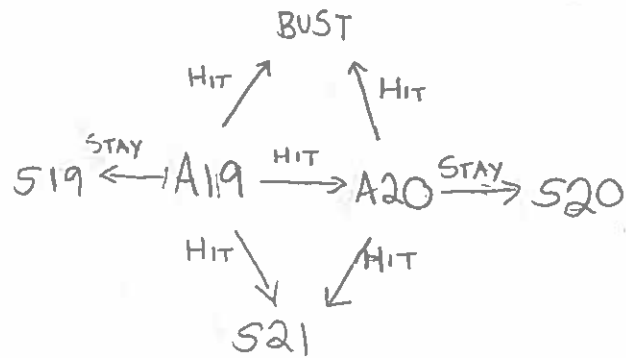# MARKOV DECISION PROCESS

1) Rather than treating blackjack as a 2-player game where one player is Fate, we could model it as the following (nondeterministic) state machine:

$$BUST$$

$$HIT \nearrow \qquad \uparrow HIT$$

$$S19 \xleftarrow{STAY} A19 \xrightarrow{HIT} A20 \xrightarrow{STAY} S20$$

$$HIT \searrow \qquad \swarrow HIT$$

$$S21$$

i.e. $M = (Q, \Sigma, \Delta, q_0, F)$ where:

- $Q = \{A19, S19, A20, S20, S21, BUST\}$
- $\Sigma = \{HIT, STAY\}$
- $\Delta = \{(A19, STAY, S19), (A20, STAY, S20), (A19, HIT, A20),$
  $(A19, HIT, S21), (A19, HIT, BUST), (A20, HIT, S21),$
  $(A20, HIT, BUST)\}$
- $q_0 = A19$
- $F = \{S19, S20, S21, BUST\}$

2) In addition, we want to specify a reward $R_t$ for reaching a state after $t$ transitions:

$$\forall t \geq 0 \qquad R_t(S20) = .58 \qquad R_t(A20) = 0$$

$$R_t(S21) = .88 \qquad R_t(BUST) = -1$$

$$R_t(S19) = .27 \qquad R_t(A19) = 0$$

③ And finally, we want to know the likelihood of each transition, given we take a particular action in a particular state.

$w(<19, STAY, S19>) = 1$

$w(<20, STAY, S20>) = 1$

$w(<19, HIT, 20>) = \frac{1}{13}$

$w(<19, HIT, S21>) = \frac{1}{13}$

$w(<19, HIT, BUST?>) = \frac{11}{13}$

$\left. \begin{array}{l} \end{array} \right\} - \text{sum } \underset{q'}{} w(<19, HIT, q'>) = 1$

$w(<20, HIT, S21>) = \frac{1}{13}$

$w(<20, HIT, BUST>) = \frac{12}{13}$

$\left. \begin{array}{l} \end{array} \right\} - \text{sum } \underset{q'}{} w(<20, HIT, q'>) = 1$

$w(<19, HIT, q'>) = P(q' \mid q=19, \sigma=HIT)$

---

④ This is called a <u>Markov Decision Process</u> (MDP). Formally an MDP is a triple $(M, R, w)$, where:

- $M = (Q, \Sigma, \Delta, q_0, F)$ is a state machine
- $R: Q \times \mathbb{N} \rightarrow \mathbb{R}$ is a "reward function"
- $w: \Delta \mapsto \mathbb{R}$ s.t. for all $q \in Q, \sigma \in \Sigma$: $\displaystyle\sum_{(q,\sigma,q') \in \Delta} w(<q, \sigma, q'>) = 1$

We assume that $R(q, t) = \gamma \cdot R(q, t-1) \; \forall q \in Q, t \geq 1$ for some "discounting factor" $\gamma$ s.t. $0 < \gamma \leq 1$.

