

# Tabular methods and dynamic programming

Optimal control problem

$$\max_{\rho} V^\rho(x) = \mathbb{E} \left[ \sum_{k=0}^{\infty} \gamma^k r(X_k, \rho(X_k)) \mid X_0 = x \right]$$

s.t.  $X_{k+1} \sim f(x_{k+1} \mid x_k, u_k)$  condition

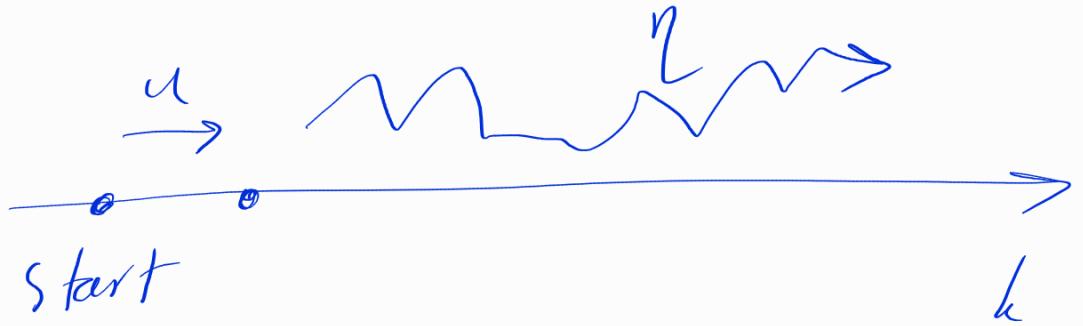
$$r_k = r(x_k, u_k)$$

The optimal policy:

$$\rho^*(x) = \arg \max_{\rho} V^\rho(x)$$

---

Let's define  $\hat{\rho}(u|\eta)$  — concatenation  
of an action  $u$  with an  
arbitrary tail policy  $\eta$



Let's denote  $X_+^u := f(x_+ | x, u)$

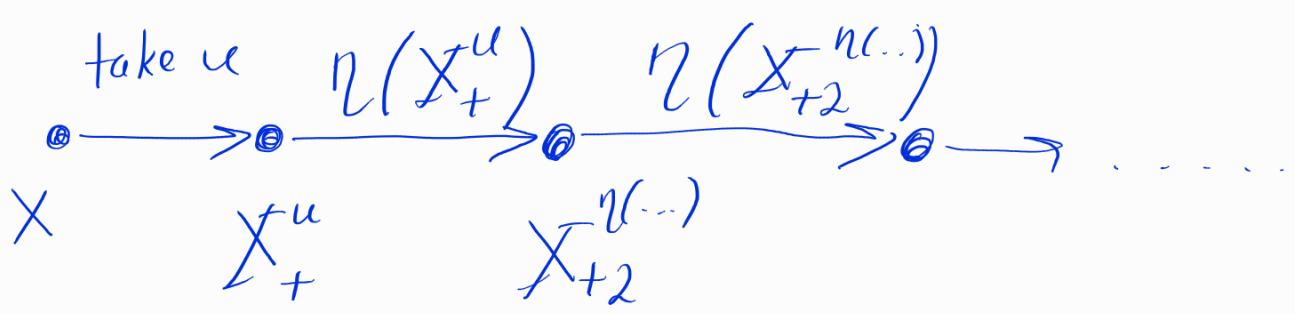
(next state upon applying the action  $u$ )

Let's take a look at the optimal value:

$$V^*(x) = \max_{\rho} V^\rho(x)$$

$$= \max_{\rho} \mathbb{E} \left[ \sum_{k=0}^{\infty} \gamma^k r(X_k, \rho(X_k)) \mid X_0 = x \right]$$

$$= \max_{\rho} \left\{ r(x, u) + \gamma \mathbb{E} [V^\rho(X_+^u)] \right\}$$



$\hat{P} = \langle$  first apply  $u$ ,  
then follow  $\eta \rangle$

---

$$V^*(x) = \max_{\hat{P}} \left\{ r(x, u) + \gamma \mathbb{E}[V^u(X^u_+)] \right\}$$

$$\leq \max_u \left\{ r(x, u) + \gamma \mathbb{E}[V^*(X^u_+)] \right\}$$

On the other hand,

$$\max_{\hat{P}} \left\{ r(x, u) + \gamma \mathbb{E}[V^u(X^u_+)] \right\} \geq r(x, v) + \gamma \mathbb{E}[V^\xi(X^v_+)]$$

