

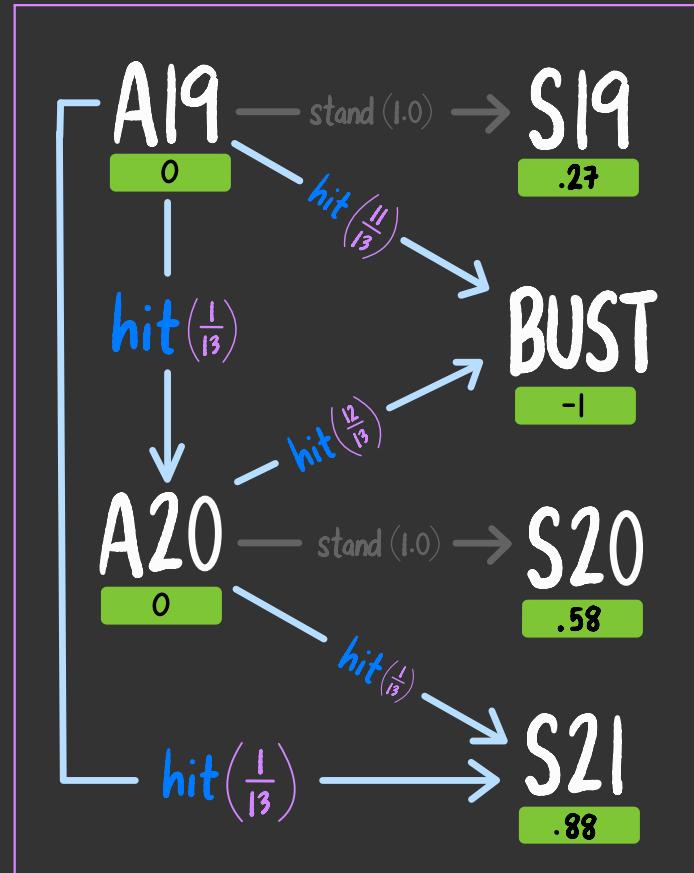
value
iteration

CSCI
373

last time, we defined the expected utility of being in state q , given policy π

$$U^\pi(q) = \sum_{\text{path}} U(q \xrightarrow{\pi(q)} \text{path}) P(q \xrightarrow{\pi(q)} \text{path})$$

$$U^\pi(A19) = -.84 \rightarrow$$

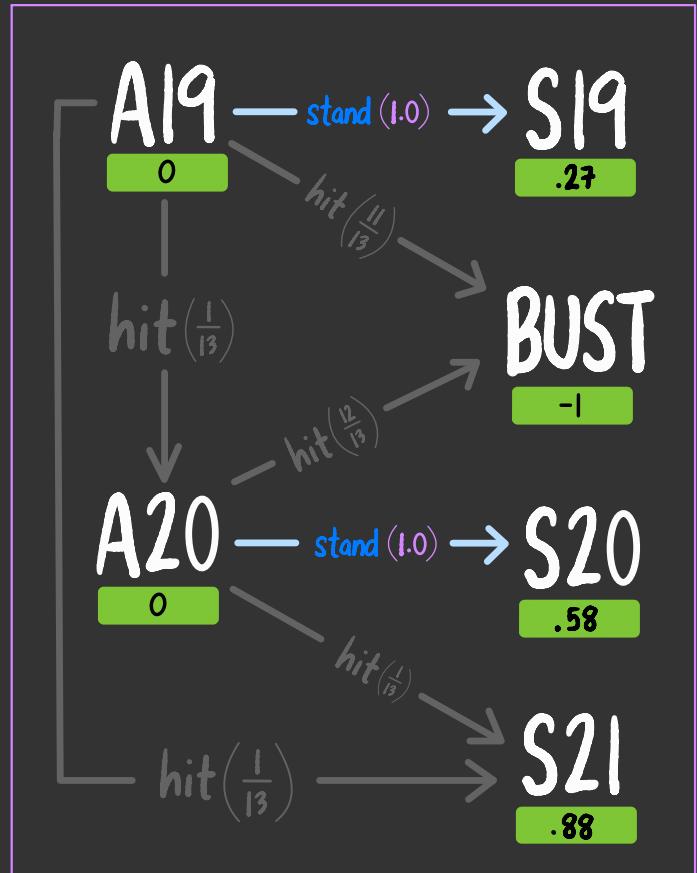


markov decision process

in order to make decisions,
we'd like to know the
policy that maximizes
our expected utility

