

Last Time

Continuous  $S, A$

LQR - Linear Dynamics Quad. Reward

This time

More Continuous  $S, A$

Reinforcement Learning

Model Based Tabular

Offline

└ ADP

Online

└ MPC

└ Sparse Tree/Progressive Widening

---

ADP Value Function Approx

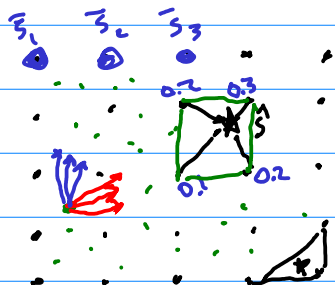
$$V(s) = \theta^T \beta(s)$$

↑  
features

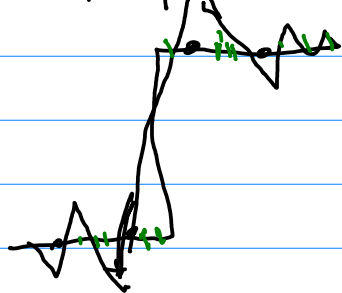
## Local

- Calculate value  $\tilde{z}$
- Interpolate Between

value at  $\tilde{z}$   
 $\Theta^T \beta(\tilde{z})$



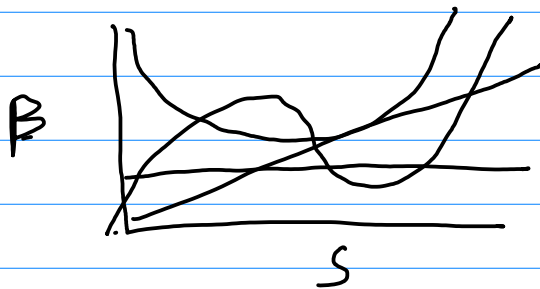
Multilinear  
 Simplex  $z^d$   
 $d+1$



## Global

$$\Theta^T \beta(s)$$

- Fourier
- Polynomial



## LSVI

$$\Theta \leftarrow 0$$

loop

for  $\tilde{s}_i$  in sampled states

$$\tilde{v}_i \leftarrow \max_a (R(\tilde{s}_i, a) + \gamma \frac{1}{m} \sum_{j=1}^m \Theta^T \beta(G(\tilde{s}_i, a, w_j)))$$

$$\Theta \leftarrow \text{Regress}(\tilde{s}, \beta, \tilde{v})$$

To get rid of Gibbs in local value approx choose  $\tilde{z}$  to be on interpolation points



## → MPC

Open Loop → Horizon  
 Differentiable Dynamic

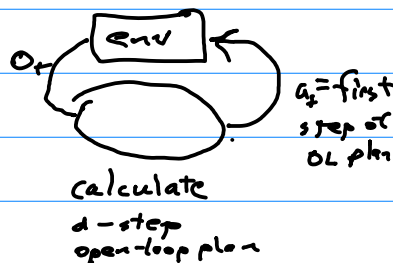
maximize  $a_{1:d} \quad \tilde{z}_{1:d}$

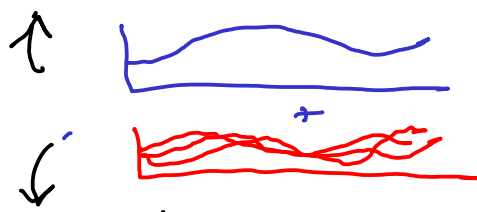
subject to

$$\sum_{t=1}^d \gamma^t R(s_t, a_t)$$

$$s_{t+1} = \text{mean}(T(s_t, a_t)) \quad \forall t$$

off the shelf solver

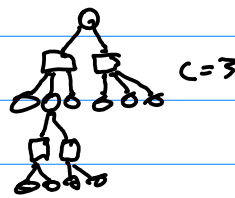




$$\begin{aligned} & \underset{a_{1:d}, s_{2:d}}{\text{maximize}} && \sum_{t=1}^d \sum_{i=1}^m \gamma^t R(s_t^i, a_t) \\ & \text{subject to} && s_{t+1}^i = G(s_t^i, a_t, w_t^i) \quad \forall t \end{aligned}$$

## Sparse Tree Search

continuous  $S$ , discrete  $A$



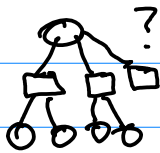
$$C = S \quad 8$$

$$C = 3$$

## Progressive Widening

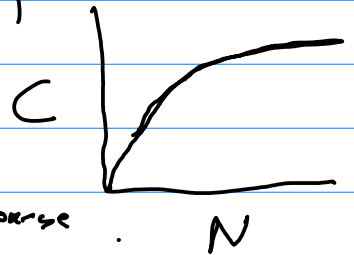
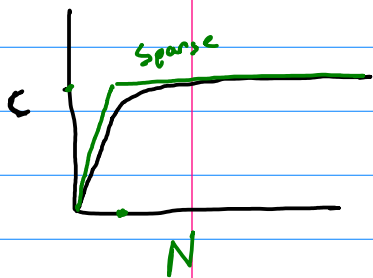
DPW  
cont.  $S, A$

