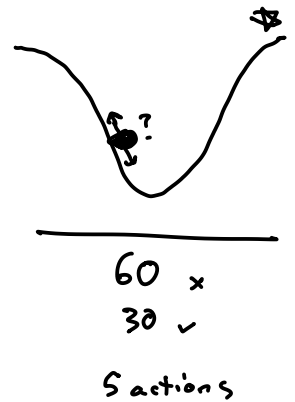


4

HW  $\leftarrow$  Traditional 10000 episodes  
Deep Q Learning.jl

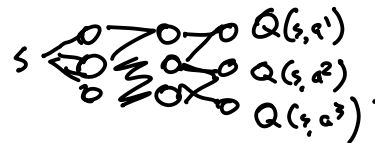
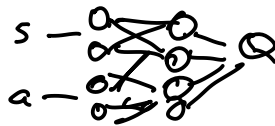
give expert  
expert 5 episode  
episode

full credit  
above 30  
45/50 above 25



DQN

$Q(s,a)$



5 actions  
[1, 0.5, 0, 0.5, 1]

~~tabular~~ Corrections

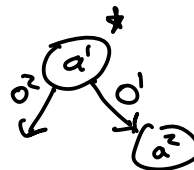
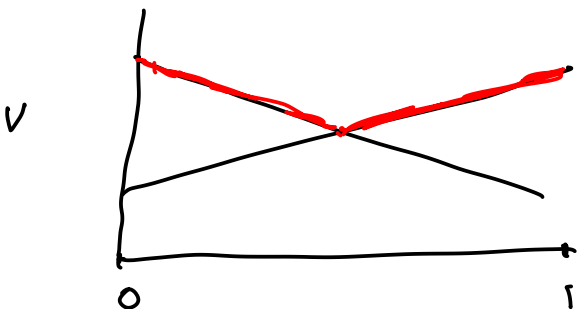
$Q(\lambda)$

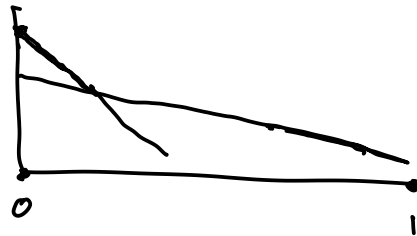
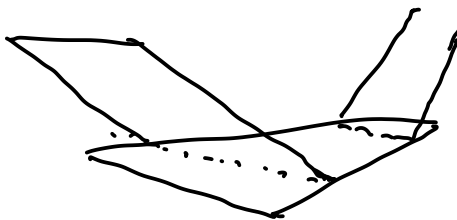
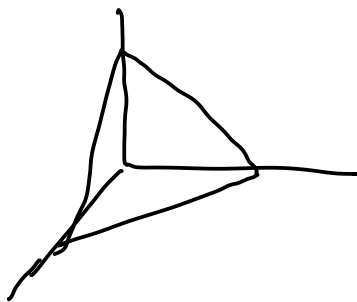
Off-policy learning with eligibility traces  
is unstable

decaying  $\alpha$  is a good idea.

POMDPs

Last Time: - Value Function  $\leftrightarrow$  belief  
- What do POMDP Policies Look like?

