efficiently computable
2. utility maximizing
3. decompose game into sum of two games
4. unique solution
5. can be extended to stochastic games (coco-Q, **NOT NON-EXPANSION**)
6. not necessarily equilibrium -- needs to be binding to reach equilibrium
7. doesn't generalize for  $n > 2$

### Mechanism Design

design payoff so that players do what we want them to do

e.g. peer-coaching: student +1 when answer a question correctly; teacher +1 when student gets what shouldn't get right right or what shouldn't get wrong wrong

e.g. King Solomon

similar to auction design, make people reveal their true v values

# RL Sutton Barto Reading Notes

## Sarsa

consider action-value function rather than state-value

1. on-policy TD control
2. transition from  $\langle s, a \rangle$  to  $\langle s', a' \rangle$  and learn  $V(s, a)$
3. on-policy methods: continually estimate  $q_\pi$  for the behavior  $\pi$ , and at the same time change  $\pi$  toward greediness with respect to  $q_\pi$  (update behavior action while learning value)
4. sarsa **converges with probability 1** to an optimal policy and action-value function as long as all  $\langle s, a \rangle$  pairs are **visited an infinite # times**; policy converges to greedy policy in limits

sidenote:

on-policy TD control convergence requirement: **GLIE** (greedy limit in exploration)

GLIE has following 2 properties:

1. If a state is visited infinitely often, then each action in that state is chosen infinitely often (with probability 1). <explore indefinitely>
2. In the limit (as  $t \rightarrow \infty$ ), the learning policy is greedy with respect to the learned Q-function (with probability 1). <policy converges to greedy>

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)]$$

### Sarsa (on-policy TD control) for estimating $Q \approx q_*$

Algorithm parameters: step size  $\alpha \in (0, 1]$ , small  $\varepsilon > 0$

Initialize  $Q(s, a)$ , for all  $s \in S^+$ ,  $a \in \mathcal{A}(s)$ , arbitrarily except that  $Q(\text{terminal}, \cdot) = 0$

Loop for each episode:

    Initialize  $S$

    Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\varepsilon$ -greedy)

    Loop for each step of episode:

        Take action  $A$ , observe  $R, S'$

        Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\varepsilon$ -greedy)

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$

$S \leftarrow S'; A \leftarrow A';$

    until  $S$  is terminal

## Q-learning

1. off-policy TD control
2. the learned action-value function  $Q$  directly approximates  $q^*$  (optimal action-value function), **independent of policy being followed**
3. difference b/w sarsa & q-learning: whether behavior action depends on learned  $q$  (sarsa  $Y$  q-learning  $N$ ); both learns optimal action-value function  $q^*$
4. converges with probability 1 to  $q^*$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \left[ R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t) \right]$$

### Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

Algorithm parameters: step size  $\alpha \in (0, 1]$ , small  $\varepsilon > 0$

Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}^+, a \in \mathcal{A}(s)$ , arbitrarily except that  $Q(\text{terminal}, \cdot) = 0$

Loop for each episode:

    Initialize  $S$

    Loop for each step of episode:

        Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\varepsilon$ -greedy)

        Take action  $A$ , observe  $R, S'$

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$

$S \leftarrow S'$

    until  $S$  is terminal

## Expected sarsa

instead of maximum over next state-action pairs (Q-learning), take into account how likely each action is under the current policy

$$\begin{aligned} Q(S_t, A_t) &\leftarrow Q(S_t, A_t) + \alpha \left[ R_{t+1} + \gamma \mathbb{E}_\pi [Q(S_{t+1}, A_{t+1}) \mid S_{t+1}] - Q(S_t, A_t) \right] \\ &\leftarrow Q(S_t, A_t) + \alpha \left[ R_{t+1} + \gamma \sum_a \pi(a \mid S_{t+1}) Q(S_{t+1}, a) - Q(S_t, A_t) \right], \end{aligned}$$

## Papers

Knows What It Knows, Li Littman Walsh

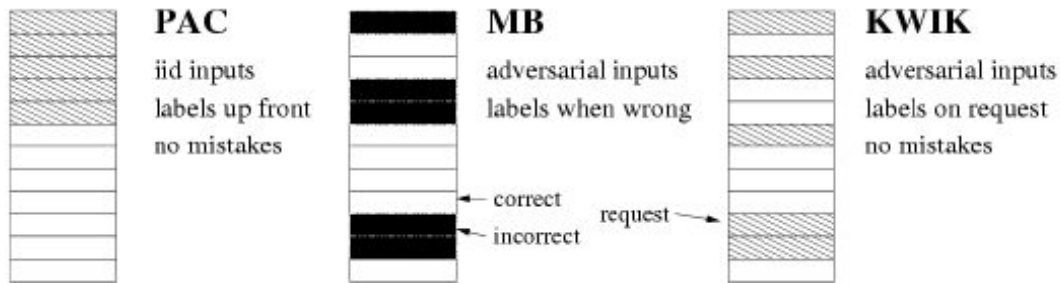
Rmax

- estimates are made separately
- guarantees polynomial bound on timestamps in which it has non-near-optimal policy

