

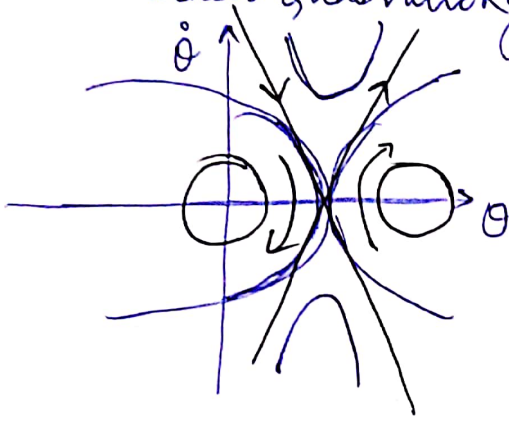
# Lecture - 6 (Dynamic Programming - III)

(Using Function Approximators like NNs)

Recap

- Tabular (Discrete State, action, time)  $\rightarrow$  Computationally expensive  
optimization over graph.

- Linear Quadratic Regular ~~state~~ (LQR) (Continuous State, action, Continuous time or discrete time)



LQR works well around the vicinity of some fixed point. But, it does not solve the problem globally. Linearization isn't valid across enough of the state space for that to be good.

Value Iteration ~~with~~ w/ Function Approximators (e.g. NNs)

$\hat{J}(x)$   
 $\nearrow$   
approximate cost-to-go

$$\begin{bmatrix} \hat{J}(s_0) \\ \hat{J}(s_1) \\ \hat{J}(s_2) \\ \vdots \end{bmatrix}$$

Fully discrete case

$$\hat{J}(x) = x^T S x \quad (\text{quadratic function})$$

(LQR case)

In General

$\hat{J}_\alpha(x)$   
 $\nearrow$   
parametrized by vector  $\alpha$

( $\alpha$  would be the weights & biases of the NNs (parameters))

## Fully discrete

$$\hat{J}(s_i) \leftarrow \min_a [l(s_i, a) + \hat{J}(f(s_i, a))]$$

take the current estimate of  $\hat{J}$  and compute for all of the points and treat it as desired target

(Minimize Cost function using ~~sub~~ gradient descent)

Function Approximator

At sample  $x_i$ ,

$$\hat{J}_\alpha(x_i) \approx \min_a [l(x_i, a) + \hat{J}_\alpha(f(x_i, a))]$$

$J^d(x_i)$

Minimize over parameter  $\alpha$

Sample points  $x_i$  Applying Gradient descent on above loss

$$\min_\alpha \sum_i [J_\alpha(x_i) - J^d(x_i)]^2$$

system dynamics

$$\Delta \alpha = \eta [J_\alpha(x_i) - J^d(x_i)] \frac{\partial \hat{J}}{\partial \alpha}$$

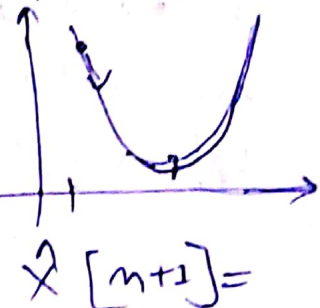
$$\alpha = \alpha + \Delta \alpha$$

When can we expect a Gradient Descent algorithm on a function approximator to actually converge to the true cost-to-go?

→ Fitted Value Iteration Algorithm

$$\min_x \|Ax - b\|_2^2$$

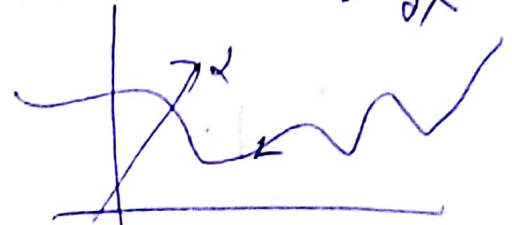
$$\frac{\partial R}{\partial x} = A^T(Ax - b) = 0$$



Really important special case -

Linear function approximators

$$\hat{J}_\alpha(x) = \underbrace{\alpha^T}_{\text{parameter}} \underbrace{\phi(x)}_{\text{basis func}} = \sum_i \alpha_i \phi_i(x)$$



If the parameters only enter linearly even if the basis functions are non-linear.

Radial basis functions

Bayesian interpolations

Deep neural networks — in the limit of  
ultra wide network  
(Neural Tangent Kernel)

→ Maybe the reason that <sup>big</sup> neural networks work ~~works~~  
even though they are non-linear in their  
parameters because they are so big and so wide  
that they start looking linear in the parameters  
again.

→ Explicitly, if we create a Multilayer perceptron, it is  
non linear in the parameters but the theoretical  
models that look at limits of arbitrarily big  
networks start to look like linear function  
approximators.

→ Some people say that why RL works when it works.

Fitted value iteration Converges.

