

function Approximation

BlOE 498/598

Value functions vs. Policies

Policy: $\pi(s) \rightarrow a$

Difficult to engineer from Scratch.

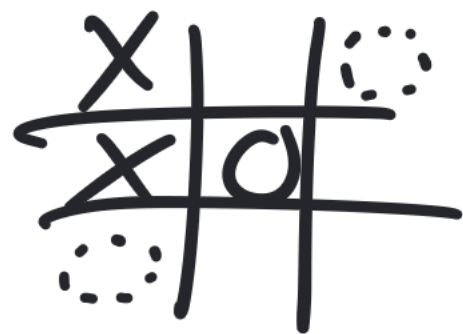
- Value functions are often more important than policies.
- We can "easily" construct π from V .
- It's easy to come up with a starting policy.

How Does AlphaGo work?

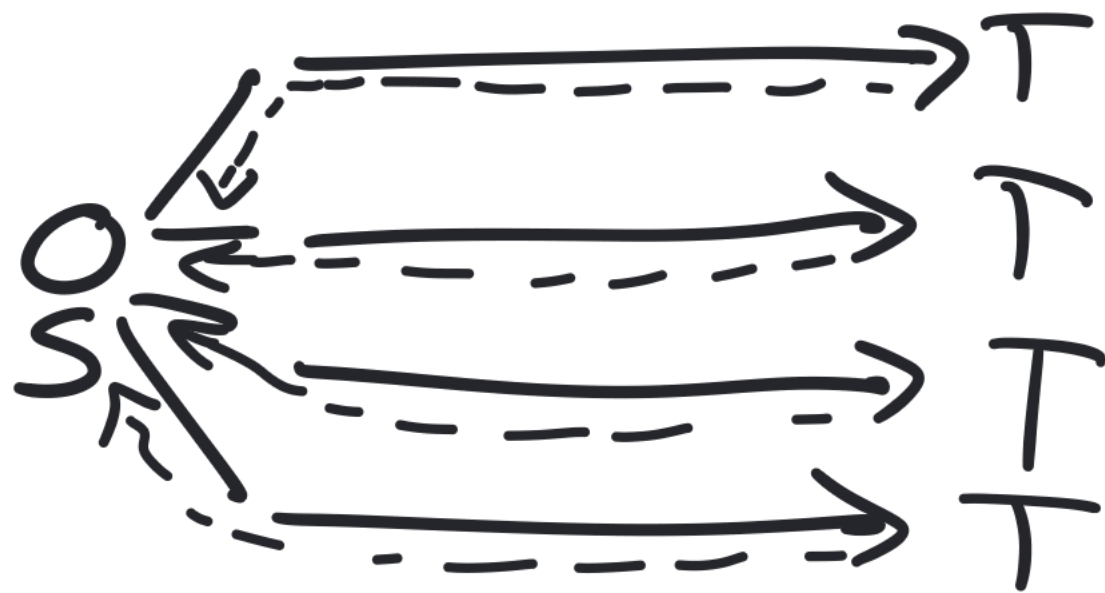
Began with a "human policy".