$$\pi^* = \arg \max_{\pi} U^\pi(q)$$

$$\pi^* = \{A19 \mapsto \text{stand}, A20 \mapsto \text{stand}\} \rightarrow$$



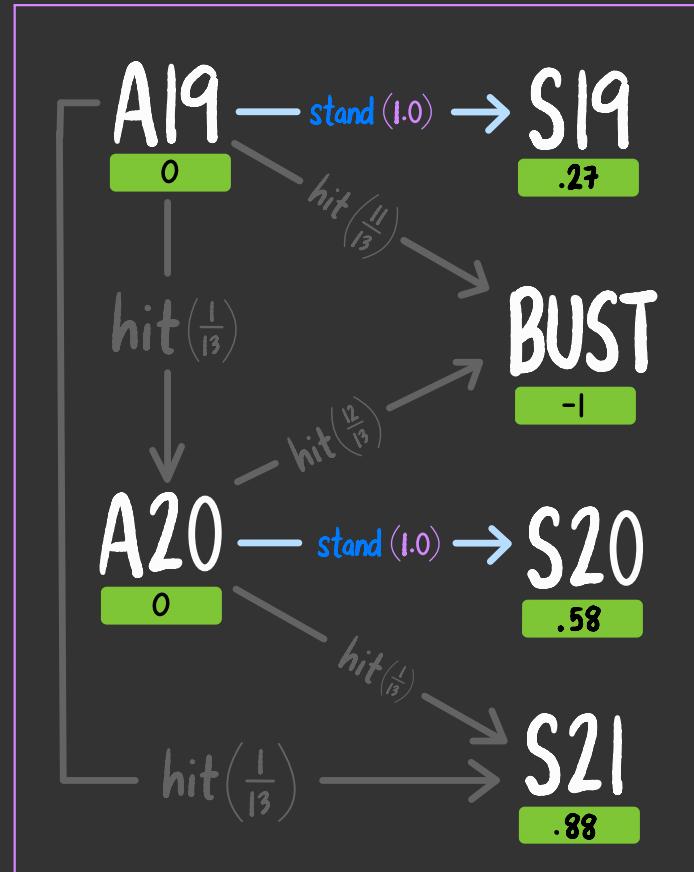
markov decision process

but first we'll focus on
how to compute the

maximum
expected utility

$$U(q) = \max_{\pi} U^{\pi}(q)$$

$$U(A19) = .27 \rightarrow$$



markov decision process

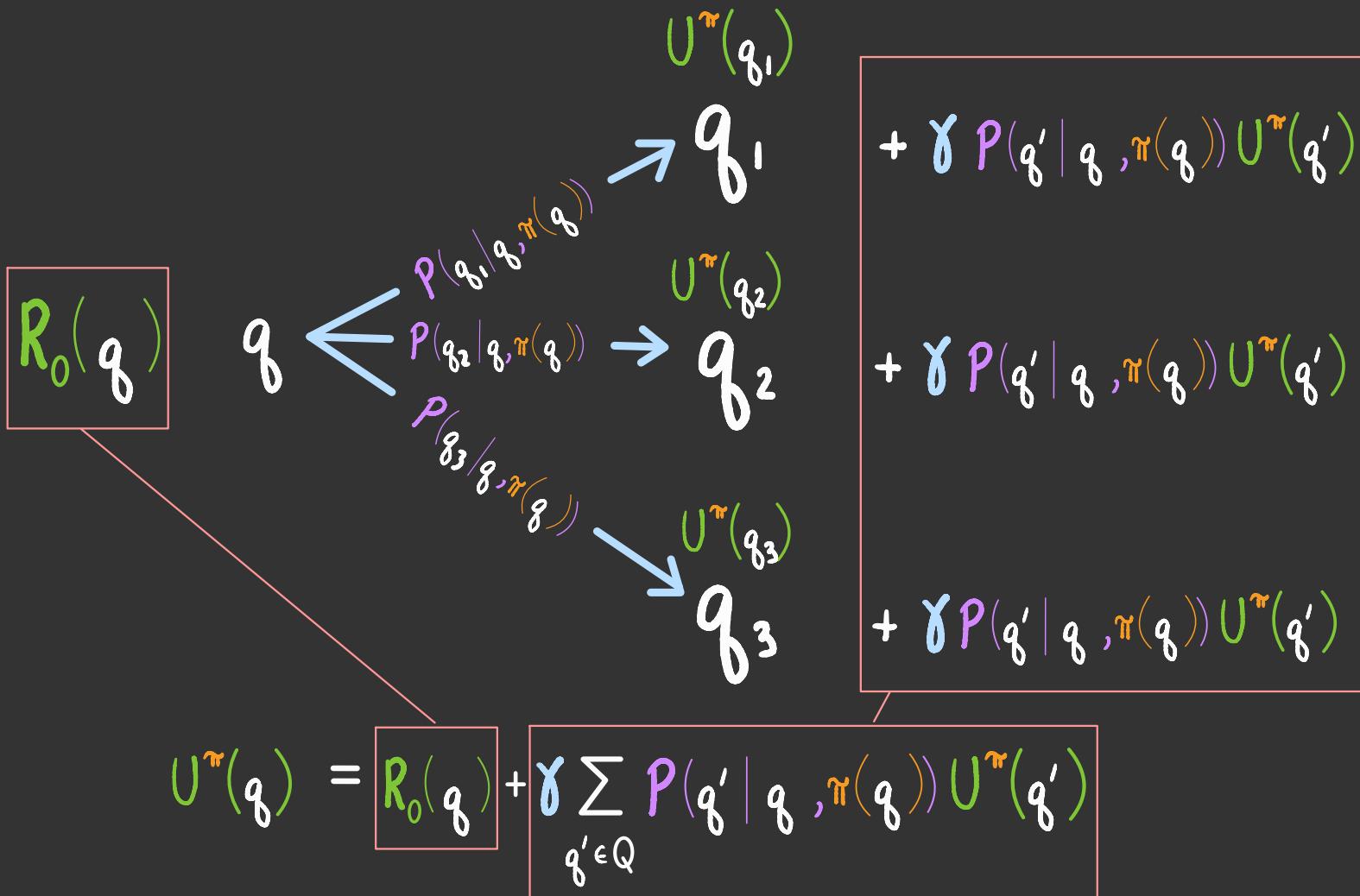
want: maximum expected utility $U(q) = \max_{\pi} U^{\pi}(q)$

last time, we observed:

$$\frac{\text{the expected utility of policy } \pi \text{ in state } q}{U^{\pi}(q)} = \frac{\text{discounting factor}}{R_0(q) + \gamma \sum_{q' \in Q} P(q' | q, \pi(q)) U^{\pi}(q')}$$

the immediate reward of being in state q

the probability of getting from state q to state q' using the action advised by the policy



$$U(q) = \max_{\pi} U^\pi(q)$$



want: maximum
expected utility $U(q) = \max_{\pi} U^\pi(q)$

$$U(q) = \max_{\pi} U^\pi(q)$$

$$= \max_{\pi} R_0(q) + \gamma \sum_{q' \in Q} P(q' | q, \pi(q)) U^\pi(q')$$

last time, we observed:

$$U^\pi(q) = R_0(q) + \gamma \sum_{q' \in Q} P(q' | q, \pi(q)) U^\pi(q')$$

$$\begin{aligned}
 U(q) &= \max_{\pi} U^\pi(q) \\
 &= \max_{\pi} R_0(q) + \gamma \sum_{q' \in Q} P(q'|q, \pi(q)) U^\pi(q') \\
 &= R_0(q) + \gamma \max_{\pi} \sum_{q' \in Q} P(q'|q, \pi(q)) U^\pi(q')
 \end{aligned}$$


 always nonnegative

$$U(q) = \max_{\pi} U^{\pi}(q)$$

$$= \max_{\pi} R_0(q) + \gamma \sum_{q' \in Q} P(q'|q, \pi(q)) U^{\pi}(q')$$

$$= R_0(q) + \gamma \max_{\pi} \sum_{q' \in Q} P(q'|q, \pi(q)) U^{\pi}(q')$$

$$= R_0(q) + \gamma \max_{\sigma} \sum_{q' \in Q} P(q'|q, \sigma) U(q')$$

markov
property

we obtain something called a
bellman equation

the maximum expected
utility of state q

discounting
factor

the maximum expected
utility of state q'

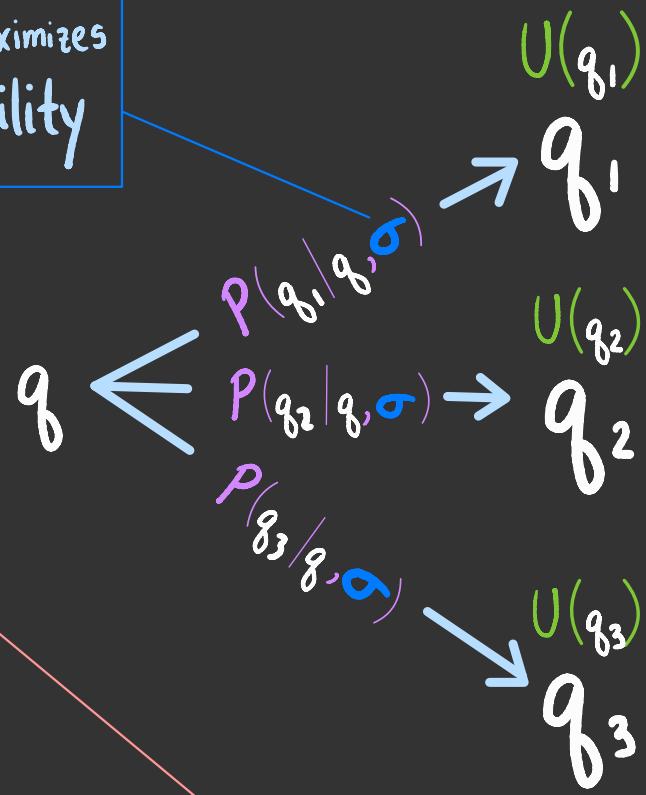
$$\overline{U(q)} = \overline{R(q) + \gamma \max_{\sigma} \sum_{q' \in Q} P(q' | q, \sigma) U(q')}$$

