

1. Let  $V_0(x) \equiv$  maximal expected return  
 if gambler has a total amt. of  $x$  & is  
 allowed  $N$  more gambling moves.

Then we have the DP algorithm as

$$(a) V_k(x) = \max_{0 \leq \alpha \leq 1} [\beta V_{k+1}(x + \alpha x) + q V_{k+1}(x - \alpha x)],$$

$$k = N-1, N-2, \dots, 0$$

with  $V_N(x) = \log x$ .

(b) When  $\beta > q$ , since  $\beta + q = 1$ ,

$$1 > 2q \Rightarrow q < \frac{1}{2}$$

$$\Rightarrow \beta > \frac{1}{2}$$

Then,

$$V_{N-1}(x) = \max_{0 \leq \alpha \leq 1} [\beta V_N(x+\alpha x) + q V_N(x-\alpha x)]$$

$$= \max_{0 \leq \alpha \leq 1} [\beta \log(x+\alpha x) + q \log(x-\alpha x)]$$

Let  $f(\alpha) = \beta \log(x+\alpha x) + q \log(x-\alpha x)$

$$= \beta \log((1+\alpha)x) + q \log((1-\alpha)x)$$

$$= \beta \log(1+\alpha) + q \log(1-\alpha)$$

$$+ (\beta + q) \log x$$

||  
/

Then  $f'(\alpha) = \frac{P}{1+\alpha} - \frac{q}{1-\alpha} = 0$

$$\Rightarrow P - q - (\beta + q)\alpha = 0$$

$$\Rightarrow \alpha = p-q$$

(Note:  $0 \leq \alpha \leq 1$ )

Now  $f''(\alpha) = -p(1+\alpha)^{-2} - q(1-\alpha)^{-2}$

$$< 0 \quad \forall \alpha \in [0, 1]$$

$\Rightarrow \alpha = p-q$  is maximizer for  $f(\alpha)$

Now

$$V_{N-1}(x) = p \log(1+p-q) + q \log(1-p+q) + \log x$$

$$= p \log(2p) + q \log(2q) + \log x$$

1/1

K

$$\therefore V_{N-1}(x) = \log x + K$$

Now

$$V_{N-2}(x) = \max_{0 \leq \alpha \leq 1} [p V_{N-1}(x+\alpha x) + q V_{N-1}(x-\alpha x)]$$

$$= \max_{0 \leq \alpha \leq 1} [p \log(x+\alpha x) + q \log(x-\alpha x)] + K$$

Again the max above is obtained at  
 $\alpha = p-q$ .

And  $V_{N-2}(x) = \log x + 2K$

⋮

Proceeding similarly,

$$V_0(x) = \log x + NK.$$

Optimal action: Bet  $(p-q)$  fraction of

the total amount at each instant.

(C) When  $p < q$ , again consider

$$f(x) = p \log(1+x) + q \log(1-x) + \log x,$$
$$0 \leq x \leq 1.$$

$$\text{So, } f'(x) = \frac{p}{1+x} - \frac{q}{1-x}.$$

For  $0 \leq x \leq 1$ ,

$$\frac{p}{1+x} < \frac{q}{1-x}$$

and  $f'(x) < 0 \quad \forall x \in [0, 1]$

$\Rightarrow f$  is a decreasing function  
on  $[0, 1]$

$\Rightarrow$  Max is attained at  $\alpha = 0$   
in interval  $[0, 1]$ .

$$\begin{aligned} & V_{N-1}(x) = p \log(1) + q \log(1) \\ & \quad + \log x \\ & \quad = \log x. \end{aligned}$$

Similarly, it can be seen that

$$V_k(x) = \log x, \quad \forall k = 0, 1, \dots, N-1$$

& the optimal strategy is to bet  
0 at any instant.

2. Note that

$$R(i, \mu_i(i)) + \gamma \sum_{j \in S} p_{ij}(\mu_i(i)) \max\{J^{\mu_1(j)}, J^{\mu_2(j)}\}$$

$$= \max_a \left( R(i, a) + \gamma \sum_{j \in S} p_{ij}(a) \max\{J^{\mu_1(j)}, J^{\mu_2(j)}\} \right)$$

$$\geq \max_a \left( R(i, a) + \gamma \sum_{j \in S} p_{ij}(a) J^{\mu_1(j)} \right)$$

$$\geq R(i, \mu_1(i)) + \gamma \sum_{j \in S} p_{ij}(\mu_1(i)) J^{\mu_1(j)}$$

$$= J^{\mu_1(i)}. \quad \longrightarrow (1)$$

Similarly, we obtain

$$\begin{aligned}
& R(i, \mu(i)) + \gamma \sum_{j \in S} p_{ij} \cdot (\mu(i)) \max \left\{ J^{M_1(j)}, \right. \\
& \quad \left. J^{M_2(j)} \right\} \\
& \geq R(i, \mu_2(i)) + \gamma \sum_{j \in S} p_{ij} \cdot (\mu_2(i)) J^{M_2(j)} \\
& = J^{M_2(i)}. \quad \longrightarrow (2)
\end{aligned}$$