$\forall$  action  $v$  and  
tail policy  $\xi$

If we take  $\mathcal{S} = \rho^*(x)$ ,  
 $\xi = \text{tail of } \rho^*$ , then

$$\max_{\hat{\rho}} \left\{ r(x, u) + \gamma \mathbb{E} \left[ \sqrt{\rho^*(x_+^u)} \right] \right\} \geq r(x, \rho^*(x)) + \gamma \mathbb{E} \left[ \sqrt{\rho^*(x_+^{\rho^*(x)})} \right] = \sqrt^*(x) !$$

Combining the two inequalities,  
that we obtained, we conclude:

$$\sqrt^*(x) = \max_u \left\{ r(x, u) + \gamma \mathbb{E} \left[ \sqrt^*(x_+^u) \right] \right\}$$



Dynamic programming principle  
"DP"

Let's introduce a DP operator:  
for an arbitrary function  $W$ ,

$$T[W](x) := \max_u \{ r(x, u) + \gamma \mathbb{E}[W(x^u)] \}$$

maps functions  
to functions

It so happens that  $T$   
satisfies the following two  
conditions:

①

$$\forall x, u, W_1, W_2$$

$$r(x, u) + \gamma \mathbb{E}[W_1(x^u)] \leq$$

$$r(x, u) + \gamma \mathbb{E}[W_2(x^u)] \text{ whenever}$$

$$W_1(x) \leq W_2(x), \forall x$$

And so

$$T[W_1] \leq T[W_2]$$

$\Rightarrow$  monotonicity condition

(the 1st Blackwell's contractive mapping condition)

② ( $T$  is discounting)

$$T[W+a](x) =$$

$$\max_u \{ r(x, u) + \gamma \mathbb{E}[W(X_t^u)] + \gamma a \} \\ = T[W] + \gamma a$$

Thus,  $T$  is contractive with modulus  $\gamma$ , and by the fixed-point theorem, successive application of  $T$  to itself converges to a unique fixed point

