

RL: Lecture 22— Approximation

Harvey Mudd College

April 20, 2020

Neil Rhodes

Wednesday Guest Lecture

Prof. Erin Talvitie

Model-based RL

Grading up-to-date

except:

A few late submittal PA 6

All have an option for P/No credit

Quiz today

on

grade scope

Overview

- Stochastic Gradient Descent

SGD

input

S

actual value

$v_{\pi}(s)$

(y)

initialized randomly

w

Feature extraction
 $x(s)$

$x(s)$

$\hat{v}(w, x(s)) = w^T x$

\hat{v}

w

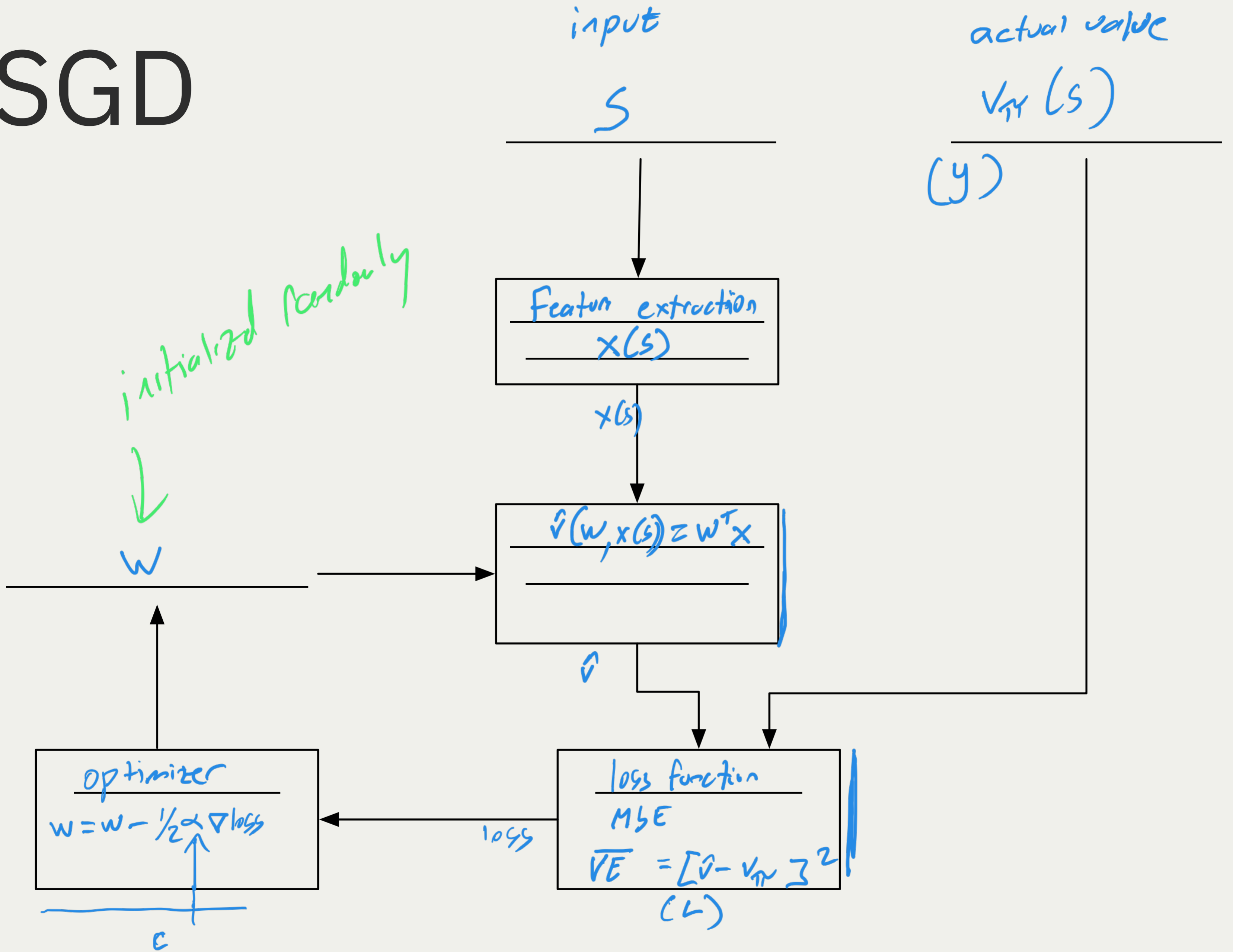
optimizer
 $w = w - \frac{1}{2} \alpha \nabla \log s$