From (1) and (2), it follows that

$$\begin{aligned}
& R(i, \mu(i)) + \gamma \sum_{j \in S} p_{ij} \cdot (\mu(i)) \max \left\{ J^{M_1(j)}, \right. \\
& \quad \left. J^{M_2(j)} \right\} \\
& \geq \max \left\{ J^{M_1(i)}, J^{M_2(i)} \right\} \\
& \quad \# i \in S.
\end{aligned}$$

In terms of the operator  $T_\mu$ , the above

is analogous to

$$T_\mu \max\left\{\frac{I_1^{\mu_1}}{I}, \frac{I_2^{\mu_2}}{I}\right\} \geq \max\left\{\frac{I_1^{\mu_1}}{I}, \frac{I_2^{\mu_2}}{I}\right\}$$

Using monotonicity of  $T_\mu$ , we have

$$\overline{I}^\mu = \lim_{n \rightarrow \infty} T_\mu^n \max\left\{\frac{I_1^{\mu_1}}{I}, \frac{I_2^{\mu_2}}{I}\right\}$$

$$\geq \dots \geq T_\mu^{k+1} \max\left\{\frac{I_1^{\mu_1}}{I}, \frac{I_2^{\mu_2}}{I}\right\}$$

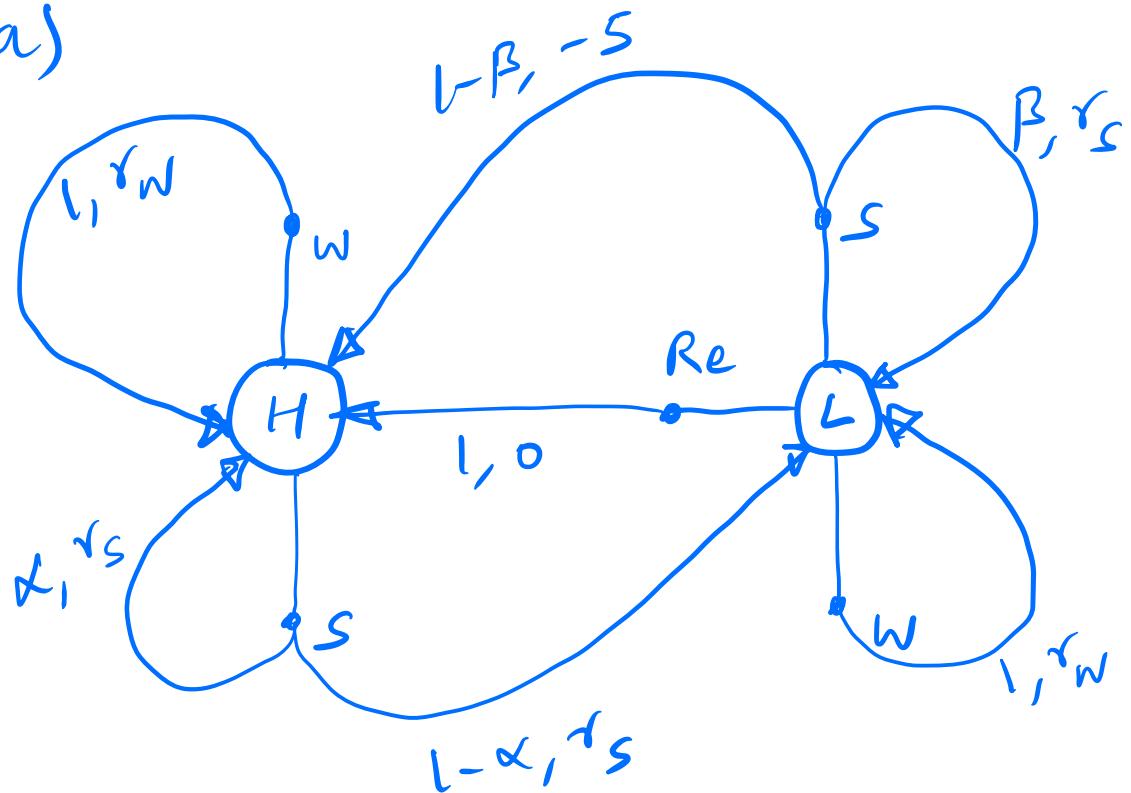
$$\geq T_\mu^k \max\left\{\frac{I_1^{\mu_1}}{I}, \frac{I_2^{\mu_2}}{I}\right\}$$