$$k=8 \quad \alpha = \frac{1}{30}$$



add new branch if  
 $C < k N^\alpha$

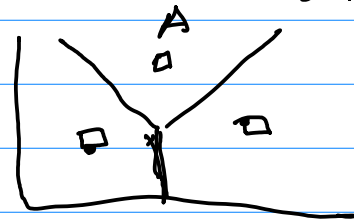
$$\alpha < 1$$



For states: ends up being close to sparse

For action: need a good heuristic

Voronoi



# Reinforcement Learning

## Course

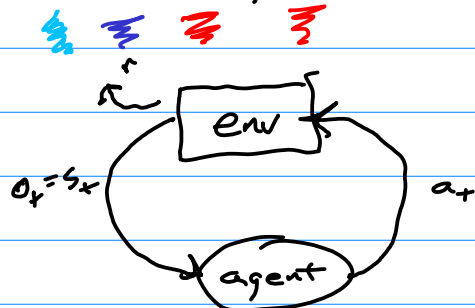
Immediate vs Future Rewards ✓

Unknown Models ←

Partial Observability

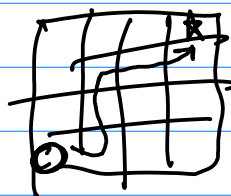
Other Agents

$(S, A, T, R, \gamma)$



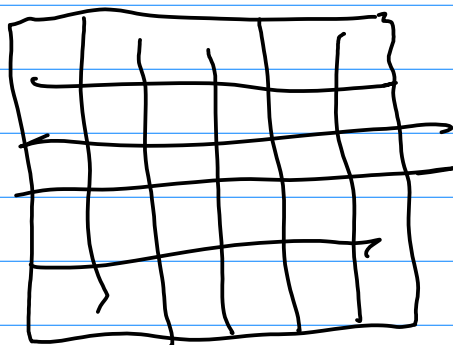
```

actions(env)
reset!(env)
s', r, done ← step!(env, a)
s', r ← G(s, a, w)
    
```



## Breakout Rooms

$S, A$



How would you solve this problem

→ Genetic Algorithm

$$\pi(s) = f(\theta, s)$$

$$\pi[s]$$

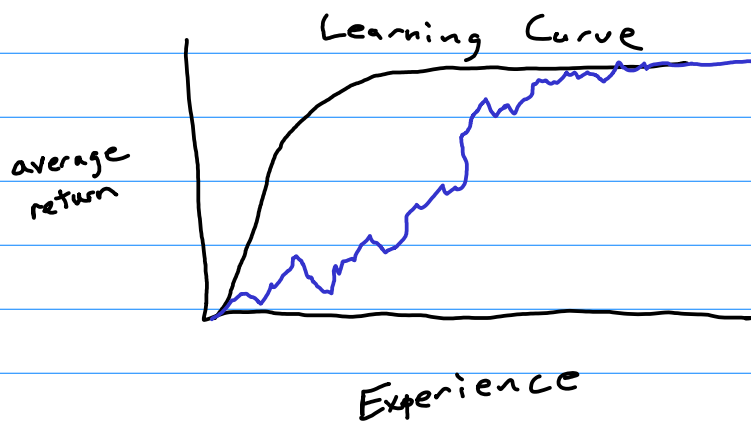
→ Learn Q values  
gradually take fewer random actions

→ Understand R  
trade off exploration/exploit

→ Q-learning  
Thompson Sampling  
T, R

# Challenges

1. Exploitation vs Explore
2. Credit Assignment
3. Generalization



Episode starting state  $\rightarrow$  terminal state

Model Based / Model Free  
 $\nwarrow$  learn  $Q, \pi$   
 $\nearrow$  learn  $T, R$

Tabular Max. Likelihood Model-Based RL

$$N[a][s, s']$$

$$N \leftarrow \text{zero}$$

$$p \leftarrow \text{zero} \quad p[a][s]$$

$$s \leftarrow s_0$$

$$\pi \leftarrow \text{random policy}$$

loop

$$a \leftarrow \text{rand}(A) \quad \text{w.p. } \epsilon, \pi(s) \text{ o.w.}$$

$$s', r \leftarrow \text{step!}(\text{env}, a)$$

$$N[a][s, s'] += 1$$

$$p[a][s] += r$$

$$T[a][s, s'] \leftarrow \frac{N[a][s, s']}{\sum_{s'} N[a][s, s']} \quad \forall s, a, s'$$

$N, T, R$

$$R[a][s] \leftarrow \frac{p[a][s]}{\sum_{s'} N[a][s, s']}$$

$\pi$

{

$$\pi \leftarrow \text{solve}(T, R)$$

$$s \leftarrow s'$$