KWIK (knows what it knows)

- make only accurate predictions
- has polynomial bound on idk response
- mistake bounded by epsilon, probability of failed run bounded by delta
- total number of steps on idk bounded by  $\text{Binomial}(\epsilon, \delta)$ , ideally polynomial in  $1/\epsilon$ ,  $1/\delta$  and parameters defining target function  $H$

comparisons:



PAC (prob. approx. correct):

- provide label for initial sequence of inputs (correct labels), after which the learner make accurate outputs

MB (mistake bounded):

- provide labels when learner makes mistake, inputs are selected adversarially
- no bound on when the last mistake might be made
- total # mistake is small,  $\text{incorrect/correct} \rightarrow 0$

KWIK (knows what it knows)

- not allowed to make mistakes
- bound on # label requests ("idk")

KWIK learnable classes

1. memorization and enumeration  
rule out all other possibilities
2. real-value functions
3. noisy observations

Learn to Predict by the Methods of Temporal Differences, Sutton88

learn to predict. example: predict weather

1. TD method: compare predictions with that made the following day
2. advantages:



- more incremental hence easier to compute (conventional: wait till label date & save during training vs. TD: update prediction with one extra observation)

- TD converges faster & produces better predictions

3. single-step & multistep:

- single-step problems: data comes out in observation-outcome pair (TD == SL)

- multi-step problems: next year economy, election result, chess game - TD, adjust prediction with more observations

4. theorem1: on multi-step prediction problems, linear TD(1) procedure produces same per-sequence weight changes as Windrow-Hoff procedure

$\Delta w_t = \alpha(z - P_t) \nabla_w P_t$  SL (back propagation)

$$z - P_t = \sum_{k=t}^m (P_{k+1} - P_k) \quad \text{where } P_{m+1} \stackrel{\text{def}}{=} z. \quad \text{TD definition}$$

$$\Delta w_t = \alpha(P_{t+1} - P_t) \sum_{k=1}^t \nabla_w P_k \quad \text{TD(1)}$$

5. TD(lambda) exponential weighting with recency

$$\Delta w_t = \alpha(P_{t+1} - P_t) \sum_{k=1}^t \lambda^{t-k} \nabla_w P_k \quad 0 \leq \lambda \leq 1$$

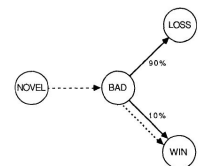
6. **TD(1)=SL**; TD(0) weight change only determined by its effect on the prediction associated with the most recent observation

7. game-play example:

if novel state always lead to win:

SL: novel is good (obs-outcome pair b/w novel & wins)

TD: novel is bad (tracks novel-bad-win & uses knowledge about bad- 90%loss)



if novel->bad->do bad thing

SL: novel is bad

TD: novel has 90% chance losing

if previous evaluation is in error: TD will be affected

if bad state is followed by defeats except when its preceded by novel state (victory):

TD may perform worse than SL.

**"TD method try to take advantage of the information provided by the temporal sequence of states, whereas SL ignore it"**

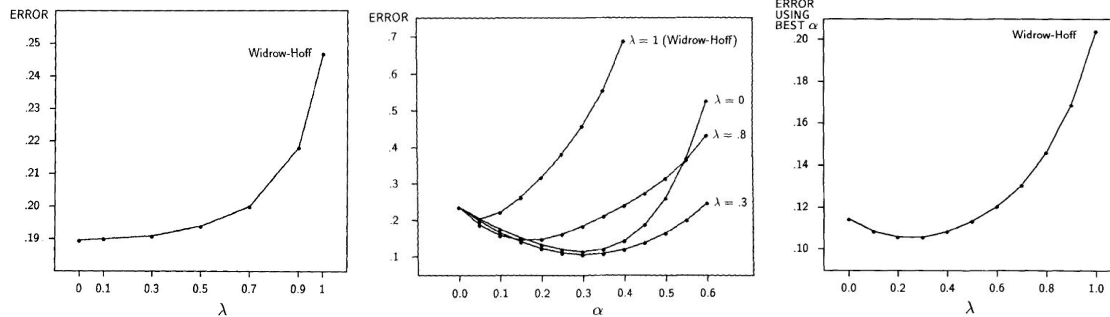
8. random-walk example

experiment1: update weight after presenting one train\_set + repeated presentation until converge (SL)

experiment2: update weight after each sequence + no repeated presentation

(fig 3,4,5) 3,4 exp1; 5 exp2





- TD(0) in exp1 works best (shocking), because TD(1) Widrow-Hoff procedure minimizes error in training but not necessarily testing. TD(0) converges to optimal estimation
- TD with  $\lambda < 1$  in exp2 perform better than SL (exp1)
- in exp2, TD(0) is not as good. because it's slowest at back propagating. with repeated presentation it's no problem since the change works back with an additional step in each rep. But single presentation means slow learning. this handicap can be avoided by working backwards thru the sequence (works for offline training)