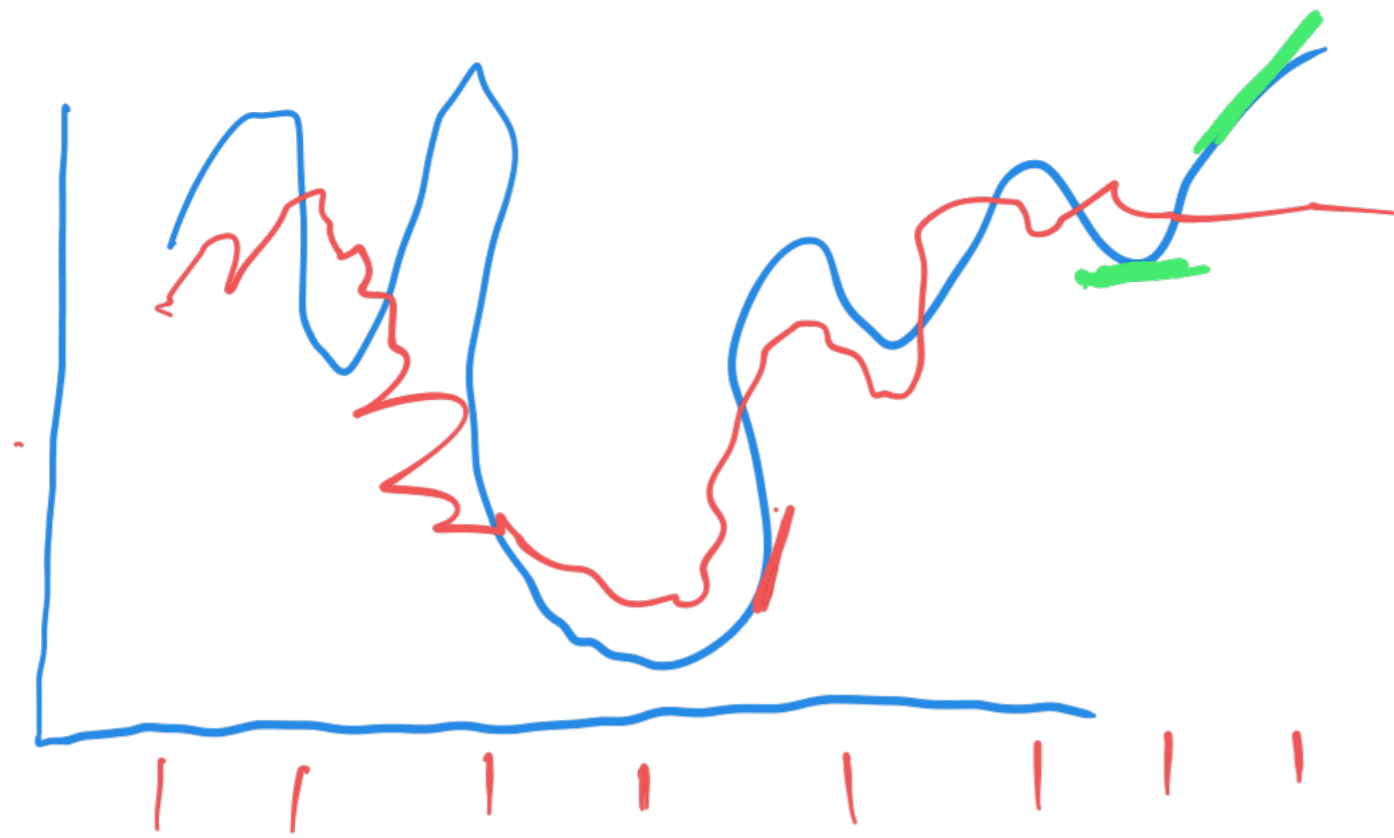
ϵ

loss function
MSE
 $\overline{VE} = [v - v_{\pi}]^2$

(L)

loss





gradient descent

one input
 quicker to calculate
 updating w same



in between
 mini-batches



all inputs (training
 example)