$$\geq \dots \geq T_\mu \max\left\{\frac{I_1^{\mu_1}}{I}, \frac{I_2^{\mu_2}}{I}\right\}$$

$$\geq \max\left\{\frac{I_1^{\mu_1}}{I}, \frac{I_2^{\mu_2}}{I}\right\}$$

$$\Rightarrow \overline{I}^\mu(i) \geq \max\left\{\frac{I_1^{\mu_1}(i)}{I}, \frac{I_2^{\mu_2}(i)}{I}\right\}, \forall i \in S$$

3. (a)



State transition Diagram

(b) Sample trajectory:

$H, w, \alpha_w, H, s, \alpha_s, L, R_e, 0, H, s, r_s, L, w, r_w, L, s, r_s, L, s, -5, H$

final state  
↓

- High state  $H$  occurs 3 times when an action is chosen & reward obtained and once as the final (terminal) state
- Low state  $L$  occurs 4 times

First visit Monte Carlo estimates for H & L

$$\hat{V}(H) \leftarrow 1 + 2 + 0 + 2 + 1 + 2 - 5 = 3$$

$$\hat{V}(L) \leftarrow 0 + 2 + 1 + 2 - 5 = 0$$

∴ First visit estimates

$$\hat{V}(H) = 3 \quad \& \quad \hat{V}(L) = 0$$

(c) Every visit Monte Carlo Estimates  
for H:

$$\hat{V}_1(H) \leftarrow 1 + 2 + 0 + 2 + 1 + 2 - 5 = 3$$

$$\hat{V}_2(H) \leftarrow 2 + 0 + 2 + 1 + 2 - 5 = 2$$

$$\hat{V}_3(H) \leftarrow 2 + 1 + 2 - 5 = 0$$

$$\Rightarrow \hat{V}(H) = \frac{3+2+0}{3} = \frac{5}{3}$$

Every visit Monte Carlo Estimates for L:

$$\hat{V}_1(L) \leftarrow 0 + 2 + 1 + 2 - 5 = 0$$

$$\hat{V}_2(L) \leftarrow 1 + 2 - 5 = -2$$

$$\hat{V}_3(L) \leftarrow 2 - 5 = -3$$

$$\hat{V}_4(L) \leftarrow -5$$

$$\Rightarrow \hat{V}(L) = \frac{0 - 2 - 3 - 5}{4} = -\frac{5}{2}.$$

-: Every visit estimates

$$\hat{V}(H) = \frac{5}{3} \text{ and } \hat{V}(L) = -\frac{5}{2}.$$

(d)  $V_0(H) = V_0(L) = 0$  (given initial values)

$\beta = 1$  (learning rate)

$\gamma = 1$  (no discounting)

## TD(0) update

$$V(s) \leftarrow V(s) + \gamma (r(s) + \gamma V(s') - V(s))$$

## steps of TD(0) update

$$\begin{aligned} V_1(H) &= V_0(H) + 1 \cdot (r_w + V_0(H) - V_0(H)) \\ &= 0 + 1 = 1 \end{aligned}$$

$$V_1(L) = V_0(L) = 0$$

$$\begin{aligned} V_2(H) &= \cancel{V_1(H)} + 1 \cdot (r_s + V_1(L) - \cancel{V_1(H)}) \\ &= 2 + 0 = 2 \end{aligned}$$

$$V_2(L) = V_1(L) = 0$$

$$V_3(L) = \cancel{V_2(L)} + 1 \cdot (r_{re} + V_2(H) - \cancel{V_2(L)})$$

$$= 0 + 2 = 2$$

$$V_3(H) = V_2(H) = 2$$

$$\begin{aligned}V_4(H) &= \cancel{V_3(H)} + 1 \cdot (r_s + V_3(L) - \cancel{V_3(H)}) \\&= 2 + 2 = 4\end{aligned}$$