the immediate reward
of being in state q

the probability of getting
from state q to state q'
if we do action σ

the action that maximizes
expected utility

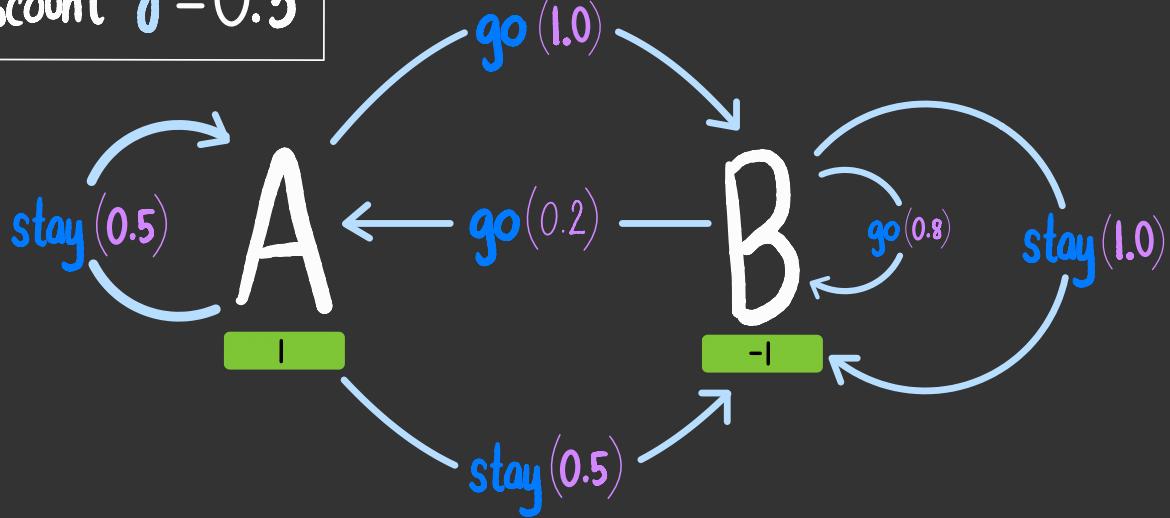
$$R_0(q)$$



$$U(q) = R_0(q) + \max_{\sigma} \gamma \sum_{q' \in Q} P(q'|q, \sigma) U(q')$$

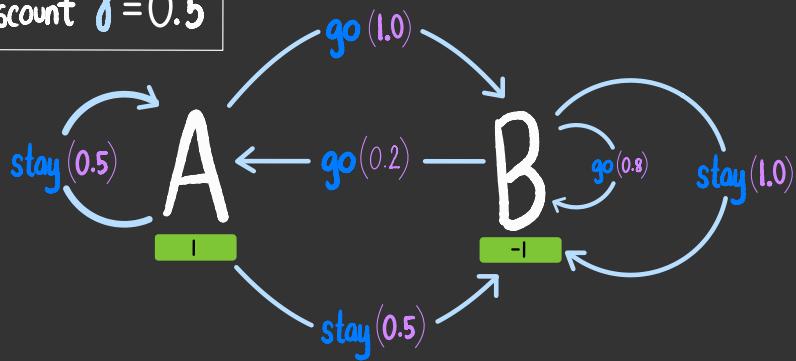
$$\begin{aligned} &+ \gamma P(q_1|q, \sigma) U(q_1) \\ &+ \gamma P(q_2|q, \sigma) U(q_2) \\ &+ \gamma P(q_3|q, \sigma) U(q_3) \end{aligned}$$

discount $\gamma = 0.5$



Consider the bellman equations for
the above markov decision process

discount $\gamma = 0.5$



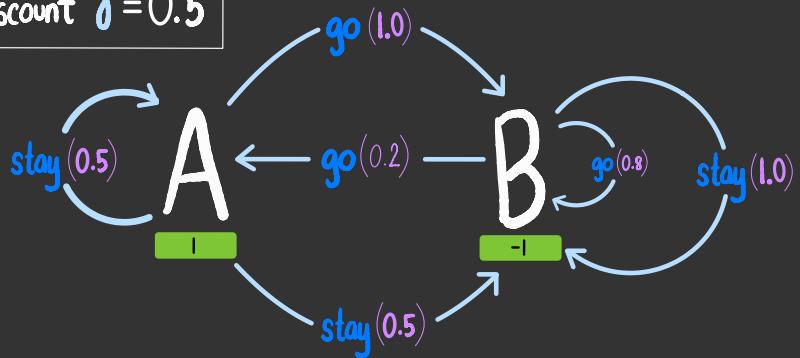
bellman equation:

$$U(q) = R(q) + \gamma \max_{\sigma} \sum_{q' \in Q} P(q'|q, \sigma) U(q')$$

$$U(A) = ?$$

$$U(B) = ?$$

discount $\gamma = 0.5$



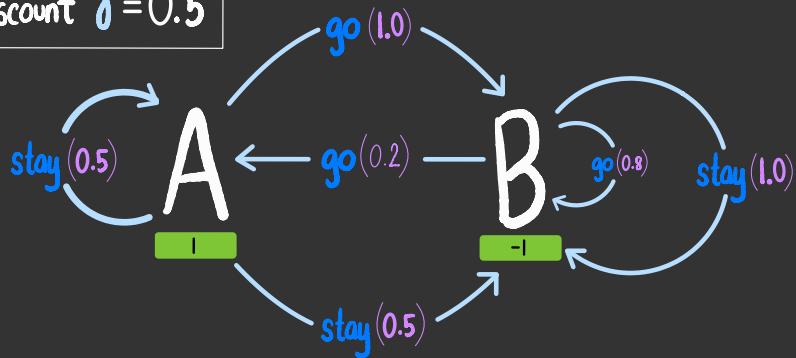
bellman equation:

$$U(q) = R(q) + \gamma \max_{\sigma} \sum_{q' \in Q} P(q'|q, \sigma) U(q')$$

$$U(A) = R(A) + \gamma \max \left\{ \begin{array}{l} P(A|A, \text{stay}) U(A) + P(B|A, \text{stay}) U(B) \\ P(A|A, \text{go}) U(A) + P(B|A, \text{go}) U(B) \end{array} \right\}$$

$$U(B) = R(B) + \gamma \max \left\{ \begin{array}{l} P(A|B, \text{stay}) U(A) + P(B|B, \text{stay}) U(B) \\ P(A|B, \text{go}) U(A) + P(B|B, \text{go}) U(B) \end{array} \right\}$$

discount $\gamma = 0.5$



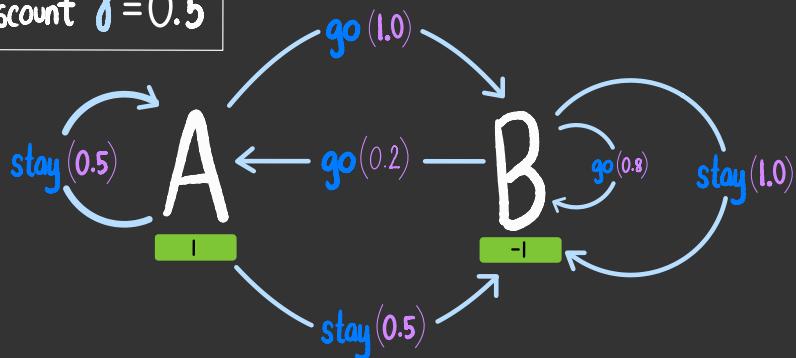
bellman equation:

$$U(q) = R(q) + \gamma \max_{\sigma} \sum_{q' \in Q} P(q'|q, \sigma) U(q')$$

$$U(A) = 1 + 0.5 \max \left\{ \begin{array}{l} 0.5 U(A) + 0.5 U(B) \\ 0.0 U(A) + 1.0 U(B) \end{array} \right\}$$

$$U(B) = -1 + 0.5 \max \left\{ \begin{array}{l} 0.0 U(A) + 1.0 U(B) \\ 0.2 U(A) + 0.8 U(B) \end{array} \right\}$$

discount $\gamma = 0.5$



bellman equation:

$$U(q) = R(q) + \gamma \max_{\sigma} \sum_{q' \in Q} P(q'|q, \sigma) U(q')$$

$$U(A) = 1 + 0.5 \max \left\{ U(B), 0.5 U(A) + 0.5 U(B) \right\}$$

$$U(B) = -1 + 0.5 \max \left\{ U(B), 0.2 U(A) + 0.8 U(B) \right\}$$

this is a system of
two equations with two unknowns

$$U(A) = | + 0.5 \max \{ U(B), 0.5 U(A) + 0.5 U(B) \}$$

$$U(B) = -| + 0.5 \max \{ U(B), 0.2 U(A) + 0.8 U(B) \}$$

this is a system of
two equations with two unknowns
but they aren't linear equations