# MARKOV DECISION PROCESS

⑤ It can be helpful to define the following shorthand for dealing with MDPs:

→ $R_t(q) \triangleq R(q, t) \qquad \forall q \in Q, t \geq 0$

→ $P(q' \mid q, \sigma) \triangleq w(\langle q, \sigma, q' \rangle) \qquad \forall \langle q, \sigma, q' \rangle \in \Delta$

→ We write path $\langle (q, \sigma_0, q_1), (q_1, \sigma_1, q_2), \ldots, (q_{n-1}, \sigma_{n-1}, q_n) \rangle$ as:

$$q \xrightarrow{\sigma_0} q_1 \xrightarrow{\sigma_1} \ldots \xrightarrow{\sigma_{n-1}} q_n$$

→ reward $\left( q \xrightarrow{\sigma_0} q_1 \xrightarrow{\sigma_1} \ldots \xrightarrow{\sigma_{n-1}} q_n \right)$

$$= R_0(q) + R_1(q_1) + \ldots + R_n(q_n)$$

→ $P( q \xrightarrow{\sigma_0} q_1 \xrightarrow{\sigma_1} \ldots \xrightarrow{\sigma_{n-1}} q_n )$

$$= P(q_1 \mid q, \sigma_0) \cdot P(q_2 \mid q_1, \sigma_1) \cdot \ldots \cdot P(q_n \mid q_{n-1}, \sigma_{n-1})$$

---

⑥ The main computational challenge, given an MDP, is to determine the best decision to make in each state.

For example:

If I have 19, should I HIT or STAY?
If I have 20, should I HIT or STAY?

We can formalize this as a function $\pi : (Q \backslash F) \to \Sigma$, which we call a policy.

e.g.
$\pi(A19) = HIT$
$\pi(A20) = HIT$

# Markov Decision Process

⑦ Given a policy $\pi$, I can compute my expected reward $U^\pi$, starting from various states:

$$U^\pi(A19) = \text{reward}\left(A19 \xrightarrow{HIT} A20 \xrightarrow{HIT} S21\right) P\left(A19 \xrightarrow{HIT} A20 \xrightarrow{HIT} S21\right)$$
$$+ \text{reward}\left(A19 \xrightarrow{HIT} A20 \xrightarrow{HIT} BUST\right) P\left(A19 \xrightarrow{HIT} A20 \xrightarrow{HIT} BUST\right)$$
$$+ \text{reward}\left(A19 \xrightarrow{HIT} S21\right) P\left(A19 \xrightarrow{HIT} S21\right)$$
$$+ \text{reward}\left(A19 \xrightarrow{HIT} BUST\right) P\left(A19 \xrightarrow{HIT} BUST\right)$$

$$U^\pi(A20) = \text{reward}\left(A20 \xrightarrow{HIT} S21\right) P\left(A20 \xrightarrow{HIT} S21\right)$$
$$+ \text{reward}\left(A20 \xrightarrow{HIT} BUST\right) P\left(A20 \xrightarrow{HIT} BUST\right)$$

---

⑧ These expected rewards (usually called <u>expected utility</u>) can be expressed in terms of each other:

$$U^\pi(A19) = \left(R_0(A19) + R_1(A20) + R_2(S21)\right) P\left(A19 \xrightarrow{HIT} A20 \xrightarrow{HIT} S21\right)$$
$$+ \left(R_0(A19) + R_1(A20) + R_2(BUST)\right) P\left(A19 \xrightarrow{HIT} A20 \xrightarrow{HIT} BUST\right)$$
$$+ \left(R_0(A19) + R_1(S21)\right) P\left(A19 \xrightarrow{HIT} S21\right)$$
$$+ \left(R_0(A19) + R_1(BUST)\right) P\left(A19 \xrightarrow{HIT} BUST\right)$$

$$= R_0(A19)\left[ P\left(A19 \xrightarrow{HIT} A20 \xrightarrow{HIT} S21\right) + P\left(A19 \xrightarrow{HIT} A20 \xrightarrow{HIT} BUST\right) \right.$$
$$\left. + P\left(A19 \xrightarrow{HIT} S21\right) + P\left(A19 \xrightarrow{HIT} BUST\right) \right]$$

$$+ \left(R_1(A20) + R_2(S21)\right) P\left(A19 \xrightarrow{H} A20 \xrightarrow{H} S21\right)$$
$$+ \left(R_1(A20) + R_2(BUST)\right) P\left(A19 \xrightarrow{H} A20 \xrightarrow{H} BUST\right)$$
$$+ R_1(S21) P\left(A19 \xrightarrow{H} S21\right)$$
$$+ R_1(BUST) P\left(A19 \xrightarrow{H} BUST\right)$$

this equals 1, b/c it is a prob. distribution.

