

Check the "Today's central concepts" from the previous sessions. My notes are available at:

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

Read through the lecture slides (these you can bring to exam, so familiarize yourself with them).

Solve the exercises we've solved here again on your own. Also try to solve the other exercises (solutions are available to all of them).

Good luck on the exam! 😊

Ex 6.5 d

Prove that

$$\mathbb{E}_{\pi_{\theta}} \{ \nabla \log \pi_{\theta}(s_t, a_t) \} = 0,$$

for  $a_t$  finite.

Solution:

We have that

$$\mathbb{E}_{\pi_{\theta}} \{ \nabla \log \pi_{\theta}(s_t, a_t) \} =$$

$$\sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \Pr \{ s_t = s, a_t = a \} \nabla \log \pi_{\theta}(s, a) =$$

$$\sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \Pr \{ a_t = a | s_t = s \} \Pr \{ s_t = s \} \nabla \log \pi_{\theta}(s, a) =$$

$$\sum_{s \in \mathcal{S}} \Pr \{ s_t = s \} \underbrace{\sum_{a \in \mathcal{A}} \Pr \{ a_t = a | s_t = s \} \nabla \log \pi_{\theta}(s, a)}_{= \pi_{\theta}(s, a)} =$$

$$\sum_{s \in \mathcal{S}} \Pr \{ s_t = s \} \sum_{a \in \mathcal{A}} \pi_{\theta}(s, a) \nabla \log \pi_{\theta}(s, a) =$$

$$\sum_{s \in \mathcal{S}} \Pr \{ s_t = s \} \sum_{a \in \mathcal{A}} \pi_{\theta}(s, a) \frac{\nabla \pi_{\theta}(s, a)}{\pi_{\theta}(s, a)} =$$

$$\sum_{s \in \mathcal{S}} \Pr \{ s_t = s \} \sum_{a \in \mathcal{A}} \nabla \pi_{\theta}(s, a) =$$

$$\sum_{s \in \mathcal{S}} \Pr \{ s_t = s \} \nabla \underbrace{\sum_{a \in \mathcal{A}} \pi_{\theta}(s, a)}_{= 1} = \left. \begin{array}{l} \text{Recall that} \\ \pi_{\theta}(s, a) = \Pr \{ a_t = a | s_t = s \}. \\ \text{"The probability that we take some} \\ \text{action is 1."} \end{array} \right\} =$$

$$\sum_{s \in \mathcal{S}} \Pr \{ s_t = s \} \nabla 1 = 0.$$

Ex 6.4 | Consider a system with  $S = \{A, B, C\}$  and  $\mathcal{A} = \{\varphi, \beta, \gamma\}$ . We observe a trajectory

$$(\dots, A, \varphi, 200, B, \varphi, 150, A, \gamma, 140, C, \gamma, \dots)$$

$s_t = 5734, a_t = 5734, r_t = 5734$

and

$$q^{(5734)}(s, a) = \begin{matrix} & \varphi & \beta & \gamma \\ \begin{matrix} A \\ B \\ C \end{matrix} & \begin{bmatrix} 100 & 200 & 80 \\ 120 & 50 & 140 \\ 90 & 200 & 100 \end{bmatrix} \end{matrix}.$$

- a, Provide updated q-tables for Q-learning.  $\alpha = \frac{1}{2} \quad \lambda = 0.1$
- b, What is the greedy policy?
- c, Provide updated q-tables for SARSA.  $\alpha = \frac{1}{2} \quad \lambda = 0.1$
- d, What is the greedy policy?
- e, with  $\epsilon$ -greedy action selection, what q-tables will SARSA and Q-learning converge to? How can we obtain the optimal policy with SARSA?