$$U(A) = | + 0.5 \max \{ U(B), 0.5 U(A) + 0.5 U(B) \}$$

$$U(B) = -| + 0.5 \max \{ U(B), 0.2 U(A) + 0.8 U(B) \}$$

an iterative strategy

- guess values for the unknowns
 - U_o^A is a guess for $U(A)$
 - U_o^B is a guess for $U(B)$
- compute new guesses using the equations
- repeat

$$U(A) = | + 0.5 \max \left\{ U_o^B, 0.5 U_o^A + 0.5 U_o^B \right\}$$

$$U(B) = -| + 0.5 \max \left\{ U_o^B, 0.2 U_o^A + 0.8 U_o^B \right\}$$

an iterative strategy

- guess values for the unknowns
 - U_0^A is a guess for $U(A)$
 - U_0^B is a guess for $U(B)$
- compute new guesses using the equations
- repeat

$$U(A) = | + 0.5 \max \left\{ U_0^B, 0.5 U_0^A + 0.5 U_0^B \right\}$$

$$U(B) = -| + 0.5 \max \left\{ U_0^B, 0.2 U_0^A + 0.8 U_0^B \right\}$$

an iterative strategy

- guess values for the unknowns
 - U_0^A is a guess for $U(A)$
 - U_0^B is a guess for $U(B)$
- compute new guesses using the equations
- repeat

$$U_i^A \quad \cancel{U(A)} \leftarrow | + 0.5 \max \left\{ U_0^B, 0.5 U_0^A + 0.5 U_0^B \right\}$$

$$U_i^B \quad \cancel{U(B)} \leftarrow -| + 0.5 \max \left\{ U_0^B, 0.2 U_0^A + 0.8 U_0^B \right\}$$

an iterative strategy

- guess values for the unknowns
 - U_0^A is a guess for $U(A)$
 - U_0^B is a guess for $U(B)$
- compute new guesses using the equations
- repeat

$$U_1^A \leftarrow | + 0.5 \max \left\{ U_0^B, 0.5 U_0^A + 0.5 U_0^B \right\}$$

$$U_1^B \leftarrow -| + 0.5 \max \left\{ U_0^B, 0.2 U_0^A + 0.8 U_0^B \right\}$$

an iterative strategy

- guess values for the unknowns
 - U_o^A is a guess for $U(A)$
 - U_o^B is a guess for $U(B)$
- compute new guesses using the equations
- repeat

$$U_{t+1}^A \leftarrow | + 0.5 \max \left\{ U_t^B, 0.5 U_t^A + 0.5 U_t^B \right\}$$

$$U_{t+1}^B \leftarrow -| + 0.5 \max \left\{ U_t^B, 0.2 U_t^A + 0.8 U_t^B \right\}$$

to the
laptop!

DISCOUNT =		0.5			
	t	U(A)	U(B)		
	0	1.2500	-1.1300		
	1	1.6250	-1.3270		
	2	1.8125	-1.3683		
	3	1.9063	-1.3661		
	4	1.9531	-1.3558		
	5	1.9766	-1.3470		
	6	1.9883	-1.3411		
	7	1.9941	-1.3376		
	8	1.9971	-1.3356		
	9	1.9985	-1.3345		
	10	1.9993	-1.3340		
	11	1.9996	-1.3337		
	12	1.9998	-1.3335		
	13	1.9999	-1.3334		
	14	2.0000	-1.3334		
	15	2.0000	-1.3334		

value iteration

- guess values for the unknowns

for each state q , U_0^q is a guess for $U(q)$

- for $t = 1$ to T :

compute new guesses using the bellman equations, i.e.

$$\text{for each state } q, U_{t+1}^q \leftarrow R(q) + \gamma \max_{\sigma} \sum_{q' \in Q} P(q'|q, \sigma) U_t^{q'}$$

Claim: if $\gamma < 1$, then value iteration converges to the maximum expected utilities, i.e. $\lim_{t \rightarrow \infty} U_t^q = U(q)$

Claim: if $\gamma < 1$, then value iteration converges to the maximum expected utilities, i.e. $\lim_{t \rightarrow \infty} U_t^\gamma = U(q)$

- assume a solution to the bellman equations exists, i.e there exists a function $U: Q \rightarrow \mathbb{R}$ such that

$$U(q) = R(q) + \gamma \max_{\sigma} \sum_{q' \in Q} P(q'|q, \sigma) U(q') \quad \text{for every state } q$$

this is the
solution to the
bellman equations

→

$U(q_0)$
$U(q_n)$

Claim: if $\gamma < 1$, then value iteration converges to the maximum expected utilities, i.e. $\lim_{t \rightarrow \infty} U_t^\gamma = U(q)$

- measure how far our guesses at iteration t are from the solution using the following distance function:

$$\text{dist} \left(\begin{bmatrix} U_t^{q_0} \\ \vdots \\ U_t^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right) = \max_{q \in Q} |U_t^q - U(q)|$$

guesses at
time t

true
utilities

$$\text{dist} \left(\begin{bmatrix} U_t^{q_0} \\ \vdots \\ U_t^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right) = \max_{q \in Q} |U_t^q - U(q)|$$

e.g.

$$\text{dist} \left(\begin{bmatrix} 2 \\ 4 \end{bmatrix}, \begin{bmatrix} 5 \\ 3 \end{bmatrix} \right) = ?$$

guesses at
 time t true
 utilities

$$\text{dist} \left(\begin{bmatrix} U_t^{q_0} \\ \vdots \\ U_t^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right) = \max_{q \in Q} |U_t^q - U(q)|$$

e.g.

$$\text{dist} \left(\begin{bmatrix} 2 \\ 4 \end{bmatrix}, \begin{bmatrix} 5 \\ 3 \end{bmatrix} \right) = \max \{ |2-5|, |4-3| \} = 3$$

guesses at
 time t true
 utilities

$$\text{dist} \left(\begin{bmatrix} U_t^{q_0} \\ \vdots \\ U_t^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right) = \max \left[\begin{array}{c} |U_t^{q_0} - U(q_0)| \\ \vdots \\ |U_t^{q_n} - U(q_n)| \end{array} \right]$$

e.g.

equals zero iff our guesses are all correct

$$\text{dist} \left(\begin{bmatrix} 2 \\ 4 \end{bmatrix}, \begin{bmatrix} 2 \\ 4 \end{bmatrix} \right) = \max \{ |2-2|, |4-4| \} = 0$$

Claim: if $\gamma < 1$, then value iteration converges to the maximum expected utilities, i.e. $\lim_{t \rightarrow \infty} U_t^\gamma = U(q_\gamma)$