$$\pi_h \rightarrow V_{\pi_h}$$

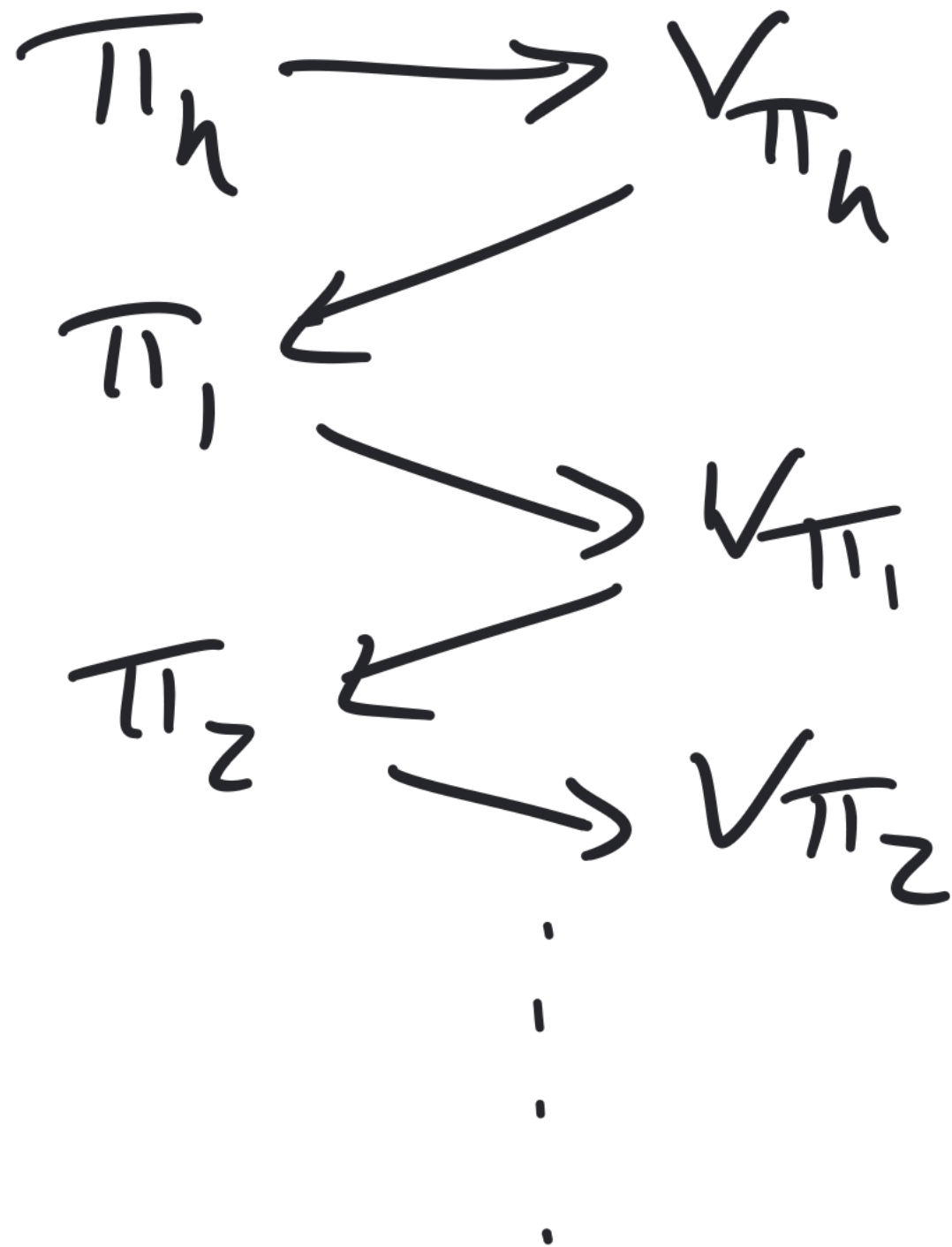
"Knowledge gradient" is very flat for Go.



Monte Carlo tree search.



RL for AlphaGo



AlphaGo Zero

No human policy

Why was Go so hard?

The state space is enormous.

Go is 19×19 grid.

Each square can have no, black, or white stones

\Rightarrow # of possible states is $3^{19^2} > \text{everything}$

$$V(s) \rightarrow \sum_i R_i$$

$\hookrightarrow \pi$ choose the best action
for every state

\Rightarrow Function Approximation.

Artificial Neural Networks

- ANN is an "overfitted" linear model with a nonlinearly transformed response.