## Residual Algorithm, Baird

1. value iteration form:  $V(x) \leftarrow \alpha R + \gamma V(x')$
2. if  $V(x)$  is separate entry in a lookup table, then  $\alpha$  is equivalent to 3 RL algorithms:  
TD(0), Q-learning, advantage learning

### 3.1. direct algorithms

- when each possible transition is experienced an infinite number of times during learning, then **lookup table** is guaranteed to converge to optimal function as learning rate decays to 0 at appropriate rate
- weight change (exactly TD(0)):

$$\Delta w = \alpha (R + \gamma V(x') - V(x)) \frac{\partial V(x)}{\partial w}$$

- doesn't converge for a general **function-approximation** system. e.g. star problem. if  $w$ 's initialized with non-zeroes,  $w_0$  will increase every fast

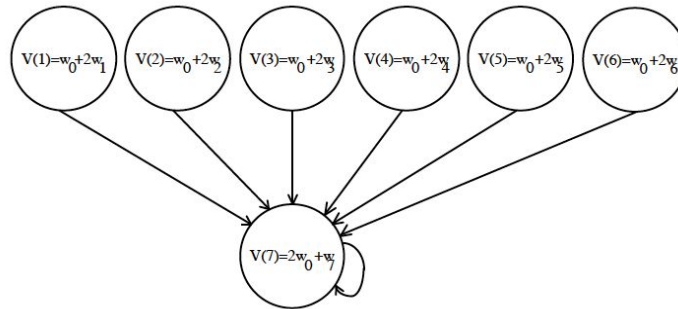


Figure 1. The star problem

### 3.2. residual gradient algorithms

- residual gradient algorithm: performing gradient descent on bellman residual
- where bellman residual:

$$E = \frac{1}{n} \sum_x [R + \gamma V(x') - V(x)]^2$$

$$\Delta w = -\alpha [R + \gamma V(x') - V(x)]$$

$$\left[ \frac{\partial}{\partial w} \gamma V(x') - \frac{\partial}{\partial w} V(x) \right]$$

- if bellman residual != 0 then suboptimal. but suboptimal reinforcement can be bounded although residual gradient algorithms have guaranteed convergence does not necessarily learn as quickly as direct algorithms

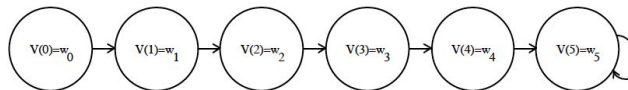


Figure 2. The hall problem.

- in the hall problem, set w5=0 w4=10, if use direct method, w5->0 w4->0; if use residual gradient, decrease both w5 & w4. info flows both ways & learn slow

### 3.3. residual algorithms

- direct algorithms: fast but unstable
- residual gradient algorithms: stable but slow
- solution? weighted average! (residual algorithms)

$$\Delta \mathbf{W}_r = (1 - \phi) \Delta \mathbf{W}_d + \phi \Delta \mathbf{W}_{rg} \quad \text{where } w_d: \text{direct}; w_{rg}: \text{residual gradient}$$

$$\Delta \mathbf{W}_d = -\alpha \sum_x [R + \gamma V(x') - V(x)] [-\nabla_w V(x)]$$

$$\Delta \mathbf{W}_{rg} = -\alpha \sum_x [R + \gamma V(x') - V(x)] \left( \nabla_w \gamma V(x') - \nabla_w V(x) \right)$$

- guaranteed to converge for appropriate phi

$$\Delta w = -\alpha [R + \gamma V(x') - V(x)]$$

$$\left[ \phi \gamma \frac{\partial}{\partial w} V(x') - \frac{\partial}{\partial w} V(x) \right]$$

- residual algorithm identical to residual gradient with one term \*phi
- how to choose phi?

on dot line, perpendicular to gradient, no change;

on the right, result in increase of E, difficult to predict convergence

closer to  $\Delta W_d$ , learns quicker

best: on left of dot line but as close to  $\Delta W_d$  as possible



Figure 3. Epoch-wise weight-change vectors for direct and residual gradient algorithms

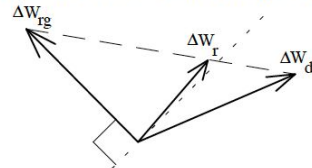


Figure 4. Weight-change vectors for direct, residual gradient, and residual algorithms.