1 problem

lots of computation

advantage is

gradient is

exactly right

TD loop:

$$g(s, a) = g(s, a) + \alpha [R + \max_{a'} Q(s', a') - \hat{g}(s, a)]$$

input

y

approximate Q with our \hat{g}

New value of $g(s, a)$

our new training example is

(s, a)
pair

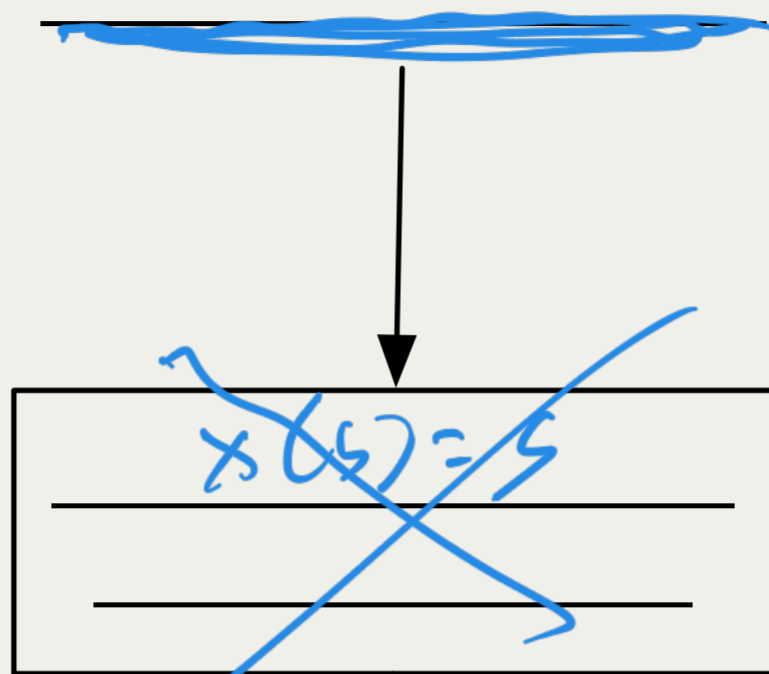
new ^{calculated} $g(s, a)$
target value (only)

HW 10

$$\hat{v}(s) = \underline{\quad}$$

input
state z

$v_{\pi}(s)$
7



$$\frac{\partial \text{loss}}{\partial w_0} = \frac{\partial \text{loss}}{\partial \hat{v}} \cdot \frac{\partial \hat{v}}{\partial w_0}$$

$$\frac{\partial \text{loss}}{\partial \hat{v}} = 2(\hat{v} - v_{\pi})$$

$$\frac{\partial \hat{v}}{\partial w_0} = 1$$

$$\frac{\partial \hat{v}}{\partial w_1} = s$$

$$\frac{\partial \hat{v}}{\partial w_2} = s^2$$

$$\text{loss} = (\hat{v} - v_{\pi})^2$$

$$\text{loss} = \frac{1}{2} (\hat{v} - v_{\pi})^2$$

$$w = \begin{bmatrix} 1 \\ 3 \\ 3 \end{bmatrix}$$

$$s = 2$$

$$\hat{v}(w, s) = w_2 s^2 + w_1 s + w_0$$

$$= 3 \cdot 4 + 3 \cdot 2 + 1$$

$$= 19$$

19

7

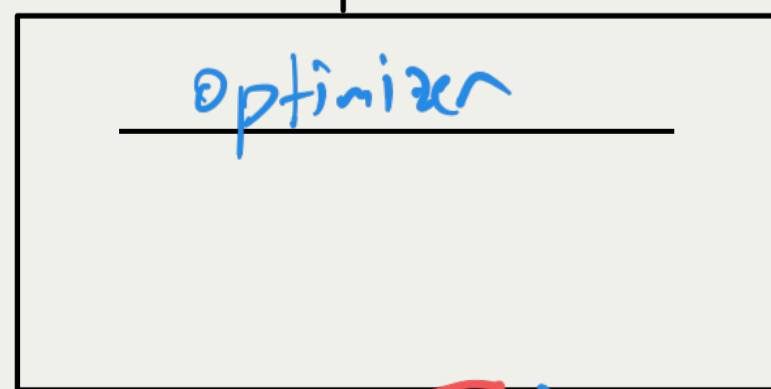
State loss

$$\text{loss} = (\hat{v} - v_{\pi})^2$$

$$= (19 - 7)^2$$

$$= 12^2 = 144$$

loss = 144



$$w_0 = w_0 - \alpha \frac{\partial \text{loss}}{\partial w_0} = 1 - \alpha \cdot 2(\hat{v} - v_{\pi}) = 1 - \alpha \cdot 2(19 - 7) = 1 - \alpha \cdot 24 = .88$$