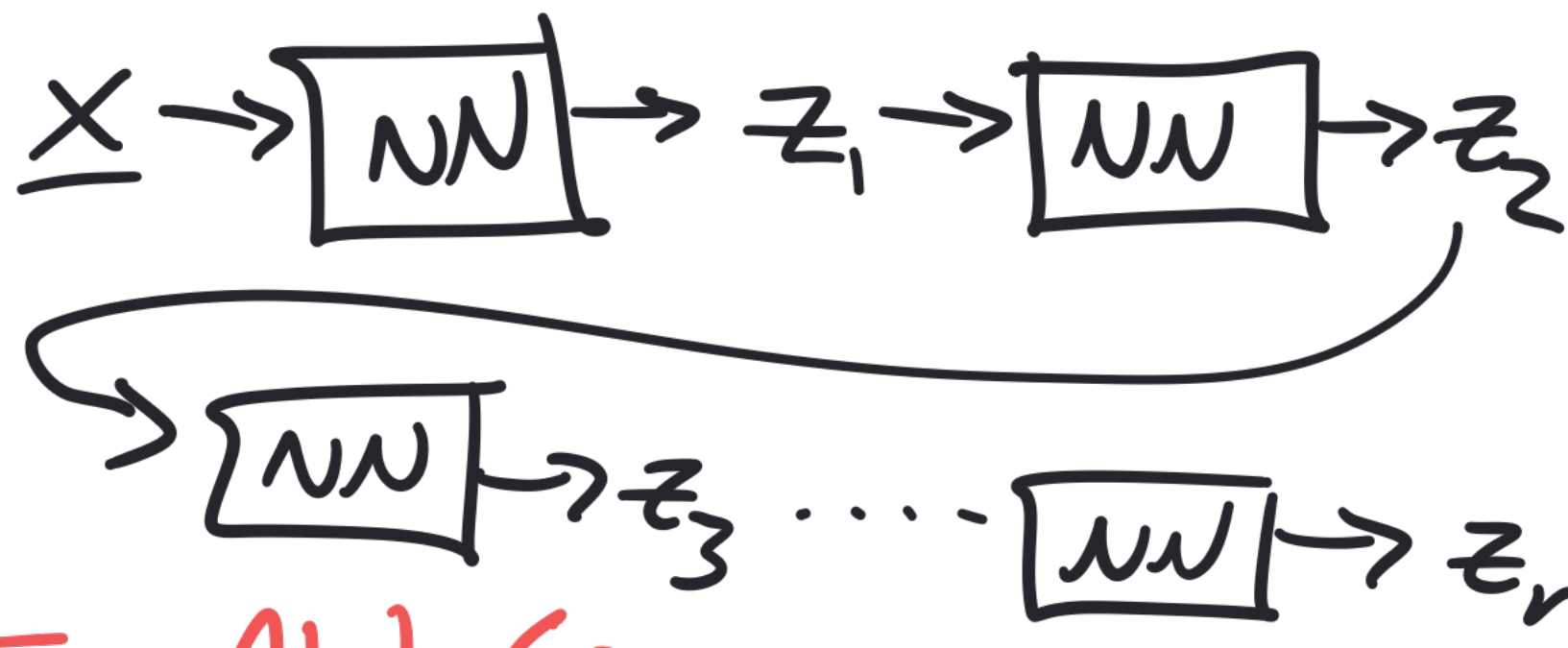
$$\underline{x} \rightarrow \underline{Ax} + \underline{b} \rightarrow \underline{y} \xrightarrow{\text{f}} \underline{z}$$

\mathbb{R}^n \mathbb{R}^p \mathbb{R}^p

$$\mathbb{R}^{p \times n}$$

$$\mathbb{R}^p$$

n inputs
 p outputs
 $p \times n + p$ parameters



In AlphaGo,

$V(s)$ is a deep NN

In AlphaGo:

$V(s)$ is a deep NN

$\pi(s)$ is also a deep NN



The value function is king.

Real-world RL.

- Self-driving cars
- Chat bots
- "Play" is expensive.
- How do we construct $V(s)$ from minimal data?