$$V_4(L) = V_3(L) = 2$$

$$\begin{aligned}V_5(L) &= V_4(L) + 1 \left( r_w + \cancel{V_4(L)} - \cancel{V_4(L)} \right) \\&= 2 + 1 = 3\end{aligned}$$

$$V_5(H) = V_4(H) = 4$$

$$\begin{aligned}V_6(L) &= V_5(L) + 1 \left( r_s + V_5(L) - \cancel{V_5(L)} \right) \\&= 3 + 2 = 5\end{aligned}$$

$$V_6(H) = V_5(H) = 4$$

$$\begin{aligned}V_7(L) &= \cancel{V_6(L)} + 1(-5 + V_6(H) - \cancel{V_6(L)}) \\&= -5 + 4 = -1\end{aligned}$$

$$V_7(H) = V_6(H) = 4$$

Final values obtained from TD(0):

$$\hat{V}(H) = 4, \quad \hat{V}(L) = -1$$

(e) Application of TD(1) =

$$V(H) = V(L) = 0 \quad (\text{Initialization})$$

$$\begin{aligned}V(H) &:= V(H) + 1(1+2+0+2+1+2 \\&\quad - 5 \\&\quad - V(H)) \\&= 3\end{aligned}$$

$$V(H) := V(\cancel{H}) + 1(2+0+2+1+2-5 - V(\cancel{H}))$$
$$= 2$$

$$V(L) := V(\cancel{L}) + 1(0+2+1+2-5 - V(\cancel{L}))$$
$$= 0$$

$$V(H) := V(\cancel{A}) + 1(2+1+2-5 - V(\cancel{A}))$$
$$= 0$$

$$V(L) := V(\cancel{J}) + 1(1+2-5 - V(\cancel{J}))$$
$$= -2$$

$$V(L) := V(\cancel{K}) + 1(2-5 - V(\cancel{L}))$$
$$= -3$$

$$V(L) := V(\cancel{J}) + 1(-5 - V(\cancel{L}))$$
$$= -5$$

Av. of the value estimates for TD(1) gives

$$\hat{V}(H) = \frac{3+2+0}{3} = 5/3$$

$$\hat{V}(L) = \frac{0-2-3-5}{4} = -\frac{10}{4} = -\frac{5}{2}$$

Note: These are same as every visit  
Monte Carlo estimates.

4. Note that  $W^P : \mathbb{R}^{|S|} \rightarrow \mathbb{R}^{|S|}$  is a contraction.  
(a)

$\Rightarrow W^P$  has a unique fixed pt.  $V^*$  such  
that  $V^* = W^P V^*$ . Applying  $W$  to both  
sides of the fixed pt. equation above giving  
 $WV^* = W^{P+1}V^* = W^P WV^*$ .

$\Rightarrow WU^*$  is also a fixed pt. of  $W^P$ .

By uniqueness of fixed pt, we have

$$U^* = WU^*.$$

(b) Since  $W^P$  is a contraction,

$\exists \alpha \in (0, 1)$  s.t.

$$\|W^P U - W^P \bar{U}\| \leq \alpha \|U - \bar{U}\|,$$

$U, \bar{U} \in \mathbb{R}^{l \times l}$ .

Thus,

$$\begin{aligned} \|W^{2P} U - W^{2P} \bar{U}\| &\leq \alpha \|W^P U - W^P \bar{U}\| \\ &\leq \alpha^2 \|U - \bar{U}\| \end{aligned}$$

$\vdots$

Proceeding in this manner, we obtain

$$\|W^{kp}U - W^{kp}\bar{U}\| \leq \alpha^k \|U - \bar{U}\|, \quad k \geq 1.$$

Let if  $\bar{U} = U^*$ , since  $W^p U^* = U^*$ ,

it follows that  $W^{kp} U^* = U^*$   
 $\forall k \geq 1$ .

Thus,

$$\|W^{kp}U - U^*\| \leq \alpha^k \|U - U^*\|, \quad \forall k \geq 1.$$

$$\Rightarrow \lim_{k \rightarrow \infty} \|W^{kp}U - U^*\| = 0$$

Since  $U$  is arbitrary, we may replace

$U$  with  $W^i U$  for  $i = 0, 1, \dots, p-1$ .

Then, we have

$$\lim_{k \rightarrow \infty} \| w^{kb+i} v - v^* \| = 0,$$

$\forall i=0, 1, \dots, b-1$

Now any  $n \in \{0, 1, 2, \dots\}$  can be written

a)  $n = kb+i$ , for some  $i \in \{0, 1, \dots, b-1\}$ ,

$k \geq 0$ . Thus, we have that

$$\lim_{n \rightarrow \infty} w^n v = v^*. \quad \square$$

Q5.