So if we do iterations like

$$V_i := T[V_{i-1}], i \in \mathbb{N}_0$$

(natural numbers)  
with zero

starting at any initial guess

$V_0$ , we converge to a limit

with the rate

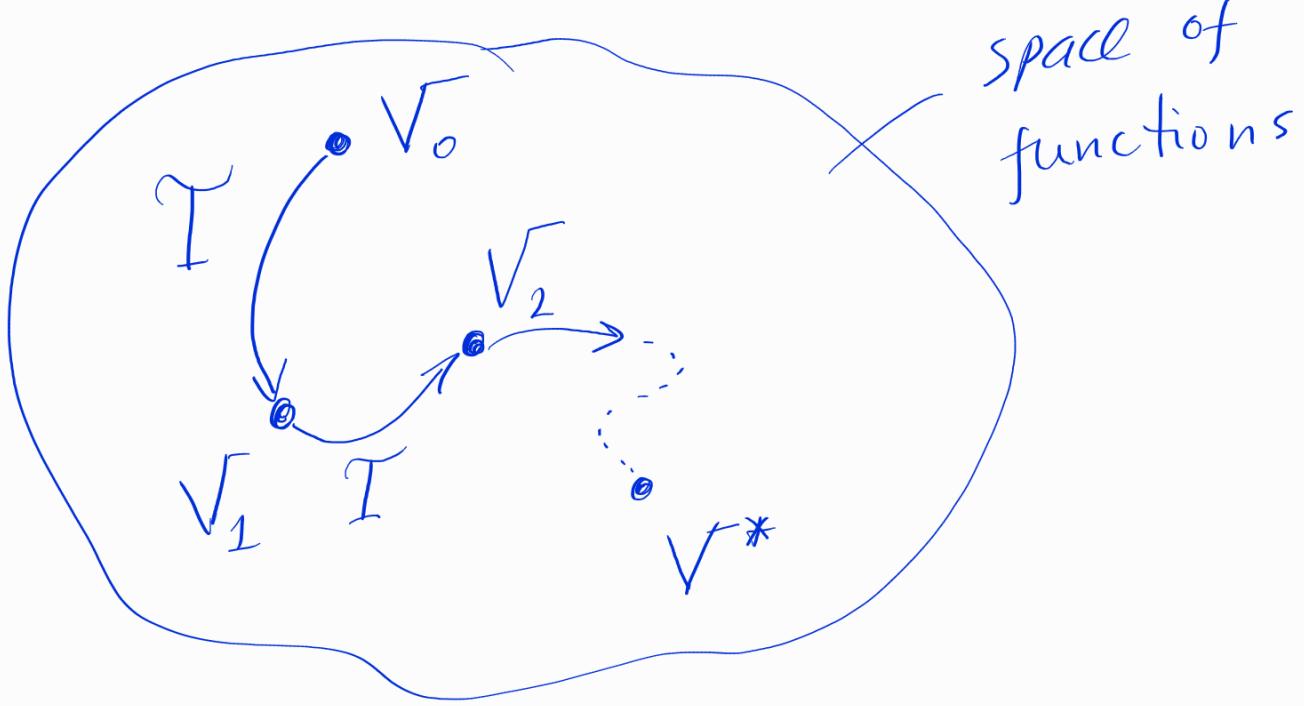
i-fold application of  $T$

$$\|T^i[V_0] - V^*\| \leq \frac{\gamma^i}{1-\gamma} \|T[V_0] - V_0\|$$



SUP-norm :

$$\|f - g\| \triangleq \sup_{x \in X} |f(x) - g(x)|$$



## Alg-m 1 ("Value iteration")

Start with a  $V_0$

apply  $T$  succ- ly

after suff - t accuracy , say,

$\|V_i - V_{i-1}\| \leq \epsilon$  is achieved,

compute a maximizer

$p_i$  of  $V_i$  and that's

appr - ly the optimal policy

## Algorithm 2 (Value iteration)

- Init value estimate  $V_0$
- do until  $\|V_i - V_{i-1}\| \leq \epsilon$   
for all states  $x$ :