Solution:

2, Recall the Q-learning update:

$$q(s, a) \leftarrow q(s, a) + \alpha [r(s, a) + \lambda \max_{a'} q(s, a') - q(s, a)]$$

Then:

$q^{(5734)}$ :

	$\rho$	$\beta$	$\delta$
A	100	200	80
B	120	50	140
C	90	200	100

$s = A$   
 $a = \rho$   
 $r = 200$   
 $s' = B$

•  $q(s, a) = q(A, \rho) \leftarrow$

$$q(s, a) + \alpha [r(s, a) + \lambda \max_{a'} q(s, a') - q(s, a)] =$$

$$q(A, \rho) + \alpha [r(A, \rho) + \lambda \max_{a'} q(B, a') - q(A, \rho)] =$$

140 for  $a' = \delta$

$$100 + \frac{1}{2} [200 + 0.1 \times 140 - 100] =$$

$$100 + \frac{1}{2} [100 + 14] = 157$$

$q^{(5735)}$ :

	$\rho$	$\beta$	$\delta$
A	157	200	80
B	120	50	140
C	90	200	100

$s = B$   
 $a = \rho$   
 $r = 150$   
 $s' = A$

•  $q(B, \rho) \leftarrow$

$$q(B, \rho) + \alpha [r(B, \rho) + \lambda \max_{a'} q(A, a') - q(B, \rho)] =$$

200 for  $a' = \beta$

$$120 + \frac{1}{2} [150 + 0.1 \times 200 - 120] =$$

$$120 + \frac{1}{2} [150 - 100] = 145$$

$q^{(5736)}$ :

	$\rho$	$\beta$	$\delta$
A	157	200	80
B	145	50	140
C	90	200	100

$s = A$   
 $a = \delta$   
 $r = 140$   
 $s' = C$

•  $q(A, \delta) \leftarrow$

$$q(A, \delta) + \alpha [r(A, \delta) + \lambda \max_{a'} q(C, a') - q(A, \delta)] =$$

200 for  $a' = \beta$

$$80 + \frac{1}{2} [140 + 0.1 \times 200 - 80] =$$

$$80 + \frac{1}{2} [140 - 60] = 120$$

5

	$\varphi$	$\beta$	$\gamma$
$q^{(5737)}$ : A	157	200	120
B	145	50	140
C	90	200	100

And then we don't have more data.

b, This corresponds to a greedy policy:

- $\pi^{(5737)}(s=A) = \arg \max_a q^{(5737)}(A, a) = \beta$
- $\pi^{(5737)}(s=B) = \varphi$
- $\pi^{(5737)}(s=C) = \beta$

5) Recall the SARSA update:

$$q(s, a) \leftarrow q(s, a) + \alpha [r(s, a) + \lambda q(s', a') - q(s, a)]$$

Then:

$q^{(5734)}$ :

	$\phi$	$\beta$	$\gamma$
A	100	200	80
B	120	50	140
C	90	200	100

•  $q(s, a) = q(A, \phi) \leftarrow$

$$q(s, a) + \alpha [r(s, a) + \lambda q(s', a') - q(s, a)] =$$

$$q(A, \phi) + \alpha [r(A, \phi) + \lambda q(B, \phi) - q(A, \phi)] =$$

$$100 + \frac{1}{2} [200 + 0.1 \times 120 - 100] =$$

$$100 + \frac{1}{2} [100 + 12] = 156$$

$s = A$   
 $a = \phi$   
 $r = 200$   
 $s' = B$   
 $a' = \phi$

$q^{(5735)}$ :

	$\phi$	$\beta$	$\gamma$
A	156	200	80
B	120	50	140
C	90	200	100

•  $q(B, \phi) \leftarrow$

$$q(B, \phi) + \alpha [r(B, \phi) + \lambda q(A, \gamma) - q(B, \phi)] =$$

$$120 + \frac{1}{2} [150 + 0.1 \cdot 80 - 120] =$$

$$120 + \frac{1}{2} [30 + 8] = 139$$

$s = B$   
 $a = \phi$   
 $r = 150$   
 $s' = A$   
 $a' = \gamma$

$q^{(5736)}$ :