TSitsiklis + Van Roy 1997

proof



# LQR

$\hat{J} = x^T \hat{S} x$  is a linear function approximator.

scalar  $\rightarrow$  linear in the parameters of  $\hat{S}$  but quadratic in  $x$

$$= \text{tr}(x^T \hat{S} x) = \text{tr}(S x x^T)$$

$\phi(x)$   $\rightarrow$  property of trace (cyclic elements)  
quadratic in  $x$  & coefficients are  $S$

## Algorithm

Given  $\hat{S}$ ,  $\langle x_i \rangle$ ,  $\langle u_i \rangle$

discrete  $u$ 's

① compute  $x_{i,j}^{\text{next}} = f(x_i, u_j) = Ax_i + Bu_j$   
next state (state, action pair)

(Applied on linear dynamics)

②  $J_{i,j}^{\text{next}} = (x_{i,j}^{\text{next}})^T \hat{S} (x_{i,j}^{\text{next}})$   
(for every state, action pair)

③  $J_i^d = \min_j [L(x_i, u_j) + J_{i,j}^{\text{next}}]$

④  $\text{loss}(L) = \sum_i [x_i^T \hat{S} x_i - J_i^d]^2$

⑤  $\Delta \hat{S} = -\eta \frac{\partial L}{\partial S}$  (gradient of  $\text{loss}$  (a scalar) w.r.t. a matrix ( $S$ ))

Loop

Is the infinite horizon cost finite?

$\rightarrow$  What we really need is ~~at least one~~ at least one the  $J_{i,j}^{\text{next}}$  to be 0. in order for the algorithm to eventually terminate

If it ~~can't~~ can't get directly to the origin, it's just going to be increasing cost forever and everytime we run the algorithm, it's going to add cost to the cost-to-go & eventually blow up.

→ for every  $X$  if we could find a  $u$  that would get us directly to the origin but if we discretize  $u$ , that's not gonna work.

### Alternative version

Continuous  $u$ 's

Given  $\hat{S}$ ,  $\langle x_i \rangle$   $\rightarrow$  bunch of  $x_i$ 's, all  $u$  available

$$J_i^d = \min_u \left[ x_i^T Q x_i + u^T R u + (A x_i + B u)^T \hat{S} (A x_i + B u) \right]$$

Even ~~if~~ when we have the wrong  $\hat{S}$ , we can still find the minimizing  $u$  exactly.

$$u_i^* = -R^{-1} B^T \hat{S} x_i$$

optimal  $u$

In the particular case of a linear function approximator when an infinite horizon solution exists, a linear function approximator is guaranteed to converge.

Note - In the first algorithm, we expected the LQR that it would work. It didn't <sup>quite</sup> work because the sampled version of LQR can't obtain zero finite cost but for the particular case of LQR



(Algorithm) we can do better and if we learn the samples but only change the actions to continuous actions (from discrete) then it resolves the problem & it converges.

## Discounted version

$$= [x[0]^T Q x[0] + u[0]^T R u[0] + J(0)]$$

$$\min_{u[n]} \sum_{k=0}^{\infty} \gamma^k (x[n]^T Q x[n] + u[n]^T R u[n])$$

$\gamma$  sol<sup>n</sup> to this can be written as  $x^T S x$  (controlled version)

$$0 < \gamma \leq 1$$

If  $\gamma$  becomes small, that means we are not worried about future, but only worried about today. discount factor

even if this doesn't go to 0, then that sum will converge to a finite sum.

LQR solution:  $\text{DARE}(A, B, Q, R) \Leftarrow$

discrete

sol<sup>n</sup> to the

undiscounted discrete-time

LQR problem

$$\text{DARE}(\sqrt{\gamma} A, B, Q, \frac{1}{\gamma} R) \Leftarrow \text{sol<sup>n</sup> to discounted LQR.}$$

→ Discounted LQR has a closed form solution and the closed form solution is equivalent to solving the original LQR problem but with a smaller  $A$  ( $\sqrt{\gamma} A$ ) which makes it more stable and it also modulates the cost ( $\frac{1}{\gamma} R$ ).

→ So, Discounted LQR is like assuming  $A$  is more stable than it actually is.

# Neural Fitted value iteration

$\hat{J}_2(x)$  as a neural network

Multi-layer Perceptron

(ReLU for activation)

3 layers ~ 255 hidden units per layer.

Discrete time:

$$\min_u [l_2(x) + u^T R u + \hat{J}^\pi(f(x, u))]$$

function approximation

result of  $f$  is inside our NNs (function approximation)

Cont. time:

$$J^d(x) \approx \min_u \left[ l(x, u) + \frac{\partial \hat{J}}{\partial x} f(x, u) \right]$$

In Cont. time, at least we don't

$l_2(x) + u^T R u + \frac{\partial \hat{J}}{\partial x} [f_1(x) + f_2(x)u]$  have to run through the NNs (function approx.). we just have to go through  $f$ .

$f_1(x) + f_2(x)u$   
(a lot of robots we care about an Control Affine).

we only have to depend on the partial derivative (slope of the cost-to-go function) not the value of the cost-to-go func. in some of the future place. That's the benefit of being instantaneous.

therefore have a closed form solver for  $u$

## Summary

### 3 Approaches to Value-Iteration

#### ① On a Grid

- everything works — it's beautiful, ~~is~~ fast
- But, grid gets really really big in higher state spaces.

#### ② LQR

- solves to arbitrary dimension basically but is restricted to the linearization of the system.

#### ③ NN

- don't have the same ~~guarantees~~ guarantees but they can often be made to work. We have to make decisions about discount factor, initial weights, choice of cost, etc.