$$\rho_i(x) := \arg \max_u \{r(x, u) + \gamma \mathbb{E}[V_i(x')]\}$$

$$V_{i+1}(x) := r(x, \rho_i(x)) + \gamma \mathbb{E}[V_i(x')]$$


---

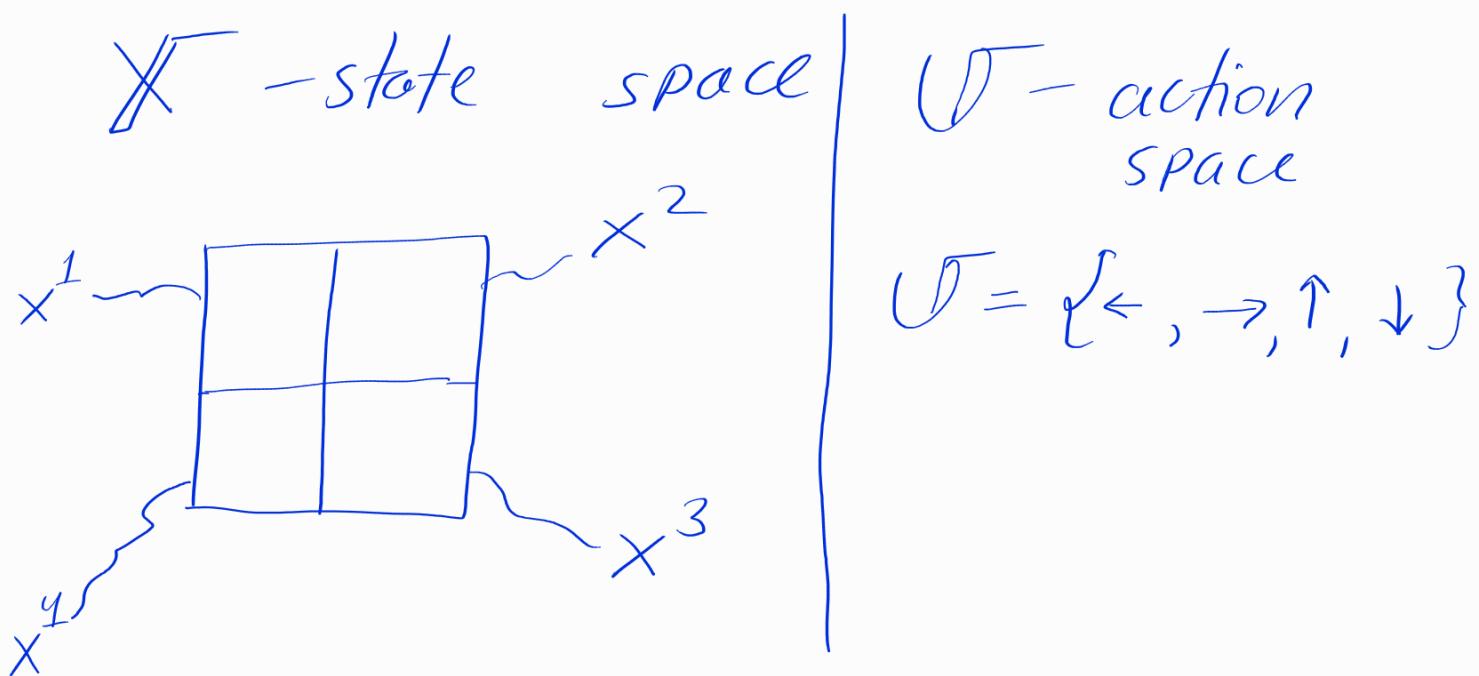
## Algorithm 3 (Policy iteration)

- Init policy  $\rho_0$
- do until  $\|V_i - V_{i-1}\| \leq \epsilon$ 
  - \*Solve (value update)
  - $V_i(x) = r(x, \rho_i(x)) + \gamma \mathbb{E}[V_i(x')]$
  - $\rho_{i+1}(x) := \arg \max_u \{r(x, u) + \gamma \mathbb{E}[V_i(x')]\}$

\* Start with a guess  $V_i^0$ , and do (for  $j$ ):

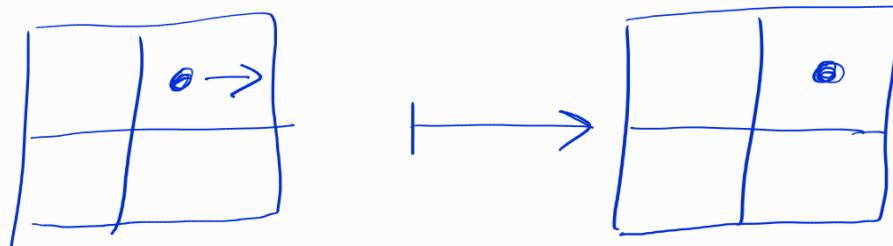
$$V_i^j(x) := r(x, p_i(x)) + \gamma E [V_i^{j-1}(x_{+}^{p_i(x)})]$$

Example: a small grid world



Environment dynamics:

in a cell, go to the next one if allowed, or stay if not



Reward

1	-100
1	100

Say  $x^3$  is terminal, i.e., any action will force the agent to stay at  $x^3$

Let's do value iteration.