	$\phi$	$\beta$	$\gamma$
A	156	200	80
B	139	50	140
C	90	200	100

•  $q(A, \gamma) \leftarrow$

$$q(A, \gamma) + \alpha [r(A, \gamma) + \lambda q(C, \gamma) - q(A, \gamma)] =$$

$$80 + \frac{1}{2} [140 + 0.1 \times 100 - 80] =$$

$$80 + \frac{1}{2} [60 + 10] = 115$$

$s = A$   
 $a = \gamma$   
 $r = 140$   
 $s' = C$   
 $a' = \gamma$

	$\varphi$	$\beta$	$\gamma$
$q^{(5737)}$ : A	156	200	115
B	139	50	140
C	90	200	100

And then there is no more data.

d, This corresponds to a greedy policy:

$$\pi^{(5737)}(s=A) = \beta$$

$$\pi^{(5737)}(s=B) = \gamma$$

$$\pi^{(5737)}(s=C) = \beta$$

e, with  $\epsilon$ -greedy action selection, SARSA will converge to the optimal policy taking into account that we do random exploration with probability  $\epsilon$ . By letting  $\epsilon \rightarrow 0$ , we will tend to the optimal policy (which is what Q-learning is converging to).

Ex 6.3 | An order is received w.p.  $p$  at each time step.

One can choose to process all orders (setup cost  $K > 0$ ), or wait (costs  $c > 0$  per order on hold). The maximum number of unfilled orders is  $n$ . There is a discount factor  $\lambda$ . Characterize an optimal processing policy.

Solution:

State-space:  $S = \{0, 1, \dots, n\}$

# of unfilled orders when starting the period

Actions:

$$A_n = \{P\}$$

process all orders

$$A_i = \{W, P\}$$

wait

process all orders

$$i = 0, 1, \dots, n-1$$

Rewards:

In this case, we deal with costs:

$$\bullet r(i, P) = -K$$

$$\bullet r(i, W) = -ci$$

$$\bullet r(n, P) = -K$$

$$i = 0, 1, \dots, n-1$$

Time-horizon and objective:

$\infty$ -horizon, discounted:

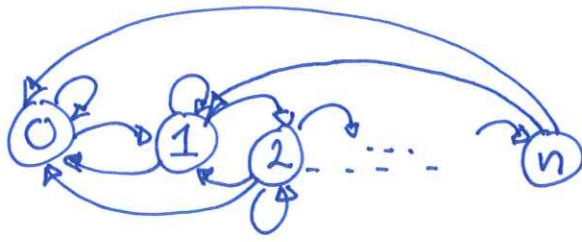
$$\mathbb{E} \left\{ \sum_{t=0}^{\infty} \lambda^t r(s_t, a_t) \right\}$$

Remark:

we can interpret  $1-\lambda$  as the probability of "going out of business" (e.g., bankruptcy).



## Transitions:



- $P(i|i, W) = 1-p$   
wait, and no new order
- $P(i+1|i, W) = p$   
wait, and a new order  
 $i = 0, 1, \dots, n-1$
- $P(0|i, P) = 1-p$   
process everything, and no new order
- $P(1|i, P) = p$   
process everything, and receive a new order
- $P(0|n, P) = 1-p$
- $P(1|n, P) = p$

## Bellman equation:

$$u^*(s) = \max_{a \in A_s} \left\{ r(s, a) + \lambda \sum_{s' \in S} p(s'|s, a) u^*(s') \right\}$$

If  $s=n$ :

$$u^*(n) = \max_{a \in A_n} \left\{ \underbrace{r(n, a)}_{= -K} + \lambda \sum_{s' \in S} p(s'|n, a) u^*(s') \right\}$$

$$= -K + \lambda [(1-p)u^*(0) + pu^*(1)]$$

(\*)

If  $s = i \wedge n$ :

$$u^*(i) = \max_{a \in A_i} \left\{ \underbrace{r(i, a)}_{= \{P, W\}} + \lambda \sum_{s' \in S} p(s' | i, a) u^*(s') \right\}$$