Suppose we can show that every step of value iteration brings our guesses closer to the maximum expected utilities

$$\text{dist} \left(\begin{bmatrix} U_{t+1}^{q_0} \\ \vdots \\ U_{t+1}^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right) < K \text{dist} \left(\begin{bmatrix} U_t^{q_0} \\ \vdots \\ U_t^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right)$$

for some fraction $K \in [0, 1]$

suppose for some fraction $K \in [0, 1]$:

$$\text{dist} \left(\begin{bmatrix} U_t^{q_0} \\ \vdots \\ U_t^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right) \leq K \text{dist} \left(\begin{bmatrix} U_{t-1}^{q_0} \\ \vdots \\ U_{t-1}^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right)$$

suppose for some fraction $K \in [0, 1]$:

$$\text{dist} \left(\begin{bmatrix} u_t^{q_0} \\ \vdots \\ u_t^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right) \leq K \text{dist} \left(\begin{bmatrix} u_{t-1}^{q_0} \\ \vdots \\ u_{t-1}^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right)$$

then:

$$\text{dist} \left(\begin{bmatrix} u_t^{q_0} \\ \vdots \\ u_t^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right) \leq K^t \text{dist} \left(\begin{bmatrix} u_0^{q_0} \\ \vdots \\ u_0^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right)$$

$$\text{dist} \left(\begin{bmatrix} u_t^{q_0} \\ \vdots \\ u_t^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right) \leq K \text{dist} \left(\begin{bmatrix} u_{t-1}^{q_0} \\ \vdots \\ u_{t-1}^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right) \leq K \cdot K \text{dist} \left(\begin{bmatrix} u_{t-2}^{q_0} \\ \vdots \\ u_{t-2}^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right) \leq \dots$$

suppose for some fraction $K \in [0, 1)$:

$$\text{dist} \left(\begin{bmatrix} u_t^{q_0} \\ \vdots \\ u_t^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right) \leq K \text{dist} \left(\begin{bmatrix} u_{t-1}^{q_0} \\ \vdots \\ u_{t-1}^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right)$$

then:

$$\text{dist} \left(\begin{bmatrix} u_t^{q_0} \\ \vdots \\ u_t^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right) \leq K^t \text{dist} \left(\begin{bmatrix} u_0^{q_0} \\ \vdots \\ u_0^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right)$$

$$\lim_{t \rightarrow \infty} \text{dist} \left(\begin{bmatrix} u_t^{q_0} \\ \vdots \\ u_t^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right) \leq \lim_{t \rightarrow \infty} K^t \text{dist} \left(\begin{bmatrix} u_0^{q_0} \\ \vdots \\ u_0^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right)$$

suppose for some fraction $K \in [0, 1]$:

$$\text{dist} \left(\begin{bmatrix} u_t^{q_0} \\ \vdots \\ u_t^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right) \leq K \text{dist} \left(\begin{bmatrix} u_{t-1}^{q_0} \\ \vdots \\ u_{t-1}^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right)$$

then:

$$\text{dist} \left(\begin{bmatrix} u_t^{q_0} \\ \vdots \\ u_t^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right) \leq K^t \text{dist} \left(\begin{bmatrix} u_0^{q_0} \\ \vdots \\ u_0^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right)$$

$$\lim_{t \rightarrow \infty} \text{dist} \left(\begin{bmatrix} u_t^{q_0} \\ \vdots \\ u_t^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right) \leq \text{dist} \left(\begin{bmatrix} u_0^{q_0} \\ \vdots \\ u_0^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right) \lim_{t \rightarrow \infty} K^t$$

doesn't depend on t

suppose for some fraction $K \in [0, 1)$:

$$\text{dist} \left(\begin{bmatrix} u_t^{q_0} \\ \vdots \\ u_t^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right) \leq K \text{dist} \left(\begin{bmatrix} u_{t-1}^{q_0} \\ \vdots \\ u_{t-1}^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right)$$

then:

$$\text{dist} \left(\begin{bmatrix} u_t^{q_0} \\ \vdots \\ u_t^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right) \leq K^t \text{dist} \left(\begin{bmatrix} u_0^{q_0} \\ \vdots \\ u_0^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right)$$

$$\lim_{t \rightarrow \infty} \text{dist} \left(\begin{bmatrix} u_t^{q_0} \\ \vdots \\ u_t^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right) \leq \text{dist} \left(\begin{bmatrix} u_0^{q_0} \\ \vdots \\ u_0^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right) \lim_{t \rightarrow \infty} K^t$$

suppose for some fraction $K \in [0, 1]$:

$$\text{dist} \left(\begin{bmatrix} u_t^{q_0} \\ \vdots \\ u_t^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right) \leq K \text{dist} \left(\begin{bmatrix} u_{t-1}^{q_0} \\ \vdots \\ u_{t-1}^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right)$$

then:

$$\text{dist} \left(\begin{bmatrix} u_t^{q_0} \\ \vdots \\ u_t^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right) \leq K^t \text{dist} \left(\begin{bmatrix} u_0^{q_0} \\ \vdots \\ u_0^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right)$$

$$\lim_{t \rightarrow \infty} \text{dist} \left(\begin{bmatrix} u_t^{q_0} \\ \vdots \\ u_t^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right) \leq 0$$

suppose for some fraction $K \in [0, 1]$:

$$\text{dist} \left(\begin{bmatrix} u_t^{q_0} \\ \vdots \\ u_t^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right) \leq K \text{dist} \left(\begin{bmatrix} u_{t-1}^{q_0} \\ \vdots \\ u_{t-1}^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right)$$

then:

$$\text{dist} \left(\begin{bmatrix} u_t^{q_0} \\ \vdots \\ u_t^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right) \leq K^t \text{dist} \left(\begin{bmatrix} u_0^{q_0} \\ \vdots \\ u_0^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right)$$

$$\lim_{t \rightarrow \infty} \text{dist} \left(\begin{bmatrix} u_t^{q_0} \\ \vdots \\ u_t^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right) \leq 0$$

always nonnegative $\left(\max_{q \in Q} |U_t^q - U(q)| \right)$

suppose for some fraction $K \in [0, 1]$:

$$\text{dist} \left(\begin{bmatrix} u_t^{q_0} \\ \vdots \\ u_t^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right) \leq K \text{dist} \left(\begin{bmatrix} u_{t-1}^{q_0} \\ \vdots \\ u_{t-1}^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right)$$

then:

$$\text{dist} \left(\begin{bmatrix} u_t^{q_0} \\ \vdots \\ u_t^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right) \leq K^t \text{dist} \left(\begin{bmatrix} u_0^{q_0} \\ \vdots \\ u_0^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right)$$

$$\lim_{t \rightarrow \infty} \text{dist} \left(\begin{bmatrix} u_t^{q_0} \\ \vdots \\ u_t^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right) = 0$$

if $\text{dist}\left(\begin{bmatrix} \mathbf{u}_t^{q_0} \\ \vdots \\ \mathbf{u}_t^{q_n} \end{bmatrix}, \begin{bmatrix} \mathbb{U}(q_0) \\ \vdots \\ \mathbb{U}(q_n) \end{bmatrix}\right) \leq K \text{dist}\left(\begin{bmatrix} \mathbf{u}_{t-1}^{q_0} \\ \vdots \\ \mathbf{u}_{t-1}^{q_n} \end{bmatrix}, \begin{bmatrix} \mathbb{U}(q_0) \\ \vdots \\ \mathbb{U}(q_n) \end{bmatrix}\right)$, then: $\lim_{t \rightarrow \infty} \text{dist}\left(\begin{bmatrix} \mathbf{u}_t^{q_0} \\ \vdots \\ \mathbf{u}_t^{q_n} \end{bmatrix}, \begin{bmatrix} \mathbb{U}(q_0) \\ \vdots \\ \mathbb{U}(q_n) \end{bmatrix}\right) = 0$

for some fraction $K \in [0, 1)$

if we can show that every step of value iteration brings our guesses closer to the maximum expected utilities

then value iteration converges to the maximum expected utilities

$$\text{if } \text{dist} \left(\begin{bmatrix} u_t^{q_0} \\ \vdots \\ u_t^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right) \leq K \text{dist} \left(\begin{bmatrix} u_{t-1}^{q_0} \\ \vdots \\ u_{t-1}^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right),$$

for some fraction $K \in [0, 1)$

$$\text{then: } \lim_{t \rightarrow \infty} \text{dist} \left(\begin{bmatrix} u_t^{q_0} \\ \vdots \\ u_t^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right) = 0$$



if we can show that every step of value iteration brings our guesses closer to the maximum expected utilities