# Markov Decision Process

(8) $U^\pi(A19) = R_0(A19)$

$\qquad + P(A20 \mid A19, HIT)\left[\begin{array}{l}(R_1(A20) + R_2(S21))P(A20 \xrightarrow{H} S21) \\ + (R_1(A20) + R_2(BUST))P(A20 \xrightarrow{H} BUST)\end{array}\right]$

$\qquad + P(S21 \mid A19, HIT)\, R_1(S21)$

$\qquad + P(BUST \mid A19, HIT)\, R_1(BUST)$

$= R_0(A19)$

$\qquad + P(A20 \mid A19, HIT)\left[\begin{array}{l}(\gamma R_0(A20) + \gamma R_1(S21))P(A20 \xrightarrow{H} S21) \\ + (\gamma R_0(A20) + \gamma R_1(BUST))P(A20 \xrightarrow{H} BUST)\end{array}\right]$

$\qquad + P(S21 \mid A19, HIT)(\gamma \cdot R_0(S21))$

$\qquad + P(BUST \mid A19, HIT)(\gamma \cdot R_0(BUST))$

$= R_0(A19)$

$\qquad + \gamma \cdot P(A20 \mid A19, HIT)\left[\begin{array}{l}\overbrace{(R_0(A20) + R_1(S21))}^{\text{reward}(A20 \xrightarrow{H} S21)}P(A20 \xrightarrow{H} S21) \\ + \underbrace{(R_0(A20) + R_1(BUST))}_{\text{reward}(A20 \xrightarrow{H} BUST)}P(A20 \xrightarrow{H} BUST)\end{array}\right]$

$\qquad + \gamma \cdot P(S21 \mid A19, HIT)\, R_0(S21)$

$\qquad + \gamma \cdot P(BUST \mid A19, HIT)\, R_0(BUST)$

$= R_0(A19)$

$\qquad + \gamma \cdot P(A20 \mid A19, HIT)\, U^\pi(A20)$

$\qquad + \gamma \cdot P(S21 \mid A19, HIT)\, U^\pi(S21)$

$\qquad + \gamma \cdot P(BUST \mid A19, HIT)\, U^\pi(BUST)$

$= R_0(A19)$

$\qquad + \gamma \cdot \left[\begin{array}{l} U^\pi(A20) \cdot P(A20 \mid A19, \pi(A19)) \\ + U^\pi(S21) \cdot P(S21 \mid A19, \pi(A19)) \\ + U^\pi(BUST) \cdot P(BUST \mid A19, \pi(A19))\end{array}\right]$

$= R_0(A19) + \gamma \cdot \displaystyle\sum_{q' \in Q} U^\pi(q') \cdot P(q' \mid A19, \pi(A19))$

# MARKOV DECISION PROCESS

⑨ This is a general result:

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

⑩ Usually we're not simply interested in computing the expected utility of a state, given some arbitrary policy $\pi$. Rather, we'd like to know how much reward we should expect if we execute the **best** policy $\pi^*$.

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

$$= \max_\pi \left[ R_o(q) + \gamma \sum_{q' \in Q} U^\pi(q') \cdot P(q' | q, \pi(q)) \right]$$

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

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

we choose some action $\sigma$ for state $q$

$$= R_o(q) + \gamma \cdot \max_\sigma \sum_{q' \in Q} \left( \max_{\pi | \pi(q) = \sigma} U^\pi(q') \right) \cdot P(q' | q, \sigma)$$

?/

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

reward from the current state

maximum expected utility of the next state, given optimal action

# MARKOV DECISION PROCESS

⑪ For our blackjack example, we get the following equations.

