HW #2: RETRIEVE THE STATE-VALUE FUNCTION OF FIGURE 3.2

OBSERVE THAT THE BELLMAN EQUATION (3.14) MAY BE WRITTEN IN MATRIX FORM AS

$$V = Av + b$$

WHERE THE ELEMENT OF A CORRESPONDING TO STATE PAIR $(s, s')$ IS

$$A_{s,s'} = \sum_{a} \pi(a|s) \sum_{s',r} p(s',r|s,a)\, \gamma$$

(FOR $\gamma \cdot 1$, THIS MATRIX IS STOCHASTIC (RIGHT), MEANING THAT EACH ROW IS A DISTRIBUTION, FOR SUCH MATRICES, ALL EIGENVALUES HAVE $|\lambda_i| \leq 1$ )

ALSO

$$b_s = \sum_{a} \pi(a|s) \sum_{s',r} p(s',r|s,a)\, r$$

ESPECIALLY FOR THE GRIDWORLD

$$A_{s,s'} = \gamma \sum_{a} \sum_{s'} \pi(a|s)\, p(s'|s,a) =$$

$$\gamma \sum_{a} \pi(a|s)\, \mathbb{1}_{s' = s'(a,s)}$$

$$b_s = \sum_{a} \pi(a|s)\, r(s,a)$$

THEREFORE:

$$V = AV + b \iff IV - AV = b \iff$$

$$V = (I-A)^{-1} b$$

THIS IS THE SOLUTION THAT THE NN MUST PROVIDE

(3.6) OPTIMAL POLICIES AND OPTIMAL VALUE FUNCTIONS

DEFINITION 1) A POLICY $\pi'$ IS BETTER THAN $\pi$ IF OR EQUAL

$$v_{\pi'}(s) \geq v_{\pi}(s) \quad \forall s \in S$$

SO THERE IS A PARTIAL ORDERING (AS OPPOSED TO A TOTAL ORDERING) AMONG POLICIES

2) A POLICY $\pi'$ IS OPTIMAL IF IT IS BETTER THAN OR EQUAL TO ALL OTHER POLICIES.

THEOREM: THERE IS ALWAYS AT LEAST ONE OPTIMAL POLICY

COMMENTS: 1) NOT OBVIOUS AT ALL

2) PROOF IS BY SHOWING THE BELLMAN OPTIMALITY CONDITION (TO BE SHOWN LATER) IS A CONTRACTION MAPPING.

WE DEFINE

$$v_*(s) \overset{\Delta}{=} \max_{\pi} v_\pi(s) = E_*[G_t \mid S_t = s] \quad \forall s \in S$$

$$q_*(s, a) \overset{\Delta}{=} \max_{\pi} q_\pi(s, a) = E_*[G_t \mid S_t = s, A_t = a]$$
$$\forall s \in S, \forall a \in A(s)$$

OPTIMAL RETURN IF WE TAKE ACTION $a$ AND THEN BEHAVE OPTIMALLY

THE TWO ARE CONNECTED:

$$q_*(s,a) = E_*[R_{t+1} + \gamma G_{t+1} \mid S_t=s, A_t=a]$$

$$= E_*[R_{t+1} \mid S_t=s, A_t=a]$$

$$+ \gamma \sum_{s',r} p(s',r \mid s,a) E_*[G_{t+1} \mid S_t=s, A_t=a, S_{t+1}=s', R_{t+1}=r]$$

$$= E_*[R_{t+1} \mid S_t=s, A_t=a]$$

$$+ \gamma \sum_{s',r} p(s',r \mid s,a) v_*(s') =$$

$$E_*[R_{t+1} \mid S_t=s, A_t=a] + E_*[\gamma v_*(S_{t+1}) \mid S_t=s, A_t=a]$$

$$\Rightarrow \boxed{q_*(s,a) = E_*[R_{t+1} + \gamma v_*(S_{t+1}) \mid S_t=s, A_t=a]}$$ Ⓐ

(INTUITIVELY CLEAR THAT THIS SHOULD HOLD)

We CAN ALSO GIVE $v_*$ IN TERMS OF $q$

$$\boxed{v_*(s) = \max_{a \in A(s)} q_*(s,a)}$$ Ⓑ

INDEED, SKETCH of proof:

$$v_*(s) = \max_\pi E_\pi[G_t \mid S_t=s] =$$

$$= \max_\pi \sum_{a \in A(s)} \pi(a \mid s) \underbrace{E_\pi[G_t \mid S_t=s, A_t=a]}_{}$$

$$\leq q_{\pi_*}(s,a),$$

ACHIEVED BY MAXIMIZING THESE AND THEN PICKING ITS COEFFICIENT UNITY

SO WE PROVED $\leq$.