$$= .76$$

$$\nabla \text{loss} = \left[\frac{\partial \text{loss}}{\partial w_0}, \frac{\partial \text{loss}}{\partial w_1}, \frac{\partial \text{loss}}{\partial w_2} \right]$$

$$= \left[(\hat{y} - y_{\text{target}}), (\hat{y} - y_{\text{target}})s, (\hat{y} - y_{\text{target}})s^2 \right]$$

$$= [12, 24, 48]$$

$$w = w - \alpha (\nabla \text{loss})^T$$

$$= \begin{bmatrix} 1 \\ 3 \\ 3 \end{bmatrix} - .01 \begin{bmatrix} 12 \\ 24 \\ 48 \end{bmatrix}$$

$$= \begin{bmatrix} 1 \\ 3 \\ 3 \end{bmatrix} - \begin{bmatrix} .12 \\ .24 \\ .48 \end{bmatrix} = \begin{bmatrix} .88 \\ 2.76 \\ 2.52 \end{bmatrix}$$

Features for Linear Methods

Polynomial

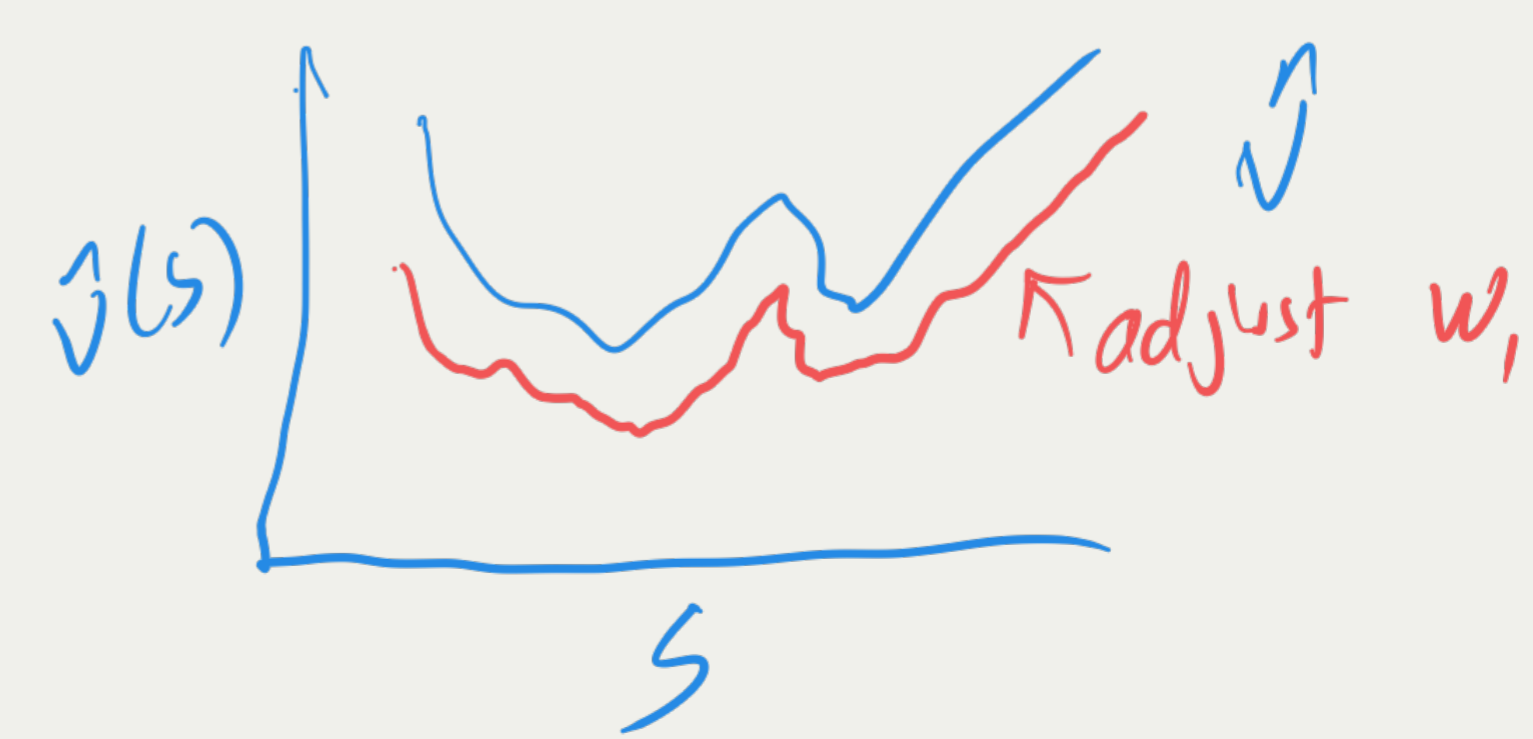
Let's say we have a state (s_1, s_2)