so
let's
show
this

then value iteration converges to the maximum expected utilities

want to
Show :

$$\text{dist} \left(\begin{bmatrix} U_t^{q_0} \\ \vdots \\ U_t^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right) \leq K \text{dist} \left(\begin{bmatrix} U_{t-1}^{q_0} \\ \vdots \\ U_{t-1}^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right) \quad \text{for some fraction } K \in [0, 1)$$


$$\text{dist} \left(\begin{bmatrix} U_t^{q_0} \\ \vdots \\ U_t^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right) = \max_{q \in Q} |U_t^q - U(q)|$$

want to
Show : $\max_{q \in Q} |U_{t+1}^q - U(q)| \leq K \max_{q \in Q} |U_t^q - U(q)|$ for some fraction $K \in [0, 1]$

$$|U_{t+1}^q - U(q)| =$$

want to
Show : $\max_{q \in Q} |U_{t+1}^q - U(q)| \leq K \max_{q \in Q} |U_t^q - U(q)|$ for some fraction $K \in [0, 1]$

$$|U_{t+1}^q - U(q)| =$$

$$U_{t+1}^q = R(q) + \gamma \max_{\sigma} \sum_{q' \in Q} P(q'|q, \sigma) U_t^{q'}$$

$$U(q) = R(q) + \gamma \max_{\sigma} \sum_{q' \in Q} P(q'|q, \sigma) U(q')$$

new guess at iteration $t+1$

bellman equation

want to show : $\max_{q \in Q} |U_{t+1}^q - U(q)| \leq K \max_{q \in Q} |U_t^q - U(q)|$ for some fraction $K \in [0, 1)$

$$|U_{t+1}^q - U(q)| = \left| \left(R(q) + \gamma \max_{\sigma} \sum_{q' \in Q} P(q'|q, \sigma) U_t^{q'} \right) - \left(R(q) + \gamma \max_{\sigma} \sum_{q' \in Q} P(q'|q, \sigma) U(q') \right) \right|$$

want to
Show : $\max_{q \in Q} |U_{t+1}^q - U(q)| \leq K \max_{q \in Q} |U_t^q - U(q)|$ for some fraction $K \in [0, 1]$

$$\begin{aligned}
 |U_{t+1}^q - U(q)| &= \left| \left(\cancel{R(q)} + \gamma \max_{\sigma} \sum_{q' \in Q} P(q'|q, \sigma) U_t^{q'} \right) - \left(\cancel{R(q)} + \gamma \max_{\sigma} \sum_{q' \in Q} P(q'|q, \sigma) U(q') \right) \right| \\
 &= \left| \gamma \max_{\sigma} \sum_{q' \in Q} P(q'|q, \sigma) U_t^{q'} - \gamma \max_{\sigma} \sum_{q' \in Q} P(q'|q, \sigma) U(q') \right|
 \end{aligned}$$

want to
Show : $\max_{q \in Q} |U_{t+1}^q - U(q)| \leq K \max_{q \in Q} |U_t^q - U(q)|$ for some fraction $K \in [0, 1)$

$$\begin{aligned}
 |U_{t+1}^q - U(q)| &= \left| \left(\cancel{R(q)} + \gamma \max_{\sigma} \sum_{q' \in Q} P(q'|q, \sigma) U_t^{q'} \right) - \left(\cancel{R(q)} + \gamma \max_{\sigma} \sum_{q' \in Q} P(q'|q, \sigma) U(q') \right) \right| \\
 &= \left| \gamma \max_{\sigma} \sum_{q' \in Q} P(q'|q, \sigma) U_t^{q'} - \gamma \max_{\sigma} \sum_{q' \in Q} P(q'|q, \sigma) U(q') \right| \\
 &= \gamma \left| \max_{\sigma} \sum_{q' \in Q} P(q'|q, \sigma) U_t^{q'} - \max_{\sigma} \sum_{q' \in Q} P(q'|q, \sigma) U(q') \right|
 \end{aligned}$$


 always nonnegative

want to show : $\max_{q \in Q} |U_{t+1}^q - U(q)| \leq K \max_{q \in Q} |U_t^q - U(q)|$ for some fraction $K \in [0, 1]$

$$|U_{t+1}^q - U(q)| = \left| \max_{\sigma} \sum_{q' \in Q} P(q'|q, \sigma) U_t^{q'} - \max_{\sigma} \sum_{q' \in Q} P(q'|q, \sigma) U(q') \right|$$

want to show : $\max_{q \in Q} |U_{t+1}^q - U(q)| \leq K \max_{q \in Q} |U_t^q - U(q)|$ for some fraction $K \in [0, 1]$

$$|U_{t+1}^q - U(q)| = \frac{\left| \max_{\sigma} \sum_{q' \in Q} P(q'|q, \sigma) U_t^{q'} - \max_{\sigma} \sum_{q' \in Q} P(q'|q, \sigma) U(q') \right|}{f(\sigma)} g(\sigma)$$

want to show : $\max_{q \in Q} |U_{t+1}^q - U(q)| \leq K \max_{q \in Q} |U_t^q - U(q)|$ for some fraction $K \in [0, 1)$

$$|U_{t+1}^q - U(q)| = \gamma \left| \max_{\sigma} f(\sigma) - \max_{\sigma} g(\sigma) \right|$$

$$\begin{aligned} f(\sigma) &= \sum_{q' \in Q} P(q'|q, \sigma) U_t^{q'} \\ g(\sigma) &= \sum_{q' \in Q} P(q'|q, \sigma) U(q') \end{aligned}$$

want to
Show : $\max_{q \in Q} |U_{t+1}^q - U(q)| \leq K \max_{q \in Q} |U_t^q - U(q)|$ for some fraction $K \in [0, 1]$

$$|U_{t+1}^q - U(q)| = \gamma \left| \max_{\sigma} f(\sigma) - \max_{\sigma} g(\sigma) \right|$$

$$f(\sigma) = \sum_{q' \in Q} P(q'|q, \sigma) U_t^{q'} \\ g(\sigma) = \sum_{q' \in Q} P(q'|q, \sigma) U(q')$$

$$\left| \max_{\sigma} f(\sigma) - \max_{\sigma} g(\sigma) \right| \stackrel{?}{=} \max_{\sigma} |f(\sigma) - g(\sigma)|$$

how are these related?

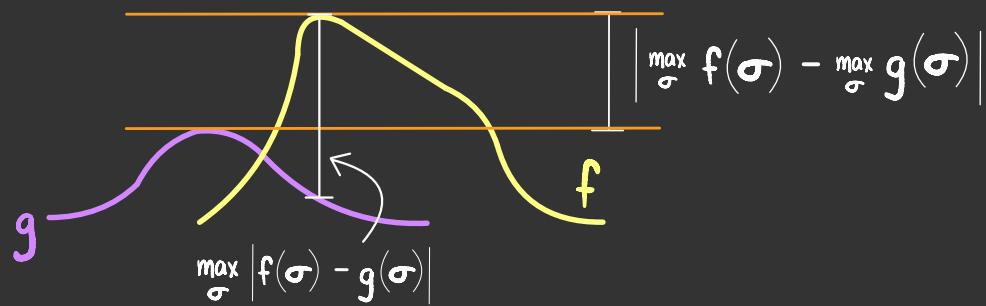
want to
Show : $\max_{q \in Q} |U_{t+1}^q - U(q)| \leq K \max_{q \in Q} |U_t^q - U(q)|$ for some fraction $K \in [0, 1)$

$$|U_{t+1}^q - U(q)| = \gamma \left| \max_{\sigma} f(\sigma) - \max_{\sigma} g(\sigma) \right|$$

$$f(\sigma) = \sum_{q' \in Q} P(q'|q, \sigma) U_t^{q'} \\ g(\sigma) = \sum_{q' \in Q} P(q'|q, \sigma) U(q')$$