WE CAN EXCLUDE THE STRICT EQUALITY BY
CONTRADICTION

## BELLMAN OPTIMALITY CONDITION FOR $v_*(s)$

$$v_*(s) = \max_{a \in A(s)} q_*(s, a)$$

$$\overset{(A)}{=} \max_{a \in A(s)} E_*\left[R_{t+1} + \gamma v_*(S_{t+1}) \mid S_t = s, A_t = a\right]$$

$$= \max_{a \in A(s)} \sum_{s', r} p(s', r \mid s, a)\left[r + \gamma v_*(s')\right]$$

$$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad (3.19)$$

$$\Rightarrow \boxed{v_*(s) = \max_{a \in A(s)} \sum_{s', r} p(s', r \mid s, a)\left[r + \gamma v_*(s')\right]}$$

(nonlin
SYSTEM) (INTUITIVELY CLEAR. COMPARE WITH (3.14) )

## BELLMAN OPTIMALITY CONDITION FOR $q_*(s, a)$

$$q_*(s, a) \overset{(A)}{=} E_*\left[R_{t+1} + \gamma v_*(S_{t+1}) \mid S_t = s, A_t = a\right]$$

$$= \sum_{s', r} p(s', r \mid s, a)\left[r + \gamma v_*(s')\right]$$

$$\overset{(B)}{=} \sum_{s', r} p(s', r \mid s, a)\left[r + \gamma \max_{a'} q(s', a')\right]$$
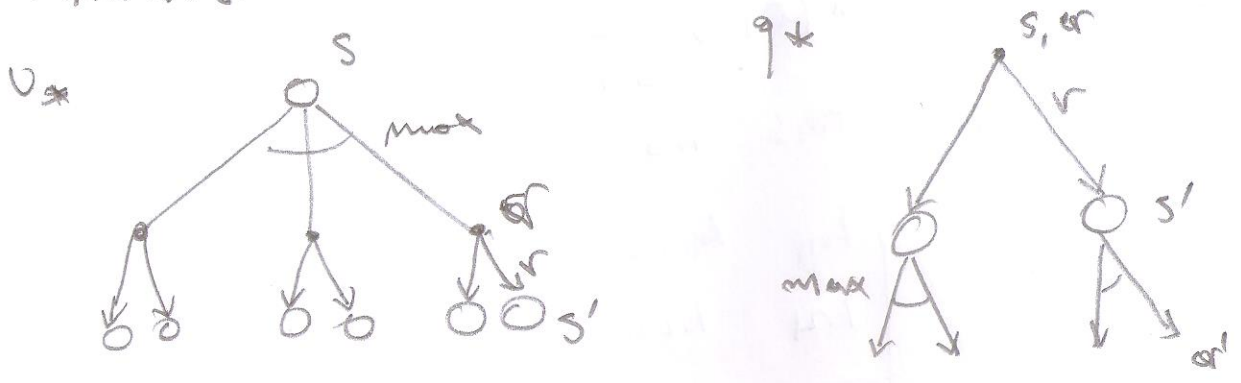
(ALSO INTUITIVELY CLEAR)

(AGAIN,
non SYSTEM)

COMMENTS: 1) IN BOTH CASES, WE HAVE AS MANY
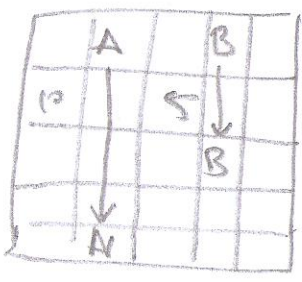EQUATIONS AS WE HAVE UNKNOWNS, BUT NOW SYSTEM
IS NONLINEAR.