Example: $x(s) = [1, s_1, s_2, s_1 s_2, s_1^2, s_2^2, s_1 s_2^2, s_1^2 s_2, s_1^2 s_2^2]$

Annotations: "linear" with arrows pointing to s_1 and s_2 ; "linear" with arrows pointing to $s_1 s_2$; a bracket under $s_1^2 s_2^2$; a blue 'x' over $s_1^2 s_2^2$.

$$x(s) = \begin{bmatrix} 1 \\ \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

(draw up card, total so far)



non-NNs:

$S \rightarrow$ feature extraction $\rightarrow x(S) \rightarrow$

function approx model

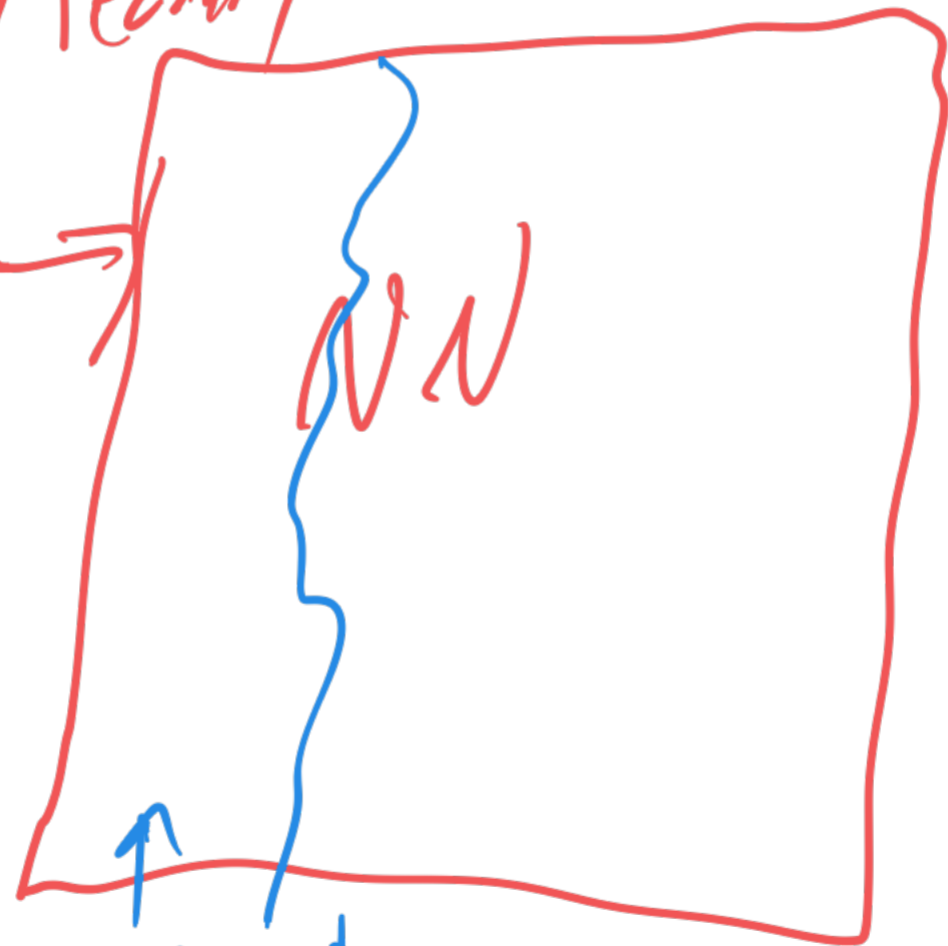


w

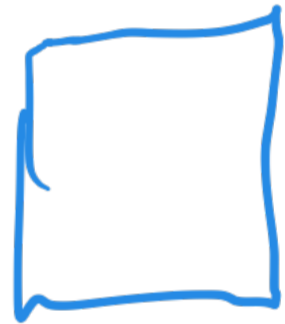
end-to-end learning
NNs: $S \rightarrow$