- (a) Given a history dependent policy  $\hat{\pi} \in \Pi_1$  and an initial state  $i \in S$ , define a markovian policy  $\pi \in \Pi_2$  such that

$$\pi = (\pi_0, \pi_1, \pi_2, \dots)$$

where

$$\pi_t(a_t = a | s_t = s) = P^{\hat{\pi}} \left\{ a_t = a \mid s_t = s, s_0 = i \right\}$$

Claim :  $P^{\hat{\pi}} \left\{ s_t = s, a_t = a \mid s_0 = i \right\} = P^\pi \left\{ s_t = s, a_t = a \mid s_0 = i \right\}$

proof: (By induction)

Base Case: for  $t=0$ ,

$$\text{LHS} = P^{\hat{\pi}} \left\{ s_0 = s, a_0 = a \mid s_0 = i \right\}$$

$$= P^{\hat{\pi}} \left\{ a_0 = a \mid s_0 = i \right\} \mathbb{1}_{\{s=i\}}$$

$$= \pi_i(a_0 = a \mid s_0 = i) \mathbb{1}_{\{s=i\}} = P^{\pi} \left\{ a_0 = a \mid s_0 = i \right\} \mathbb{1}_{\{s=i\}}$$

$$= P^{\pi} \left\{ s_0 = s, a_0 = a \mid s_0 = i \right\}$$

Induction : Assume the relation holds  $\forall t = 1, 2, \dots, n-1$

Now,  $P^{\hat{\pi}} \{ S_n = s \mid S_0 = i \}$

$$= \sum_{k \in S} \sum_{u \in A} P^{\hat{\pi}} \{ S_{n-1} = k, a_{n-1} = u \mid S_0 = i \} \underbrace{P(s \mid k, u)}_{\text{Transition prob. for the MDP.}}$$

$$= \sum_{k \in S} \sum_{u \in A} P^{\pi} \{ S_{n-1} = k, a_{n-1} = u \mid S_0 = i \} \underbrace{P(s \mid k, u)}_{\text{(from inductive step)}}$$

$$= P^\pi \{ S_n = s \mid S_0 = i \}$$

Now,  $P^{\hat{\pi}} \{ S_n = s, a_n = a \mid S_0 = i \}$

$$= P^{\hat{\pi}} \{ a_n = a \mid S_n = s, S_0 = i \} P^\pi \{ S_n = s \mid S_0 = i \}$$

$$= \pi_n(a_n = a \mid S_n = s) P^\pi \{ S_n = s \mid S_0 = i \}$$

$$= P^\pi \{ a_n = a \mid S_n = s, S_0 = i \} P^\pi \{ S_n = s \mid S_0 = i \}$$

$$= P^\pi \{ S_n = s, a_n = a \mid S_0 = i \}$$

This completes the proof for (a).

(b)

$$J_{\hat{\pi}}(i) = \mathbb{E}_{\hat{\pi}} \left( \sum_{t=0}^T c(s_t, a_t, s_{t+1}) \mid s_0 = i \right)$$

$$= \sum_{t=0}^T \sum_{s_t, a_t, s_{t+1}} c(s_t, a_t, s_{t+1}) \hat{P}^{\hat{\pi}} \{ s_t, a_t, s_{t+1} \mid s_0 = i \}$$

$$= \sum_{t=0}^T \sum_{s_t, a_t, s_{t+1}} c(s_t, a_t, s_{t+1}) \underbrace{\hat{P}^{\hat{\pi}} \{ s_{t+1} \mid s_t, a_t, s_0 = i \}}_{\text{Depends only on the transition prob. of the MDP and not on the policy}} \underbrace{\hat{P}^{\hat{\pi}} \{ s_t, a_t \mid s_0 = i \}}_{P^{\pi} \{ s_t, a_t \mid s_0 = i \}}$$

Depends only on  
the transition  
prob. of the MDP  
and not on the  
policy

$$P(s_{t+1} \mid s_t, a_t)$$

from part (a)

$$= \sum_{t=0}^T \sum_{s_t, a_t, s_{t+1}} c(s_t, a_t, s_{t+1}) P^\pi \{ s_t, a_t, s_{t+1} \mid s_0 = i \}$$

$$= \mathbb{E}_\pi \left( \sum_{t=0}^T c(s_t, a_t, s_{t+1}) \mid s_0 = i \right)$$