

Quioto Forest Legends

A B C



time up → loss
 0 eggs → loss
 100 eggs → win
 must get 10 every 4 skirmishes

eggs ticks

20	24
25	23
27	21
32	18
32	16
32	15
32	14
32	13
32	12
<u>DEFENSE WINS</u>	



0 1 2 3

$$Q((1, 1, 1, 10), 0) = .9$$

Images by eBay user furret_inc <https://www.ebay.com/itm/OLD-Original-Vintage-Pokemon-10-Card-LOT-1st-Edition-Holo-Rare-Played/233176095585?hash=item364a60db61:g:30QAAOSwG7NclWhm>
 Used under Fair Use Doctrine

OFFENSE

Hitmonchan
 3 ticks except
 as noted

Machop Magnetron Mewtwo

	1-4	5	6-10	11-12	13-20	Hitmonchan	Machop	Magnetron	Mewtwo	
Chansey	1-4	5	6-10	11-12	13-20	2	5	19	0	4 ticks
	5	8	5	-1	3	12	-100	4	4	
	6-10	5	-100	47	0	5	22	-100	1 tick	
	11-12	-1	0	47	0	-100	47	47	otherwise 2 ticks	
	13-20	3	0	0	0	0	0	0	0	
Charizard	1-4	5	6-10	11-12	13-20	1	4	11	48	4 ticks
	5	14	2	-4	2	14	32	-100	-9	
	6-10	2	5	14	0	2	5	14	41	
	11-12	-4	0	-9	0	-4	0	-9	-100	
	13-20	2	0	0	0	2	0	0	0	
Clefairy	1-4	5	6-10	11-12	13-20	1	21	-7	53	4 ticks
	5	4	2	1	-3	4	-4	8	-100	
	6-10	2	0	46	5	2	0	15	-11	
	11-12	1	46	-8	18	1	46	-8	-100	
	13-20	-3	5	18	0	-3	5	18	0	

states encapsulate

(yards to go, downs lefts, yards to go to 1st down, ticks left)
eggs needed to win, skirmishes left to get 10 eggs, eggs needed for current set of 4 skirmishes

possible values 101, 5, 99, 25
total # of states ~1.25 million

$$V(s) = P(\text{offense wins in state } (s))$$

matrix for each state

$$\begin{matrix} & A & B & C \\ \begin{matrix} 0 \\ 1 \\ 2 \\ B \end{matrix} & \begin{pmatrix} - & - & - \\ - & - & - \\ - & - & - \\ - & - & - \end{pmatrix} \end{matrix}$$

An arrow points to the element in row 2, column B, which is circled in green.

$$\sum_{s \rightarrow s'} P(s \rightarrow s') \cdot V(s')$$

value of equilibrium is $V(s)$

compute $V(s)$ in order of ↑ time left

Q Learning

$Q(s, a)$ = value of taking action a in state s
(expected reward from taking action a in state s
+ expected discounted future reward from resulting state)

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

initialize $Q(s, a) = 0$ for all s, a (viewing reward as earned on transition into winning states)

while not done

$s \leftarrow s_0$ initial state

$1-\epsilon$ greedy: pick a to maximize $Q(s, a)$
 ϵ pick random action

while s not terminal

choose action a

ϵ -greedy

observed reward

observe transition (s, a, r, s')

update $Q(s, a) \leftarrow Q(s, a) + \alpha (r + \gamma \cdot \max_a Q(s', a) - Q(s, a))$

error

learning rate (should \downarrow as episodes \uparrow)

episode

Function Approximators

— instead of learning $Q(s,a)$ for every (s,a) , learn f_{π} to approximate them

Linear Approximator

Define features of states or (state, action) pairs

time left	↑ fns of state	$f_1(s,a)$	on pace to earn upper bonus	1.0 have already earned 0.0 no chance
yards left / ticks left		$f_2(s,a)$	is chance unused	
yards-to-first / downs left		$f_3(s,a)$	both LS, SS unused	let $f(s)$ be fn from state to $[0, 1]$
good features to start with		\vdots		$f_a(s,a) = \begin{cases} f(s) & \text{if } a=a \\ 0 & \text{otherwise} \end{cases}$
		$f_4(s,a)$	γ unused	

$$\hat{Q}(s,a) = \underbrace{w_1 \cdot f_1(s,a) + w_2 \cdot f_2(s,a) + \dots + w_n \cdot f_n(s,a)}_{\text{learn } w_i \text{ instead of } Q(s,a) \text{ directly}} \quad \text{for nonterminal } s \quad + \underbrace{(w_{n+1} \cdot 1)}_{\text{constant term}}$$