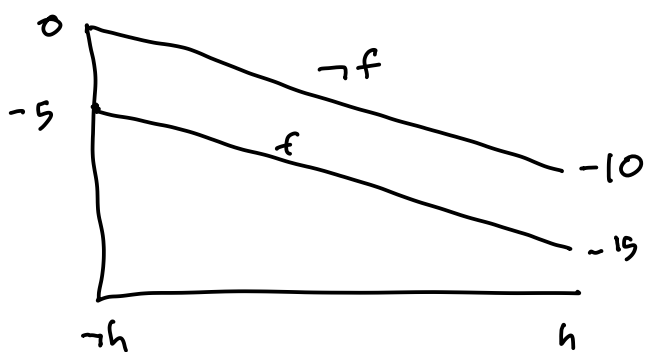




$$S = \{h, \neg h\} \quad \begin{matrix} -10 & h \\ -5 & f \end{matrix}$$

$$A = \{f, \neg f\}$$

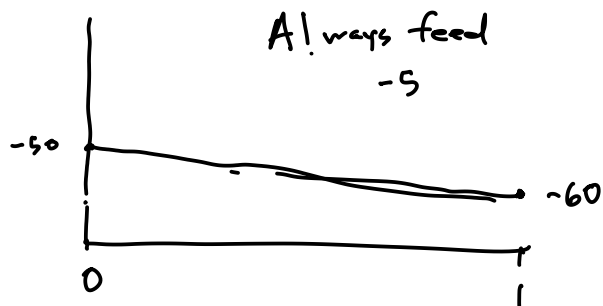
$$O = \{c, \neg c\}$$



(f)

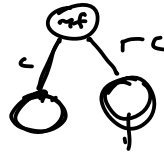
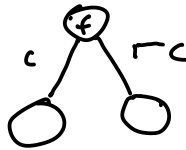
(\neg f)

$$\frac{-5}{1-\gamma} = 50$$



$$\alpha[s] = R(s, p_0) + \gamma \sum_{s'} T(s'|s, a) \sum_o Z(o|s, s') \alpha[s']^{p(o)} \leftarrow \text{action after } o$$

$\Gamma$  = set of alpha vectors       $p_0$  initial action



Bellman Backup  $(b, \Gamma)$

for  $a \in A$

for  $o \in O$

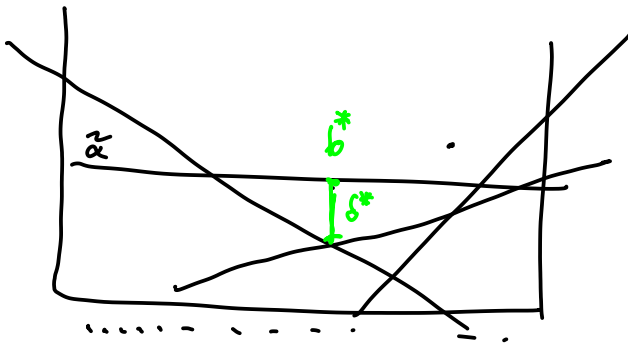
$$b' = \tau(b, a, o)$$

$$\alpha_{ao} = \arg \max_{\alpha \in \Gamma} b'^T \alpha$$

for  $s \in S$

$$\alpha_a[s] = R(s, a) + \gamma \sum_{s', o} Z(o|s', a) T(s'|s, a) \alpha_{ao}(s')$$

$$\Gamma = \Gamma \cup \{\alpha\}$$



maximize  $\delta$   
 $\delta, b$

Subject to

$$b^T \alpha \geq \delta + b^T \alpha' \quad \forall \alpha' \in \Gamma$$

$$\begin{aligned} b \text{ is prob} \quad & - b^T 1 = 1 \\ & - b \geq 0 \end{aligned}$$

if  $\delta^* > 0$

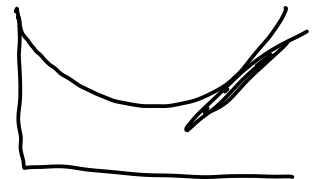
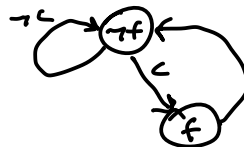
$b^*$  is a "witness"

Witness algorithm (1996)

Incremental Pruning (1997)

Baby POMDP

Feed When Crying

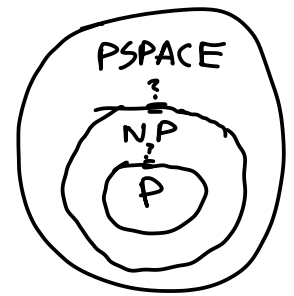


Computational Complexity of POMDPs

"NP-Hard" - at least as hard as all problems in NP

"NP-Complete" - NP-Hard and in NP

PSPACE - can be solved with polynomial memory space



QSAT — known to be PSPACE Complete  
 $\forall x \exists y \exists z ((x \vee z) \wedge y)$

QSAT can be transformed into a finite horizon POMDP

$\therefore$  POMDPs are PSPACE Complete

---

Approximations

- Numerical

- Optimization objective / Formulation