Init  $V_0$

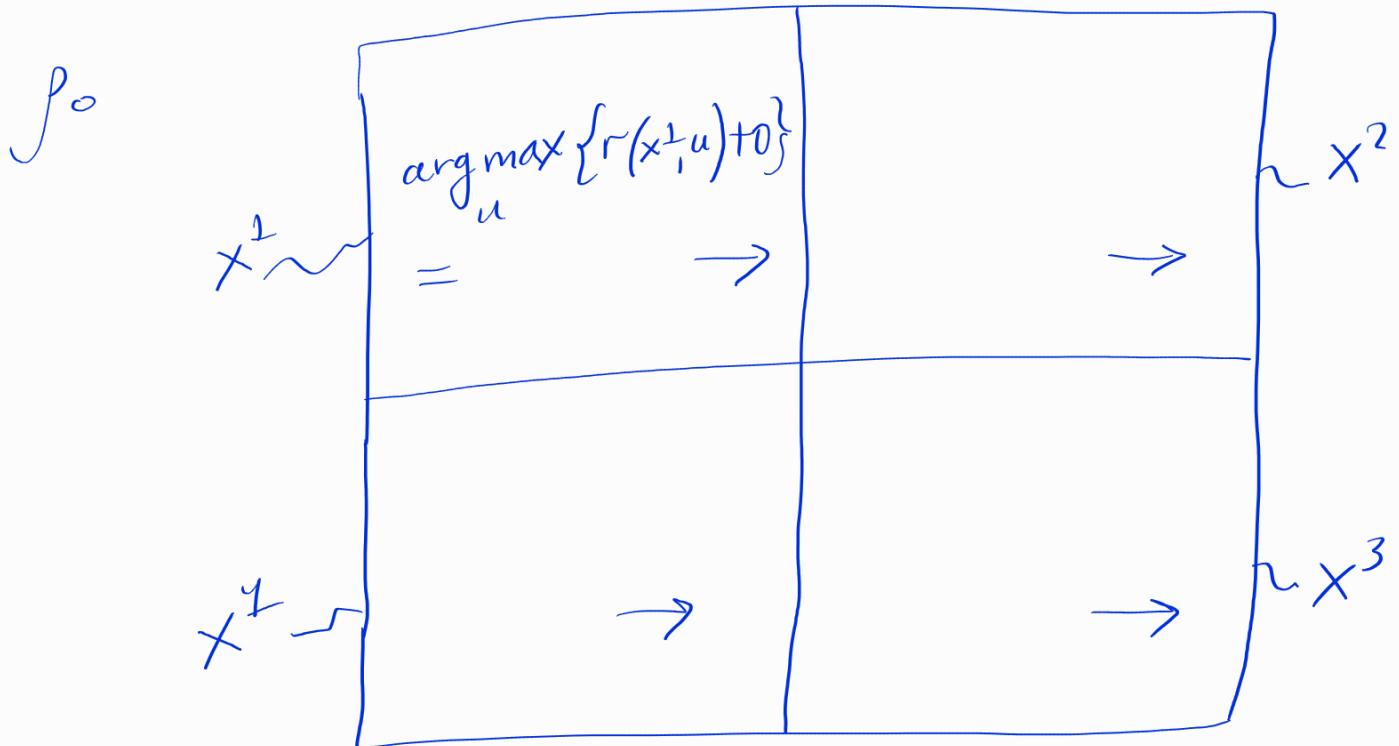
0	0
0	0

Let's calculate  $\rho_0$ :

algorithm recap

$$\rho_i(x) := \arg \max_u \{ r(x, u) + \gamma \mathbb{E}[V_i(x')] \}$$

Take  $\gamma = 1$ .  $E$  can be dropped.