2) IF WE SOLVE THE EQUATIONS, (EITHER SET), THE
OPTIMAL POLICY IS TRIVIAL TO FIND.

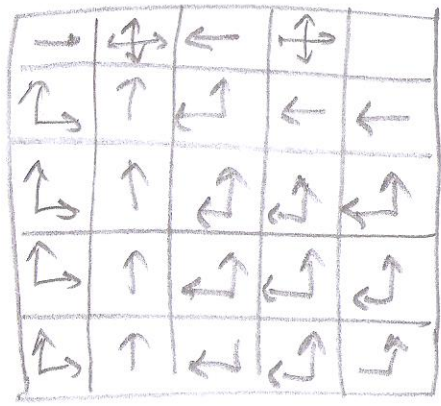3) BOTH SETS CAN BE DESCRIBED IN TERMS OF BACKUP
DIAGRAMS.



EXAMPLE 3.8    GRIDWORLD



| 22 | 24.4 | 22 | 19.4 | 17.5 |
|------|------|------|------|------|
| 19.8 | 22 | 19.8 | 17.8 | 16. |
| 17.8 | 19.8 | 17.8 | 16 | 14.4 |
| 16 | 17.8 | 16.0 | 14.4 | 13 |
| 14.4 | 16 | 14.4 | 13 | 11.7 |

OBSERVE THAT THESE
NUMBERS INDEED
SATISFY
BELLMAN'S OPTIMALITY
   CONDITION

# CHAPTER 4: DYNAMIC PROGRAMMING

## 4.1 POLICY EVALUATION

PROBLEM: GIVEN A POLICY $\pi$, FIND ITS VALUE FUNCTION $v_\pi(s)$ FOR ALL $s \in S$

WE HAVE BELLMAN'S EQUATION

$$v_\pi(s) = \sum_a \pi(a|s) \left[ \sum_{s',r} p(s',r|s,a) \left[ r + \gamma v_\pi(s') \right] \right]$$

WE CAN WRITE THIS AS:

$$\begin{pmatrix} v_\pi(1) \\ v_\pi(2) \\ \vdots \\ v_\pi(m) \end{pmatrix} = \begin{pmatrix} A_{11} & A_{12} & \cdots & A_{1m} \\ A_{21} & A_{22} & \cdots & A_{2m} \\ & & \cdots & \\ A_{m1} & A_{m2} & \cdots & A_{mm} \end{pmatrix} \begin{pmatrix} v_\pi(1) \\ v_\pi(2) \\ \vdots \\ v_\pi(m) \end{pmatrix} + \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{pmatrix}$$

$(\Rightarrow)$

$$v = Av + b \iff Iv - Av = b \overset{(s)}{\iff} (I-A)v = b$$

$(\Rightarrow)$

$$\boxed{v = (I-A)^{-1} b}$$

SO WE CAN FIND THE VALUE FUNCTION BY SOLVING A LINEAR SYSTEM.

THERE IS ANOTHER WAY, WHICH WORKS BETTER AND IS MORE GENERALIZABLE:

## ITERATIVE POLICY EVALUATION ALGORITHM

1) WE SET $V_0(s)$ ARBITRARILY

( BUT IF THERE IS A TERMINAL STATE $S_T$, THEN WE SET $V_0(S_T) = 0$, BECAUSE THIS IS THE RIGHT VALUE, AND THE ALGORITHM WILL NOT UPDATE IT)

2)     SET     $v_{t+1}(s) = \sum_{a} \pi(a|s) \sum_{s',r} p(s',r|s,a)\left[r + \gamma v_t(s')\right]$

