

Last Time
Imitation Learning
IRL

match
 π^*
IRL
Given: $S, A, T, \{Z\}$
Find: R
RL
 S, A, T, R
 π^*

This Time
Peer Review
Transfer and Meta-Learning

Transfer Learning: Use experience from one set of tasks for faster learning and better performance on a new task,

In RL task = MDP

"source" domain \rightarrow "target" domain

"shot" = attempt in "target domain"

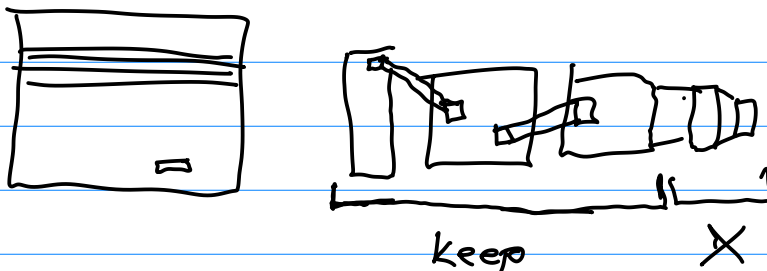
"0-shot" = run policy in target domain

"1-shot" = try task once

"few-shot"

How is prior knowledge stored?

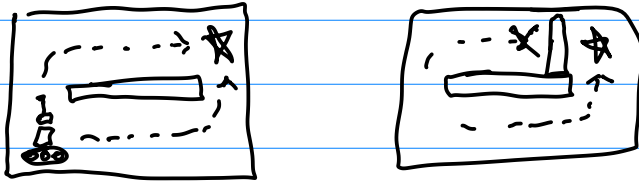
- Features \leftarrow
- Policy (Bad)
- Q-function
- Model (Physics)



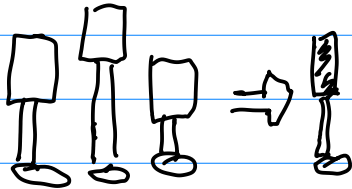
Approaches

Pre-training + Finetuning

Pretrain with robustness and diversity



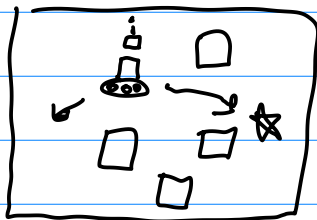
- Soft-Q-learning
- EPDpt



- CADZRL

Multi-Task RL

Make environment part of the state



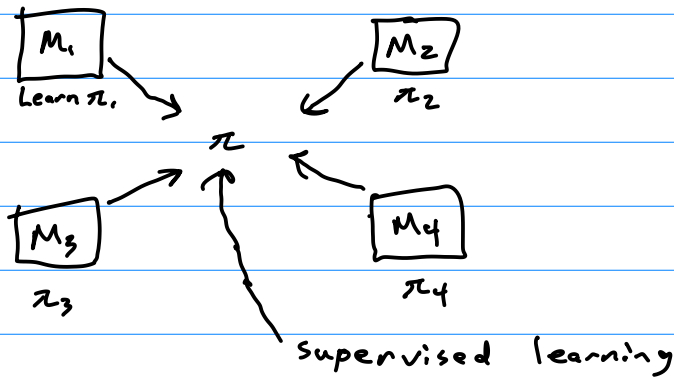
$$s = (r_x, r_y, g_x, g_y, o_{11}, o_{12}, \dots)$$

Problems

- Gradient Interference
- Winner-take-all

In practice, often negative transfer

Actor - Mimic



Successor Features - Reason about multiple reward functions at once

Problems all have same S, A, T, γ
different R

Recall $V^\pi = (I - \gamma T^\pi)^{-1} R^\pi$ $|S|$ vector

$$Q^\pi = (I - \gamma \bar{T}^\pi)^{-1} \bar{R}$$

$|S||A|$ vectors

Let $R(s,a) = \theta^T \phi(s,a)$

↑
features

$\bar{\phi}$ feature matrix $|S||A|$ columns, N rows

$$\bar{\phi}_{i,k} = \phi_i(s,a)$$

↑
 $k = s_a + a$

→ Successor Feature: $\Psi^\pi(s,a) = E \left[\sum_{t=0}^{\infty} \gamma^t \phi(s,a) \right]$

$$Q^\pi(s,a) = \theta^T \Psi^\pi(s,a)$$

$$\Psi^\pi(s,a) = \phi(s,a) + \gamma E \left[\Psi(s',a') \right]$$

$$\bar{\Psi}^\pi = (I - \gamma \bar{T}^\pi)^{-1} \bar{\phi}$$

$s' \sim T^\pi(s,a)$
 $a' = \pi(s')$

$$R'(s,a) = \Theta'^T \phi(s,a)$$

$$Q^\pi(s,a) = \Theta'^T \psi^\pi(s,a)$$

Meta Learning : Learning to Learn in Target domain

RL

$$\Theta^* = \operatorname{argmax}_{\Theta} E_{\pi_{\Theta}}[R(\tau)]$$

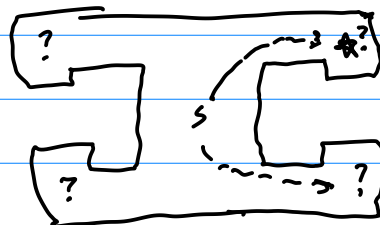
$$= f_{RL}(M)$$

Meta RL

$$\Theta^* = \operatorname{argmax}_{\Theta} \sum_i E_{\pi_{\phi_i}}[R(\tau)]$$

where $\phi_i = f_{\Theta}(M_i)$

Approach 1: This is a POMDP



RL²

Approach 2: Gradient Based Meta RL MAML

RL (policy gradient)

$$\Theta^{k+1} \leftarrow \Theta^k + \alpha \nabla_{\Theta^k} J(\Theta^k)$$

G.B.
Meta RL

$$f_{\Theta}(M_i) = \Theta + \alpha \nabla_{\Theta} J_i(\Theta)$$

Meta