Reward

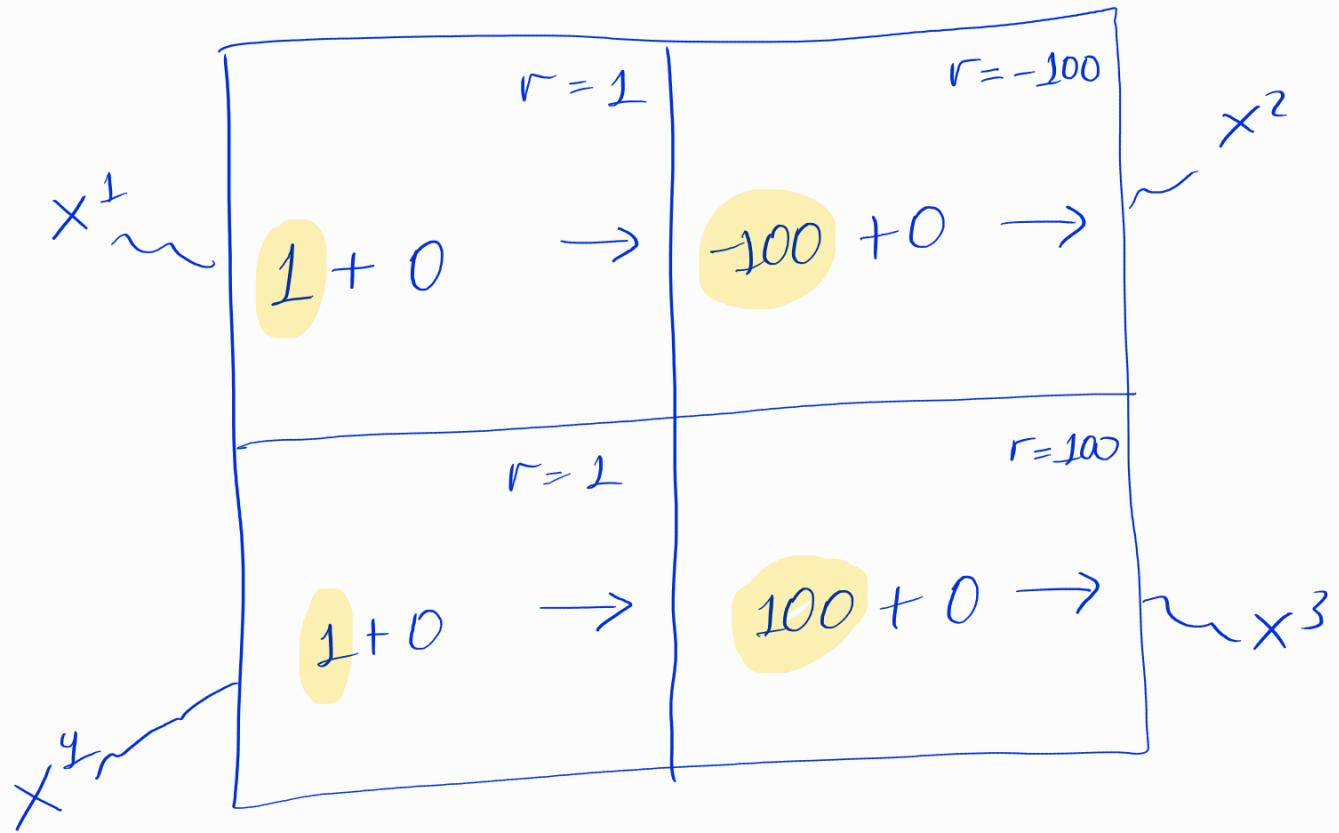
1	-100
1	100

Step 1

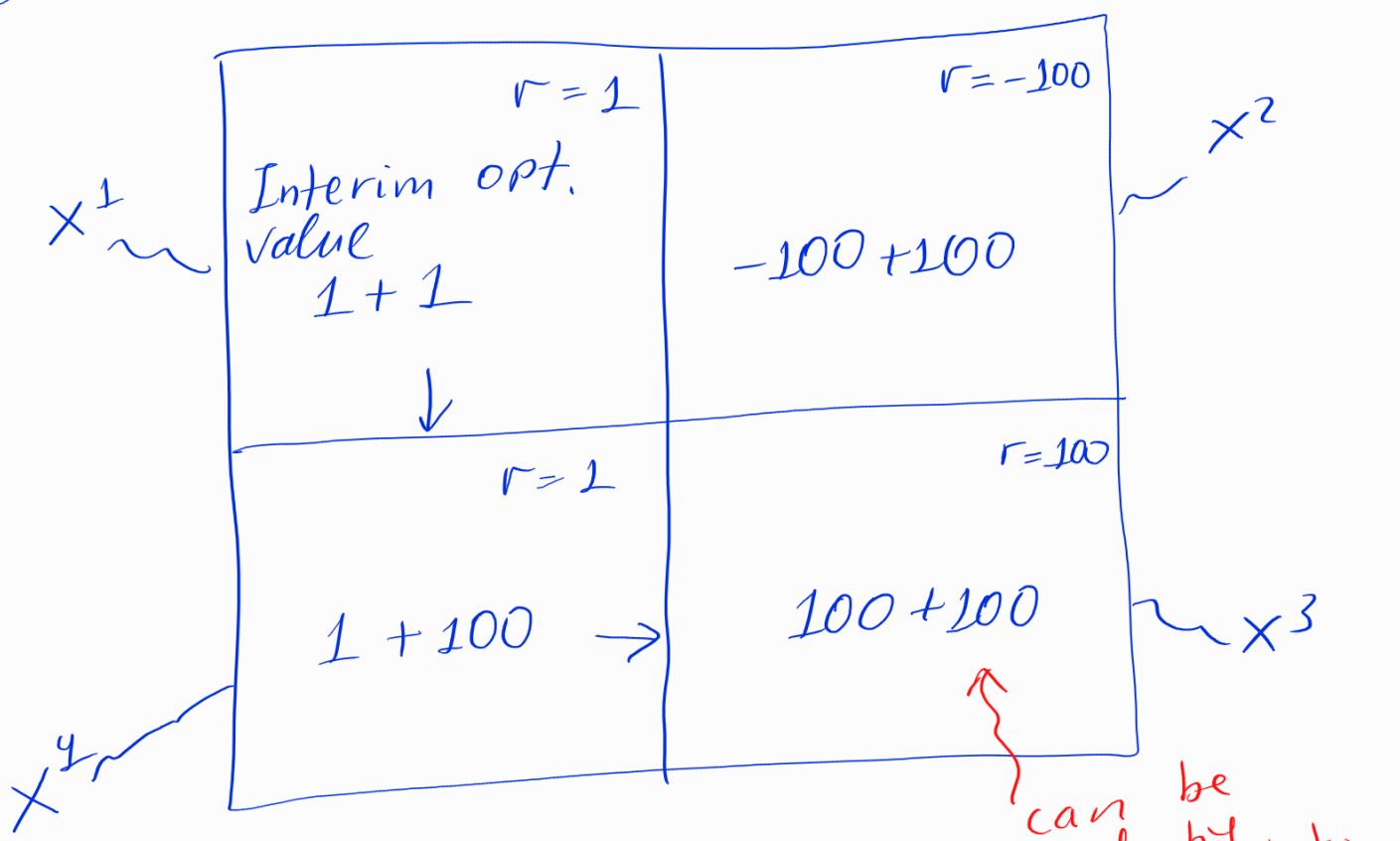
$V_1$

Algorithm recap

$$V_{i+1}(x) := r(x, p_i(x)) + \gamma E[V_i(x^{p_i(x)})]$$



$\beta_1$



algorithm recap

can be fixed by teleporting  $x^s$  into  $x^s$  with  $r^s = 0$

$$\beta_i(x) := \arg \max_u \{r(x, u) + \gamma E[V_i(x^u)]\}$$

$\checkmark^*$

$r=1$	$r=-100$
102	2
$r=1$	$r=100$
101	100

If game  
stops after  
3 steps

$\checkmark^*$

$r=-1$	$r=-100$
98	0
$r=-1$	$r=100$
99	100

Let's pretend the reward table  
was

-1	-100
-1	100