$$U(A19) = \max \begin{cases} U(S19) \cdot P(S19 \mid A19, STAY), \\[8pt] U(A20) \cdot P(A20 \mid A19, HIT) \\ \quad + U(S21) \cdot P(S21 \mid A19, HIT) \\ \quad + U(BUST) \cdot P(BUST \mid A19, HIT) \end{cases}$$

$$U(A20) = \max \begin{cases} U(S20) \cdot P(S20 \mid A20, STAY), \\[8pt] U(S21) \cdot P(S21 \mid A20, HIT) \\ \quad + U(BUST) \cdot P(BUST \mid A20, HIT) \end{cases}$$

$$U(S19) = R_0(S19)$$
$$U(S20) = R_0(S20)$$
$$U(S21) = R_0(S21)$$
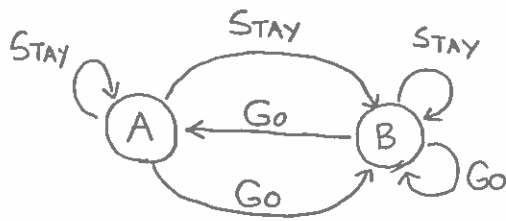$$U(BUST) = R_0(BUST)$$

We have six equations for six unknowns ($U(A19), U(A20), U(S19), U(S20), U(S21), U(BUST)$).

⑫ Importantly, these equations are not linear, so we can't use linear algebra techniques.

# MARKOV DECISION PROCESS

(13) To motivate our equation-solving technique, let's use a smaller MDP:



where:

$$P(A \mid A, \text{STAY}) = \tfrac{1}{2}$$
$$P(B \mid A, \text{STAY}) = \tfrac{1}{2}$$
$$P(B \mid A, \text{Go}) = 1$$

$$P(B \mid B, \text{STAY}) = 1$$
$$P(A \mid B, \text{Go}) = \tfrac{1}{5}$$
$$P(B \mid B, \text{Go}) = \tfrac{4}{5}$$

and

$$R_t(A) = \gamma^t \cdot 1$$
$$R_t(B) = \gamma^t \cdot (-1)$$

where $\gamma = \tfrac{1}{2}$.

(14) We get the following equations:

$$U(A) = R_0(A) + \gamma \cdot \max \left\{ U(B), \; U(A)P(A \mid A, \text{STAY}) + U(B)P(B \mid A, \text{STAY}) \right\}$$

$$= 1 + \tfrac{1}{2} \max \left\{ U(B), \; \tfrac{1}{2}U(A) + \tfrac{1}{2}U(B) \right\}$$

$$U(B) = R_0(B) + \gamma \cdot \max \left\{ U(B), \; U(A)P(A \mid B, \text{Go}) + U(B)P(B \mid B, \text{Go}) \right\}$$

$$= -1 + \tfrac{1}{2} \max \left\{ U(B), \; \tfrac{1}{5}U(A) + \tfrac{4}{5}U(B) \right\}$$

# MARKOV DECISION PROCESS

(15) Suppose we guess the values of $U(A)$ and $U(B)$:

$U_0^A$ is our guess at $U(A)$

$U_0^B$ is our guess at $U(B)$

Consider the following iterative algorithm:

for $i = 1$ to $\infty$:

$$\text{let } U_i^A = 1 + \frac{1}{2} \max\left\{ U_{i-1}^B, \frac{1}{2} U_{i-1}^A + \frac{1}{2} U_{i-1}^B \right\}$$

$$\text{let } U_i^B = -1 + \frac{1}{2} \max\left\{ U_{i-1}^B, \frac{1}{5} U_{i-1}^A + \frac{4}{5} U_{i-1}^B \right\}$$

At each iteration, we assume our guesses for $U(A)$ and $U(B)$ from the previous iteration <u>are correct</u>, and we compute new guesses using our equations.

---

(16) Why would this ever work? Well, it seems to converge...

| $i$ | $U_i^A$ | $U_i^B$ |
|-----|---------|---------|
| 0 | 0 | 0 |
| 1 | 1 | -1 |
| 2 | 1 | -1.3 |
| 3 | .925 | -1.42 |
| 4 | .876 | -1.476 |
| 5 | .850 | -1.503 |
| 6 | .836 | -1.516 |
| 7 | .830 | -1.523 |

(17) But can we prove it? Assume for the moment that a solution exists, i.e. there's a vector

$$U^* = \begin{bmatrix} U^A \\ U^B \end{bmatrix}$$

such that:

$$U^A = 1 + \frac{1}{2} \max \left\{ U^B, \frac{1}{2} U^A + \frac{1}{2} U^B \right\}$$

$$U^B = -1 + \frac{1}{2} \max \left\{ U^B, \frac{1}{5} U^A + \frac{4}{5} U^B \right\}$$

---

(18) Let's also measure how bad our guesses are.