Evaluate the actions separately:

If  $a = P$ :

$$\underbrace{r(i, P)}_{=-K} + \lambda \sum_{s' \in S} p(s' | i, P) u^*(s') =$$

$$-K + \lambda [(1-p)u^*(0) + pu^*(1)] \stackrel{\text{def.}}{=} \delta$$

If  $a = W$ :

$$\underbrace{r(i, W)}_{=-ci} + \lambda \sum_{s' \in S} p(s' | i, W) u^*(s') =$$

$$-ci + \lambda [(1-p)u^*(i) + pu^*(i+1)]$$

So that:

$$u^*(i) = \max \left\{ \underset{P}{\delta}, \underset{W}{-ci + \lambda [(1-p)u^*(i) + pu^*(i+1)]} \right\}. \quad (**)$$

It is reasonable to assume that the optimal policy is of threshold type. Can we prove that it is?

A threshold policy is such that

$$a^*(i) = \tau \Rightarrow a^*(i+1) = \tau.$$

This is equivalent to

$$\delta > -c_i + \lambda [(1-p)u^*(i) + pu^*(i+1)] \Rightarrow$$

$$\delta > -c_{(i+1)} + \lambda [(1-p)u^*(i+1) + pu^*(i+2)]$$

Note that this holds if

$$u^*(i) \geq u^*(i+1), \quad (\square)$$

since  $c > 0$ .

Because then:

$$\delta > -c_i + \lambda [(1-p)u^*(i) + pu^*(i+1)] \quad \leftarrow \text{since } c > 0$$

$$> -c_{(i+1)} + \lambda [(1-p)u^*(i) + pu^*(i+1)]$$

$$\geq -c_{(i+1)} + \lambda [(1-p)\underbrace{u^*(i)}_{\geq u^*(i+1)} + p\underbrace{u^*(i+1)}_{\geq u^*(i+2)}]$$

Can we prove that  $(\square)$  holds?

We will use value iteration to prove  $(\square)$ .

Base case:

Let  $u_0(s) = 0$  for all  $s \in S$ .

Then  $u_0(i) \geq u_0(i+1)$  holds trivially.

Induction:

Assume that  $u_k(i) \geq u_k(i+1)$ .

(we will show that then  $u_{k+1}(i) \geq u_{k+1}(i+1)$  holds.)

If  $s = i+1 \in N$ : (see p. 15 for an alternative proof.)

It is clear that the following holds:

$$\begin{aligned}
 & \underbrace{-c(i+1)}_{\substack{\leq -c_i \\ \text{since} \\ c > 0}} + \lambda \left[ \underbrace{(1-p)u_k(i+1)}_{\substack{\leq u_k(i) \\ \text{by assumption}}} + \underbrace{pu_k(i+2)}_{\substack{\leq u_k(i+1) \\ \text{by assumption}}} \right] \leq -c_i + \lambda \left[ (1-p)u_k(i) + pu_k(i+1) \right] \quad (\Delta)
 \end{aligned}$$

Let  $F_k(x)$  be defined as

$$F_k(x) = \max \left\{ -K + \lambda \left[ (1-p)u_k(0) + pu_k(1) \right], x \right\}.$$

This function is increasing in x.

One iteration of value iteration is (compare the computations that led to  $(**)$ ):

$$\begin{aligned}
 u_{k+1}(i+1) &= \max \left\{ -K + \lambda \left[ (1-p)u_k(0) + pu_k(1) \right], \right. \\
 & \quad \left. -c(i+1) + \lambda \left[ (1-p)u_k(i+1) + pu_k(i+2) \right] \right\} \\
 &= F_k(-c(i+1) + \lambda \left[ (1-p)u_k(i+1) + pu_k(i+2) \right])
 \end{aligned}$$

$$\leq F_k(-ci + \lambda [(1-p)u_k(i) + pu_k(i+1)])$$

by  
(Δ)  
and that  
 $F_k(x) \uparrow x$

$$= \max \left\{ -K + \lambda [(1-p)u_k(0) + pu_k(1)], \right. \\ \left. -ci + \lambda [(1-p)u_k(i) + pu_k(i+1)] \right\}$$

$$= u_{k+1}(i).$$

by definition of one  
iteration of value iteration,  
compare (\*\*).