$$\left| \max_{\sigma} f(\sigma) - \max_{\sigma} g(\sigma) \right|$$

$$\leq \max_{\sigma} |f(\sigma) - g(\sigma)|$$



want to show : $\max_{q \in Q} |U_{t+1}^q - U(q)| \leq K \max_{q \in Q} |U_t^q - U(q)|$ for some fraction $K \in [0, 1)$

$$|U_{t+1}^q - U(q)| = \gamma \left| \max_{\sigma} f(\sigma) - \max_{\sigma} g(\sigma) \right|$$

$$\leq \gamma \max_{\sigma} \left| f(\sigma) - g(\sigma) \right|$$

$$f(\sigma) = \sum_{q' \in Q} P(q'|q, \sigma) U_t^{q'}$$

$$g(\sigma) = \sum_{q' \in Q} P(q'|q, \sigma) U(q')$$

want to show : $\max_{q \in Q} |U_{t+1}^q - U(q)| \leq K \max_{q \in Q} |U_t^q - U(q)|$ for some fraction $K \in [0, 1)$

$$|U_{t+1}^q - U(q)| = \gamma \left| \max_{\sigma} f(\sigma) - \max_{\sigma} g(\sigma) \right|$$

$$f(\sigma) = \sum_{q' \in Q} P(q'|q, \sigma) U_t^{q'} \\ g(\sigma) = \sum_{q' \in Q} P(q'|q, \sigma) U(q')$$

$$\leq \gamma \max_{\sigma} \left| f(\sigma) - g(\sigma) \right|$$

$$= \gamma \max_{\sigma} \left| \sum_{q' \in Q} P(q'|q, \sigma) U_t^{q'} - \sum_{q' \in Q} P(q'|q, \sigma) U(q') \right|$$

want to show : $\max_{q \in Q} |U_{t+1}^q - U(q)| \leq K \max_{q \in Q} |U_t^q - U(q)|$ for some fraction $K \in [0, 1)$

$$|U_{t+1}^q - U(q)| = \gamma \left| \max_{\sigma} f(\sigma) - \max_{\sigma} g(\sigma) \right|$$

$$f(\sigma) = \sum_{q' \in Q} P(q'|q, \sigma) U_t^{q'} \\ g(\sigma) = \sum_{q' \in Q} P(q'|q, \sigma) U(q')$$

$$\leq \gamma \max_{\sigma} \left| f(\sigma) - g(\sigma) \right|$$

$$= \gamma \max_{\sigma} \left| \sum_{q' \in Q} P(q'|q, \sigma) U_t^{q'} - \sum_{q' \in Q} P(q'|q, \sigma) U(q') \right|$$

$$= \gamma \max_{\sigma} \left| \sum_{q' \in Q} P(q'|q, \sigma) U_t^{q'} - P(q'|q, \sigma) U(q') \right|$$

want to show : $\max_{q \in Q} |U_{t+1}^q - U(q)| \leq K \max_{q \in Q} |U_t^q - U(q)|$ for some fraction $K \in [0, 1)$

$$|U_{t+1}^q - U(q)| = \gamma \left| \max_{\sigma} f(\sigma) - \max_{\sigma} g(\sigma) \right|$$

$$f(\sigma) = \sum_{q' \in Q} P(q'|q, \sigma) U_t^{q'} \\ g(\sigma) = \sum_{q' \in Q} P(q'|q, \sigma) U(q')$$

$$\leq \gamma \max_{\sigma} \left| f(\sigma) - g(\sigma) \right|$$

$$= \gamma \max_{\sigma} \left| \sum_{q' \in Q} P(q'|q, \sigma) U_t^{q'} - \sum_{q' \in Q} P(q'|q, \sigma) U(q') \right|$$

$$= \gamma \max_{\sigma} \left| \sum_{q' \in Q} P(q'|q, \sigma) (U_t^{q'} - U(q')) \right|$$

want to show : $\max_{q \in Q} |U_{t+1}^q - U(q)| \leq K \max_{q \in Q} |U_t^q - U(q)|$ for some fraction $K \in [0, 1]$

$$|U_{t+1}^q - U(q)| \leq \gamma \max_{\sigma} \left| \sum_{q' \in Q} P(q'|q, \sigma) (U_t^{q'} - U(q')) \right|$$

want to show : $\max_{q \in Q} |U_{t+1}^q - U(q)| \leq K \max_{q \in Q} |U_t^q - U(q)|$ for some fraction $K \in [0, 1]$

$$|U_{t+1}^q - U(q)| \leq \gamma \max_{\sigma} \left| \sum_{q' \in Q} P(q'|q, \sigma) \frac{(U_t^{q'} - U(q'))}{f(q')} \right|$$

want to
Show : $\max_{q \in Q} |U_{t+1}^q - U(q)| \leq K \max_{q \in Q} |U_t^q - U(q)|$ for some fraction $K \in [0, 1]$

$$|U_{t+1}^q - U(q)| \leq \gamma \max_{\sigma} \left| \sum_{q' \in Q} f(q') \right|$$

$$f(q') = P(q'|q, \sigma) (U_t^{q'} - U(q'))$$

want to
Show : $\max_{q \in Q} |U_{t+1}^q - U(q)| \leq K \max_{q \in Q} |U_t^q - U(q)|$ for some fraction $K \in [0, 1]$

$$|U_{t+1}^q - U(q)| \leq \gamma \max_{\sigma} \left| \sum_{q' \in Q} f(q') \right|$$

$$f(q') = P(q'|q, \sigma) (U_t^{q'} - U(q'))$$

$$\overline{\left| \sum_{q' \in Q} f(q') \right|}$$

$$\left| \sum_{q' \in Q} f(q') \right| \stackrel{?}{=} \sum_{q' \in Q} |f(q')|$$

how are these related?

want to
Show : $\max_{q \in Q} |U_{t+1}^q - U(q)| \leq K \max_{q \in Q} |U_t^q - U(q)|$ for some fraction $K \in [0, 1]$

$$|U_{t+1}^q - U(q)| \leq \gamma \max_{q' \in Q} \left| \sum_{q' \in Q} f(q') \right|$$

$$f(q') = P(q'|q, \sigma) (U_t^{q'} - U(q'))$$

$$\left| \sum_{q' \in Q} f(q') \right| \leq \sum_{q' \in Q} |f(q')|$$

e.g.

$$|2 + (-3) + 4| \leq |2| + |-3| + |4|$$

want to show : $\max_{q \in Q} |U_{t+1}^q - U(q)| \leq K \max_{q \in Q} |U_t^q - U(q)|$ for some fraction $K \in [0, 1]$

$$|U_{t+1}^q - U(q)| \leq \gamma \max_{\sigma} \left| \sum_{q' \in Q} f(q') \right|$$

$$f(q') = P(q'|q, \sigma) (U_t^{q'} - U(q'))$$

$$\leq \gamma \max_{\sigma} \sum_{q' \in Q} |f(q')|$$

$$= \gamma \max_{\sigma} \sum_{q' \in Q} \left| P(q'|q, \sigma) (U_t^{q'} - U(q')) \right|$$

want to show : $\max_{q \in Q} |U_{t+1}^q - U(q)| \leq K \max_{q \in Q} |U_t^q - U(q)|$ for some fraction $K \in [0, 1]$

$$|U_{t+1}^q - U(q)| \leq \gamma \max_{\sigma} \sum_{q' \in Q} |P(q'|q, \sigma) (U_t^{q'} - U(q'))|$$

want to show : $\max_{q \in Q} |U_{t+1}^q - U(q)| \leq K \max_{q \in Q} |U_t^q - U(q)|$ for some fraction $K \in [0, 1]$