__THEOREM__   IF   $\gamma < 1$   OR   EVENTUAL TERMINATION IS GUARANTEED (FOR EPISODIC TASKS) THEN ALGORITHM IS WILL CONVERGE TO SOLUTION
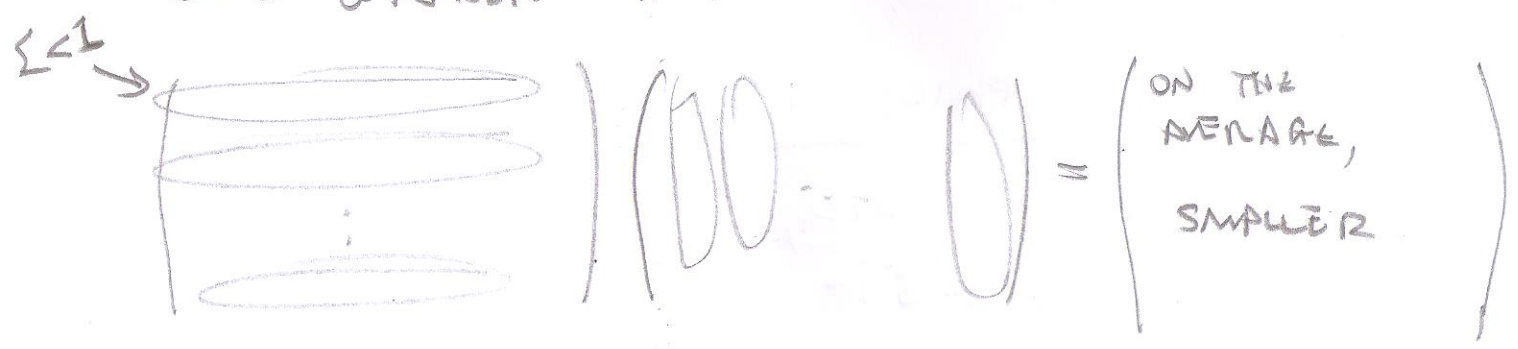
SKETCH OF PROOF FOR $\gamma < 1$ CASE:

$$V_1 = A V_0 + b, \quad V_2 = A(A V_0 + b) + b = A^2 V_0 + Ab + b$$

$$V_3 = A(A V_0 + b) + b = A^3 V_0 + A^2 b + Ab + b \dots$$

$$V_t = (A^{t-1} + A^{t-2} + \dots + A + I) b + A^t V_0$$

HOWEVER,   $\gamma < 1 \Rightarrow$   NORM   $|A| < 1$

INDEED CONSIDER   $A^L$:

$\{ < 1 \rightarrow$



$$\left( \begin{array}{} \\ \\ \\ \end{array} \right) \left( \begin{array}{} 0 & 0 & \dots & 0 \end{array} \right) = \left( \begin{array}{} \text{ON THE} \\ \text{AVERAGE,} \\ \text{SMALLER} \end{array} \right)$$

IT FOLLOWS THAT

$$V_{\infty} = \left( \sum_{t=0}^{\infty} A^t \right) b .$$

HOWEVER, WE KNOW THAT:

$$\sum_{t=0}^{\infty} A^t = (I - A)^{-1} \quad (\text{NEUMANN SERIES})$$

WHICH CAN BE PROVEN SIMILARLY TO

$$\sum_{t=0}^{\infty} r^t = \frac{1}{1-r}$$

WE AN IMPROVE SPEED OF CONVERGENCE IF WE USE VALUES THE MOMENT WE COMPUTE THEM:

## ITERATIVE POLICY EVALUATION (pg. 75)

INPUT: $\pi$, $\theta$ (FOR TERMINATION CONDITION),

$V(s)$ (ARBITRARY, BUT $V(TERMINAL) = 0$)

LOOP:

$\Delta \leftarrow 0$
LOOP FOR EACH $s \in S$:

$v \leftarrow V(s)$

$V(s) \leftarrow \sum\limits_{\alpha} \pi(\alpha|s) \sum\limits_{s',r} p(s',r|s,\alpha)[r + \gamma V(s')]$

$\Delta \leftarrow \max(\Delta, |v - V(s)|)$

UNTIL $\Delta < \theta$

## 4.2 POLICY IMPROVEMENT

GIVEN A SPECIFIC POLICY, WE KNOW THAT:

$$\boxed{\begin{aligned} v_\pi(s) &\doteq E_\pi[G_t | S_t = s] = E_\pi[R_{t+1} + \gamma G_{t+1} | S_t = s] \\ &= \sum_\alpha \pi(\alpha|s) \sum_{s',r} p(s',r|s,\alpha)[r + \gamma v_\pi(s')] \end{aligned}} \quad \text{①}$$