Suppose the solution is $\begin{bmatrix} 2 \\ 4 \end{bmatrix}$ and our guess is $\begin{bmatrix} 5 \\ 3 \end{bmatrix}$.

We'll measure the distance of our guesses to the solution as the absolute difference between our worst guess and the solution, i.e.

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

⑲ Now what if we could show that our guesses get better every iteration? i.e. that

$$\text{dist}\left(\begin{bmatrix} U^A_{i+1} \\ U^B_{i+1} \end{bmatrix}, \begin{bmatrix} U^A \\ U^B \end{bmatrix}\right) \leq K \cdot \text{dist}\left(\begin{bmatrix} U^A_i \\ U^B_i \end{bmatrix}, \begin{bmatrix} U^A \\ U^B \end{bmatrix}\right)$$

for $0 \leq K < 1$.

That would mean:

$$\lim_{i \to \infty} \text{dist}\left(\begin{bmatrix} U^A_i \\ U^B_i \end{bmatrix}, \begin{bmatrix} U^A \\ U^B \end{bmatrix}\right)$$

$$\leq \lim_{i \to \infty} K^i \, \text{dist}\left(\begin{bmatrix} U^A_0 \\ U^B_0 \end{bmatrix}, \begin{bmatrix} U^A \\ U^B \end{bmatrix}\right)$$

$$= \lim_{i \to \infty} K^i$$

$$= 0$$

Since $\text{dist}\left(\begin{bmatrix} U^A_i \\ U^B_i \end{bmatrix}, \begin{bmatrix} U^A \\ U^B \end{bmatrix}\right) \geq 0$, thus $\lim_{i \to \infty} \text{dist}\left(\begin{bmatrix} U^A_i \\ U^B_i \end{bmatrix}, \begin{bmatrix} U^A \\ U^B \end{bmatrix}\right) = 0$

So our guesses would converge to the solution.

⑳ So then let's show it.

$$\text{dist}\left(\begin{bmatrix} U^A_{i+1} \\ U^B_{i+1} \end{bmatrix}, \begin{bmatrix} U^A \\ U^B \end{bmatrix}\right) = \max_{q \in \{A,B\}} |U^q_{i+1} - U^q|$$

If we simplify $|U^q_{i+1} - U^q|$, we get:

$$|U^q_{i+1} - U^q| = \left| \gamma \cdot \left[ \max_\sigma \sum_{q'} U^{q'}_i P(q'|q,\sigma) - \max_\sigma \sum_{q'} U^{q'} P(q'|q,\sigma) \right] \right|$$

because $\forall f, g$

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

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

$$\leq \gamma \cdot \max_\sigma \left| \sum_{q'} U^{q'}_i P(q'|q,\sigma) - \sum_{q'} U^{q'} P(q'|q,\sigma) \right|$$

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

$$\leq \gamma \cdot \max_\sigma \sum_{q'} |U^{q'}_i - U^{q'}| P(q'|q,\sigma)$$

$$\leq \gamma \cdot \max_\sigma \max_{q'} |U^{q'}_i - U^{q'}|$$

(because the weighted average of a set of numbers is at most the max)

$$= \gamma \cdot \max_{q'} |U^{q'}_i - U^{q'}|$$

Thus:

$$\text{dist}\left(\begin{bmatrix} U^A_{i+1} \\ U^B_{i+1} \end{bmatrix}, \begin{bmatrix} U^A \\ U^B \end{bmatrix}\right) \leq \max_q \gamma \cdot \max_{q'} |U^{q'}_i - U^{q'}|$$

$$= \gamma \max_{q'} |U^{q'}_i - U^{q'}|$$

$$= \gamma \cdot \text{dist}\left(\begin{bmatrix} U^A_i \\ U^B_i \end{bmatrix}, \begin{bmatrix} U^A \\ U^B \end{bmatrix}\right)$$

# MARKOV DECISION PROCESS

21) So as long as the discounting factor is between 0 and 1 (not-including 1), then our iterative technique (called "value iteration") will converge to the correct solution.