$$\begin{aligned}
 |U_{t+1}^q - U(q)| &\leq \gamma \max_{\sigma} \sum_{q' \in Q} \left| P(q'|q, \sigma) (U_t^{q'} - U(q')) \right| \\
 &= \gamma \max_{\sigma} \sum_{q' \in Q} \frac{P(q'|q, \sigma)}{\text{always nonnegative}} \left| U_t^{q'} - U(q') \right|
 \end{aligned}$$

want to
Show : $\max_{q \in Q} |U_{t+1}^q - U(q)| \leq K \max_{q \in Q} |U_t^q - U(q)|$ for some fraction $K \in [0, 1]$

$$|U_{t+1}^q - U(q)| \leq \gamma \max_{\sigma} \sum_{q' \in Q} |P(q'|q, \sigma) (U_t^{q'} - U(q'))|$$

$$= \gamma \max_{\sigma} \sum_{q' \in Q} P(q'|q, \sigma) |U_t^{q'} - U(q')|$$



the expected value of

$$|U_t^{q'} - U(q')|$$

with respect to $P(q'|q, \sigma)$

want to
Show : $\max_{q \in Q} |U_{t+1}^q - U(q)| \leq K \max_{q \in Q} |U_t^q - U(q)|$ for some fraction $K \in [0, 1]$

$$|U_{t+1}^q - U(q)| \leq \gamma \max_{\sigma} \sum_{q' \in Q} |P(q'|q, \sigma) (U_t^{q'} - U(q'))|$$

$$= \gamma \max_{\sigma} \sum_{q' \in Q} P(q'|q, \sigma) |U_t^{q'} - U(q')|$$



the expected value of

$$|U_t^{q'} - U(q')| \leq$$

the maximum value of

$$|U_t^{q'} - U(q')|$$

with respect to $P(q'|q, \sigma)$

want to
Show : $\max_{q \in Q} |U_{t+1}^q - U(q)| \leq K \max_{q \in Q} |U_t^q - U(q)|$ for some fraction $K \in [0, 1)$

$$|U_{t+1}^q - U(q)| \leq \gamma \max_{\sigma} \sum_{q' \in Q} |P(q'|q, \sigma) (U_t^{q'} - U(q'))|$$

$$= \gamma \max_{\sigma} \sum_{q' \in Q} P(q'|q, \sigma) |U_t^{q'} - U(q')|$$



the expected value of

$$|U_t^{q'} - U(q')| \leq \max_{q'} |U_t^{q'} - U(q')|$$

with respect to $P(q'|q, \sigma)$

want to show : $\max_{q \in Q} |U_{t+1}^q - U(q)| \leq K \max_{q \in Q} |U_t^q - U(q)|$ for some fraction $K \in [0, 1]$

$$\begin{aligned}
 |U_{t+1}^q - U(q)| &\leq \gamma \max_{\sigma} \sum_{q' \in Q} \left| P(q'|q, \sigma) (U_t^{q'} - U(q')) \right| \\
 &= \gamma \max_{\sigma} \sum_{q' \in Q} P(q'|q, \sigma) \left| U_t^{q'} - U(q') \right| \\
 &\leq \gamma \max_{\sigma} \max_{q'} \left| U_t^{q'} - U(q') \right|
 \end{aligned}$$

want to show : $\max_{q \in Q} |U_{t+1}^q - U(q)| \leq K \max_{q \in Q} |U_t^q - U(q)|$ for some fraction $K \in [0, 1]$

$$\begin{aligned}
 |U_{t+1}^q - U(q)| &\leq \gamma \max_{\sigma} \sum_{q' \in Q} |P(q'|q, \sigma) (U_t^{q'} - U(q'))| \\
 &= \gamma \max_{\sigma} \sum_{q' \in Q} P(q'|q, \sigma) |U_t^{q'} - U(q')| \\
 &\leq \gamma \max_{\sigma} \max_{q'} |U_t^{q'} - U(q')| \\
 &\leq \gamma \max_{q'} |U_t^{q'} - U(q')|
 \end{aligned}$$



want to
Show : $\max_{q \in Q} |U_{t+1}^q - U(q)| \leq K \max_{q \in Q} |U_t^q - U(q)|$ for some fraction $K \in [0, 1]$

for every state $q \in Q$:

$$|U_{t+1}^q - U(q)| \leq \gamma \max_{q' \in Q} |U_t^{q'} - U(q')|$$

want to show : $\max_{q \in Q} |U_{t+1}^q - U(q)| \leq K \max_{q \in Q} |U_t^q - U(q)|$ for some fraction $K \in [0, 1)$

for every state $q \in Q$:

$$|U_{t+1}^q - U(q)| \leq \gamma \max_{q' \in Q} |U_t^{q'} - U(q')|$$

therefore :

$$\max_{q \in Q} |U_{t+1}^q - U(q)| \leq \gamma \max_{q' \in Q} |U_t^{q'} - U(q')|$$

want to show : $\max_{q \in Q} |U_{t+1}^q - U(q)| \leq K \max_{q \in Q} |U_t^q - U(q)|$ for some fraction $K \in [0, 1]$

for every state $q \in Q$:

$$|U_{t+1}^q - U(q)| \leq \gamma \max_{q' \in Q} |U_t^{q'} - U(q')|$$

therefore :

$$\max_{q \in Q} |U_{t+1}^q - U(q)| \leq \gamma \max_{q \in Q} |U_t^q - U(q)|$$

if $\text{dist} \left(\begin{bmatrix} u_t^{q_0} \\ \vdots \\ u_t^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right) \leq \gamma \text{dist} \left(\begin{bmatrix} u_{t-1}^{q_0} \\ \vdots \\ u_{t-1}^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right)$, then: $\lim_{t \rightarrow \infty} \text{dist} \left(\begin{bmatrix} u_t^{q_0} \\ \vdots \\ u_t^{q_n} \end{bmatrix}, \begin{bmatrix} U(q_0) \\ \vdots \\ U(q_n) \end{bmatrix} \right) = 0$
 for some fraction $K \in [0, 1]$

~~if we showed~~
~~if we can show~~ that every
 step of value iteration brings
 our guesses closer to the
 maximum expected utilities

if discount $\gamma < 1$

so
~~then~~ value iteration
 converges to the
 maximum expected
 utilities

t	U(A)	U(B)
0	1.2500	-1.1300
1	1.6250	-1.3270
2	1.8125	-1.3683
3	1.9063	-1.3661
4	1.9531	-1.3558
5	1.9766	-1.3470
6	1.9883	-1.3411
7	1.9941	-1.3376
8	1.9971	-1.3356
9	1.9985	-1.3345
10	1.9993	-1.3340
11	1.9996	-1.3337
12	1.9998	-1.3335
13	1.9999	-1.3334
14	2.0000	-1.3334
15	2.0000	-1.3334

value iteration

- guess values for the unknowns

for each state q , U_t^q is a guess for $U(q)$

- for $t = 1$ to T :

compute new guesses using the bellman equations, i.e.

$$\text{for each state } q, U_{t+1}^q = R(q) + \gamma \max_{\sigma} \sum_{q' \in Q} P(q'|q, \sigma) U_t^{q'}$$

1 measure the distance between $U: Q \rightarrow \mathbb{R}$ and our guesses U_t^q at iteration t using:

$$\max_{q \in Q} |U_{t+1}^q - U(q)|$$

2 show:

$$\max_{q \in Q} |U_{t+1}^q - U(q)|$$

$$\leq \gamma \max_{q \in Q} |U_t^q - U(q)|$$

3 this implies:

$$\lim_{t \rightarrow \infty} U_t^q = U(q)$$

proof

if $\gamma < 1$, then value iteration converges to the maximum expected utilities, i.e. $\lim_{t \rightarrow \infty} U_t^q = U(q)$