If  $s = n$ :

Again, one iteration of value iteration (compare  
computation of (\*)) yields:

$$u_{k+1}(n) = -K + \lambda [(1-p)u_k(0) + pu_k(1)]$$

$$\leq \max \left\{ -K + \lambda [(1-p)u_k(0) + pu_k(1)], \right. \\ \left. -c(n-1) + \lambda [(1-p)u_k(n-1) + pu_k(n)] \right\}$$

obviously  
smaller than  
the max of itself  
and something  
else

$$= u_{k+1}(n-1).$$

compare  
(\*\*)

In summary, we have now shown that

$$u_k(i) \geq u_k(i+1) \Rightarrow u_{k+1}(i) \geq u_{k+1}(i+1)$$

for all  $i = 0, \dots, n-1$ .

Taking the limit:

We have shown that

- $u_0(i) \geq u_0(i+1)$
- $u_k(i) \geq u_k(i+1) \Rightarrow u_{k+1}(i) \geq u_{k+1}(i+1)$ .

Since it is well known that

$$u^*(s) = \lim_{k \rightarrow \infty} u_k(s),$$

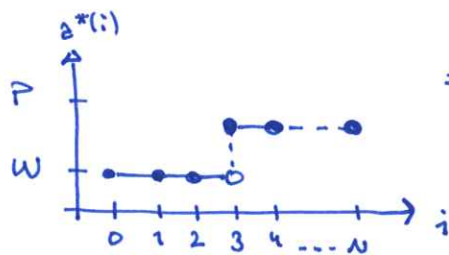
when  $u_k$  is computed using value iteration for discounted problems, we conclude that

$$u^*(i) \geq u^*(i+1),$$

which is  $\square$ .

See, e.g.,  
p. 84 in Vol. 2  
of Bertsekas'  
"Dynamic Program.  
and Optimal Control".

In summary, this proves that  $a^*(i) = P \Rightarrow a^*(i+1) = P$



$\Rightarrow$  To find the optimal policy, we just need to check which of the  $N+1$  possible threshold policies is the best.

( This entails performing the policy evaluation step of policy iteration  $N+1$  times. )

If  $s = i+1 < n$ : (Alternative)

One iteration of value iteration is:

$$u_{k+1}(i+1) = \max_{a \in \mathcal{A}_{i+1}} \left\{ r(i+1, a) + \sum_{s' \in \mathcal{S}} p(s'|i+1, a) u_k(s') \right\}$$

= { compare calculations for (\*\*) }

$$= \max \left\{ \underbrace{-K + \lambda [(1-p)u_k(0) + pu_k(1)]}_{\stackrel{\text{def.}}{=} \delta_k^p, \text{ constant in } i}, \right.$$

$$\left. -c(i+1) + \lambda [(1-p)u_k(i+1) + pu_k(i+2)] \right\}$$

$$= \max \left\{ \delta_k^p, -c(i+1) + \lambda [(1-p)u_k(i+1) + pu_k(i+2)] \right\}$$

= { Note:  
 $\max \{ \delta_k^p, x \}$  is increasing in  $x$  }

$$\stackrel{(\Delta)}{\leq} \max \left\{ \delta_k^p, -c i + \lambda [(1-p)u_k(i) + pu_k(i+1)] \right\}$$

$$= \max \left\{ -K + \lambda [(1-p)u_k(0) + pu_k(1)], \right.$$

$$\left. -c i + \lambda [(1-p)u_k(i) + pu_k(i+1)] \right\}$$

$$= u_{k+1}(i).$$