$\hat{Q}(s,a) = 0$ if s is terminal

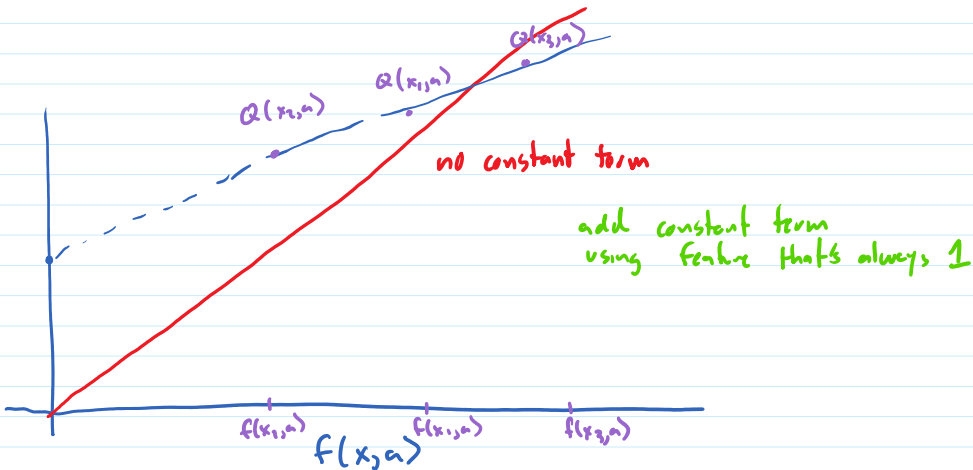
while not done
 $s \leftarrow s'$
 while not terminal
 In state s

Choose action a

Observe transition (s, a, r, s')

Update $Q(s,a) \leftarrow Q(s,a) + \alpha (\max_{a'} Q(s',a') - Q(s,a))$

for each feature $w_i \leftarrow w_i + \alpha (r + \gamma \max_{a'} Q(s',a') - Q(s,a)) \cdot f_i(s,a)$



quantized features: $f_{\text{short-yardage}}(s) = \begin{cases} 1 & \text{if yards-to-first/downs-left} \leq 2 \\ 0 & \text{otherwise} \end{cases}$

$f_1(s) = \begin{cases} 0 & \text{if yards-to-first/downs-left} \leq 2 \\ 1 & \text{if between 2 and 8} \\ 2 & \text{if " " } \geq 8 \end{cases}$

2 it

28

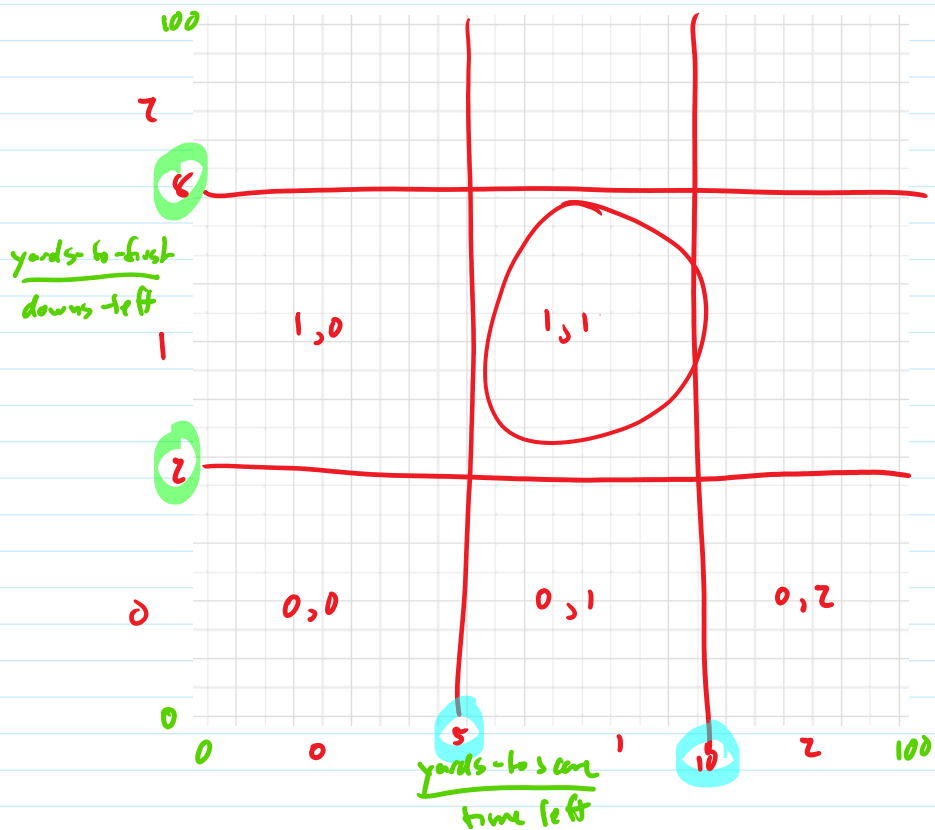
$$f_2(s) = \begin{cases} 0 \\ \frac{1}{2} \end{cases}$$

yards-to-score/time left ≤ 5

210

buckets

Buckets



$$f_{0,0} = \begin{cases} 1 & \text{if } r_1 \in \mathbb{Z} \text{ and } r_2 \leq 5 \\ 0 & \text{otherwise} \end{cases}$$

$$f_{1,1} = \begin{cases} 1 & \text{if } 2 < r_1 < 8 \text{ and } 5 < r_2 < 10 \\ 0 & \text{otherwise} \end{cases}$$

⋮

equivalent to Q-learning
except using superstates instead
of states

(terminal states in their own
bucket)

while not done

$s \leftarrow s_0$

while s not terminal

choose action a

observe transition (s, a, r, s')

compute superstate (bucket) \hat{s} for s , \hat{s}' for s'

update $Q(\hat{s}, a) \leftarrow Q(\hat{s}, a) + \alpha (r + \gamma \cdot \max_a Q(\hat{s}', a) - Q(\hat{s}, a))$