chess board

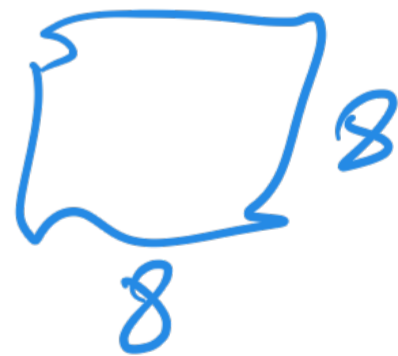
8x8 matrix
w/ number @
each entry
representing piece
 \uparrow



deep part
learns to do
feature extraction



8x8 for
white
pawns



white
rooks

Cartpole problem

[Actual Cartpole video](#) [Computer Cartpole Video](#)

State:

- Cart Position $[-2.4, 2.4]$ \mathbb{R} from left to right
 - Cart velocity: \mathbb{R}
 - Pole Angle: $[-41.8^\circ, 41.8^\circ]$ \mathbb{R}
 - Pole tip velocity: \mathbb{R} speed of pole
- (cart pos, cart v, pole angle, pole velocity)

Action:

- Left
- Right

2.39 2.3899999995

\mathbb{R} from left to right

\mathbb{R}
speed of pole

(cart pos, cart v,
pole angle, pole velocity)

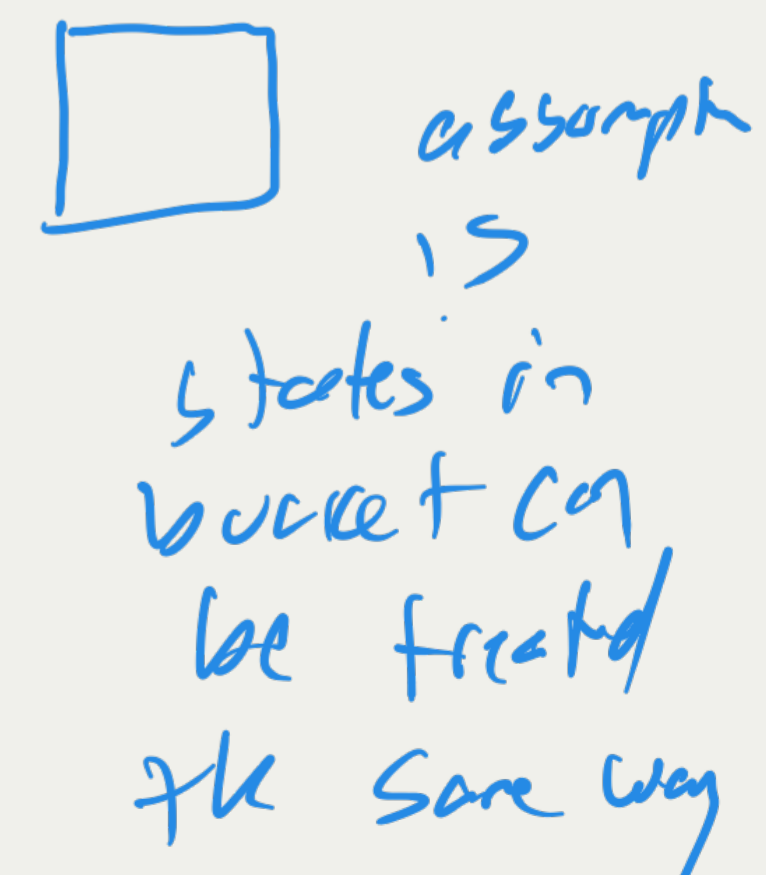
\mathbb{R}^4

Cartpole problem

Reward: 1 for every step taken, including termination step

Episode Termination:

- Pole angle more than $\pm 12^\circ$
- Cart position more than ± 24
- Episode length ≥ 500



why approximation
infinite # of states
SBD deals w/ approximation
Linear model:
how many features?
 $\{s_1, s_2, s_3, s_4\}$ s_3, s_4