$q_\pi(s,\alpha) \doteq E_\pi[G_t | S_t = s, A_t = \alpha]$

$= E_\pi[R_{t+1} + \gamma G_{t+1} | S_t = s, A_t = \alpha]$

$= E_\pi[R_{t+1} | S_t = s, A_t = \alpha] + \gamma E_\pi[G_{t+1} | S_t = s, A_t = \alpha]$

$$= E[R_{t+1} | S_t = s, A_t = \alpha]$$

$$+ \gamma \sum_{s', r} p(s', r | s, \alpha) \left[ E_\pi \left[ G_{t+1} | S_t = s, A_t = \alpha, S_{t+1} = s' \right] \right]_{R_{t+1} = r}$$

$$= E[R_{t+1} | S_t = s, A_t = \alpha]$$

$$+ \gamma \sum_{s', r} p(s', r | s, \alpha) v_\pi(s') =$$

$$= E[R_{t+1} | S_t = s, A_t = \alpha] + \gamma E\left[ v_\pi(S_{t+1}) | S_t = s, A_t = \alpha \right]$$

$$= E[R_{t+1} + \gamma v_\pi(S_{t+1}) | S_t = s, A_t = \alpha]$$

$$\Rightarrow \boxed{\begin{aligned} q_\pi(s, \alpha) &= E[R_{t+1} + \gamma v_\pi(S_{t+1}) | S_t = s, A_t = \alpha] \\ &= \sum_{s', r} p(s', r | s, \alpha) [r + \gamma v_\pi(s')] \end{aligned}} \quad ②$$

LET US WRITE ① AND ② FOR **DETERMINISTIC** POLICIES

$$v_\pi(s) = \sum_{s', r} p(s', r | s) [r + \gamma v_\pi(s')] \quad ①$$

NEW NOTATION →

$$q_\pi(s, \pi(s)) = \sum_{s', r} p(s', r | s, \pi(s)) [r + \gamma v_\pi(s')] \quad ②'$$

ALSO OBSERVE THAT THEN

$$v_\pi(s) = q_\pi(s, \pi(s))$$

POLICY IMPROVEMENT THEOREM (SPECIAL CASE)

LET $\pi, \pi'$ DETERMINISTIC POLICIES SUCH THAT:

$$\forall s \in S \quad q_\pi(s, \pi'(s)) \geq v_\pi(s) \quad \text{(A)}$$

THEN IT FOLLOWS THAT

$$v_{\pi'}(s) \geq v_\pi(s) \quad \text{(B)}$$

MOREOVER, IF (A) IS STRICT FOR SOME $s$, THEN IT IS ALSO STRICT FOR THIS $s$.

(OBSERVE THAT IT MAKES SENSE!)

(COMPARE MMN)

PROOF:

$$v_\pi(s) \leq q_\pi(s, \pi'(s)) = E\left[R_{t+1} + \gamma v_\pi(S_{t+1}) \mid S_t = s, A_t = \pi'(s)\right]$$

$$= \sum_{s', r} P_{S_{t+1}, R_{t+1} \mid S_t, A_t}(s', r \mid s, \pi'(s)) \left[r + \gamma v_\pi(s')\right]$$

$$= E_{\pi'}\left[R_{t+1} + \gamma v_\pi(S_{t+1}) \mid S_t = s\right]$$

$$\leq E_{\pi'}\left[R_{t+1} + \gamma q_\pi(S_{t+1}, \pi'(S_{t+1})) \mid S_t = s\right]$$

$$= E_{\pi'}\left[R_{t+1} + \gamma E\left[R_{t+2} + \gamma v_\pi(S_{t+2}) \mid S_{t+1}, A_{t+1} = \pi'(S_{t+1})\right] \mid S_t = s\right]$$

$$= E_{\pi'}\left[E\left[R_{t+1} + \gamma R_{t+2} + \gamma^2 v_\pi(S_{t+2}) \mid S_{t+1}, A_{t+1} = \pi'(S_{t+1}), S_t = s\right]\right]$$

(NOTE:

$$E\left[X + E[Y \mid Z]\right] = E\left[E[E[X \mid Z] + E[Y \mid Z]]\right] = E\left[E[X + Y \mid Z]\right] \; )$$